

Robust Q -learning Algorithm for Markov Decision Processes under Wasserstein Uncertainty [★]

Ariel Neufeld ^a, Julian Sester ^b

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

^b*National University of Singapore, Department of Mathematics, 21 Lower Kent Ridge Road, 119077*

Abstract

We present a novel Q -learning algorithm tailored to solve distributionally robust Markov decision problems where the corresponding ambiguity set of transition probabilities for the underlying Markov decision process is a Wasserstein ball around a (possibly estimated) reference measure. We prove convergence of the presented algorithm and provide several examples also using real data to illustrate both the tractability of our algorithm as well as the benefits of considering distributional robustness when solving stochastic optimal control problems, in particular when the estimated distributions turn out to be misspecified in practice.

Key words: Q -learning; Markov Decision Process; Wasserstein Uncertainty; Distributionally Robust Optimization; Reinforcement Learning.

1 Introduction

The among practitioners popular and widely applied Q -learning algorithm provides a tractable reinforcement learning methodology to solve Markov decision problems (MDP). The Q -learning algorithm learns an optimal policy online via observing at each time the current state of the underlying process as well as the reward depending on the current (and possibly next) state when acting according to a (not necessarily) optimal policy and by assuming to act optimally after the next state. The observed rewards determine a function Q depending on a state-action pair that describes the *quality* of the chosen action when being in the observed state. After a sufficient amount of observations the function Q then allows in each state to decide which actions possess the most *quality*. In this way the Q -learning algorithm determines an optimal policy.

The Q -learning algorithm was initially proposed in Watkins' PhD thesis ([57]). [27] and [58] then provided a rigorous mathematical proof of the convergence of the Q -learning algorithm to the optimal Q -value function using results from stochastic approximation theory

(see e.g. [16] and [41]). The design of the Q -learning algorithm as well as the proof of its convergence to the optimal Q -value both rely on the dynamic programming principle of the corresponding Markov decision problem, which allows to find an optimal policy for the involved infinite horizon stochastic optimal control problem by solving a one time-step optimization problem. We refer to [1], [2], [3], [11], [12], [24], [25], [28], [29], [35], [38], and [55] for various successful applications of the Q -learning algorithm.

Recently, there has been a huge focus in the literature starting from the viewpoint that one might have an estimate of the correct transition probability of the underlying Markov decision process, for example through the empirical measure derived from past observed data, but one faces the risk of misspecifying the correct distribution and hence would like to consider a distributionally robust Markov decision process (compare [5], [6], [13], [17], [23], [30], [31], [32], [37], [39], [47], [48], [52], [56], [59], [61], [62], [64], and [66]), also called *Markov decision process under model uncertainty*, where one maximizes over the worst-case scenario among all probability measures of an ambiguity set of transition probabilities. We also refer to, e.g. the following related distributionally robust stochastic control problems [13], [14], [22], [52], [53], [60], and [63] beyond the MDP setting. Indeed, as discussed in [31], there is a common risk in practice that

[★] Corresponding author A. Neufeld

Email addresses: ariel.neufeld@ntu.edu.sg (Ariel Neufeld), jul_ses@nus.edu.sg (Julian Sester).

one cannot fully capture the probabilities of the real-world environment due to its complexity and hence the corresponding reinforcement learning algorithm will be trained based on misspecified probabilities. In addition, there is the risk that the environment shifts between the training period and the testing period. This situation can often be observed in practice as the future evolution of random processes rarely behaves *exactly* according to, for example, the observed historical evolution. One may think as a prime example of financial markets, where several financial crises revealed repeatedly that used models were strongly misspecified. We refer to [31] for further examples, e.g. in robotics, and a further general discussion on the need of considering distributionally robust Markov decision processes and corresponding reinforcement learning based algorithms.

While there has been a lot of contributions in the literature on distributionally robust Markov decision problems, only very recently, to the best of our knowledge, there has been a first Q -learning algorithm developed in [31] to solve distributionally robust Markov decision problems. More precisely, in [31] the authors recently introduced a Q -learning algorithm tailored for distributionally robust Markov decision problems where the corresponding ambiguity set of transition probabilities consists of all probability measures which are ε -close to a reference measure with respect to the Kullback-Leibler (KL) divergence, and prove its convergence to the optimal robust Q -value function.

The goal of this paper is to provide a Q -learning algorithm which can solve distributionally robust Markov decision problems where the corresponding ambiguity set of transition probabilities for the underlying Markov decision process is a Wasserstein ball around a (possibly estimated) reference measure. We obtain theoretical guarantees of convergence of our Q -learning algorithm to the corresponding optimal robust Q -value function (see also (12)). The design of our Q -learning algorithm combines the dynamic programming principle of the corresponding Markov decision process under model uncertainty (see, e.g., [37]) and a convex duality result for worst-case expectations with respect to a Wasserstein ball (see [4], [9], [19], [34], and [65]).

From an application point of view, considering the Wasserstein distance has the crucial advantage that a corresponding Wasserstein-ball consists of probability measures which do not necessarily share the same support as the reference measure, compared to the KL-divergence, where by definition probability measures within a certain fixed distance to the reference measure all need to have a corresponding support included in the support of the reference measure. We highlight that from a structural point of view, our Q -learning algorithm is different than the one in [31], which roughly speaking comes from the fact that the dual optimization problem with respect to the Wasserstein distance has

a different structure than the corresponding one with respect to the KL-divergence.

We demonstrate in several examples also using real data that our *robust* Q -learning algorithm determines *robust* policies that outperform non-robust policies, determined by the classical Q -learning algorithm, given that the probabilities for the underlying Markov decision process turn out to be misspecified.

The remainder of the paper is as follows. In Section 2 we introduce the underlying setting of the corresponding Markov decision process under model uncertainty. In Section 3 we present our new Q -learning algorithm and provide our main result: the convergence of this algorithm to the optimal robust Q -value function. Numerical examples demonstrating the applicability as well as the benefits of our Q -learning algorithm compared to the classical Q -learning algorithm are provided in Section 4. All proofs and auxiliary results are provided in Appendix A.1 and A.2, respectively

2 Setting and Preliminaries

In this section we provide the setting and define necessary quantities to define our Q -learning algorithm for distributionally robust stochastic optimization problems under Wasserstein uncertainty.

2.1 Setting

Optimal control problems are defined on a state space containing all the states an underlying stochastic process can attain. We model this state space as a finite subset $\mathcal{X} \subset \mathbb{R}^d$ where $d \in \mathbb{N}$ refers to the dimension of the state space. We consider the robust control problem over an infinite time horizon, hence the space of all attainable states in this horizon is given by the infinite Cartesian product $\Omega := \mathcal{X}^{\mathbb{N}_0} = \mathcal{X} \times \mathcal{X} \times \dots$, with the corresponding σ -algebra $\mathcal{F} := 2^{\mathcal{X}} \otimes 2^{\mathcal{X}} \otimes \dots$. On Ω we consider a stochastic process that describes the states that are attained over time. To this end, we let $(X_t)_{t \in \mathbb{N}_0}$ be the canonical process on Ω , that is defined by $X_t(x_0, x_1, \dots, x_t, \dots) := x_t$ for each $(x_0, x_1, \dots, x_t, \dots) \in \Omega$, $t \in \mathbb{N}_0$.

Given a realization X_t of the underlying stochastic process at some time $t \in \mathbb{N}_0$, the outcome of the next state X_{t+1} can be influenced through actions that are executed in dependence of the current state X_t . At any time the set of possible actions is given by a finite set $A \subseteq \mathbb{R}^m$, where $m \in \mathbb{N}$ is the dimension of the action space (also referred to as control space). The set of admissible policies \mathcal{A} over the entire time horizon contains all sequences of actions that depend at any time only on the current

observation of the state process $(X_t)_{t \in \mathbb{N}_0}$ formalized by

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

The current state and the chosen action influence the outcome of the next state by influencing the probability distribution with which the subsequent state is realized. As we take into account model uncertainty we assume that the correct probability kernel is unknown and hence, for each given state x and action a , we consider an ambiguity set of probability distributions representing the set of possible probability laws for the next state. We denote by $\mathcal{M}_1(\Omega)$ and $\mathcal{M}_1(\mathcal{X})$ the set of probability measures on (Ω, \mathcal{F}) and $(\mathcal{X}, 2^{\mathcal{X}})$ respectively, and we assume that an ambiguity set of probability measures is modelled by a set-valued map

$$\mathcal{X} \times A \ni (x, a) \mapsto \mathcal{P}(x, a) \subseteq \mathcal{M}_1(\mathcal{X}). \quad (1)$$

Hence, if at time $t \in \mathbb{N}_0$ the process X_t attains the value $x \in \mathcal{X}$, and the agent decides to execute action $a \in A$, then $\mathcal{P}(x, a)$ describes the set of possible probability distributions with which the next state X_{t+1} is realized. If $\mathcal{P}(x, a)$ is single-valued, then the state-action pair (x, a) determines unambiguously the transition probability, and the setting coincides with the usual setting used for classical (i.e., non-robust) Markov decision processes, compare e.g. [7].

The ambiguity set of admissible probability distributions on Ω depends therefore on the initial state $x \in \mathcal{X}$ and the chosen policy $\mathbf{a} \in \mathcal{A}$. We define for every initial state $x \in \mathcal{X}$ and every policy $\mathbf{a} \in \mathcal{A}$ the set of admissible underlying probability distributions of $(X_t)_{t \in \mathbb{N}_0}$ by

$$\mathfrak{P}_{x, \mathbf{a}} := \left\{ \delta_x \otimes P_0 \otimes P_1 \otimes \cdots \mid \begin{aligned} &\text{for all } t \in \mathbb{N}_0 : \\ &P_t : \mathcal{X} \rightarrow \mathcal{M}_1(\mathcal{X}) \text{ Borel-measurable,} \\ &\text{and } P_t(x_t) \in \mathcal{P}(x_t, a_t(x_t)) \text{ for all } x_t \in \mathcal{X} \end{aligned} \right\},$$

where the notation $P = \delta_x \otimes P_0 \otimes P_1 \otimes \cdots \in \mathfrak{P}_{x, \mathbf{a}}$ abbreviates

$$\begin{aligned} P(B) &:= \sum_{x_0 \in \mathcal{X}} \cdot \sum_{x_t \in \mathcal{X}} \cdots \mathbf{1}_B((x_t)_{t \in \mathbb{N}_0}) \cdots P_{t-1}(x_{t-1}; \{x_t\}) \\ &\quad \cdots P_0(x_0; \{x_1\}) \delta_x(\{x_0\}), \quad B \in \mathcal{F}. \end{aligned}$$

Remark 1 In the literature of robust Markov decision processes one refers to $\mathfrak{P}_{x, \mathbf{a}}$ as being (s, a) -rectangular, see, e.g., [26], [45], [59]. This is a common assumption which turns out to be crucial to obtain a dynamic programming principle (see, e.g., [37, Theorem 2.7] and [43]) and therefore to enable efficient and tractable computations. Indeed, if one weakens this assumption the problem becomes computationally more expensive (see, e.g. [8, Section 2]), or can be provably intractable (compare [30]) and therefore cannot be solved by dynamic programming methods. Several approaches to solve robust MDPs w.r.t. non-rectangular ambiguity sets using methods other than dynamic programming however have recently been proposed, and are described in [21], [30], and [50].

To determine *optimal* policies we reward actions in dependence of the current state-action pair and the subsequent realized state. To this end, let $r : \mathcal{X} \times A \times \mathcal{X} \rightarrow \mathbb{R}$ be some *reward function*, and let $\alpha \in \mathbb{R}$ be a *discount factor* fulfilling

$$0 < \alpha < 1. \quad (2)$$

Then, our *robust* optimization problem consists, for every initial value $x \in \mathcal{X}$, 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 policies $\mathbf{a} \in \mathcal{A}$. More precisely, we aim for every $x \in \mathcal{X}$ to maximize $\inf_{\mathbf{P} \in \mathfrak{P}_{x, \mathbf{a}}} \left(\mathbb{E}_{\mathbf{P}} \left[\sum_{t=0}^{\infty} \alpha^t r(X_t, a_t, X_{t+1}) \right] \right)$ among all policies $\mathbf{a} \in \mathcal{A}$. The value function given by

$$\mathcal{X} \ni x \mapsto V(x) := \sup_{\mathbf{a} \in \mathcal{A}} \inf_{\mathbf{P} \in \mathfrak{P}_{x, \mathbf{a}}} \mathbb{E}_{\mathbf{P}} \left[\sum_{t=0}^{\infty} \alpha^t r(X_t, a_t, X_{t+1}) \right] \quad (3)$$

then describes the expectation of $\sum_{t=0}^{\infty} \alpha^t r(X_t, a_t, X_{t+1})$ under the worst case measure from $\mathfrak{P}_{x, \mathbf{a}}$ and under the optimal policy from $\mathbf{a} \in \mathcal{A}$ in dependence of the initial value.

2.2 Specification of the Ambiguity Sets

To specify the ambiguity set $\mathcal{P}(x, a)$ for each $(x, a) \in \mathcal{X} \times A$, we first consider for each $(x, a) \in \mathcal{X} \times A$ a reference probability measure. In applications, this reference measure may be derived from observed data. Considering an ambiguity set related to this reference measure then allows to respect deviations from the historic behavior in the future and leads therefore to a more *robust* optimal control problem that allows to take into account adverse scenarios, compare also [37]. To that end, let

$$\mathcal{X} \times A \ni (x, a) \mapsto \hat{\mathbf{P}}(x, a) \in \mathcal{M}_1(\mathcal{X}). \quad (4)$$

be a probability kernel, where $\hat{\mathbf{P}}(x, a)$ acts as reference probability measure for each $(x, a) \in \mathcal{X} \times A$. Then, for

every $(x_0, \mathbf{a}) \in \mathcal{X} \times \mathcal{A}$ we denote by

$$\hat{P}_{x_0, \mathbf{a}} := \delta_{x_0} \otimes \hat{P}(\cdot, a_0(\cdot)) \otimes \hat{P}(\cdot, a_1(\cdot)) \otimes \cdots \in \mathcal{M}_1(\Omega) \quad (5)$$

the corresponding probability measure on Ω that determines the distribution of $(X_t)_{t \in \mathbb{N}_0}$ in dependence of initial value $x_0 \in \mathcal{X}$ and the policy $\mathbf{a} \in \mathcal{A}$, i.e., we have for any $B \in \mathcal{F}$ that

$$\begin{aligned} \hat{P}_{x_0, \mathbf{a}}(B) := & \sum_{x_0 \in \mathcal{X}} \cdots \sum_{x_t \in \mathcal{X}} \cdots \mathbf{1}_B((x_t)_{t \in \mathbb{N}_0}) \cdots \\ & \cdot \hat{P}(x_{t-1}, a_{t-1}(x_{t-1}); \{x_t\}) \\ & \cdots \hat{P}(x_0, a_0(x_0); \{x_1\}) \delta_x(\{x_0\}). \end{aligned}$$

We provide two specifications of ambiguity sets of probability measures $\mathcal{P}(x, a)$, $(x, a) \in \mathcal{X} \times A$, as defined in (1). Both ambiguity sets rely on the assumption that for each given $(x, a) \in \mathcal{X} \times A$ the uncertainty with respect to the underlying probability distribution is modelled through a Wasserstein-ball around the reference probability measure $\hat{P}(x, a)$ on \mathcal{X} .

To that end, for any $q \in \mathbb{N}$, and any $P_1, P_2 \in \mathcal{M}_1(\mathcal{X})$, consider the q -Wasserstein-distance

$$W_q(P_1, P_2) := \left(\inf_{\pi \in \Pi(P_1, P_2)} \int_{\mathcal{X} \times \mathcal{X}} \|x - y\|^q d\pi(x, y) \right)^{1/q},$$

where $\|\cdot\|$ denotes the Euclidean norm on \mathbb{R}^d and where $\Pi(P_1, P_2) \subset \mathcal{M}_1(\mathcal{X} \times \mathcal{X})$ denotes the set of joint distributions of P_1 and P_2 . Since we consider probability measures on a finite space we have a representation of the form

$$P_i = \sum_{x \in \mathcal{X}} a_{i,x} \delta_x, \text{ with } \sum_{x \in \mathcal{X}} a_{i,x} = 1, \quad a_{i,x} \geq 0$$

for all $x \in \mathcal{X}$ for $i = 1, 2$, where δ_x denotes the Dirac-measure at point $x \in \mathcal{X}$. Hence, the q -Wasserstein-distance can also be written as

$$W_q(P_1, P_2) := \left(\min_{\pi_{x,y} \in \tilde{\Pi}(P_1, P_2)} \sum_{x,y \in \mathcal{X}} \|x - y\|^q \cdot \pi_{x,y} \right)^{1/q},$$

where

$$\begin{aligned} \tilde{\Pi}(P_1, P_2) := & \left\{ (\pi_{x,y})_{x,y \in \mathcal{X}} \subseteq [0, 1] \mid \sum_{x' \in \mathcal{X}} \pi_{x',y} = a_{2,y}, \right. \\ & \left. \sum_{y' \in \mathcal{X}} \pi_{x,y'} = a_{1,x} \text{ for all } x, y \in \mathcal{X} \right\}. \end{aligned}$$

Relying on the above introduced Wasserstein-distance we define two ambiguity sets of probability measures.

Setting 1.) The ambiguity set $\mathcal{P}_1^{(q,\varepsilon)}$

We consider for any fixed $\varepsilon > 0$ and $q \in \mathbb{N}$ the ambiguity set

$$\mathcal{X} \times A \ni (x, a) \mapsto \mathcal{P}_1^{(q,\varepsilon)}(x, a) := \left\{ P \in \mathcal{M}_1(\mathcal{X}) \text{ s.t. } W_q(P, \hat{P}(x, a)) \leq \varepsilon \right\} \quad (6)$$

being the q -Wasserstein ball with radius ε around the reference measure $\hat{P}(x, a)$, defined in (4). For each $(x, a) \in \mathcal{X} \times A$ the ambiguity set $\mathcal{P}_1^{(q,\varepsilon)}(x, a)$ contains all probability measures that are close to $\hat{P}(x, a)$ with respect to the q -Wasserstein distance. In particular, $\mathcal{P}_1^{(q,\varepsilon)}(x, a)$ contains also measures that are not necessarily dominated by the reference measure $\hat{P}(x, a)$.

Setting 2.) The ambiguity set $\mathcal{P}_2^{(q,\varepsilon)}$

We next define an ambiguity set that can particularly be applied when autocorrelated time-series are considered. In this case we assume that the past $h \in \mathbb{N} \cap [2, \infty)$ values of a time series $(Y_t)_{t=-h+1, -h+2, \dots}$ may have an influence on the subsequent value of the state process. Then, at time $t \in \mathbb{N}_0$ the state vector is given by

$$X_t = (Y_{t-h+1}, \dots, Y_t) \in \mathcal{X} := \mathcal{Y}^h \subset \mathbb{R}^{D \cdot h}, \quad (7)$$

with $\mathcal{Y} \subset \mathbb{R}^D$ finite, where $D \in \mathbb{N}$ describes the dimension of each value $Y_t \in \mathcal{Y} \subset \mathbb{R}^D$.

An example is given by financial time series of financial assets, where not only the current state, but also past realizations may influence the subsequent evolution of the assets and can therefore be modelled to be a part of the state vector, compare also the presentation in [37, Section 4.3].

Note that at each time $t \in \mathbb{N}_0$ the part $(Y_{t-h+2}, \dots, Y_t) \in \mathbb{R}^{D \cdot (h-1)}$ of the state vector X_{t+1} that relates to past information can be derived once the current state $X_t = (Y_{t-h+1}, Y_{t-h+2}, \dots, Y_t)$ is known. Only the realization of Y_{t+1} is subject to uncertainty. Conditionally on X_t the distribution of X_{t+1} should therefore be of the form $\delta_{(Y_{t-h+2}, \dots, Y_t)} \otimes \tilde{P} \in \mathcal{M}_1(\mathcal{X})$ for some probability measure $\tilde{P} \in \mathcal{M}_1(\mathcal{Y})$.

We write, given some $x = (x_1, \dots, x_h) \in \mathcal{X}$,

$$\pi(x) := (x_2, \dots, x_h) \in \mathcal{Y}^{h-1} \quad (8)$$

such that $x = (x_1, \pi(x)) \in \mathcal{X}$ and such that $\pi(X_t) = (Y_{t-h+2}, \dots, Y_t)$. The vector $\pi(x)$ denotes the projection

of x onto the last $h - 1$ components and represents the part of the state $x \in \mathcal{X}$ that is carried over to the subsequent state and is therefore not subject to any uncertainty. To reflect the fact that the first $h - 1$ components can be deterministically derived once the previous state is known, we impose now the assumption that the reference kernel is of the form

$$\mathcal{X} \times A \ni (x, a) \mapsto \hat{P}(x, a) = \delta_{\pi(x)} \otimes \hat{\tilde{P}}(x, a) \in \mathcal{M}_1(\mathcal{X}), \quad (9)$$

where $\hat{\tilde{P}}$ is a probability kernel defined by $\mathcal{X} \times A \ni (x, a) \mapsto \hat{\tilde{P}}(x, a) \in \mathcal{M}_1(\mathcal{Y})$. This allows us to define for any fixed $\varepsilon > 0$ and $q \in \mathbb{N}$ the ambiguity set¹

$$\mathcal{X} \times A \ni (x, a) \mapsto \mathcal{P}_2^{(q, \varepsilon)}(x, a) := \left\{ P \in \mathcal{M}_1(\mathcal{X}) \text{ s.t. } P = \delta_{\pi(x)} \otimes \tilde{P} \text{ for } \tilde{P} \in \mathcal{M}_1(\mathcal{Y}) \text{ with } W_q(\tilde{P}, \hat{\tilde{P}}(x, a)) \leq \varepsilon \right\}, \quad (10)$$

i.e., for each $(x, a) \in \mathcal{X} \times A$ we consider all measures of the form $\delta_{\pi(x)} \otimes \tilde{P}$ for \tilde{P} being close in the q -Wasserstein distance to $\hat{\tilde{P}}(x, a)$.

From now on, the ambiguity set of probability measures $\mathcal{P}(x, a)$, $(x, a) \in \mathcal{X} \times A$, either corresponds to $\mathcal{P}_1^{(q, \varepsilon)}(x, a)$, $(x, a) \in \mathcal{X} \times A$, defined by (6), or to $\mathcal{P}_2^{(q, \varepsilon)}(x, a)$, $(x, a) \in \mathcal{X} \times A$, defined by (10).

Remark 2 In various applications such as for example portfolio optimization in finance ([37, Section 4]), an agent would like to choose at each time t an action a_t not only based on the current observation of the state process but also based on some historical observations. To be able to cover such a scenario also in the context of Markov Decision Problems, it is a well-known procedure to extend the state space to be able to also include historical observations into the current state. The ambiguity set $\mathcal{P}_2^{(q, \varepsilon)}$ can therefore be seen as the natural extension of $\mathcal{P}_1^{(q, \varepsilon)}$ tailored exactly for that scenario described above. We highlight that in that case, given an agent observes $X_t = (Y_{t-h+1}, \dots, Y_t)$ at time t , the only uncertainty on $X_{t+1} = (Y_{t-h+2}, \dots, Y_t, Y_{t+1})$ lies in the last component Y_{t+1} , and not in the whole vector X_{t+1} , as the other components are observed through X_t . This explains the structure of the corresponding measures in $\mathcal{P}_2^{(q, \varepsilon)}$ involving Dirac measures.

Remark 3 Using the Wasserstein distance for capturing distributional uncertainty differs significantly from employing the Kullback-Leibler distance, which was used,

¹ By abuse of notation W_q here denotes the q -Wasserstein distance on $\mathcal{M}_1(\mathcal{Y})$.

e.g., in [31]. By using an ambiguity set defined via the Wasserstein distance, one can consider all probability distributions that are in proximity to a reference measure, even if they are not necessarily absolutely continuous with respect to it. This becomes important when the reference measure is estimated from historical data and contains point masses at the observed values, but one does not want to restrict future values to those observed in the past. In contrast, if one is confident about the support of the underlying transition kernel, it can be advantageous to use an ambiguity set defined using a distance such as the Kullback-Leibler distance which only considers probability measures with the same support (or smaller) as the reference measure.

2.3 Definition of Operators

We consider the following single time step optimization problem

$$\mathcal{TV}(x) := \sup_{a \in A} \inf_{P \in \mathcal{P}(x, a)} \mathbb{E}_P[r(x, a, X_1) + \alpha V(X_1)], \quad x \in \mathcal{X}, \quad (11)$$

where $\mathcal{X} \ni x \mapsto V(x)$ is the value function defined in (3), and we define the optimal robust Q -value function by

$$\mathcal{X} \times A \ni (x, a) \mapsto Q^*(x, a) := \inf_{P \in \mathcal{P}(x, a)} \mathbb{E}_P[r(x, a, X_1) + \alpha V(X_1)]. \quad (12)$$

Note that if (2) holds and $\mathcal{P}(x, a)$ is either $\mathcal{P}_1^{(q, \varepsilon)}(x, a)$ or $\mathcal{P}_2^{(q, \varepsilon)}(x, a)$ for all $(x, a) \in \mathcal{X} \times A$, then the values of Q^* are finite, since for all $(x, a) \in \mathcal{X} \times A$ we have

$$\begin{aligned} |Q^*(x, a)| &\leq \inf_{P \in \mathcal{P}(x, a)} \mathbb{E}_P[|r(x, a, X_1)| + \alpha |V(X_1)|] \\ &\leq \sup_{y \in \mathcal{X}} |r(x, a, y)| + \alpha \sup_{y \in \mathcal{X}} |V(y)| < \infty, \end{aligned} \quad (13)$$

where the finiteness of V follows from [37, Theorem 2.7]. Then we obtain as a consequence of the main result from [37, Theorem 3.1] the following proposition showing that the infinite time horizon distributionally robust optimization problem defined in (3) can be solved by the consideration of a suitable one time-step fixed point equation, which is the key result that allows to derive Q -learning type of algorithms.

Proposition 4 Assume that (2) holds and that the ambiguity set $\mathcal{P}(x, a)$ is either given by $\mathcal{P}_1^{(q, \varepsilon)}(x, a)$ or $\mathcal{P}_2^{(q, \varepsilon)}(x, a)$ for all $(x, a) \in \mathcal{X} \times A$. Then for all $x \in \mathcal{X}$ we have $\sup_{a \in A} Q^*(x, a) = \mathcal{TV}(x) = V(x)$, where $\mathcal{X} \ni x \mapsto V(x)$ corresponds to the value function of the robust stochastic optimal control problem defined in (3).

3 The Robust Q -learning Algorithm

In this section we present a novel robust Q -learning algorithm for the corresponding distributionally robust stochastic optimization problem (3) and prove its convergence.

A robust Q -learning algorithm intends to approximate $Q^*(x, a) = \inf_{P \in \mathcal{P}(x, a)} \mathbb{E}_P[r(x, a, X_1) + \alpha V(X_1)]$ which involves the minimization over an infinite amount of probability measures. Due to the particular choice of ambiguity sets (6) and (10) w.r.t. the Wasserstein-distance, we can transform this minimization problem into a tractable problem using a duality from, e.g., [4].

To this end, for a function $f : \mathcal{X} \rightarrow \mathbb{R}$ we define, as in [4, Section 2] or [54, Section 5] its λc -transform.

Definition 5 (λc -transform) *Let $f : \mathcal{X} \rightarrow \mathbb{R}$, let $\lambda \geq 0$, and let $c : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$. Then the λc -transform of f is defined by $\mathcal{X} \ni x \mapsto (f)^{\lambda c}(x) := \sup_{y \in \mathcal{X}} \{f(y) - \lambda \cdot c(x, y)\}$.*

Indeed, the λc -transform now allows to rephrase the optimization problem involved in the definition of Q^* in more tractable terms involving only an expectation with respect to the reference kernel, compare also Proposition 16. We use this representation to define our robust Q -learning algorithm which is summarized in Algorithm 1.

The update rule from (16) in Algorithm 1 means that for all $(x, a) \in \mathcal{X} \times A$, $t \in \mathbb{N}_0$, we have $Q_{t+1}(x, a) = Q_t(x, a) + \tilde{\gamma}_{\text{visits}(x, a)} \left(-(-f_{t, (x, a)})^{\lambda c}(X_{t+1}) - \varepsilon^q \lambda_t - Q_t(x, a) \right)$ if $(x, a) = (X_t, a_t(X_t))$ and $Q_{t+1}(x, a) = Q_t(x, a)$ else, i.e., the update of Q_{t+1} only takes that state-action pair into account which was realized by the process $(X_t)_{t \in \mathbb{N}}$. Further, note that Algorithm 1 assumes for each time $t \in \mathbb{N}_0$ the existence of some $\lambda_t \in [0, \infty)$ such that (15) holds. The following result ensures that this requirement is indeed fulfilled.

Lemma 6 *Let $(x, a) \in \mathcal{X} \times A$, $t \in \mathbb{N}_0$, let $\hat{P} \in \mathcal{M}_1(\mathcal{X})$ and recall $\mathcal{X} \ni y \mapsto f_{t, (x, a)}(y)$ defined in (14). Further let $\mathcal{X} \times \mathcal{X} \ni (x, y) \mapsto c(x, y) \in [0, \infty]$ satisfy $\min_{y \in \mathcal{X}} c(x, y) = 0$ for all $x \in \mathcal{X}$. Then, there exists some $\lambda^* \in [0, \infty)$ such that $\mathbb{E}_{\hat{P}}[-(-f_{t, (x, a)})^{\lambda^* c}(X_1) - \varepsilon^q \lambda^*] = \sup_{\lambda \geq 0} \left(\mathbb{E}_{\hat{P}}[-(-f_{t, (x, a)})^{\lambda c}(X_1) - \varepsilon^q \lambda] \right)$.*

The following main result now shows that the function $(Q_t)_{t \in \mathbb{N}_0}$ obtained as the output of Algorithm 1 converges indeed against the optimal robust Q -value function Q^* defined in (12).

Algorithm 1 Robust Q -learning

Input State space $\mathcal{X} \subset \mathbb{R}^d$; Control space $A \subset \mathbb{R}^m$; Reward function r ; Discount factor $\alpha \in (0, 1)$; Kernel \hat{P} ; Starting point x_0 ; Policy $\mathbf{a} \in \mathcal{A}$; Cost function c of the λc -transform; Ambiguity parameter $\varepsilon > 0$; Parameter $q \in \mathbb{N}$ related to the Wasserstein-distance; Sequence of learning rates $(\tilde{\gamma}_t)_{t \in \mathbb{N}_0} \subseteq [0, 1]$;

- 1: Initialize $Q_0(x, a)$ for all $(x, a) \in \mathcal{X} \times A$ to an arbitrary real value;
- 2: Initialize $\text{visits}(x, a) \leftarrow 0$ for all $(x, a) \in \mathcal{X} \times A$;
- 3: **for** $t = 0, 1, \dots$ **do**
- 4: Set for all $(x, a) \in \mathcal{X} \times A$:
- 5:

$$\text{visits}(x, a) \leftarrow \begin{cases} \text{visits}(x, a) + 1 & \text{if } (x, a) = (X_t, a_t(X_t)), \\ \text{visits}(x, a) & \text{else;} \end{cases}$$

- 6: Define the map:

$$\begin{aligned} \gamma_t : \mathcal{X} \times A \times \mathcal{X} &\rightarrow \mathbb{R}, \\ (x, a, x') &\mapsto \gamma_t(x, a, x') := \tilde{\gamma}_{\text{visits}(x, a)} \mathbb{I}_{\{(x', a_t(x')) = (x, a)\}}; \end{aligned}$$

- 7: For every $(x, a) \in \mathcal{X} \times A$ we set:

$$\begin{aligned} f_{t, (x, a)} : \mathcal{X} &\rightarrow \mathbb{R}, \\ y &\mapsto r(x, a, y) + \alpha \max_{b \in A} Q_t(y, b); \end{aligned} \quad (14)$$

- 8: Choose $\lambda_t \in [0, \infty)$ which satisfies:

$$\begin{aligned} &\mathbb{E}_{\hat{P}(X_t, a_t(X_t))} [-(-f_t(X_t, a_t(X_t)))^{\lambda_t c}(X_{t+1}) - \varepsilon^q \lambda_t] \\ &= \sup_{\lambda \geq 0} \mathbb{E}_{\hat{P}(X_t, a_t(X_t))} [-(-f_t(X_t, a_t(X_t)))^{\lambda c}(X_{t+1}) - \varepsilon^q \lambda]; \end{aligned} \quad (15)$$

- 9: For all $(x, a) \in \mathcal{X} \times A$ we define the following update rule:

$$\begin{aligned} Q_{t+1}(x, a) &:= Q_t(x, a) \\ &\quad + \gamma_t(x, a, X_t) \cdot \left(-(-f_{t, (x, a)})^{\lambda_t c}(X_{t+1}) - \varepsilon^q \lambda_t - Q_t(x, a) \right); \end{aligned} \quad (16)$$

- 10: **end for**

Output A sequence $(Q_t(x, a))_{t \in \mathbb{N}_0, x \in \mathcal{X}, a \in A}$

Theorem 7 *Assume that (2) holds, and let $(x_0, \mathbf{a}) \in \mathcal{X} \times \mathcal{A}$ such that*

$$\sum_{t=1}^{\infty} \gamma_t(x, a, X_t) = \infty, \quad \sum_{t=1}^{\infty} \gamma_t^2(x, a, X_t) < \infty \quad (17)$$

for all $(x, a) \in \mathcal{X} \times A$ $\hat{P}_{x_0, \mathbf{a}}$ - almost surely.

(i) *Let the ambiguity set be given by $\mathcal{P}(x, a) = \mathcal{P}_1^{(q, \varepsilon)}(x, a)$*

for all $(x, a) \in \mathcal{X} \times A$ for some $\varepsilon > 0$ and $q \in \mathbb{N}$, and consider² $c_1 : \mathcal{X} \times \mathcal{X} \ni (x, y) \mapsto \|x - y\|^q$. Then, we have for all $(x, a) \in \mathcal{X} \times A$ that

$$\lim_{t \rightarrow \infty} Q_t(x, a) = Q^*(x, a) \quad \hat{P}_{x_0, \mathbf{a}} - \text{almost surely.}^3$$

(ii) Let $\mathcal{X} = T^h$ for some $h \in \mathbb{N} \cap [2, \infty)$ and $T \subset \mathbb{R}^D$ finite for some $D \in \mathbb{N}$, let the ambiguity set be given by $\mathcal{P}(x, a) = \mathcal{P}_2^{(q, \varepsilon)}(x, a)$ for all $(x, a) \in \mathcal{X} \times A$ for some $\varepsilon > 0$ and $q \in \mathbb{N}$, and consider⁴ $c_2 : \mathcal{X} \times \mathcal{X} \ni (x, y) \mapsto \infty \cdot \mathbf{1}_{\{(x_1, \dots, x_{h-1}) \neq (y_1, \dots, y_{h-1})\}}(x, y) + \|x_h - y_h\|^q$, where $(x, y) = ((x_1, \dots, x_h), (y_1, \dots, y_h))$. Then, we have for all $(x, a) \in \mathcal{X} \times A$ that

$$\lim_{t \rightarrow \infty} Q_t(x, a) = Q^*(x, a) \quad \hat{P}_{x_0, \mathbf{a}} - \text{almost surely.}$$

Remark 8 Note that condition (17) can be ensured by considering a sequence of learning rates $(\tilde{\gamma}_t)_{t \in \mathbb{N}_0} \subseteq [0, 1]$ satisfying

$$\sum_{t=0}^{\infty} \tilde{\gamma}_t = \infty, \quad \sum_{t=0}^{\infty} \tilde{\gamma}_t^2 < \infty, \quad (18)$$

and $(X_t)_{t \in \mathbb{N}_0}$ is a (positive) recurrent irreducible Markov decision process under $\hat{P}_{x_0, \mathbf{a}}$.

Remark 9 Note that in the non-robust case it has been empirically shown that an efficient choice for $\mathbf{a} \in \mathcal{A}$ when applying Q -learning is given by the so called $\tilde{\varepsilon}$ -greedy policy, see e.g. [15, Chapter 9], [33], or [51]. The $\tilde{\varepsilon}$ -greedy policy $\mathbf{a} := (a_1, a_2, \dots) \in \mathcal{A}$ is, for $\tilde{\varepsilon} > 0$, $t \in \mathbb{N}_0$, defined by

$$\mathcal{X} \ni x \mapsto a_t(x) := \begin{cases} \operatorname{argmax}_{b \in B} Q_t(x, b) & \text{prob. } 1 - \tilde{\varepsilon}, \\ a \sim \mathcal{U}(A) & \text{prob. } \tilde{\varepsilon}, \end{cases}$$

where $a \sim \mathcal{U}(A)$ means that a random action a is chosen uniformly at random from the finite set A . A popular modification of the $\tilde{\varepsilon}$ -greedy policy is to start with a relatively large $\tilde{\varepsilon}$ and to decrease the value of $\tilde{\varepsilon}$ over time, see, e.g., [33].

Remark 10 Note that from the optimal Q -value function one can infer $\mathcal{X} \ni x \mapsto a_{\text{loc}}^*(x) := \operatorname{argmax}_{a \in A} Q^*(x, a)$ and $\mathbf{a}^* := (a_{\text{loc}}^*(X_0), a_{\text{loc}}^*(X_1), \dots) \in \mathcal{A}$ which solves the robust stochastic optimal control problem (3), compare Proposition 4 and [37, Theorem 2.7]. Analogously,

² The function c_1 is used to determine the λc -transform in the algorithm, see (15) and (16).

³ For the definition of the probability measure $\hat{P}_{x_0, \mathbf{a}}$ we refer to (5).

⁴ The function c_2 is used to determine the λc -transform in the algorithm, see (15) and (16).

by considering $\mathcal{X} \ni x \mapsto \operatorname{argmax}_{a \in A} Q_t(x, a)$ for a sufficiently large $t \in \mathbb{N}$, we can derive an approximation of the optimal action.

The following result based on [36] shows that whenever an agent possesses a good enough guess about the true (but to her unknown) probability kernel $P^{\text{true}}(x, a)$ so that it is contained in the ambiguity set, one can bound the difference of the values of the robust and non-robust Markov decision problems. This is important since $\lim_{t \rightarrow \infty} Q_t(x, a) = Q^*(x, a) \hat{P}_{x_0, \mathbf{a}} - \text{a.s.}$ and $\sup_{a \in A} Q^*(x, a) = V(x)$, hence the following result also provides an upper bound on the sub-optimality of the performance of our robust Q -learning algorithm. We see that it can be controlled to be arbitrarily small when $\varepsilon \rightarrow 0$, as long as the agent possesses a good enough guess for $P^{\text{true}}(x, a)$ as discussed above. Note that compared to [36], no regularity assumptions on the map $(x, a) \mapsto P^{\text{true}}(x, a)$ nor on the reward function are necessary due to the finiteness of both the state and action space.

Proposition 11 Let $\varepsilon > 0$, $q \in \mathbb{N}$, and let

$$\mathcal{X} \ni x \mapsto V^{\text{true}}(x) := \sup_{\mathbf{a} \in \mathcal{A}} \left(\mathbb{E}_{P_{x, \mathbf{a}}^{\text{true}}} \left[\sum_{t=0}^{\infty} \alpha^t r(X_t, a_t, X_{t+1}) \right] \right), \quad (19)$$

with

$$P_{x, \mathbf{a}}^{\text{true}} := \delta_x \otimes P^{\text{true}} \otimes P^{\text{true}} \otimes P^{\text{true}} \otimes P^{\text{true}} \dots \in \mathcal{M}_1(\Omega),$$

where $\mathcal{X} \times A \ni (x, a) \mapsto P^{\text{true}}(x, a) \in \mathcal{P}_i^{(q, \varepsilon)}(x, a)$, $i \in \{1, 2\}$, for all $x \in \mathcal{X}$, $a \in A$. Moreover, assume that the discount factor satisfies (2.2) as well as $\alpha L_P < 1$, where

$$L_P := \sup_{\substack{(x, a), (x', a') \in \mathcal{X} \times A: \\ (x, a) \neq (x', a')}} \frac{W_q(P^{\text{true}}(x, a), P^{\text{true}}(x', a'))}{\|x - x'\| + \|a - a'\|}. \quad (20)$$

Then for any $x \in \mathcal{X}$ we have

$$0 \leq V^{\text{true}}(x) - V(x) \leq 2L_r \varepsilon (1 + \alpha) \sum_{i=0}^{\infty} \alpha^i \sum_{j=0}^i (L_P)^j < \infty, \quad (21)$$

where

$$L_r := \sup_{\substack{(x_0, a, x_1) \in \mathcal{X} \times A \times \mathcal{X}, \\ (x'_0, a', x'_1) \in \mathcal{X} \times A \times \mathcal{X}: \\ (x_0, a, x_1) \neq (x'_0, a', x'_1)}} \frac{|r(x_0, a, x_1) - r(x'_0, a', x'_1)|}{\|x_0 - x'_0\| + \|a - a'\| + \|x_1 - x'_1\|}. \quad (22)$$

4 Numerical Examples

In this section we provide three numerical examples that illustrate how the robust Q -learning Algorithm 1 can

be applied to specific problems. The examples highlight that a distributionally robust approach can outperform non-robust approaches whenever the assumed underlying distribution of the non-robust Markov Q -learning approach turns out to be misspecified during the testing period.

The selection of examples in this section is intended to give a small impression on the broad range of different applications of Q -learning algorithms for stochastic optimization problems. We refer to [7], [15], and [24] for an overview on several applications in finance and to [49] for a range of applications outside the world of finance.

4.1 On the Implementation

To apply the numerical method from Algorithm 1, we use for all of the following examples a discount factor of $\alpha = 0.45$, an $\tilde{\varepsilon}$ -greedy policy with $\tilde{\varepsilon} = 0.1$ (compare Remark 9), $q = 1$, and as a sequence of learning rates we use $\tilde{\gamma}_t = \frac{1}{1+t}$ for $t \in \mathbb{N}_0$. Moreover, we train all implementations with 50 000 iterations. The parameter λ_t from (15) is determined by maximizing the right-hand-side of (15) with a numerical solver relying on the Broyden—Fletcher—Goldfarb—Shanno (BFGS) algorithm ([10], [18], [20], [44]). Further details of the implementation can be found under <https://github.com/juliansester/Wasserstein-Q-learning>.

4.2 Examples

Example 12 (Coin Toss) We consider an agent playing the following game: At each time $t \in \mathbb{N}_0$ the agent observes the result of 10 coins that either show heads (encoded by 1) or tails (encoded by 0). The state X_t at time $t \in \mathbb{N}_0$ is then given by the sum of the heads observed in the 10 coins, i.e., we have $\mathcal{X} := \{0, \dots, 10\}$. At each time t the agent can bet whether the sum of the heads of the next throw strictly exceeds the previous sum (i.e. $X_{t+1} > X_t$), or whether it is strictly smaller (i.e. $X_{t+1} < X_t$).

If the agent is correct, she gets 1 \$, if the agent is wrong she has to pay 1 \$. The agent also has the possibility not to play. We model this by considering the reward function: $\mathcal{X} \times A \times \mathcal{X} \ni (x, a, x') \mapsto r(x, a, x') := a \mathbf{1}_{\{x < x'\}} - a \mathbf{1}_{\{x > x'\}} - |a| \mathbf{1}_{\{x = x'\}}$, where the possible actions are given by $A := \{-1, 0, 1\}$, where for example $a = 1$ corresponds to betting $X_{t+1} > X_t$. We then rely on Setting 1.) from Section 2.2 and consider as a reference measure a binomial distribution with $n = 10, p = 0.5$, i.e., $\mathcal{X} \times A \ni (x, a) \mapsto \hat{P}(x, a) := \text{Bin}(10, 0.5)$. We then define, according to Setting 1.) from Section 2.2, an ambiguity set, in dependence of $\varepsilon > 0$, by

$$\mathcal{P}(x, a) := \left\{ P \in \mathcal{M}_1(\mathcal{X}) \text{ s.t. } W_1\left(P, \text{Bin}(10, 0.5)\right) \leq \varepsilon \right\} \quad (23)$$

for every $(x, a) \in \mathcal{X} \times A$. Let $p \in [0, 1]$. Then, we denote the cumulative distribution function of a $B(10, p)$ -distributed random variable by $F_{10,p}$. Then we compute for the 1-Wasserstein distance that

$$\begin{aligned} W_1\left(\text{Bin}(10, 0.5), \text{Bin}(10, p)\right) &= \int_{\mathbb{R}} |F_{10,0.5}(x) - F_{10,p}(x)| dx \\ &= \int_0^\infty F_{10,\min\{p, 0.5\}}(x) - F_{10,\max\{0.5, p\}}(x) dx \\ &= \int_0^\infty (1 - F_{10,\max\{0.5, p\}}(x)) dx \\ &\quad - \int_0^\infty (1 - F_{10,\min\{0.5, p\}}(x)) dx \\ &= 10 \cdot \max\{0.5, p\} - 10 \cdot \min\{0.5, p\} = 10 \cdot |0.5 - p|, \end{aligned} \quad (24)$$

where the first equality of (24) follows e.g. from [42, Equation (3.5)] and the second equality of (24) follows since $F_{10,\min\{0.5, p\}}(x) \geq F_{10,\max\{0.5, p\}}(x)$ for all $x \in \mathbb{R}$. This means that all binomial distributions $\text{Bin}(10, p)$ with $p \in [0.5 - \frac{\varepsilon}{10}, 0.5 + \frac{\varepsilon}{10}]$ are contained in the ambiguity set⁵. The calculation from (24) gives a good indication how choosing a different value of ε may influence the measures contained in the ambiguity set. We then train actions $\mathbf{a}^{\text{robust}, \varepsilon} = (a_t^{\text{robust}, \varepsilon})_{t \in \mathbb{N}_0} \in \mathcal{A}$ according to the robust Q -learning approach proposed in Algorithm 1 for different values of ε , compare also Remark 10. Additionally we train an action $\mathbf{a}^{\text{non-robust}} = (a_t^{\text{non-robust}})_{t \in \mathbb{N}_0} \in \mathcal{A}$ according to the classical non-robust Q -learning approach, see, e.g., [57], where we assume that the underlying process $(X_t)_{t \in \mathbb{N}_0}$ develops according to the reference measure \hat{P} . We obtain after applying Algorithm 1 the strategies depicted in Table 1. In particular, we see that in compari-

| X_t | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|----|----|----|----|----|
| $a_t^{\text{non-robust}}(X_t)$ | 1 | 1 | 1 | 1 | 1 | 0 | -1 | -1 | -1 | -1 | -1 |
| $a_t^{\text{robust}, \varepsilon=0.5}(X_t)$ | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | -1 | -1 | -1 |
| $a_t^{\text{robust}, \varepsilon=1}(X_t)$ | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | -1 | -1 |
| $a_t^{\text{robust}, \varepsilon=2}(X_t)$ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Table 1

The trained actions $a_t^{\text{robust}, \varepsilon=0.5}(X_t)$, $a_t^{\text{robust}, \varepsilon=1}(X_t)$, $a_t^{\text{robust}, \varepsilon=2}(X_t)$, and $a_t^{\text{non-robust}}(X_t)$ in dependence of the realized state X_t at time $t \in \mathbb{N}_0$.

son with the non-robust action $\mathbf{a}^{\text{non-robust}}$, the robust actions $\mathbf{a}^{\text{robust}, \varepsilon}$ behave more carefully where a larger value of ε corresponds to a more careful behavior, which can be clearly seen for $\varepsilon = 2$, in which case the agent decides not to play for every realization of the state.

Then, we test the profit of the resultant actions $\mathbf{a}^{\text{robust}, \varepsilon}$

⁵ We highlight that of course the ambiguity set not only contains binomial distributions.

and $\mathbf{a}^{\text{non-robust}}$ by playing 100 000 rounds of the game according to these actions. For simulating the 100 000 rounds we assume an underlying binomial distribution $P_{\text{true}} = \text{Bin}(10, p_{\text{true}})$ with a fix probability p_{true} for heads which we vary from 0.1 to 0.9. We depict the cumulated profits of the considered actions in Table 2. We observe

| p_{true} | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 | 0.6 | 0.7 | 0.8 | 0.9 |
|-----------------------------|----------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|----------|
| Non-Robust | -31386 | -18438 | -1567 | 22892 | 35082 | 22956 | -656 | -18374 | -31091 |
| Robust, $\varepsilon = 0.5$ | -24728 | 4554 | 16491 | 13323 | 9920 | 13170 | 16825 | 4451 | -24427 |
| Robust, $\varepsilon = 1$ | -8174 | 15201 | 11091 | 4387 | 2050 | 4373 | 11139 | 15276 | -7611 |
| Robust, $\varepsilon = 2$ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Table 2

Overall Profit of the game described in Example 12 in dependence of different trained strategies and of the probability distribution $P_{\text{true}} = \text{Bin}(10, p_{\text{true}})$ of the underlying process. The best performing strategy in each case is indicated with bold characters.

that if the probability for heads p_{true} is similar as probability for heads in the reference measure ($p = 0.5$), then the non-robust approach (w.r.t. $\hat{P}(x, a) := \text{Bin}(10, 0.5)$) outperforms the robust approaches. If however the model with which the non-robust action was trained was clearly misspecified then $\mathbf{a}^{\text{robust}, \varepsilon}$ outperforms $\mathbf{a}^{\text{non-robust}}$. More precisely, the larger the degree of misspecification the more favorable it becomes to choose a larger ε . This can be well explained by the choice of the ambiguity set that covers, according to (24), the more measures under which we test, the larger we choose ε .

This simple example showcases that if in practice one is uncertain about the correct law according to which the state process evolves and one faces the risk of misspecifying the probabilities, then it can be advantageous to rely on a distributionally robust approach, whereas the choice of the radius of the Wasserstein-ball is extremely important as it corresponds to the degree of misspecification one wants to be robust against.

Example 13 (Comparison with KL-Uncertainty)

We reconsider an example of a supply-chain model provided in [31, Section 4]. In this example we have for some $n \in \mathbb{N}$ the state space $\mathcal{X} = \{0, 1, \dots, n\}$ representing the possible goods in the inventory and the action space $A = \mathcal{X}$ representing the possible goods we can order. The reward function is defined as the negative of the costs that are composed of holding costs and fixed ordering costs depending on parameters $h, p, k \in \mathbb{R}$ and on the demand which is, for the reference measure, uniformly distributed on \mathcal{X} , see [31, Section 4] for more details.

In the setting described in [31, Section 4], the optimal non-robust strategy (w.r.t. the reference measure) given current number of goods x is $a_t^{\text{non-robust}}(x) = (8 - x) \cdot \mathbf{1}_{\{x \leq 2\}}$ while we compute for a Wasserstein-uncertainty parameter $\varepsilon = 1$ an optimal robust strategy $a_t^{\text{Wasserstein}}(x) = (8 - x) \cdot \mathbf{1}_{\{x \leq 1\}} + 5 \mathbf{1}_{\{x \in \{2, 3\}\}}$. The robust strategy computed in [31, Section 4] that takes uncertainty w.r.t. Kullback-Leibler distance in account is given by $a_t^{\text{KL}}(x) = (7 - x) \cdot \mathbf{1}_{\{x \leq 4\}}$.

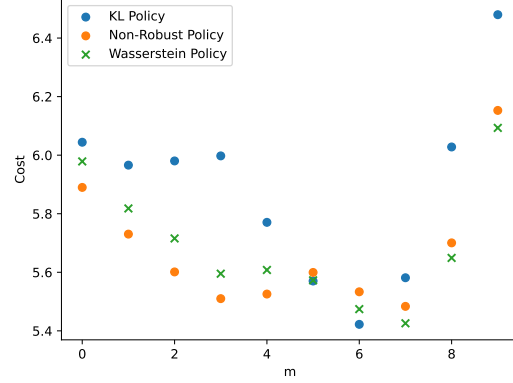


Fig. 1. Total Costs for $b = 1$ after 100000 iterations in the setting of Example 13, compare also [31, Figure 1].

As in [31, Figure 1], we evaluate the strategies on a distribution which does not coincide with the reference measure. To this end, we follow the example from [31, Section 4] and consider a perturbed uniform distribution depending on parameters m and b . With parameter $b = 1$ we compute after evaluation on 100 000 iterations the costs depicted in Figure 1, in dependence of the parameter m . The figure shows that for this particular example the Wasserstein approach leads for all values that are considered, except for $m \in \{5, 6\}$, to smaller costs than the approach provided in [31, Section 4]. Moreover, since the true distribution does not coincide with the reference distribution, the robust strategies can outperform the non-robust ones (defined w.r.t. the reference distribution).

Example 14 (Stock Movement Prediction)

We study the problem of predicting the movement of stock prices. We aim to predict whether in the next time step the return of an underlying stock is strongly negative (encoded by -2), slightly negative (encoded by -1), slightly positive (encoded by 1), or strongly positive (encoded by 2). Hence the space of the numerically encoded returns is given by $T := \{-2, -1, 1, 2\}$. We want to rely our prediction for the movement of the next return on the last $h = 5$ values. Hence, we consider, in line with the setting outlined in (7) $\mathcal{X} := T^h = \{-2, -1, 1, 2\}^5$. The space of actions is modelled by $A := \{-2, -1, 1, 2\} = T$ as the actions correspond to the future returns that are predicted. To construct a reference measure, we consider the historic evolution of the (numerically encoded) returns of the underlying stock. This time series is denoted by $(\mathcal{R}_j)_{j=1, \dots, N} \subset T^N$ for some $N \in \mathbb{N}$, see also Figure 2 for an illustration.

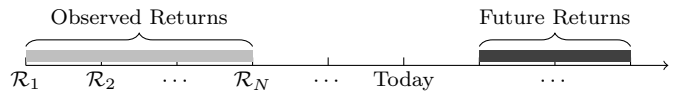


Fig. 2. Illustration of the time relation between the time series of observed returns $(\mathcal{R}_j)_{j=1, \dots, N}$ and the future returns which we want to predict.

We then define for some small⁶ $\gamma > 0$ the set-valued map $\mathcal{X} \times A \ni (x, a) \mapsto \widehat{\mathbf{P}}(x, a) := \sum_{i \in T} p_i(x) \cdot \delta_{\{i\}} \in \mathcal{M}_1(T)$ where for $x \in \mathcal{X}, i \in T$ we define

$$p_i(x) := \frac{\tilde{p}_i(x) + \frac{\gamma}{4}}{\gamma + \sum_{j \in T} \tilde{p}_j(x)} \in [0, 1], \quad (25)$$

as well as⁷

$$\tilde{p}_i(x) := \sum_{j=1}^{N-h+1} \mathbf{1}_{\{(\pi(x), i) = (\mathcal{R}_j, \dots, \mathcal{R}_{j+h-1})\}}. \quad (26)$$

This means the construction of $\widehat{\mathbf{P}}(x, a)$ relies, according to (26), on the relative frequency of the sequence $(\pi(x), i)$ in the time series of past realized returns $(\mathcal{R}_j)_{j=1, \dots, N}$. Equation (25) is then applied to convert the frequencies to probabilities. Then, as a reference measure we consider, as in (9), the set-valued map

$$\mathcal{X} \times A \ni (x, a) \mapsto \widehat{\mathbf{P}}(x, a) = \delta_{\pi(x)} \otimes \widehat{\mathbf{P}}(x, a) \in \mathcal{M}_1(\mathcal{X}). \quad (27)$$

Moreover, as a reward function we consider⁸ $\mathcal{X} \times A \times \mathcal{X} \ni (x, a, x') \mapsto r(x, a, x') := \mathbf{1}_{\{x'_h = a\}}$, i.e., we reward only correct predictions. We apply the setting described above to real data. To this end, we consider as series of realized returns $(\mathcal{R}_j)_{j=1, \dots, N}$ the daily returns of the stock of Apple in the time horizon from 1 January 2010 until 28 September 2018 and hence we take into account $N = 2200$ daily returns. To encode the observed market returns to values in T , we distinguish between small returns and large returns by saying that a daily return is strongly positive if it is larger than 0.01. Analogously a daily return is strongly negative if smaller than -0.01 . This leads to the distribution of returns as depicted in Table 3.

| Type of Encoded Return (Numerical Value) | Total Amount |
|--|--------------|
| Strongly Negative Returns (-2) | 404 |
| Slightly Negative Returns (-1) | 637 |
| Slightly Positive Returns (1) | 627 |
| Strongly Positive Returns (2) | 532 |

Table 3

The distribution of the numerically encoded daily returns of Apple between January 2010 and September 2018. The threshold to distinguish slightly positive (negative) returns from strongly positive returns is 0.01 (-0.01).

We then train a non-robust action $\mathbf{a}^{\text{non-robust}} = (a_t^{\text{non-robust}})_{t \in N_0} \in \mathcal{A}$ according to the classical non-robust Q -learning algorithm ([58]) as well as robust

actions $\mathbf{a}^{\text{robust}} = (a_t^{\text{robust}})_{t \in N_0} \in \mathcal{A}$ according to Algorithm 1 that takes into account an ambiguity set defined in (10) with $\varepsilon = 0.1$. Moreover, for comparison, we consider a trivial action $\mathbf{a}^{\text{trivial}} = (a_t^{\text{trivial}})_{t \in N_0} \in \mathcal{A}$ which always, independent of the state-action pair, predicts -1 since, according to Table 3, -1 is the most frequent appearing value in the time series $(\mathcal{R}_j)_{j=1, \dots, N}$.

We then evaluate the trained actions, in a small back-testing study, on realized daily returns of Apple that occurred after the end of the training period. To this end, we consider an evaluation period from 1 October 2018 until 26 February 2019 consisting of 100 daily returns that are distributed according to Table 4. We observe that in the

| Type of Encoded Return (Numerical Value) | Total Amount |
|--|--------------|
| Strongly Negative Returns (-2) | 29 |
| Slightly Negative Returns (-1) | 21 |
| Slightly Positive Returns (1) | 22 |
| Strongly Positive Returns (2) | 28 |

Table 4

The distribution of the numerically encoded daily returns of Apple between 1 October 2018 and 26 February 2019.

evaluation period, in contrast to the training period, the large negative returns impose the largest class of appearing returns. Overall the distribution is significantly different from the distribution of the classes on the training data. We illustrate in Table 5 the results of predictions of the actions $\mathbf{a}^{\text{trivial}}$, $\mathbf{a}^{\text{non-robust}}$, $\mathbf{a}^{\text{robust}}$ evaluated in the evaluation period, and we observe that indeed the robust action $\mathbf{a}^{\text{robust}}$ outperforms the other two actions clearly in this period where the distribution of returns significantly differs from the distributions of the returns on which the actions were trained. This showcases again that if there

| Action | Share of Correct Predictions |
|----------------------------------|------------------------------|
| $\mathbf{a}^{\text{non-robust}}$ | 23.40% |
| $\mathbf{a}^{\text{robust}}$ | 28.72% |
| $\mathbf{a}^{\text{trivial}}$ | 21.27% |

Table 5

The proportion of correct stock movement predictions in the evaluation period between 1 October 2018 and 10 January 2019

is the risk that the underlying distribution on which the actions were trained turns out to be misspecified, then it can be advantageous to use a robust approach.

Acknowledgements

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.

⁶ Note that γ is only introduced to avoid a division by 0.

⁷ Note that π is defined in (8).

⁸ Here $x' := (x'_1, \dots, x'_h) \in \mathcal{X}$, and hence, x'_h denotes the last component of x' .

References

- [1] Asma Al-Tamimi, Frank L Lewis, and Murad Abu-Khalaf. Model-free Q-learning designs for linear discrete-time zero-sum games with application to h-infinity control. *Automatica*, 43(3):473–481, 2007.
- [2] 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.
- [3] Andrea Angiuli, Jean-Pierre Fouque, and Mathieu Lauriere. Reinforcement learning for mean field games, with applications to economics. *arXiv preprint arXiv:2106.13755*, 2021.
- [4] Daniel Bartl, Samuel Drapeau, and Ludovic Tangpi. Computational aspects of robust optimized certainty equivalents and option pricing. *Mathematical Finance*, 30(1):287–309, 2020.
- [5] Nicole Bäuerle and Alexander Glauner. Distributionally robust Markov decision processes and their connection to risk measures. *Mathematics of Operations Research*, 2021.
- [6] 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.
- [7] Nicole Bäuerle and Ulrich Rieder. *Markov decision processes with applications to finance*. Springer Science & Business Media, 2011.
- [8] Bahram Behzadian, Marek Petrik, and Chin Pang Ho. Fast algorithms for l_∞ -constrained s-rectangular robust mdps. *Advances in Neural Information Processing Systems*, 34:25982–25992, 2021.
- [9] Jose Blanchet and Karthyek Murthy. Quantifying distributional model risk via optimal transport. *Mathematics of Operations Research*, 44(2):565–600, 2019.
- [10] Charles George Broyden. The convergence of a class of double-rank minimization algorithms 1. general considerations. *IMA Journal of Applied Mathematics*, 6(1):76–90, 1970.
- [11] 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.
- [12] Arthur Charpentier, Romuald Elie, and Carl Remlinger. Reinforcement learning in economics and finance. *Computational Economics*, pages 1–38, 2021.
- [13] Zhi Chen, Pengqian Yu, and William B Haskell. Distributionally robust optimization for sequential decision-making. *Optimization*, 68(12):2397–2426, 2019.
- [14] Jeremy Coulson, John Lygeros, and Florian Dörfler. Regularized and distributionally robust data-enabled predictive control. In *2019 IEEE 58th Conference on Decision and Control (CDC)*, pages 2696–2701. IEEE, 2019.
- [15] Matthew F Dixon, Igor Halperin, and Paul Bilokon. *Machine learning in Finance*, volume 1170. Springer, 2020.
- [16] Aryeh Dvoretzky. *On stochastic approximation*. University of California Press, 1956.
- [17] Laurent El Ghaoui and Arnab Nilim. Robust solutions to Markov decision problems with uncertain transition matrices. *Operations Research*, 53(5):780–798, 2005.
- [18] Roger Fletcher. A new approach to variable metric algorithms. *The computer journal*, 13(3):317–322, 1970.
- [19] Rui Gao and Anton Kleywegt. Distributionally robust stochastic optimization with Wasserstein distance. *Mathematics of Operations Research*, 2022.
- [20] Donald Goldfarb. A family of variable metric updates derived by variational means. *mathematics of computing*. 1970.
- [21] Vineet Goyal and Julien Grand-Clement. Robust markov decision processes: Beyond rectangularity. *Mathematics of Operations Research*, 48(1):203–226, 2023.
- [22] Yi Guo, Kyri Baker, Emiliano Dall’Anese, Zechun Hu, and Tyler Holt Summers. Data-based distributionally robust stochastic optimal power flow—part i: Methodologies. *IEEE Transactions on Power Systems*, 34(2):1483–1492, 2018.
- [23] Astghik Hakobyan and Insoon Yang. Distributionally robust differential dynamic programming with wasserstein distance. *IEEE Control Systems Letters*, 2023.
- [24] Ben Hambly, Renyuan Xu, and Huining Yang. Recent advances in reinforcement learning in finance. *arXiv preprint arXiv:2112.04553*, 2021.
- [25] Jian Huang, Junyi Chai, and Stella Cho. Deep learning in finance and banking: A literature review and classification. *Frontiers of Business Research in China*, 14(1):1–24, 2020.
- [26] Garud N Iyengar. Robust dynamic programming. *Mathematics of Operations Research*, 30(2):257–280, 2005.
- [27] Tommi Jaakkola, Michael I Jordan, and Satinder P Singh. On the convergence of stochastic iterative dynamic programming algorithms. *Neural computation*, 6(6):1185–1201, 1994.
- [28] Gyeun Jeong and Ha Young Kim. Improving financial trading decisions using deep Q-learning: Predicting the number of shares, action strategies, and transfer learning. *Expert Systems with Applications*, 117:125–138, 2019.
- [29] Petter N Kolm and Gordon Ritter. Modern perspectives on reinforcement learning in finance. *Modern Perspectives on Reinforcement Learning in Finance (September 6, 2019)*. *The Journal of Machine Learning in Finance*, 1(1), 2020.
- [30] Mengmeng Li, Tobias Sutter, and Daniel Kuhn. Policy gradient algorithms for robust mdps with non-rectangular uncertainty sets. *arXiv preprint arXiv:2305.19004*, 2023.
- [31] Zijian Liu, Qinxun Bai, Jose Blanchet, Perry Dong, Wei Xu, Zhengqing Zhou, and Zhengyuan Zhou. Distributionally robust Q-learning. In *International Conference on Machine Learning*, pages 13623–13643. PMLR, 2022.
- [32] Shie Mannor, Ofir Mebel, and Huan Xu. Robust MDPs with k-rectangular uncertainty. *Mathematics of Operations Research*, 41(4):1484–1509, 2016.
- [33] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- [34] Peyman Mohajerin Esfahani and Daniel Kuhn. Data-driven distributionally robust optimization using the Wasserstein metric: Performance guarantees and tractable reformulations. *Mathematical Programming*, 171(1):115–166, 2018.
- [35] Mohammad Bagher Naghibi-Sistani, MR Akbarzadeh-Tootoonchi, MH Javidi-Dashte Bayaz, and Habib Rajabi-Mashhadi. Application of Q-learning with temperature variation for bidding strategies in market based power systems. *Energy Conversion and Management*, 47(11-12):1529–1538, 2006.
- [36] Ariel Neufeld and Julian Sester. Bounding the difference between the values of robust and non-robust markov decision problems. *arXiv preprint arXiv:2308.05520*, 2023.
- [37] Ariel Neufeld, Julian Sester, and Mario Šikić. Markov decision processes under model uncertainty. *Mathematical Finance*, 33(3):618–665, 2023.

- [38] Brian Ning, Franco Ho Ting Lin, and Sebastian Jaimungal. Double deep Q-learning for optimal execution. *Applied Mathematical Finance*, 28(4):361–380, 2021.
- [39] Kishan Panaganti and Dileep Kalathil. Sample complexity of robust reinforcement learning with a generative model. In *International Conference on Artificial Intelligence and Statistics*, pages 9582–9602. PMLR, 2022.
- [40] Tiberiu Popoviciu. Sur les équations algébriques ayant toutes leurs racines réelles. *Mathematica*, 9(129-145):20, 1935.
- [41] Herbert Robbins and Sutton Monro. A stochastic approximation method. *The annals of mathematical statistics*, pages 400–407, 1951.
- [42] Ludger Rüschendorf. Monge-Kantorovich transportation problem and optimal couplings. *Jahresbericht der DMV*, 3:113–137, 2007.
- [43] Andrzej Ruszczyński and Alexander Shapiro. Conditional risk mappings. *Mathematics of operations research*, 31(3):544–561, 2006.
- [44] David F Shanno. Conditioning of quasi-Newton methods for function minimization. *Mathematics of computation*, 24(111):647–656, 1970.
- [45] Alexander Shapiro. Rectangular sets of probability measures. *Operations Research*, 64(2):528–541, 2016.
- [46] Rajesh Sharma, Madhu Gupta, and Girish Kapoor. Some better bounds on the variance with applications. *Journal of Mathematical Inequalities*, 4(3):355–363, 2010.
- [47] Nian Si, Fan Zhang, Zhengyuan Zhou, and Jose Blanchet. Distributional robust batch contextual bandits. *arXiv preprint arXiv:2006.05630*, 2020.
- [48] Nian Si, Fan Zhang, Zhengyuan Zhou, and Jose Blanchet. Distributionally robust policy evaluation and learning in offline contextual bandits. In *International Conference on Machine Learning*, pages 8884–8894. PMLR, 2020.
- [49] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- [50] Andrea Tirinzoni, Marek Petrik, Xiangli Chen, and Brian Ziebart. Policy-conditioned uncertainty sets for robust markov decision processes. *Advances in neural information processing systems*, 31, 2018.
- [51] Michel Tokic and Günther Palm. Value-difference based exploration: adaptive control between epsilon-greedy and softmax. In *Annual conference on artificial intelligence*, pages 335–346. Springer, 2011.
- [52] Kerem Ugurlu. Robust optimal control using conditional risk mappings in infinite horizon. *Journal of Computational and Applied Mathematics*, 344:275–287, 2018.
- [53] Bart PG Van Parys, Daniel Kuhn, Paul J Goulart, and Manfred Morari. Distributionally robust control of constrained stochastic systems. *IEEE Transactions on Automatic Control*, 61(2):430–442, 2015.
- [54] Cédric Villani. *Optimal transport: old and new*, volume 338. Springer, 2008.
- [55] Yue Wang and Shaofeng Zou. Online robust reinforcement learning with model uncertainty. *Advances in Neural Information Processing Systems*, 34:7193–7206, 2021.
- [56] Yue Wang and Shaofeng Zou. Policy gradient method for robust reinforcement learning. *arXiv preprint arXiv:2205.07344*, 2022.
- [57] Christopher JCH Watkins. Learning from delayed rewards. *Ph. D. thesis, King’s College, University of Cambridge*, 1989.
- [58] Christopher JCH Watkins and Peter Dayan. Q-learning. *Machine learning*, 8(3-4):279–292, 1992.
- [59] Wolfram Wiesemann, Daniel Kuhn, and Berç Rustem. Robust Markov decision processes. *Mathematics of Operations Research*, 38(1):153–183, 2013.
- [60] Wolfram Wiesemann, Daniel Kuhn, and Melvyn Sim. Distributionally robust convex optimization. *Operations research*, 62(6):1358–1376, 2014.
- [61] Huan Xu and Shie Mannor. Distributionally robust Markov decision processes. *Mathematics of Operations Research*, 37(2):288–300, 2012.
- [62] Insoon Yang. A convex optimization approach to distributionally robust markov decision processes with wasserstein distance. *IEEE control systems letters*, 1(1):164–169, 2017.
- [63] Insoon Yang. Wasserstein distributionally robust stochastic control: A data-driven approach. *IEEE Transactions on Automatic Control*, 66(8):3863–3870, 2020.
- [64] Wenhao Yang, Liangyu Zhang, and Zhihua Zhang. Towards theoretical understandings of robust markov decision processes: Sample complexity and asymptotics. *arXiv preprint arXiv:2105.03863*, 2021.
- [65] Chaoyue Zhao and Yongpei Guan. Data-driven risk-averse stochastic optimization with Wasserstein metric. *Operations Research Letters*, 46(2):262–267, 2018.
- [66] Zhengqing Zhou, Zhengyuan Zhou, Qinxun Bai, Linhai Qiu, Jose Blanchet, and Peter Glynn. Finite-sample regret bound for distributionally robust offline tabular reinforcement learning. In *International Conference on Artificial Intelligence and Statistics*, pages 3331–3339. PMLR, 2021.

A Auxiliary Results and Proofs

In Section A.1 we provide several useful results which then allow in Section A.2 to prove the main result from Section 3.

A.1 Auxiliary Results

To establish convergence of our Q -learning algorithm that was presented in Section 3 we will make use of the following auxiliary result from stochastic approximation theory which was developed to prove the convergence of the classical Q -learning algorithm. We refer to [27, Section 3] for a discussion of the advantage of the following result compared to classical results from stochastic approximation such as, e.g., [16]. Note that for any $f : \mathcal{X} \times A \rightarrow \mathbb{R}$, we write

$$\|f\|_{\infty} := \sup_{x \in \mathcal{X}} \sup_{a \in A} |f(x, a)|. \quad (\text{A.1})$$

Lemma 15 ([27], Theorem 1) *Let $P_0 \in \mathcal{M}_1(\Omega)$ be a probability measure on (Ω, \mathcal{F}) , and consider a family of stochastic processes $(\gamma_t(x, a), F_t(x, a), \Delta_t(x, a))_{t \in \mathbb{N}_0}$, $(x, a) \in \mathcal{X} \times A$, satisfying for all $t \in \mathbb{N}_0$*

$$\Delta_{t+1}(x, a) = (1 - \gamma_t(x, a)) \Delta_t(x, a) + \gamma_t(x, a) F_t(x, a)$$

P_0 -almost surely for all $(x, a) \in \mathcal{X} \times A$. Let $(\mathcal{G}_t)_{t \in \mathbb{N}_0} \subseteq \mathcal{F}$ be a sequence of increasing σ -algebras such that for all

$(x, a) \in \mathcal{X} \times A$ the random variables $\Delta_0(x, a)$ and $\gamma_0(x, a)$ are \mathcal{G}_0 -measurable and such that $\Delta_t(x, a)$, $\gamma_t(x, a)$, and $F_{t-1}(x, a)$ are \mathcal{G}_t -measurable for all $t \in \mathbb{N}$. Further assume that the following conditions hold.

- (i) $0 \leq \gamma_t(x, a) \leq 1$, $\sum_{t=0}^{\infty} \gamma_t(x, a) = \infty$, $\sum_{t=0}^{\infty} \gamma_t^2(x, a) < \infty$ P_0 -almost surely for all $(x, a) \in \mathcal{X} \times A$, $t \in \mathbb{N}_0$.
- (ii) There exists $\delta \in (0, 1)$ such that $\|E_{P_0}[F_t(\cdot, \cdot) | \mathcal{G}_t]\|_{\infty} \leq \delta \|\Delta_t\|_{\infty}$ P_0 -almost surely for all $t \in \mathbb{N}_0$.
- (iii) There exists $C > 0$ such that $\|\text{Var}_{P_0}(F_t(\cdot, \cdot) | \mathcal{G}_t)\|_{\infty} \leq C(1 + \|\Delta_t\|_{\infty})^2$ P_0 -almost surely for all $t \in \mathbb{N}_0$.

Then, $\lim_{t \rightarrow \infty} \Delta_t(x, a) = 0$ P_0 -almost surely for all $(x, a) \in \mathcal{X} \times A$.

Next, as the following proposition shows, the λc -transform allows to compute worst case expectations with respect to probability measures contained in the Wasserstein-ball by computing its dual which solely depends on the center of the Wasserstein-ball.

Proposition 16 ([4], Theorem 2.4) Let $f : \mathcal{X} \rightarrow \mathbb{R}$, let $\varepsilon > 0$ and $q \in \mathbb{N}$, let $\mathcal{P}_1^{(q, \varepsilon)}$, $\mathcal{P}_2^{(q, \varepsilon)}$ be the ambiguity sets of probability measures defined in (6) and (10), and let $c_1, c_2 : \mathcal{X} \times \mathcal{X} \rightarrow [0, \infty]$ be defined as in Theorem 7.

(i) Then, we have for every $(x, a) \in \mathcal{X} \times A$ that

$$\inf_{P \in \mathcal{P}_1^{(q, \varepsilon)}(x, a)} E_P[f] = \sup_{\lambda \geq 0} \left(E_{\hat{P}(x, a)}[(-f)^{\lambda c_1}] - \varepsilon^q \lambda \right).$$

- (ii) In addition, let $\mathcal{X} = T^h$ for some $h \in \mathbb{N} \cap [2, \infty)$, and $T \subset \mathbb{R}^D$ finite for some $D \in \mathbb{N}$. Moreover, assume that there exists some probability kernel $\mathcal{X} \times A \ni (x, a) \mapsto \hat{P}(x, a) \in \mathcal{M}_1(T)$ such for all $(x, a) \in \mathcal{X} \times A$ we have $\hat{P}(x, a) = \delta_{\pi(x)} \otimes \hat{P}(x, a)$. Then, we have for every $(x, a) \in \mathcal{X} \times A$ that $\inf_{P \in \mathcal{P}_2^{(q, \varepsilon)}(x, a)} E_P[f] = \sup_{\lambda \geq 0} \left(E_{\delta_{\pi(x)} \otimes \hat{P}(x, a)}[(-f)^{\lambda c_2}] - \varepsilon^q \lambda \right)$.

Proof of Proposition 16

In case (i), the assertion follows by an application of the duality result from [4, Theorem 2.4] (with the specifications $d_{c_1}(\cdot, \cdot) := W_q(\cdot, \cdot)^q$, $\varphi(\cdot) := \infty \mathbf{1}_{(\varepsilon, \infty]}(\cdot)$ in the notation of [4, Theorem 2.4], see also [4, Example 2.5]). More precisely, by [4, Theorem 2.4], [4, Example 2.5] and by the definition of $\mathcal{P}_1^{(q, \varepsilon)}$ we have for all $(x, a) \in \mathcal{X} \times A$ that

$$\begin{aligned} \inf_{P \in \mathcal{P}_1^{(q, \varepsilon)}(x, a)} E_P[f] &= - \sup_{\{P \in \mathcal{M}_1(\mathcal{X}) \mid W_q(P, \hat{P}(x, a)) \leq \varepsilon\}} E_P[-f] \\ &= - \left(\inf_{\lambda \geq 0} \left\{ E_{\hat{P}(x, a)}[(-f)^{\lambda c_1}] + \varepsilon^q \lambda \right\} \right) \\ &= \sup_{\lambda \geq 0} \left\{ E_{\hat{P}(x, a)}[(-f)^{\lambda c_1}] - \varepsilon^q \lambda \right\}. \end{aligned}$$

To show (ii), we observe that in the notation of [4], we have for $P, P' \in \mathcal{M}_1(\mathcal{X})$ that

$$\begin{aligned} d_{c_2}(P, P') &:= \inf_{Q \in \Pi(P, P')} \int_{\mathcal{X} \times \mathcal{X}} c_2(x, y) dQ(x, y) \\ &= \inf_{Q \in \Pi(P, P')} \int_{\mathcal{X} \times \mathcal{X}} (\infty \cdot \mathbf{1}_{\{(x_1, \dots, x_{h-1}) \neq (y_1, \dots, y_{h-1})\}} \\ &\quad + \|x_h - y_h\|^q) dQ((x_1, \dots, x_h), (y_1, \dots, y_h)). \end{aligned}$$

Hence, we have $d_{c_2}(P, \delta_{\pi(x)} \otimes \hat{P}(x, a)) \leq \varepsilon^q$ if and only if $P = \delta_{\pi(x)} \otimes \hat{P}$ for some $\hat{P} \in \mathcal{M}_1(T)$ with $W_q(\tilde{P}, \hat{P}(x, a)) \leq \varepsilon$. Moreover, we see that c_2 is indeed a cost function in the sense of [4]. This implies by [4, Theorem 2.4] and by the definition of $\mathcal{P}_2^{(q, \varepsilon)}$ that for all $(x, a) \in \mathcal{X} \times A$ we have

$$\begin{aligned} \inf_{P \in \mathcal{P}_2^{(q, \varepsilon)}(x, a)} E_P[f] &= - \sup_{\{P \in \mathcal{M}_1(\mathcal{X}) \mid d_{c_2}(P, \delta_{\pi(x)} \otimes \hat{P}(x, a)) \leq \varepsilon^q\}} E_P[-f] \\ &= - \left(\inf_{\lambda \geq 0} \left\{ E_{\delta_{\pi(x)} \otimes \hat{P}(x, a)}[(-f)^{\lambda c_2}] + \varepsilon^q \lambda \right\} \right) \\ &= \sup_{\lambda \geq 0} \left\{ E_{\delta_{\pi(x)} \otimes \hat{P}(x, a)}[(-f)^{\lambda c_2}] - \varepsilon^q \lambda \right\}. \quad \square \end{aligned}$$

Next, consider the operator H which is defined for any $q : \mathcal{X} \times A \rightarrow \mathbb{R}$ by

$$\mathcal{X} \times A \ni (x, a) \mapsto (Hq)(x, a) := \inf_{P \in \mathcal{P}(x, a)} E_P[r(x, a, X_1) + \alpha \max_{b \in A} q(X_1, b)]. \quad (\text{A.2})$$

We derive for H the following form of the Bellman-equation.

Lemma 17 Assume that (2) holds and let the ambiguity set \mathcal{P} be either $\mathcal{P}_1^{(q, \varepsilon)}$ or $\mathcal{P}_2^{(q, \varepsilon)}$, defined in (6) and (10). Then the following equation holds true for the optimal Q -value function defined in (12): $HQ^*(x, a) = Q^*(x, a)$ for all $(x, a) \in \mathcal{X} \times A$.

Proof of Lemma 17 This follows directly by definition of Q^* and by Proposition 4. Indeed, let $(x, a) \in \mathcal{X} \times A$. Then, we have

$$\begin{aligned} (HQ^*)(x, a) &= \inf_{P \in \mathcal{P}(x, a)} E_P[r(x, a, X_1) + \alpha \sup_{b \in A} Q^*(X_1, b)] \\ &= \inf_{P \in \mathcal{P}(x, a)} E_P[r(x, a, X_1) + \alpha V(X_1)] \\ &= Q^*(x, a). \quad \square \end{aligned}$$

Moreover, we observe that the operator H is a contraction with respect to the supremum norm defined in

(A.1).

Lemma 18 For any maps $q_i : \mathcal{X} \times A \rightarrow \mathbb{R}$, $i = 1, 2$, we have $\|Hq_1 - Hq_2\|_\infty \leq \alpha \|q_1 - q_2\|_\infty$.

Proof of Lemma 18 Consider any maps $q_i : \mathcal{X} \times A \rightarrow \mathbb{R}$, $i = 1, 2$. Then, we have for all $(x, a) \in \mathcal{X} \times A$ that

$$\begin{aligned} & |(Hq_1)(x, a) - (Hq_2)(x, a)| \\ &= \left| \inf_{P \in \mathcal{P}(x, a)} \mathbb{E}_P [r(x, a, X_1) + \alpha \sup_{b \in A} q_1(X_1, b)] \right. \\ &\quad \left. - \inf_{P \in \mathcal{P}(x, a)} \mathbb{E}_P [r(x, a, X_1) + \alpha \sup_{b \in A} q_2(X_1, b)] \right| \\ &\leq \sup_{P \in \mathcal{P}(x, a)} \left| \mathbb{E}_P \left[r(x, a, X_1) + \alpha \sup_{b \in A} q_2(X_1, b) \right. \right. \\ &\quad \left. \left. - r(x, a, X_1) - \alpha \sup_{b \in A} q_1(X_1, b) \right] \right| \\ &\leq \alpha \sup_{P \in \mathcal{P}(x, a)} \mathbb{E}_P \left[\sup_{b \in A} |q_2(X_1, b) - q_1(X_1, b)| \right] \\ &\leq \alpha \sup_{(y, b) \in \mathcal{X} \times A} |q_2(y, b) - q_1(y, b)| = \alpha \|q_1 - q_2\|_\infty, \end{aligned}$$

which implies the assertion by taking the supremum with respect to the arguments of $Hq_1(\cdot, \cdot) - Hq_2(\cdot, \cdot)$. \square

A.2 Proofs

In this section we provide the proofs of the results from Section 2 and Section 3.

Proof of Proposition 4

The first equality $\sup_{a \in A} Q^*(x, a) = \mathcal{TV}(x)$ follows by definition of \mathcal{TV} . For the second equality $\mathcal{TV}(x) = V(x)$ we want to check that [37, Assumption 2.2] and [37, Assumption 2.4] hold true to be able to apply [37, Theorem 3.1]. [37, Assumption 2.2] is fulfilled (for $p = 0$ and $C_P = 1$ in the notation of [37, Assumption 2.2]) according to [37, Proposition 3.1] in the case $\mathcal{P} = \mathcal{P}_1^{(q, \varepsilon)}$, and according to [37, Proposition 3.3] in the case $\mathcal{P} = \mathcal{P}_2^{(q, \varepsilon)}$. To verify [37, Assumption 2.4] (i), note that $\mathcal{X} \times A \times \mathcal{X} \ni (x_0, a, x_1) \mapsto r(x_0, a, x_1)$ is continuous since \mathcal{X} and A are finite (endowed with the discrete topology). To show [37, Assumption 2.4] (ii) note that for all $x_0, x'_0, x_1 \in \mathcal{X}$ and $a, a' \in A$ we have $|r(x_0, a, x_1) - r(x'_0, a', x_1)| \leq L \cdot (\|x_0 - x'_0\| + \|a - a'\|)$. with $L := \left(\max_{\substack{y_0, y'_0 \in \mathcal{X}, b, b' \in A \\ (y_0, b) \neq (y'_0, b')}} \frac{|r(y_0, b, x_1) - r(y'_0, b', x_1)|}{\|y_0 - y'_0\| + \|b - b'\|} \right)$. Similarly, to show [37, Assumption 2.4] (iii), we observe that for all $x_0, x_1 \in \mathcal{X}$ and all $a \in A$ we have

$$|r(x_0, a, x_1)| \leq \max_{y_0, y_1 \in \mathcal{X}, b \in A} |r(y_0, b, y_1)|,$$

i.e., in the notation of [37, Assumption 2.4] we have $C_r := \max\{1, \max_{y_0, y_1 \in \mathcal{X}, b \in A} |r(y_0, b, y_1)|\}$. To verify

[37, Assumption 2.4] (iv) we see that, since $p = 0$, we can choose $C_P := 1$ in the notation of [37, Assumption 2.2] (ii) and hence with (2) we get $0 < \alpha < 1 = \frac{1}{C_P}$ as required. Hence, the result follows from [37, Theorem 3.1]. \square

Proof of Lemma 6 For any $\lambda \geq 0$ we have by definition of the λc -transform

$$\begin{aligned} & \mathbb{E}_{\hat{P}(x, a)} \left[-(-f_{t, (x, a)})^{\lambda c}(X_1) - \varepsilon^q \lambda \right] \\ &= \mathbb{E}_{\hat{P}(x, a)} \left[-\max_{y \in \mathcal{X}} \{-f_{t, (x, a)}(y) - \lambda c(X_1, y)\} - \varepsilon^q \lambda \right]. \end{aligned}$$

Therefore, since \mathcal{X} is finite, the map $[0, \infty) \ni \lambda \mapsto G(\lambda) := \mathbb{E}_{\hat{P}(x, a)} \left[-(-f_{t, (x, a)})^{\lambda c}(X_1) - \varepsilon^q \lambda \right]$ is continuous. Hence, the assertion of Lemma 6 follows once we have shown that $\lim_{\lambda \rightarrow \infty} G(\lambda) = -\infty$. To that end, note that as, by assumption, $\min_{y \in \mathcal{X}} c(x, y) = 0$ for all $x \in \mathcal{X}$, we have that

$$\begin{aligned} & \limsup_{\lambda \rightarrow \infty} G(\lambda) \\ &= \limsup_{\lambda \rightarrow \infty} \mathbb{E}_{\hat{P}(x, a)} \left[-\max_{y \in \mathcal{X}} \{-f_{t, (x, a)}(y) - \lambda c(X_1, y)\} - \varepsilon^q \lambda \right] \\ &\leq \limsup_{\lambda \rightarrow \infty} \mathbb{E}_{\hat{P}(x, a)} \left[\max_{z \in \mathcal{X}} f_{t, (x, a)}(z) - \max_{y \in \mathcal{X}} \{-\lambda c(X_1, y)\} - \varepsilon^q \lambda \right] \\ &= \limsup_{\lambda \rightarrow \infty} \left(\max_{z \in \mathcal{X}} f_{t, (x, a)}(z) + \lambda \mathbb{E}_{\hat{P}(x, a)} \left[\min_{y \in \mathcal{X}} c(X_1, y) \right] - \varepsilon^q \lambda \right) \\ &= \max_{z \in \mathcal{X}} f_{t, (x, a)}(z) + \limsup_{\lambda \rightarrow \infty} (-\varepsilon^q \lambda) = -\infty. \quad \square \end{aligned}$$

Proof of Theorem 7 Let $(x_0, \mathbf{a}) \in \mathcal{X} \times \mathcal{A}$.

Assume that either $\mathcal{P} = \mathcal{P}_1^{(q, \varepsilon)}$ and $c = c_1$, or $\mathcal{P} = \mathcal{P}_2^{(q, \varepsilon)}$ and $c = c_2$. Then, we show for all $(x, a) \in \mathcal{X} \times A$ $\lim_{t \rightarrow \infty} Q_t(x, a) = Q^*(x, a)$ $\hat{P}_{x_0, \mathbf{a}}$ -almost surely, which shows simultaneously both (i) and (ii). To that end, let $(x, a) \in \mathcal{X} \times A$ be fixed. Then we rearrange the terms in (16) and write

$$\begin{aligned} Q_{t+1}(x, a) &= (1 - \gamma_t(x, a, X_t)) Q_t(x, a) \\ &\quad + \gamma_t(x, a, X_t) \left(-(-f_{t, (x, a)})^{\lambda c}(X_{t+1}) - \varepsilon^q \lambda_t \right), \end{aligned} \tag{A.3}$$

where λ_t is as defined in (15), see also Lemma 6 for its existence. Note that by construction $Q_t(x, a) \in \mathbb{R}$ for all $(x, a) \in \mathcal{X} \times A$. We define for every $t \in \mathbb{N}_0$ the map $\mathcal{X} \times A \ni (x, a) \mapsto \Delta_t(x, a) := Q_t(x, a) - Q^*(x, a) \in \mathbb{R}$. Note that indeed, as for all $(x, a) \in \mathcal{X} \times A$ we have $Q_t(x, a)$ as well as $Q^*(x, a)$ is finite (compare (13)), we directly conclude the finiteness of $\Delta_t(x, a)$ for all $(x, a) \in \mathcal{X} \times A$. Moreover, we obtain by (A.3) and by using the relation $\gamma_t(x, a, X_t) = \tilde{\gamma}_t \mathbf{1}_{\{(X_t, a_t(X_t)) = (x, a)\}}$ that

$$\begin{aligned} \Delta_{t+1}(x, a) &= (1 - \gamma_t(x, a, X_t))\Delta_t(x, a) \\ &\quad + \gamma_t(x, a, X_t) \left(-(-f_{t,(X_t, a_t(X_t))})^{\lambda_t c}(X_{t+1}) \right. \\ &\quad \left. - \varepsilon^q \lambda_t - Q^*(x, a) \right) \\ &= (1 - \gamma_t(x, a, X_t))\Delta_t(x, a) \\ &\quad + \gamma_t(x, a, X_t) \mathbf{1}_{\{(X_t, a_t(X_t)) = (x, a)\}} \\ &\quad \cdot \left(-(-f_{t,(X_t, a_t(X_t))})^{\lambda_t c}(X_{t+1}) \right. \\ &\quad \left. - \varepsilon^q \lambda_t - Q^*(X_t, a_t(X_t)) \right). \end{aligned} \tag{A.4}$$

Next, we define for every $t \in \mathbb{N}_0$ the random variable $F_t(x, a) := \mathbf{1}_{\{(X_t, a_t(X_t)) = (x, a)\}} \left(-(-f_{t,(X_t, a_t(X_t))})^{\lambda_t c}(X_{t+1}) - \varepsilon^q \lambda_t - Q^*(X_t, a_t(X_t)) \right)$, which by (13) is finite for all $(x, a) \in \mathcal{X} \times A$. We consider the filtration $(\mathcal{G}_t)_{t \in \mathbb{N}_0}$ with $\mathcal{G}_t := \sigma(\{X_1, \dots, X_t\})$, $t \in \mathbb{N}$, and $\mathcal{G}_0 = \{\emptyset, \Omega\}$ being the trivial sigma-algebra. Note that, in particular, $\Delta_t(x, a)$, $\gamma_t(x, a)$ and $F_{t-1}(x, a)$ are \mathcal{G}_t -measurable for all $t \in \mathbb{N}$. Moreover, we have by (5) and by Proposition 16 that $\hat{\mathbb{P}}_{x_0, \mathbf{a}}$ - almost surely

$$\begin{aligned} &\left| \mathbb{E}_{\hat{\mathbb{P}}_{x_0, \mathbf{a}}} [F_t(x, a) \mid \mathcal{G}_t] \right| \\ &= \left| \mathbf{1}_{\{(X_t, a_t(X_t)) = (x, a)\}} \right. \\ &\quad \cdot \mathbb{E}_{\hat{\mathbb{P}}_{(X_t, a_t(X_t))}} \left[-(-f_{t,(X_t, a_t(X_t))})^{\lambda_t c}(X_{t+1}) \right. \\ &\quad \left. - \varepsilon^q \lambda_t - Q^*(X_t, a_t(X_t)) \right] \left| \right| \end{aligned}$$

$$\begin{aligned} &= \mathbf{1}_{\{(X_t, a_t(X_t)) = (x, a)\}} \\ &\quad \cdot \left| \sup_{\lambda \geq 0} \left(\mathbb{E}_{\hat{\mathbb{P}}_{(X_t, a_t(X_t))}} \left[-(-f_{t,(X_t, a_t(X_t))})^{\lambda c}(X_{t+1}) \right. \right. \right. \\ &\quad \left. \left. - \varepsilon^q \lambda \right] \right) - Q^*(X_t, a_t(X_t)) \right| \\ &= \mathbf{1}_{\{(X_t, a_t(X_t)) = (x, a)\}} \\ &\quad \cdot \left| \inf_{\mathbb{P} \in \mathcal{P}(X_t, a_t(X_t))} \mathbb{E}_{\mathbb{P}} \left[(f_{t,(X_t, a_t(X_t))})(X_{t+1}) \right] \right. \\ &\quad \left. - Q^*(X_t, a_t(X_t)) \right|. \end{aligned}$$

Thus, (14), (A.2), and Lemma 17 show that $\hat{\mathbb{P}}_{x_0, \mathbf{a}}$ - almost surely

$$\begin{aligned} &\left| \mathbb{E}_{\hat{\mathbb{P}}_{x_0, \mathbf{a}}} [F_t(x, a) \mid \mathcal{G}_t] \right| \\ &= \mathbf{1}_{\{(X_t, a_t(X_t)) = (x, a)\}} \\ &\quad \cdot \left| \inf_{\mathbb{P} \in \mathcal{P}(X_t, a_t(X_t))} \mathbb{E}_{\mathbb{P}} \left[(f_{t,(X_t, a_t(X_t))})(X_{t+1}) \right] \right. \\ &\quad \left. - Q^*(X_t, a_t(X_t)) \right| \\ &= \mathbf{1}_{\{(X_t, a_t(X_t)) = (x, a)\}} \cdot \left| (HQ_t)(X_t, a_t(X_t)) \right. \\ &\quad \left. - Q^*(X_t, a_t(X_t)) \right| \\ &= \mathbf{1}_{\{(X_t, a_t(X_t)) = (x, a)\}} \cdot \left| (HQ_t)(X_t, a_t(X_t)) \right. \\ &\quad \left. - (HQ^*)(X_t, a_t(X_t)) \right|. \end{aligned} \tag{A.5}$$

Hence it follows with Lemma 18 that $\hat{\mathbb{P}}_{x_0, \mathbf{a}}$ - almost surely

$$\begin{aligned} \left\| \mathbb{E}_{\hat{\mathbb{P}}_{x_0, \mathbf{a}}} [F_t(\cdot, \cdot) \mid \mathcal{G}_t] \right\|_{\infty} &\leq \|(HQ_t) - (HQ^*)\|_{\infty} \\ &\leq \alpha \|Q_t - Q^*\|_{\infty} = \alpha \|\Delta_t\|_{\infty}, \end{aligned}$$

where the norm $\|\cdot\|$ is defined in (A.1). Next, recall that $C_r := \max \{1, \max_{y_0, y_1 \in \mathcal{X}, b \in A} |r(y_0, b, y_1)|\}$. Note that by (14), by the λc -transform from Definition 5, and since $\inf_{x, y \in \mathcal{X}} c(x, y) = 0$, we have for all $t \in \mathbb{N}_0$ that

$$\begin{aligned} &(-f_{t,(X_t, a_t(X_t))})^{\lambda_t c}(X_{t+1}) + \alpha \min_{y' \in \mathcal{X}} \max_{b' \in A} Q^*(y', b') \\ &= \left(-r(X_t, a_t(X_t), \cdot) - \alpha \max_{b \in A} Q_t(\cdot, b) \right)^{\lambda_t c} (X_{t+1}) \\ &\quad + \alpha \min_{y' \in \mathcal{X}} \max_{b' \in A} Q^*(y', b') \end{aligned}$$

$$\begin{aligned}
&\leq \left(C_r - \alpha \max_{b \in A} Q_t(\cdot, b) \right)^{\lambda_t c} (X_{t+1}) \\
&\quad + \alpha \min_{y' \in \mathcal{X}} \max_{b' \in A} Q^*(y', b') \\
&\leq \max_{z \in \mathcal{X}} \left(C_r - \alpha \max_{b \in A} Q_t(\cdot, b) \right)^{\lambda_t c} (z) \\
&\quad + \alpha \min_{y' \in \mathcal{X}} \max_{b' \in A} Q^*(y', b')
\end{aligned}$$

The latter expression coincides with

$$\begin{aligned}
&\max_{z, y \in \mathcal{X}} \left(C_r + \alpha \min_{y' \in \mathcal{X}} \max_{b' \in A} Q^*(y', b') \right. \\
&\quad \left. - \alpha \max_{b \in A} Q_t(y, b) - \lambda_t c(z, y) \right),
\end{aligned}$$

which implies

$$\begin{aligned}
&(-f_{t, (X_t, a_t(X_t))})^{\lambda_t c} (X_{t+1}) + \alpha \min_{y' \in \mathcal{X}} \max_{b' \in A} Q^*(y', b') \\
&\leq \max_{z, y \in \mathcal{X}} \left(C_r + \alpha \max_{b' \in A} Q^*(y, b') \right. \\
&\quad \left. - \alpha \max_{b \in A} Q_t(y, b) - \lambda_t c(z, y) \right) \\
&\leq \max_{z, y \in \mathcal{X}} \left(C_r + \alpha \max_{b \in A} \{Q^*(y, b) - Q_t(y, b)\} - \lambda_t c(z, y) \right) \\
&\leq \max_{z, y \in \mathcal{X}} (C_r + \alpha \|\Delta_t\|_\infty - \lambda_t c(z, y)) \\
&= C_r + \alpha \|\Delta_t\|_\infty =: M \in \mathbb{R},
\end{aligned} \tag{A.6}$$

and similarly, since $c(z, z) = 0$ for all $z \in \mathcal{X}$,

$$\begin{aligned}
&(-f_{t, (X_t, a_t(X_t))})^{\lambda_t c} (X_{t+1}) + \alpha \min_{y' \in \mathcal{X}} \max_{b' \in A} Q^*(y', b') \\
&\geq \min_{z \in \mathcal{X}} \max_{y \in \mathcal{X}} \left(-C_r + \alpha \min_{y' \in \mathcal{X}} \max_{b' \in A} Q^*(y', b') \right. \\
&\quad \left. - \alpha \max_{b \in A} Q_t(y, b) - \lambda_t c(z, y) \right) \\
&\geq \min_{z \in \mathcal{X}} \max_{y \in \mathcal{X}} \left(-C_r + \alpha \min_{y' \in \mathcal{X}} \min_{b \in A} (Q^*(y', b) - Q_t(y, b)) \right. \\
&\quad \left. - \lambda_t c(z, y) \right) \\
&\geq \min_{z \in \mathcal{X}} \max_{y \in \mathcal{X}} \left(-C_r - \alpha \max_{y' \in \mathcal{X}, b \in A} |Q_t(y, b) - Q^*(y', b)| \right. \\
&\quad \left. - \lambda_t c(z, y) \right) \\
&\geq \min_{z \in \mathcal{X}} \left(-C_r - \alpha \max_{y' \in \mathcal{X}, b \in A} |Q_t(z, b) - Q^*(y', b)| \right) \\
&\geq -C_r - \alpha \max_{z, y' \in \mathcal{X}, b \in A} \left(|Q_t(z, b) - Q^*(z, b)| \right. \\
&\quad \left. + |Q^*(z, b) - Q^*(y', b)| \right)
\end{aligned}$$

$$\begin{aligned}
&\geq -C_r - \alpha \|\Delta_t\|_\infty \\
&\quad - \alpha \max_{z, y' \in \mathcal{X}, b \in A} |Q^*(z, b) - Q^*(y', b)| =: m \in \mathbb{R}.
\end{aligned} \tag{A.7}$$

We define $C := (4\alpha^2 + (2C_r + \alpha \max_{z, y' \in \mathcal{X}, b \in A} |Q^*(z, b) - Q^*(y', b)|)^2) < \infty$. Then, by using Popoviciu's inequality on variances⁹ applied to the bounds m, M computed in (A.6) and (A.7), and by using the inequality $(a+b)^2 \leq 2(a^2 + b^2)$ which holds for all $a, b \in \mathbb{R}$, we see for every $(x, a) \in \mathcal{X} \times A$ that $\hat{\mathbf{P}}_{x_0, \mathbf{a}}$ - almost surely

$$\begin{aligned}
&\text{Var}_{\hat{\mathbf{P}}_{x_0, \mathbf{a}}} (F_t(x, a) \mid \mathcal{G}_t) \\
&= \mathbf{1}_{\{(X_t, a_t(X_t)) = (x, a)\}} \\
&\quad \cdot \text{Var}_{\hat{\mathbf{P}}_{(X_t, a_t(X_t))}} ((-f_{t, (X_t, a_t(X_t))})^{\lambda_t c} (X_{t+1})) \\
&= \mathbf{1}_{\{(X_t, a_t(X_t)) = (x, a)\}} \\
&\quad \cdot \text{Var}_{\hat{\mathbf{P}}_{(X_t, a_t(X_t))}} \left((-f_{t, (X_t, a_t(X_t))})^{\lambda_t c} (X_{t+1}) \right. \\
&\quad \left. + \alpha \min_{y' \in \mathcal{X}} \max_{b' \in A} Q^*(y', b') \right) \\
&\leq \frac{1}{4} (M - m)^2 \\
&= \frac{1}{4} \left(2\alpha \|\Delta_t\|_\infty + 2C_r \right. \\
&\quad \left. + \alpha \max_{z, y' \in \mathcal{X}, b \in A} |Q^*(z, b) - Q^*(y', b)| \right)^2 \\
&\leq \frac{1}{2} \left(4\alpha^2 \|\Delta_t\|_\infty^2 \right. \\
&\quad \left. + (2C_r + \alpha \max_{z, y' \in \mathcal{X}, b \in A} |Q^*(z, b) - Q^*(y', b)|)^2 \right) \\
&\leq \left(4\alpha^2 + \left(2C_r + \alpha \max_{z, y' \in \mathcal{X}, b \in A} |Q^*(z, b) - Q^*(y', b)| \right)^2 \right) \\
&\quad \cdot (1 + \|\Delta_t\|_\infty^2) \leq C \cdot (1 + \|\Delta_t\|_\infty)^2.
\end{aligned}$$

This means the assumptions of Lemma 15 are fulfilled, and we obtain that $\Delta_t(x, a) \rightarrow 0$ for $t \rightarrow \infty$ $\hat{\mathbf{P}}_{x_0, \mathbf{a}}$ -almost surely, which implies, by definition of Δ_t , that $Q_t(x, a) \rightarrow Q^*(x, a)$ for $t \rightarrow \infty$ $\hat{\mathbf{P}}_{x_0, \mathbf{a}}$ -almost surely. \square

Proof of Proposition 11 The conditions of [36, Assumption 2.1-2.4] are satisfied w.r.t. L_P and L_r defined in (20) and (22) since here both the state and action space are finite. Hence the result follows from [36, Theorem 3.1]. \square

⁹ Popoviciu's inequality (see [40] or [46]) states that for all random variables Z on a probability space $(\Omega, \mathcal{F}, \mathbf{P})$ satisfying $m \leq Z \leq M$ for some $-\infty < m \leq M < \infty$ we have $\text{Var}_{\mathbf{P}}(Z) \leq \frac{1}{4} (M - m)^2$.