# Statistics and Samples in Distributional Reinforcement Learning

**Mark Rowland** [1]  **Robert Dadashi** [2]  **Saurabh Kumar** [2]  **Rémi Munos** [1]  **Marc G. Bellemare** [2]  **Will Dabney** [1]

## Abstract

We present a unifying framework for designing and analysing distributional reinforcement learning (DRL) algorithms in terms of recursively estimating statistics of the return distribution. Our key insight is that DRL algorithms can be decomposed as the combination of some statistical estimator and a method for imputing a return distribution consistent with that set of statistics. With this new understanding, we are able to provide improved analyses of existing DRL algorithms as well as construct a new algorithm (EDRL) based upon estimation of the *expectiles* of the return distribution. We compare EDRL with existing methods on a variety of MDPs to illustrate concrete aspects of our analysis, and develop a deep RL variant of the algorithm, ER-DQN, which we evaluate on the Atari-57 suite of games.

## 1. Introduction

In reinforcement learning (RL), a central notion is the *return*, the sum of discounted rewards. Typically, the average of these returns is estimated by a value function and used for policy improvement. Recently, however, approaches that attempt to learn the distribution of the return have been shown to be surprisingly effective (Morimura et al., 2010a;b; Bellemare et al., 2017; Dabney et al., 2017; 2018; Gruslys et al., 2018); we refer to the general approach of learning return distributions as *distributional RL* (DRL).

Despite impressive experimental performance (Bellemare et al., 2017; Barth-Maron et al., 2018; Dabney et al., 2018) and fundamental theoretical results (Rowland et al., 2018; Qu et al., 2018), it remains challenging to develop and analyse DRL algorithms. In this paper, we propose to address these challenges by phrasing DRL algorithms in terms of recursive estimation of sets of statistics on the return distribution. We observe that DRL algorithms can be viewed as combining a statistical estimator with a procedure we refer

to as an *imputation strategy*, which generates a return distribution consistent with the set of statistical estimates. This highly general approach (see Figure 1) requires a precise treatment of the differing roles of statistics and samples in distributional RL.

Using this framework we are able to provide new theoretical results for existing DRL algorithms as well as demonstrate the derivation of a new algorithm based on the expectiles of the return distribution. More importantly, our novel approach immediately applies to a large class of statistics and imputation strategies, suggesting several avenues for future research. Specifically, we are able to provide answers to the following questions:

 (i) Can we describe existing DRL algorithms in a unifying framework, and could such a framework be used to develop new algorithms?
 (ii) What return distribution statistics can be learnt *exactly* through Bellman updates?
(iii) If certain statistics cannot be learnt exactly, how can we estimate them in a principled manner, and give guarantees on their approximation error relative to the true values of these statistics?

After reviewing relevant background material, we begin with (i) by presenting a new framework for understanding DRL, that is, in terms of a set of *statistics* to be learnt, and an *imputation strategy* for specifying a dynamic programming update. We then formalise (ii) by introducing the notion of *Bellman closedness* for collections of statistics, and show that in a wide class of statistics, the only properties of return distributions that can be learnt exactly through Bellman updates are moments. Interestingly, this rules out statistics such as quantiles that have formed the basis of successful existing DRL algorithms. However, we then address (iii) by showing that the framework allows us to give guarantees on the approximation error introduced in learning these statistics, through the notion of *approximate Bellman closedness*. We apply the framework developed in answering these questions to the case of *expectile* statistics to develop a new distributional RL algorithm, which we term Expectile Distributional RL (EDRL). Finally, we test these new insights on a variety of MDPs and larger-scale environments to illustrate and expand on the theoretical contributions developed earlier in the paper.

---

[1]DeepMind [2]Google Brain. Correspondence to: Mark Rowland <markrowland@google.com>.
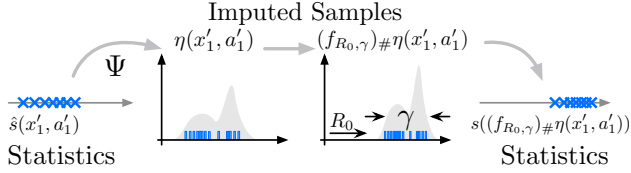
Figure 1. Illustration of learning with imputed samples from sets of statistics. Left: A distribution is imputed from the current statistical estimate. Middle: The distributional Bellman operator is applied to the imputed distribution. Right: New statistics are estimated based upon samples from the imputed distribution.

## 2. Background

Consider a Markov decision process $(\mathcal{X}, \mathcal{A}, p, \gamma, \mathcal{R})$ with finite state space $\mathcal{X}$, finite action space $\mathcal{A}$, transition kernel $p : \mathcal{X} \times \mathcal{A} \to \mathscr{P}(\mathcal{X})$, discount rate $\gamma \in [0, 1)$, and reward distributions $\mathcal{R}(x, a) \in \mathscr{P}(\mathbb{R})$ for each $(x, a) \in \mathcal{X} \times \mathcal{A}$. Thus, if an agent is at state $X_t \in \mathcal{X}$ at time $t \in \mathbb{N}_0$, and an action $A_t \in \mathcal{A}$ is taken, the agent transitions to a state $X_{t+1} \sim p(\cdot|X_t, A_t)$ and receives a reward $R_t \sim \mathcal{R}(X_t, A_t)$. We now briefly review two principal goals in reinforcement learning.

Firstly, given a Markov policy $\pi : \mathcal{X} \to \mathscr{P}(\mathcal{A})$, *evaluation* of $\pi$ consists of computing the expected returns $Q^\pi(x, a) = \mathbb{E}_\pi \left[ \sum_{t=0}^\infty \gamma^t R_t | X_0 = x, A_0 = a \right]$, where $\mathbb{E}_\pi$ indicates that at each time step $t \in \mathbb{N}$, the agent's action $A_t$ is sampled from $\pi(\cdot|X_t)$. Secondly, the task of *control* consists of finding a policy $\pi : \mathcal{X} \to \mathscr{P}(\mathcal{A})$ for which the expected returns are maximised.

### 2.1. Bellman equations

The classical *Bellman equation* (Bellman, 1957) relates expected returns at each state-action pair $(x, a) \in \mathcal{X} \times \mathcal{A}$ to the expected returns at possible next states in the MDP by:

$$Q^\pi(x, a) = \mathbb{E}_\pi[R_0 + \gamma Q^\pi(X_1, A_1)|X_0 = x, A_0 = a]. \quad (1)$$

This gives rise to the following fixed-point iteration scheme

$$Q(x, a) \leftarrow \mathbb{E}_\pi[R_0 + \gamma Q(X_1, A_1)|X_0 = x, A_0 = a], \quad (2)$$

for updating a collection of *approximations* $(Q(x, a)|(x, a) \in \mathcal{X} \times \mathcal{A})$ towards their true values. This fundamental algorithm, together with techniques from approximate dynamic programming and stochastic approximation, allows expected returns in an MDP to be learnt and improved upon, forming the basis of all value-based RL (Sutton & Barto, 2018).

The *distributional Bellman equation* describes a similar relationship to Equation (1) at the level of probability distributions (Morimura et al., 2010a;b; Bellemare et al., 2017). Letting $\eta_\pi(x, a) \in \mathscr{P}(\mathbb{R})$ be the *distribution* of the random return $\sum_{t=0}^\infty \gamma^t R_t | X_0 = x, A_0 = a$ when actions are

selected according to $\pi$, we have

$$\begin{aligned} \eta_\pi(x, a) &= (\mathcal{T}^\pi \eta_\pi)(x, a), \quad (3) \\ &= \mathbb{E}_\pi[(f_{R_0, \gamma})_\# \eta_\pi(X_1, A_1)|X_0 = x, A_0 = a], \end{aligned}$$

where the expectation gives a mixture distribution over next-states, $f_{r,\gamma} : \mathbb{R} \to \mathbb{R}$ is defined by $f_{r,\gamma}(x) = r + \gamma x$, and $g_\# \mu \in \mathscr{P}(\mathbb{R})$ is the pushforward of the measure $\mu$ through the function $g$, so that for all Borel subsets $A \subseteq \mathbb{R}$, we have $g_\# \mu(A) = \mu(g^{-1}(A))$ (Rowland et al., 2018).

Stated in terms of the random return $Z^\pi(x, a)$, distributed according to $\eta_\pi(x, a)$, this takes a more familiar form with

$$Z^\pi(x, a) \overset{D}{=} R_0 + \gamma Z^\pi(X_1, A_1).$$

In analogy with Expression (2), an update operation could be defined from Equation (3) to move a collection of approximate distributions $(\eta(x, a)|(x, a) \in \mathcal{X} \times \mathcal{A})$ towards the true return distributions. However, since the space of distributions $\mathscr{P}(\mathbb{R})$ is infinite-dimensional, it is typically impossible to work directly with the distributional Bellman equation, and existing approaches to distributional RL generally rely on parametric approximations to this equation; we briefly review some important examples of these approaches below.

### 2.2. Categorical and quantile distributional RL

To date, the main approaches to DRL employed at scale have included learning discrete *categorical* distributions (Bellemare et al., 2017; Barth-Maron et al., 2018; Qu et al., 2018), and learning distribution *quantiles* (Dabney et al., 2017; 2018; Zhang et al., 2019); we refer to these approaches as CDRL and QDRL respectively. We give brief accounts of the dynamic programming versions of these algorithms here, with full descriptions of stochastic versions, related results, and visualisations given in Appendix Section A for completeness. We note also that other approaches, such as learning mixtures of Gaussians, have been explored (Barth-Maron et al., 2018).

**CDRL.** CDRL assumes a *categorical* form for return distributions, taking $\eta(x, a) = \sum_{k=1}^K p_k(x, a)\delta_{z_k}$, where $\delta_z$ denotes the Dirac distribution at location $z$. The values $z_1 < \cdots < z_K$ are an evenly spaced, fixed set of supports, and the probability parameters $p_{1:K}(x, a)$ are learnt. The corresponding Bellman update takes the form

$$\eta(x, a) \leftarrow (\Pi_\mathcal{C} \mathcal{T}^\pi \eta)(x, a),$$

where $\Pi_\mathcal{C} : \mathscr{P}(\mathbb{R}) \to \mathscr{P}(\{z_1, \ldots, z_K\})$ is a projection operator which ensures the right-hand side of the expression above is a distribution supported only on $\{z_1, \ldots, z_K\}$; full details are reviewed in Appendix Section A.

**QDRL.** In contrast, QDRL assumes a parametric form for return distributions $\eta(x,a) = \frac{1}{K}\sum_{k=1}^{K}\delta_{z_k(x,a)}$, where now $z_{1:K}(x,a)$ are learnable parameters. The Bellman update is given by moving the atom location $z_k(x,a)$ in $\eta(x,a)$ to the $\tau_k$-quantile (where $\tau_k = \frac{2k-1}{2K}$) of the target distribution $\mu := (\mathcal{T}^\pi\eta)(x,a)$, defined as the minimiser $q^* \in \mathbb{R}$ of the quantile regression loss

$$\text{QR}(q;\mu,\tau_k) = \mathbb{E}_{Z\sim\mu}[[\tau_k\mathbb{1}_{Z>q} + (1-\tau_k)\mathbb{1}_{Z\leq q}]\,|Z-q|]. \tag{4}$$

## 3. The role of statistics in distributional RL

In this section, we describe a new perspective on existing distributional RL algorithms, with a focus on learning sets of statistics, rather than approximate distributions. We begin with a precise definition.

**Definition 3.1** (**Statistics**). *A* statistic *is a function* $s : \mathscr{P}(\mathbb{R}) \to \mathbb{R}$. *We also allow statistics to be defined on subsets of* $\mathscr{P}(\mathbb{R})$*, in situations where an assumption (such as finite moments) is required for the statistic to be defined.*

The QDRL update described in Section 2.2 is readily interpreted from the perspective of learning statistics; the update extracts the values of a finite set of quantile statistics from the target distribution, and all other information about the target is lost. It is less obvious whether the CDRL update can also be interpreted as keeping track of a finite set of statistics, but the following lemma shows that this is indeed the case.

**Lemma 3.2.** *CDRL updates, with distributions supported on* $z_1 < \ldots < z_K$*, can be interpreted as learning the values of the following statistics of return distributions:*

$$s_{z_k,z_{k+1}}(\mu) = \mathbb{E}_{Z\sim\mu}\big[h_{z_k,z_{k+1}}(Z)\big] \text{ for } k=1,\ldots,K-1\,,$$

*where for* $a < b$*,* $h_{a,b} : \mathbb{R} \to \mathbb{R}$ *is a piecewise linear function defined so that* $h_{a,b}(x)$ *is equal to 1 for* $x \leq a$*, equal to 0 for* $x \geq b$*, and linearly interpolating between* $h_{a,b}(a)$ *and* $h_{a,b}(b)$ *for* $x \in [a,b]$*.*

Although viewing distributional RL as approximating the return distribution with some parameterisation is intuitive from an algorithmic standpoint, there are advantages to thinking in terms of sets of statistics and their recursive estimation; this perspective allows us to precisely quantify what information is being passed through successive distributional Bellman updates. This in turn leads to new insights in the development and analysis of DRL algorithms. Before addressing these points, we first consider a motivating example where a lack of precision could lead us astray.

### 3.1. Expectiles

Motivated by the success of QDRL, we consider learning *expectiles* of return distributions, a family of statistics introduced by Newey & Powell (1987). Expectiles generalise the mean in analogy with how quantiles generalise the median. As the goal of RL is to maximise mean returns, we conjectured that expectiles, in particular, might lead to successful DRL algorithms. We begin with a formal definition.

**Definition 3.3** (**Expectiles**). *Given a distribution* $\mu \in \mathscr{P}(\mathbb{R})$ *with finite second moment, and* $\tau \in [0,1]$*, the* $\tau$*-expectile of* $\mu$ *is defined to be the minimiser* $q^* \in \mathbb{R}$ *of the expectile regression loss* $\text{ER}(q;\mu,\tau)$*, given by*

$$\text{ER}(q;\mu,\tau) = \mathbb{E}_{Z\sim\mu}\big[[\tau\mathbb{1}_{Z>q} + (1-\tau)\mathbb{1}_{Z\leq q}]\,(Z-q)^2\big]\,.$$

*For each* $\tau \in [0,1]$*, we denote the* $\tau$*-expectile of* $\mu$ *by* $e_\tau(\mu)$*.*

We remark that: (i) the expectile regression loss is an asymmetric version of the squared loss, just as the quantile regression loss is an asymmetric version of the absolute value loss; and (ii) the $1/2$-expectile of $\mu$ is simply its mean. Because of this, we can attempt to derive an algorithm by replacing the quantile regression loss in QDRL with the expectile regression loss in Definition 3.3, so as to learn the expectiles corresponding to $\tau_1, \ldots, \tau_K \in [0,1]$.

Following this logic, we again take approximate distributions of the form $\eta(x,a) = \frac{1}{K}\sum_{k=1}^{K}\delta_{z_k(x,a)}$, and we perform updates according to

$$z_k(x,a) \leftarrow \underset{q\in\mathbb{R}}{\arg\min}\,\text{ER}(q;\mu,\tau_k)\,, \tag{5}$$

where $\mu = (\mathcal{T}^\pi\eta)(x,a)$ is the target distribution.

In practice, however, this algorithm does not perform as we might expect, and in fact the variance of the learnt distributions collapses as training proceeds, indicating that the algorithm does not approximate the true expectiles in any reasonable sense. In Figure 2, we illustrate this point by comparing the learnt statistics for this "naive" approach with those of CDRL and our proposed algorithm EDRL (introduced in Section 3.3). All methods accurately approximate the immediate reward distribution (right), but as successive Bellman updates are applied the different algorithms show characteristic approximation errors. The CDRL algorithm overestimates the variance of the return distribution due to the projection $\Pi_\mathcal{C}$ splitting probability mass across the discrete support. By contrast, the naive expectile approach underestimates the true variance, quickly converging to a single Dirac.

We observe that there is a *"type error"* present in Expression (5); the parameter being updated, $z_k(x,a)$, has the semantics of a *statistic*, as the minimiser of the ER loss, whilst the parameters appearing in the target distribution $(\mathcal{T}^\pi\eta)(x,a)$ have the semantics of *outcomes/samples*. A crucial message of this paper is the need to distinguish between statistics and samples in distributional RL; in the next section, we describe a general framework for achieving this.
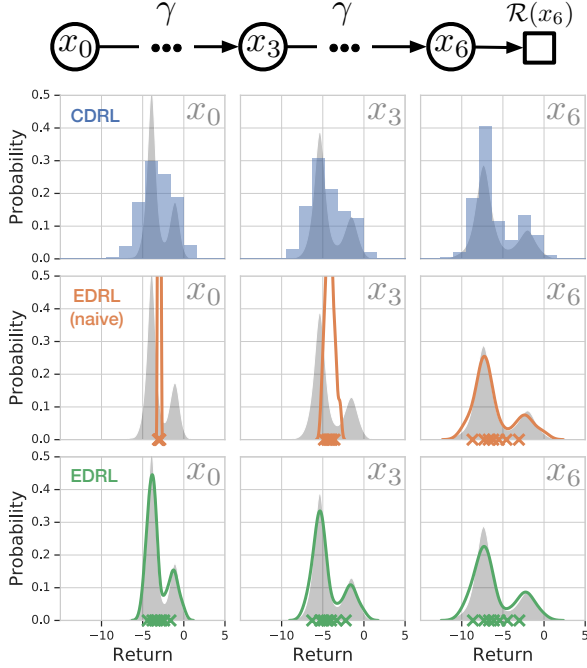
*Figure 2.* Chain MDP, one action, with bimodal reward distribution at absorbing state $x_6$ and $\gamma = 0.9$. CDRL (top, blue) fits the true return distribution (grey) well, but overestimates the variance. A naive approach to EDRL (middle, orange) accurately fits the immediate reward distribution at $x_6$, but quickly collapses to zero variance with successive Bellman updates. Our proposed approach, EDRL, using imputation strategies (bottom, green) provides an accurate approximation through many Bellman updates.

## 3.2. Imputation strategies

If we had access to full return distribution estimates $\eta(x', a')$ at each possible next state-action pair $(x', a')$, we would be able to avoid the conflation between samples and statistics described in the previous section. Denoting the approximation to the value of a statistic $s_k$ at a state-action pair $(x, a) \in \mathcal{X} \times \mathcal{A}$ by $\hat{s}_k(x, a)$, we would like to update according to:

$$\hat{s}_k(x, a) \leftarrow s_k\left((\mathcal{T}^\pi \eta)(x, a)\right). \tag{6}$$

Thus, a principled way in which to design DRL algorithms for collections of statistics is to include an *additional* step in the algorithm in which for any state-action pair $(x', a')$ that we would like to backup from, the estimated statistics $\hat{s}_{1:K}(x', a')$ are converted into a consistent distribution $\eta(x', a')$. This would then allow backups of the form in Expression (6) to be carried out. This notion is formalised in the following definition.

**Definition 3.4** (**Imputation strategies**). *Given a set of statistics* $\{s_1, \ldots, s_K\}$, *an* imputation strategy *is a function* $\Psi : \mathbb{R}^K \to \mathscr{P}(\mathbb{R})$ *that maps each vector of statistic values to a distribution that has those statistics. Mathematically,*

$\Psi$ *is such that* $s_i(\Psi(\sigma_{1:K})) = \sigma_i$, *for each* $i \in \{1, \ldots, K\}$ *and each collection of statistic values* $\sigma_{1:K} \in \mathbb{R}^K$.

Thus, an imputation strategy is simply a function that takes in a collection of values for certain statistics, and returns a probability distribution with those statistic values; in some sense, it is a pseudo-inverse of $s_{1:K}$.

**Example 3.5** (**Imputation strategies in CDRL and QDRL**). *In QDRL, the imputation strategy is given by* $\Psi(\sigma_{1:K}) = \frac{1}{K}\sum_{k=1}^{K} \delta_{\sigma_k}$. *In CDRL, given approximate statistics* $\hat{s}_{z_k, z_{k+1}}(x, a)$ *for* $k = 1, \ldots, K - 1$, *the imputation strategy is given by selecting the distribution* $\sum_{k=1}^{K} p_k \delta_{z_k}$ *such that* $p_1 = \hat{s}_{z_1, z_2}(x, a)$, $p_k = \hat{s}_{z_k, z_{k+1}}(x, a) - \hat{s}_{z_{k-1}, z_k}(x, a)$ *for* $k = 2, \ldots, K - 1$, *and* $p_K = 1 - \sum_{k<K} p_k$.

We now have a general framework for defining principled distributional RL algorithms: (i) select a family of statistics to learn; (ii) select an imputation strategy; (iii) perform (or approximate) updates of the form in Expression (6). We summarise this in Algorithm 1.

---

**Algorithm 1** Generic DRL update algorithm.

**Require:** Statistic estimates $\hat{s}_{1:K}(x, a)\ \forall (x, a) \in \mathcal{X} \times \mathcal{A}$
    and $k = 1, \ldots, K$, imputation strategy $\Psi$.
    Select state-action pair $(x, a) \in \mathcal{X} \times \mathcal{A}$ to update.
    Impute distribution at each possible next state-action pair:
      $\eta(x', a') = \Psi(\hat{s}_{1:K}(x', a')), \quad \forall (x', a') \in \mathcal{X} \times \mathcal{A}$.
    Update statistics at $(x, a) \in \mathcal{X} \times \mathcal{A}$:
      $\hat{s}_k(x, a) \leftarrow s_k\left((\mathcal{T}^\pi \eta)(x, a)\right)$.

---

## 3.3. Expectile distributional reinforcement learning

We now apply the general framework of statistics and imputation strategies developed in Section 3.2 to the specific case of *expectiles*, introduced in Section 3.1. We will define an imputation strategy so that updates of the form given in Expression (6) can be applied to learn expectiles.

The imputation strategy has the task of accepting as input a collection of expectile values $\epsilon_1, \ldots, \epsilon_K$, corresponding to $\tau_1, \ldots, \tau_K \in (0, 1)$, and computing a probability distribution $\mu$ such that $e_{\tau_i}(\mu) = \epsilon_i$ for $i = 1, \ldots, K$. Since $\text{ER}(q; \mu, \tau)$ is strictly convex as a function of $q$, this can be restated as finding a probability distribution $\mu$ satisfying the first-order optimality conditions

$$\nabla_q \text{ER}(q; \mu, \tau_i)\big|_{q=\epsilon_i} = 0 \quad \forall i \in [K]. \tag{7}$$

This defines a root-finding problem, but may equivalently be formulated as a minimisation problem, with objective

$$\sum_{i=1}^{K} \left(\nabla_q \text{ER}(q; \mu, \tau_i)\big|_{q=\epsilon_i}\right)^2. \tag{8}$$

By constraining the distribution $\mu$ to be of the form $\frac{1}{N} \sum_{n=1}^{N} \delta_{z_n}$ and viewing the minimisation objective above as a function of $z_{1:N}$, it is straightforwardly verifiable that this minimisation problem is convex. The imputation strategy is thus defined implicitly, by stating that $\Psi(\epsilon_{1:K})$ is given by a minimiser of (8) of the form $\frac{1}{N} \sum_{n=1}^{N} \delta_{z_n}$. We remark that other parametric choices for $\mu$ are possible, but the mixture of Dirac deltas described above leads to a particular tractable optimisation problem.

Having established an imputation strategy $\Psi$, Algorithm 1 now yields a full DRL algorithm for learning expectiles, which we term EDRL. Returning to Figure 2, we observe that EDRL (bottom row) is able to accurately represent the true return distribution, even after many Bellman updates through the chain, and does not exhibit the collapse observed with the naive approach in Section 3.1.

### 3.4. Stochastic approximation

Practically speaking, it is often not possible to compute the updates in Expression (6), owing to MDP dynamics being unknown and/or intractable to integrate over. Because of this, it is often necessary to apply stochastic approximation. Let $(r, x', a')$ be a sample of the random variables $(R_0, X_1, A_1)$, obtained by direct interaction with the environment. Then, we update $\hat{s}_k(x, a)$ using the gradient of a loss function $L_k : \mathbb{R} \times \mathscr{P}(\mathbb{R}) \to \mathbb{R}$:

$$\nabla_{\hat{s}_k(x,a)} L_k(\hat{s}_k(x, a); (f_{r,\gamma})_\# \eta(x', a')) . \qquad (9)$$

For EDRL, a natural such loss function for the estimated statistic $\hat{s}_k(x, a)$ is the expectile regression loss of Definition 3.3 at $\tau_k$; this yields a stochastic version of EDRL, described in Algorithm 2.

---

**Algorithm 2** Stochastic EDRL update algorithm.

**Require:** Expectile estimates $\hat{s}_k(x, a)$ for each $(x, a) \in \mathcal{X} \times \mathcal{A}$ and $k = 1, \ldots, K$.
  Collect sample $(x, a, r, x', a')$.
  Impute distribution $\frac{1}{K} \sum_{k=1}^{K} \delta_{z_k}$ from target expectiles $\hat{s}_{1:K}(x', a')$ by solving (7) or minimising (8).
  Scale/translate samples $z_i \leftarrow r + \gamma z_i \; \forall i$.
  Update estimated expectiles at $(x, a) \in \mathcal{X} \times \mathcal{A}$ by computing the gradients

$$\nabla_{\hat{s}_k(x,a)} \sum_{k=1}^{K} \mathrm{ER}(\hat{s}_k(x, a); \tfrac{1}{N} \sum_{n=1}^{N} \delta_{z_n}, \tau_k)$$

  for each $k = 1, \ldots, K$.

---

To ensure convergence of these stochastic gradient updates to the correct statistic, it should be the case that the expectation of the (sub-)gradient (9) at the true value of the statistics is equal to 0. It can be verified that this is the case whenever (i) the true statistic $q^*$ of a distribution $\mu$ satisfies

$q^* = \mathrm{argmin}_{q \in \mathbb{R}} L_k(q; \mu)$, (ii) the loss $L_k$ is *affine* in the probability distribution argument. M-estimator losses and their associated statistics (Huber & Ronchetti, 2009) satisfy these conditions, and thus represent a large family of statistics to which this approach to DRL could immediately be applied; the statistics in CDRL, QDRL and EDRL are all special cases of M-estimators.

## 4. Analysing distributional RL

We now use the framework of statistics and imputations strategies developed in Section 3 to build a deeper understanding of the accuracy with which statistics in distributional RL may be learnt via Bellman updates.

### 4.1. Bellman closedness

The classical Bellman equation (1) shows that there is a closed-form relationship between expected returns at each state-action pair of an MDP; if the goal is to learn expected returns, we are not required to keep track of any other statistics of the return distributions. This well-known observation, together with the new interpretation of DRL algorithms as learning collections of statistics of return distributions, motivates a more general question:

"Given a set of statistics $\{s_1, \ldots, s_K\}$, if we want to learn the values $s_{1:K}(\eta_\pi(x, a))$ for all $(x, a) \in \mathcal{X} \times \mathcal{A}$ via dynamic programming, is it sufficient to keep track of *only* these statistics?"

The following definition formalises this question.

**Definition 4.1** (**Bellman closedness**). *A set of statistics* $\{s_1, \ldots, s_K\}$ *is* Bellman closed *if for each* $(x, a) \in \mathcal{X} \times \mathcal{A}$, *the statistics* $s_{1:K}(\eta_\pi(x, a))$ *can be expressed, in an MDP-independent manner, in terms of the random variables $R_0$ and $s_{1:K}(\eta_\pi(X_1, A_1))|X_0 = x, A_0 = a$, and the discount factor $\gamma$. We refer to any such expression for a set of Bellman closed set of statistics as a* Bellman equation, *and write $\mathcal{T}^\pi : (\mathbb{R}^K)^{\mathcal{X} \times \mathcal{A}} \to (\mathbb{R}^K)^{\mathcal{X} \times \mathcal{A}}$ for the corresponding operator such that the Bellman equation can be written*

$$\mathbf{s}^\pi = \mathcal{T}^\pi \mathbf{s}^\pi , \qquad (10)$$

*where* $\mathbf{s}^\pi = (s_{1:K}(\eta_\pi(x, a))|(x, a) \in \mathcal{X} \times \mathcal{A})$.

Thus, the singleton set consisting of the mean statistic is Bellman closed; the corresponding Bellman equation is Equation (1). It is also known that the set consisting of the mean and variance statistics are Bellman closed (Sobel, 1982). In principle, given a Bellman closed set of statistics $\{s_1, \ldots, s_K\}$, the corresponding statistics of the return distributions can be found by solving a fixed-point equation corresponding to the relevant Bellman operator, $\mathcal{T}^\pi$. Further, if $\mathcal{T}^\pi$ is a contraction in some metric, then it is possible

to find the true statistics for the MDP via a fixed-point iteration scheme based on the operator $\mathcal{T}^\pi$. In contrast, if a collection of statistics $s_{1:K}$ is *not* Bellman closed, there is no Bellman equation relating the statistics of the return distributions, and consequently it is not possible to learn the statistics *exactly* using dynamic programming in a self-contained way; the set of statistics must either be enlarged to make it Bellman closed, or an imputation strategy can be used to perform backups as described in Section 3.2.

An important class of Bellman closed sets of statistics are given in the following result (Sobel, 1982; Lattimore & Hutter, 2012).

**Lemma 4.2.** *For each $K \in \mathbb{N}$, the set of statistics consisting of the first $K$ moments is Bellman closed.*

The next result shows that across a wide range of statistics, collections of moments are effectively the only finite sets of statistics that are Bellman closed; the proof relies on a result of Engert (1970) which characterises finite-dimensional vector spaces of measurable functions closed under translation.

**Theorem 4.3.** *The only finite sets of statistics of the form $s(\mu) = \mathbb{E}_{Z \sim \mu}[h(Z)]$ that are Bellman closed are given by collections of statistics $s_1, \ldots, s_K : \mathscr{P}(\mathbb{R}) \to \mathbb{R}$ with the property that the linear span $\{\sum_{k=0}^{K} \alpha_k s_k | \alpha_k \in \mathbb{R} \, \forall k\}$ is equal to the linear span of the set of moment functionals $\{\mu \mapsto \mathbb{E}_{Z \sim \mu}[Z^l] | l = 0, \ldots, L\}$, for some $L \le K$, where $s_0$ is the constant functional equal to $1$.*

We believe this to be an important novel result, which helps to highlight how rare it is for statistics to be Bellman closed. One important corollary of Theorem 4.3, given the characterisation of CDRL as learning expectations of return distributions in Lemma 3.2, is that the sets of statistics learnt in CDRL are *not* Bellman closed. A similar result holds for QDRL, and we record these facts in the following result.

**Lemma 4.4.** *The sets of statistics learnt under (i) CDRL, and (ii) QDRL, are not Bellman closed.*

The immediate upshot of this is that in general, the *learnt* values of statistics in distributional RL algorithms need not correspond exactly to the true underlying values for the MDP (even in tabular settings), as the statistics propagated through DRL dynamic programming updates are not sufficient to determine the statistics we seek to learn. This inexactness was noted specifically for CDRL and QDRL in the original papers (Bellemare et al., 2017; Dabney et al., 2017). In this paper, our analysis and experiments confirm that these artefacts arise even with tabular agents in fully-observed domains, thus representing intrinsic properties of the distributional RL algorithms concerned. However, empirically the distributions learnt by these algorithms are often accurate. In the next section, we provide theoretical guarantees that describe this phenomenon quantitatively.

## 4.2. Approximate Bellman closedness

In light of the results on Bellman closedness in Section 4.1, we might ask in what sense the values of the statistics learnt by DRL algorithms relate to the corresponding true underlying values for the MDP concerned. A key task in this analysis is to formalise the notion of *low approximation error* in DRL algorithms that seek to learn collections of statistics that are not Bellman closed. Perhaps surprisingly, in general it is not possible to simultaneously achieve low approximation error on all statistics in a non-Bellman closed set; we give several examples for CDRL and QDRL to this end in Appendix Section C.

Due to the fact that it is in general not possible to learn statistics uniformly well, we formalise the notion of approximate closedness in terms of the *average* approximation error across a collection of statistics, as described below.

**Definition 4.5** (**Approximate Bellman closedness**). *A collection of statistics $s_1, \ldots, s_K$, together with an imputation strategy $\Psi$, are said to be $\varepsilon$-approximately Bellman closed for a class $\mathcal{M}$ of MDPs if, for each MDP $M = (\mathcal{X}, \mathcal{A}, p, \gamma, \mathcal{R})$ in $\mathcal{M}$ and every policy $\pi \in \mathscr{P}(\mathcal{A})^{\mathcal{X}}$, we have*

$$\sup_{(x,a) \in \mathcal{X} \times \mathcal{A}} \frac{1}{K} \sum_{k=1}^{K} |s_k(\eta_\pi(x,a)) - \hat{s}_k(x,a)| \le \varepsilon \,,$$

*where $\hat{s}_k(x,a)$ denotes the learnt value of the statistic $s_k$ for the return distribution at the state-action pair $(x,a) \in \mathcal{X} \times \mathcal{A}$.*

We can now study the approximation errors of CDRL and QDRL in light of this new concept. Whilst the analysis in Section 4.1 shows that CDRL and QDRL necessarily induce some approximation error due to lack of Bellman closedness, the following results reassuringly show that the approximation error can be made arbitrarily small by increasing the number of learnt statistics.

**Theorem 4.6.** *Consider the class $\mathcal{M}$ of MDPs with a fixed discount factor $\gamma \in [0, 1)$, and immediate reward distributions supported on $[-R_{max}, R_{max}]$. The set of statistics and imputation strategy corresponding to CDRL with evenly spaced bin locations at $-R_{max}/(1 - \gamma) = z_1 < \cdots < z_K = R_{max}/(1 - \gamma)$ is $\varepsilon$-approximately Bellman closed for $\mathcal{M}$, where $\varepsilon = \frac{\gamma}{2(1-\gamma)(K-1)}$.*

**Theorem 4.7.** *Consider the class of MDPs $\mathcal{M}$ with a fixed discount factor $\gamma \in [0, 1)$, and immediate reward distributions supported on $[-R_{max}, R_{max}]$. Then the collection of quantile statistics $s_k(\mu) = F_\mu^{-1}(\frac{2k-1}{2K})$ for $k = 1, \ldots, K$, together with the standard QDRL imputation strategy, is $\varepsilon$-approximately Bellman closed for $\mathcal{M}$, where $\varepsilon = \frac{2R_{max}(5-2\gamma)}{(1-\gamma)^2 K}$.*
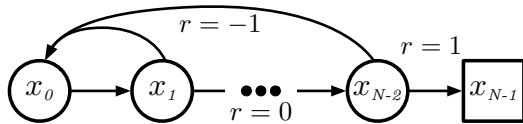
*Figure 3.* An illustration of the $N$-Chain environment.



*Figure 4.* Expectiles for state $x_0$ of the 15-Chain under policy $\pi^*$.



*Figure 5.* Expectile estimation error for varying numbers of learned expectiles and different $N$-Chain lengths.

Both of these extend existing analyses for CDRL and QDRL. In particular, Theorem 4.6 improves on the bound of Rowland et al. (2018), and Theorem 4.7 is the first approximation result for QDRL; existing results dealt solely with contraction mappings under $W_\infty$ (Dabney et al., 2017).

### 4.3. Mean consistency

So far, our discussion has been focused around *evaluation*. For *control*, it is important to correctly estimate *expected returns*, so that accurate policy improvement can be performed. We analyse to what extent expected returns are correctly learnt in existing DRL algorithms in the following result. The result for CDRL has been shown previously (Rowland et al., 2018; Lyle et al., 2019), but our proof here gives a new perspective in terms of statistics.

**Lemma 4.8.** *(i) Under CDRL updates using support locations $z_1 < \cdots < z_K$, if all approximate reward distributions have support bounded in $[z_1, z_K]$, expected returns are exactly learnt. (ii) Under QDRL updates, expected returns are not exactly learnt.*

Importantly, for EDRL, as long as the $1/2$-expectile (i.e. the mean) is included in the set of statistics, expected returns are learnt exactly; we return to this point in Section 5.2.

## 5. Experimental results

We first present results with a tabular version of EDRL to illustrate and expand upon the theoretical results presented in Sections 3 and 4. We then combine the EDRL update with a DQN-style architecture to create a novel deep RL algorithm (ER-DQN), and evaluate performance on the Atari-57 environments. We give full details of the architectures used in experiments in Appendix Section D.1.

There are several ways in which the root-finding/optimisation problems (7) and (8) may be solved in practice. In our experiments, we use a SciPy optimisation routine (Jones et al., 2001).

### 5.1. Tabular policy evaluation

We empirically validate that EDRL, which uses a sample imputation strategy, better approximates the true expectiles of a policy's return distribution as compared to the naive approach described in Section 3.1. We then show that the same is true for a variant of QDRL.
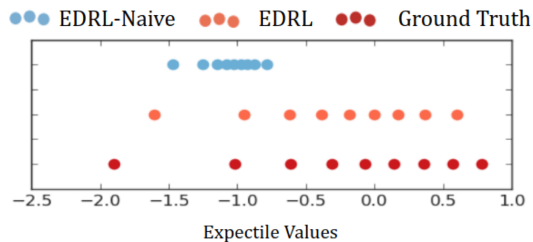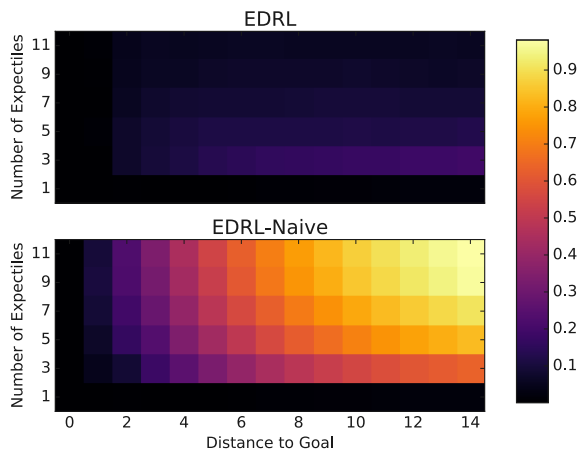
We use a variant of the classic $N$-Chain domain (see Figure 3). This environment is a one-dimensional chain of length $N$ with two possible actions at each state: (i) forward, which moves the agent right by one step with probability 0.95 and to $x_0$ with probability 0.05, and backward, which moves the agent to $x_0$ with probability 0.95 and one step to the right with probability 0.05. The reward is $-1$ when transitioning to the leftmost state, $+1$ when transitioning to the rightmost state, and zero elsewhere. Episodes begin in the leftmost state and terminate when the rightmost state is reached. The discount factor is $\gamma = 0.99$. For an $N$-Chain with length 15, we compute the return distribution of the optimal policy $\pi^*$ which selects the forward action at each state. This environment formulation induces an increasingly multimodal return distribution under the policy as the distance from the goal state increases. We compute the ground truth start state expectiles from the empirical distribution of 1,000 Monte Carlo rollouts under the policy $\pi^*$.

**EDRL.** We ran two DRL algorithms on this $N$-Chain environment: (i) EDRL, using a SciPy optimisation routine to impute target samples at each step; and (ii) *EDRL-Naive*, using the update described in Section 3.1. We learned $\{1, 3, 5, 7, 9\}$ expectiles, set the learning rate to $\alpha = 0.05$, and performed 30,000 training steps.
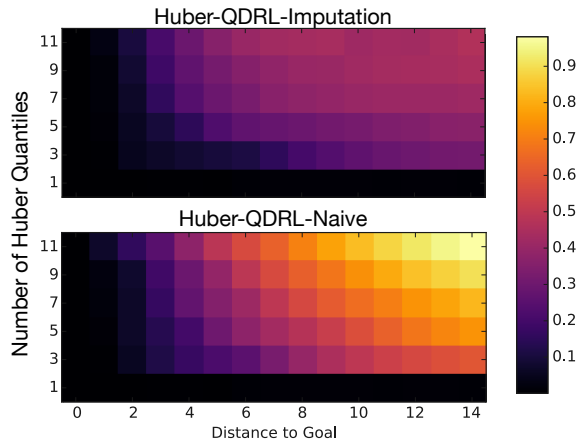
Figure 6. Huber quantile estimation error for varying numbers of learned Huber quantiles at different distances to the goal state. The environment is an $N$-Chain with $N = 15$.

In Figure 4 we illustrate the collapse of the start state expectiles learned by the EDRL-Naive algorithm with 9 expectiles, which leads to high expectile estimation error, measured as in Definition 4.5. In Figure 5, we show that this error grows as both the distance to the goal state and number of expectiles learned increase. In contrast, under EDRL these errors are much lower this error remains relatively low for varying numbers of expectiles and distances to the goal with EDRL. In Appendix E, we illustrate that this observation generalises to other return distributions in the $N$-Chain.

**QDRL.** In practical implementations, QDRL often minimises the Huber-quantile loss

$$\operatorname*{argmin}_{q \in \mathbb{R}} \mathbb{E}_{Z \sim \mu}[(\tau \mathbb{1}_{Z>q} + (1-\tau)\mathbb{1}_{Z<q}) H_\kappa(Z-q)], \quad (11)$$

rather than the quantile loss (4) for numerical stability, where $H_\kappa$ is the Huber loss function with width parameter $\kappa$, as in Dabney et al. (2017) (we set $\kappa = 1$). As with naive EDRL, simply replacing the quantile regression loss in QDRL with Expression (11) conflates samples and statistics, leading to worse approximation of the distribution. We propose a new algorithm for learning Huber quantiles, *Huber-QDRL-Imputation*, that incorporates an imputation strategy by solving an optimisation problem analogous to (8) in the case of the Huber quantile loss. In Figure 6, we compare this to *Huber-QDRL-Naive*, the standard algorithm for learning Huber quantiles, on the $N$-chain environment. As in the case of expectiles, the Huber quantile estimation error is vastly reduced when using an imputation strategy.

### 5.2. Tabular control

In Section 4.3 we argued for the importance of mean consistency. In Figure 7a we give a simple, five state, MDP in which the learned control policy is directly affected by mean
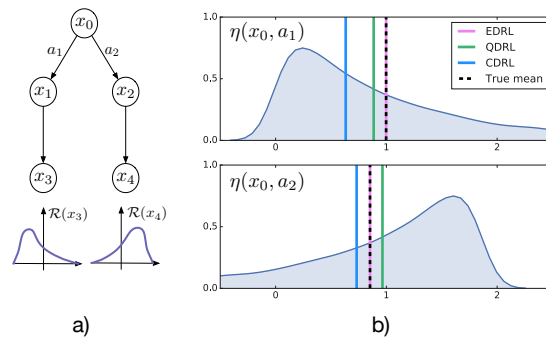


Figure 7. (a) 5-state MDP, reward is zero everywhere except at the terminal states $x_3$ and $x_4$ which have stochastic rewards. (b) We show the true return distributions $\eta(x_0, a_1)$ and $\eta(x_0, a_2)$, and the expected returns estimated by CDRL, QDRL, and EDRL.

consistency. At start state $x_0$ the agent has the choice of two actions, leading down two paths and culminating in two different reward distributions. The rewards at terminal states $x_3$ and $x_4$ are sampled from (shifted) exponential distributions with densities $e^{-\lambda}$ $(\lambda \geq 0)$ and $e^{\lambda+1.85}$ $(\lambda \leq 1.85)$, respectively. Transitions are deterministic, and $\gamma = 1$. For CDRL, we take bin locations at $(z_1, z_2, z_3) = (0, 1, 2)$.

Figure 7b shows the true return distributions, their expectations, and the means estimated by CDRL, QDRL and EDRL. Due to a lack of mean consistency both CDRL and QDRL learn a sub-optimal greedy policy. For CDRL, this is due to the true return distributions having support outside $[0, 2]$, and for QDRL, this is due to the quantiles not capturing tail behaviour. In contrast, EDRL correctly learns the means of both return distributions, and so is able to act optimally.

### 5.3. Expectile regression DQN

To demonstrate the effectiveness of EDRL at scale, we combine the EDRL update in Algorithm 2 with the architecture of QR-DQN to obtain a new deep RL agent, expectile regression DQN (ER-DQN). Precise details of the architecture, training algorithm, and environments are given in Appendix Section D. We evaluate ER-DQN on a suite of 57 Atari games using the Arcade Learning Environment (Bellemare et al., 2013). In Figure 8, we plot mean and median human normalised scores for ER-DQN with 11 atoms, and compare against DQN, QR-DQN (which learns 200 Huber quantile statistics), and a naive implementation of ER-DQN that doesn't use an imputation strategy, learning 201 expectiles. All methods were re-run for this paper, and results were averaged over 3 seeds. In practice, we found that with 11 expectiles, ER-DQN already offers strong performance relative to these other approaches, and that with this number of statistics, the additional training overhead due to the SciPy optimiser calls is low.
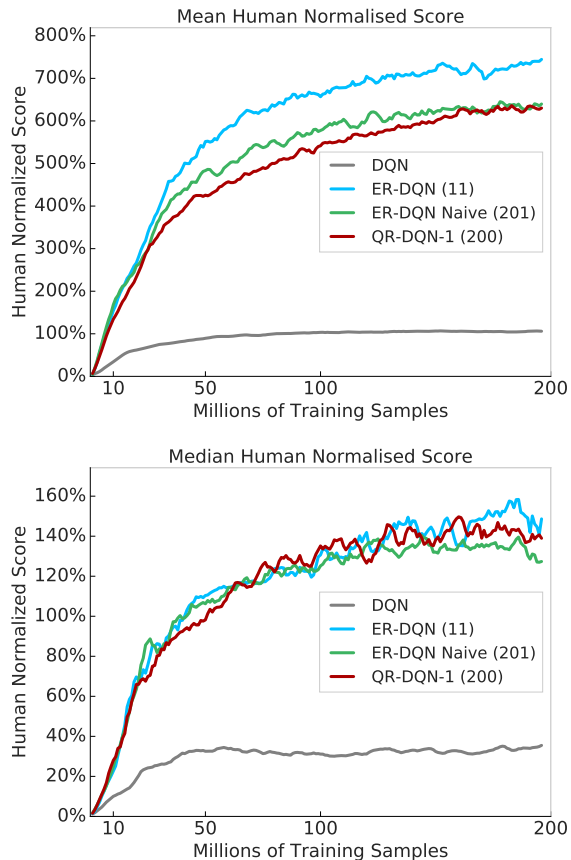
Figure 8. Mean and median human normalised scores across all 57 Atari games. Number of statistics learnt for each algorithm indicated in parentheses.

In terms of mean human normalised score, ER-DQN represents a substantial improvement over both QR-DQN and the naive version of ER-DQN that does not use an imputation strategy. We hypothesise that the mean consistency of EDRL (in contrast to other DRL methods; see Section 4.3) is partially responsible for these improvements, and leave further investigation of the role of mean consistency in DRL as a direction for future work. We also remark that the performance of ER-DQN shows that there may be significant practical value in applying the framework developed in this paper to other families of statistics. It remains to be seen if the presence of partial observability may induce non-trivial distributions, which could also explain ER-DQN's improved performance in some games. Investigation into the robustness of ER-DQN with regards to the precise imputation strategy used is also a natural question for future work.

## 6. Conclusion

We have developed a unifying framework for DRL in terms of statistical estimators and imputation strategies. Through this framework, we have developed a new algorithm, EDRL, as well as proposing algorithmic adjustments to an existing approach. We have also used this framework to define the notion of Bellman closedness, and provided new approximation guarantees for existing algorithms.

This paper also opens up several avenues for future research. Firstly, the framework of imputation strategies has the potential to be applied to a wide range of collections of statistics, opening up a large space of new algorithms to explore. Secondly, our analysis has shown that a lack of Bellman closedness necessarily introduces a source of approximation error into many DRL algorithms; it will be interesting to see how this interacts with errors introduced by function approximation. Finally, we have focused on DRL algorithms that can be interpreted as learning a finite collection of statistics in this paper. One notable alternative is implicit quantile networks (Dabney et al., 2018), which attempt to learn an uncountable collection of quantiles with a finite-capacity function approximator; it will also be interesting to extend our analysis to this setting.

## Acknowledgements

# References

Barth-Maron, G., Hoffman, M. W., Budden, D., Dabney, W., Horgan, D., TB, D., Muldal, A., Heess, N., and Lillicrap, T. Distributional policy gradients. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2018.

Bellemare, M. G., Naddaf, Y., Veness, J., and Bowling, M. The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279, 2013.

Bellemare, M. G., Dabney, W., and Munos, R. A distributional perspective on reinforcement learning. *Proceedings of the International Conference on Machine Learning (ICML)*, 2017.

Bellman, R. *Dynamic Programming*. Princeton University Press, 1st edition, 1957.

Dabney, W., Rowland, M., Bellemare, M. G., and Munos, R. Distributional reinforcement learning with quantile regression. *Proceedings of the AAAI Conference on Artificial Intelligence*, 2017.

Dabney, W., Ostrovski, G., Silver, D., and Munos, R. Implicit quantile networks for distributional reinforcement learning. *Proceedings of the International Conference on Machine Learning (ICML)*, 2018.

Engert, M. Finite dimensional translation invariant subspaces. *Pacific Journal of Mathematics*, 32(2):333–343, 1970.

Gruslys, A., Dabney, W., Azar, M. G., Piot, B., Bellemare, M., and Munos, R. The reactor: A fast and sample-efficient actor-critic agent for reinforcement learning. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2018.

Huber, P. J. and Ronchetti, E. *Robust Statistics*. Wiley New York, 2nd edition, 2009.

Jones, E., Oliphant, T., and Peterson, P. SciPy: Open source scientific tools for Python, 2001. URL http://www.scipy.org/.

Lattimore, T. and Hutter, M. PAC bounds for discounted MDPs. In *International Conference on Algorithmic Learning Theory (ALT)*, 2012.

Lyle, C., Castro, P. S., and Bellemare, M. G. A comparative analysis of expected and distributional reinforcement learning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 2019.

Morimura, T., Sugiyama, M., Kashima, H., Hachiya, H., and Tanaka, T. Nonparametric return distribution approximation for reinforcement learning. *Proceedings of the International Conference on Machine Learning (ICML)*, 2010a.

Morimura, T., Sugiyama, M., Kashima, H., Hachiya, H., and Tanaka, T. Parametric return density estimation for reinforcement learning. *Proceedings of the Conference on Uncertainty in Artificial Intelligence (UAI)*, 2010b.

Newey, W. K. and Powell, J. L. Asymmetric least squares estimation and testing. *Econometrica: Journal of the Econometric Society*, pp. 819–847, 1987.

Qu, C., Mannor, S., and Xu, H. Nonlinear distributional gradient temporal-difference learning. *arXiv preprint arXiv:1805.07732*, 2018.

Rowland, M., Bellemare, M. G., Dabney, W., Munos, R., and Teh, Y. W. An analysis of categorical distributional reinforcement learning. *Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2018.

Sobel, M. J. The variance of discounted Markov decision processes. *Journal of Applied Probability*, 19(4):794–802, 1982.

Sutton, R. S. and Barto, A. G. *Reinforcement Learning: An Introduction*. MIT Press, 2018.

Zhang, S., Mavrin, B., Yao, H., Kong, L., and Liu, B. QUOTA: The quantile option architecture for reinforcement learning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 2019.

# Appendices

# A. Distributional reinforcement learning algorithms

For completeness, we give full descriptions of CDRL and QDRL algorithms in this section, complementing the details given in Section 2.2. We also summarise CDRL, QDRL, the exact approach to distributional RL, and our proposed algorithm EDRL, in Figure 10 at the end of this section.

## A.1. The distributional Bellman operator

In accordance with the distributional Bellman equation (3), the distributional Bellman operator $\mathcal{T}^\pi : \mathscr{P}(\mathbb{R})^{\mathcal{X} \times \mathcal{A}} \to \mathscr{P}(\mathbb{R})^{\mathcal{X} \times \mathcal{A}}$ is defined by Bellemare et al. (2017) as

$$(\mathcal{T}^\pi \eta)(x, a) = \mathbb{E}_\pi \left[ (f_{R_0, \gamma})_\# \eta(X_1, A_1) | X_0 = x, A_0 = a \right] ,$$

for all $\eta \in \mathscr{P}(\mathbb{R})^{\mathcal{X} \times \mathcal{A}}$.

## A.2. Categorical distributional reinforcement learning

As described in Section 2.2, CDRL algorithms are an approach to distributional RL that restrict approximate distributions to the parametric family of the form $\{\sum_{k=1}^{K} p_k \delta_{z_k} | \sum_{k=1}^{K} p_k = 1, \ p_k \geq 0 \forall k\} \subseteq \mathscr{P}(\mathbb{R})$, where $z_1 < \cdots < z_K$ are an evenly spaced, fixed set of supports. For evaluation of a policy $\pi : \mathcal{X} \to \mathscr{P}(\mathcal{A})$, given a collection of approximations $(\eta(x, a) | (x, a) \in \mathcal{X} \times \mathcal{A})$, the approximation at $(x, a) \in \mathcal{X} \times \mathcal{A}$ is updated according to:

$$\eta(x, a) \leftarrow \Pi_{\mathcal{C}} \mathbb{E}_\pi \left[ (f_{R_0, \gamma})_\# \eta(X_1, A_1) | X_0 = x, A_0 = a \right] .$$

Here, $\Pi_{\mathcal{C}} : \mathscr{P}(\mathbb{R}) \to \mathscr{P}(\{z_1, \ldots, z_K\})$ is a projection operator defined for a single Dirac delta as

$$\Pi_{\mathcal{C}}(\delta_w) = \begin{cases} \delta_{z_1} & w \leq z_1 \\ \frac{w - z_{k+1}}{z_k - z_{k+1}} \delta_{z_k} + \frac{z_k - w}{z_k - z_{k+1}} \delta_{k+1} & z_k \leq w \leq z_{k+1} \\ \delta_{z_K} & w \geq z_K , \end{cases} \tag{12}$$

and extended affinely and continuously. In the language of operators, the CDRL update may be neatly described as $\eta \leftarrow \Pi_{\mathcal{C}} \mathcal{T}^\pi \eta$, where we abuse notation by interpreting $\Pi_{\mathcal{C}}$ as an operator on collections of distributions indexed by state-action pairs, applying the transformation in Expression (12) to each distribution. The supremum-Cramér distance is defined as

$$\bar{\ell}_2(\eta_1, \eta_2) = \sup_{(x,a) \in \mathcal{X} \times \mathcal{A}} \ell_2(\eta_1(x, a), \eta_2(x, a)) = \sup_{(x,a) \in \mathcal{X} \times \mathcal{A}} \left( \int_{\mathbb{R}} |F_{\eta_1(x,a)}(t) - F_{\eta_2(x,a)}(t)|^2 \mathrm{d}t \right)^{\frac{1}{2}} .$$

for all $\eta_1, \eta_2 \in \mathscr{P}(\mathbb{R})^{\mathcal{X} \times \mathcal{A}}$, where for any $\mu \in \mathscr{P}(\mathbb{R})$, $F_\mu$ denotes the CDF of $\mu$. The operator $\Pi_{\mathcal{C}} \mathcal{T}^\pi$ is a $\sqrt{\gamma}$-contraction in the supremum-Cramér distance, and so by the contraction mapping theorem, repeated CDRL updates converge to a unique limit point, regardless of the initial approximate distributions. For more details on these results and further background, see Bellemare et al. (2017); Rowland et al. (2018).

**Stochastic approximation.** The update $\eta \leftarrow \Pi_{\mathcal{C}} \mathcal{T}^\pi \eta$ is typically not computable in practice, due to unknown/intractable dynamics. An unbiased approximation to $(\mathcal{T}^\pi \eta)(x, a)$ may be obtained by interacting with the environment to obtain a transition $(x, a, r, x', a')$, and computing the target

$$(f_{r, \gamma})_\# \eta(x', a') .$$

It can be shown (Rowland et al., 2018) that the following is an unbiased estimator for the CDRL update $(\Pi_{\mathcal{C}} \mathcal{T}^\pi \eta)(x, a)$:

$$\Pi_{\mathcal{C}}(f_{r, \gamma})_\# \eta(x', a') .$$

Finally, the current estimate $\eta(x, a)$ can be moved towards the stochastic target by following the (semi-)gradient of some loss, in analogy with semi-gradient methods in classical RL. Bellemare et al. (2017) consider the KL loss

$$\mathrm{KL}(\Pi_{\mathcal{C}}(f_{r, \gamma})_\# \eta(x', a') \| \eta(x, a)) ,$$

and update $\eta(x, a)$ by taking the gradient of the loss through the second argument with respect to the parameters $p_{1:K}(x, a)$. Other losses, such as the Cramér distance, may also be considered (Rowland et al., 2018).

**Control.** All variants of CDRL for evaluation may be modified to become control algorithms. This is achieved by adjusting the distribution of the action $A_1$ in the backup in an analogous way to classical RL algorithms. Instead of having $A_1 \sim \pi(\cdot|X_1)$, we instead select $A_1$ based on the currently estimated expected returns for each of the actions at the state $X_1$. For Q-learning-style algorithms, the action corresponding to the highest estimated expected return is selected:

$$A_1 = \arg\max_{a \in \mathcal{A}} \mathbb{E}_{Z \sim \eta(X_1, a)}[Z] \ .$$

However, other choices are possible, such as SARSA-style $\varepsilon$-greedy action selection.

## A.3. Quantile distributional reinforcement learning

As described in Section 2.2, QDRL algorithms are an approach to distributional RL that restrict approximate distributions to the parametric family of the form $\{\frac{1}{K} \sum_{k=1}^{K} \delta_{z_k} | z_{1:K} \in \mathbb{R}^K\} \subseteq \mathscr{P}(\mathbb{R})$. For evaluation of a policy $\pi : \mathcal{X} \to \mathscr{P}(\mathcal{A})$, given a collection of approximations $(\eta(x, a)|(x, a) \in \mathcal{X} \times \mathcal{A})$, the approximation at $(x, a) \in \mathcal{X} \times \mathcal{A}$ is updated according to:

$$\eta(x, a) \leftarrow \Pi_{W_1} \mathbb{E}_\pi \left[ (f_{R_0, \gamma})_\# \eta(X_1, A_1) | X_0 = x, A_0 = a \right] \ , .$$

Here, $\Pi_{W_1} : \mathscr{P}(\mathbb{R}) \to \mathscr{P}(\mathbb{R})$ is a projection operator defined by

$$\Pi_{\mathcal{C}}(\mu) = \frac{1}{K} \sum_{k=1}^{K} \delta_{F_\mu^{-1}(\tau_k)} \ ,$$

where $\tau_k = \frac{2k-1}{2K}$, and $F_\mu$ is the CDF of of $\mu$. As noted in Section 2.2, $F_\mu^{-1}(\tau)$ may also be characterised as the minimiser (over $q \in \mathbb{R}$) of the quantile regression loss $\mathrm{QR}(q; \mu, \tau) = \mathbb{E}_{Z \sim \mu} [[\tau \mathbb{1}_{Z > q} + (1 - \tau) \mathbb{1}_{Z \leq q}] |Z - q|]$; this perspective turns out to be crucial in deriving a stochastic approximation version of the algorithm.

**Stochastic approximation.** As for CDRL, the update $\eta \leftarrow \Pi_{W_1} \mathcal{T}^\pi \eta$ is typically not computable in practice, due to unknown/intractable dynamics. Instead, a stochastic target may be computed by using a transition $(x, a, r, x', a')$, and updating each atom location $z_k(x, a)$ at the current state-action pair $(x, a)$ by following the gradient of the QR loss:

$$\nabla_q \mathrm{QR}(q; (f_{r, \gamma})_\# \eta(x', a'), \tau_k)\big|_{q = z_k(x, a)} \ .$$

Because the QR loss is affine in its second argument, this yields an unbiased estimator of the true gradient

$$\nabla_q \mathrm{QR}(q; (\mathcal{T}^\pi \eta)(x, a), \tau_k)\big|_{q = z_k(x, a)} \ .$$

**Control.** The methods for evaluation described above may be modified to yield control methods in exactly the same as described for CDRL in Section A.2.

## A.4. Quantiles versus expectiles

Quantiles of a distribution are given by the inverse of the cumulative distribution function. As such, they fundamentally represent threshold values for the cumulative probabilities. That is, the quantile at $\tau$, $q_\tau$, is greater than or equal to $\tau \times 100\%$ of the outcome values. In contrast, expectiles also take into account the *magnitude* of outcomes; the expectile at $\tau$, $e_\tau$, is such that the expectation of the deviations below $e_\tau$ of the random variable $Z$ is equal to $\frac{\tau}{1 - \tau}$ of the expectation of the deivations above $e_\tau$. We illustrate these points in Figure 9.
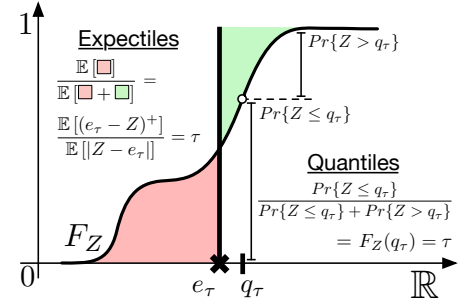


*Figure 9.* Diagram illustrating the similarities and differences of quantiles and expectiles.
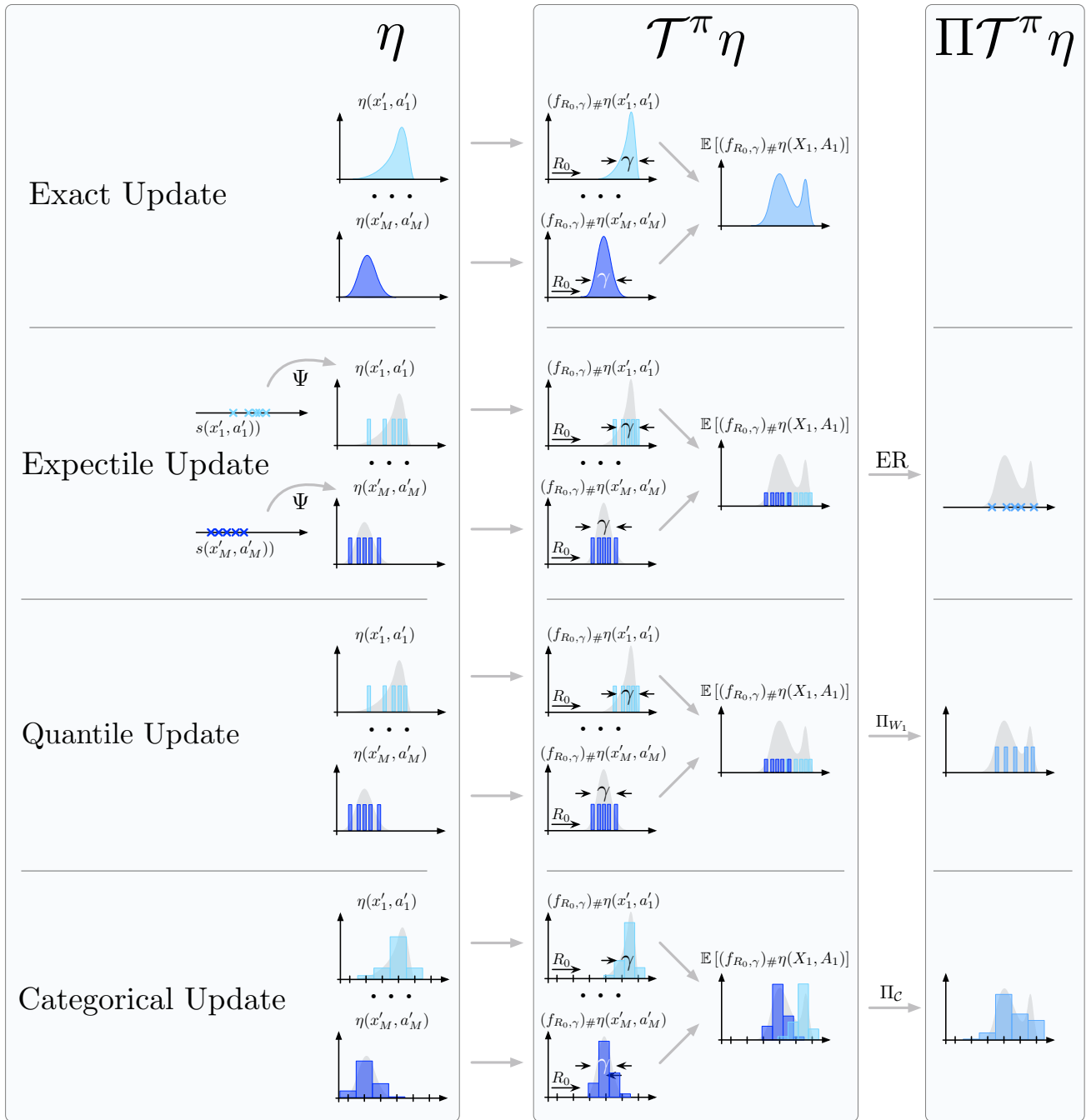
*Figure 10.* Illustration of distributional RL, with exact updates, expectile updates (EDRL), quantile updates (QDRL), and categorical updates (CDRL).

# B. Proofs

## B.1. Proofs of results from Section 3

**Lemma 3.2.** *CDRL updates, with distributions supported on $z_1 < \ldots < z_K$, can be interpreted as learning the values of the following statistics of return distributions:*

$$s_{z_k, z_{k+1}}(\mu) = \mathbb{E}_{Z \sim \mu}\left[h_{z_k, z_{k+1}}(Z)\right] \ for \ k = 1, \ldots, K-1,$$

*where for $a < b$, $h_{a,b} : \mathbb{R} \to \mathbb{R}$ is a piecewise linear function defined so that $h_{a,b}(x)$ is equal to 1 for $x \leq a$, equal to 0 for $x \geq b$, and linearly interpolating between $h_{a,b}(a)$ and $h_{a,b}(b)$ for $x \in [a, b]$.*

*Proof.* We first observe that the projection operator $\Pi_{\mathcal{C}}$, defined in Section A.2, preserves each of the statistics $s_{z_1,z_2}, \ldots, s_{z_{K-1},z_K}$, in the sense that for any distribution $\mu$, we have $s_{z_k,z_{k+1}}(\mu) = s_{z_k,z_{k+1}}(\Pi_{\mathcal{C}}\mu)$ for all $k = 1, \ldots, K$. Secondly, we observe that that the map $\{\sum_{k=1}^K p_k \delta_{z_k} | \sum_{k=1}^K p_k = 1, \ p_k \geq 0 \forall k\} \ni \mu \mapsto (s_{z_1,z_2}(\mu), \ldots, s_{z_{K-1},z_K}(\mu)) \in \mathbb{R}^{K-1}$ is injective; each distribution has a unique vector of statistics. Thus, CDRL can indeed be interpreted as learning precisely the set of statistics $s_{z_1,z_2}, \ldots, s_{z_{K-1},z_K}$. $\qquad\square$

## B.2. Proofs of results from Section 4.1

**Lemma 4.2.** *For each $K \in \mathbb{N}$, the set of statistics consisting of the first $K$ moments is Bellman closed.*

*Proof.* We begin by introducing notation. Let $s_k : \mu \mapsto \mathbb{E}_{Z \sim \mu}\left[Z^k\right]$ be the $k^{\text{th}}$ moment functional, for $k = 1, \ldots, K$. We now compute

$$\begin{aligned}
s_k(\eta_\pi(x,a)) &= \mathbb{E}_{Z \sim \eta_\pi(x,a)}\left[Z^k\right] \\
&= \sum_{(x',a') \in \mathcal{X} \times \mathcal{A}} \int_{\mathbb{R}} \mathcal{R}(\mathrm{d}r|x,a)p(x'|x,a)\pi(a'|x')\mathbb{E}_{Z \sim \eta_\pi(x',a')}\left[(r + \gamma Z)^k\right] \\
&= \sum_{(x',a') \in \mathcal{X} \times \mathcal{A}} \int_{\mathbb{R}} \mathcal{R}(\mathrm{d}r|x,a)p(x'|x,a)\pi(a'|x') \sum_{m=0}^k \binom{k}{m} \gamma^{k-m}\mathbb{E}_{Z \sim \eta_\pi(x',a')}\left[Z^{k-m}\right] r^m \\
&= \sum_{(x',a') \in \mathcal{X} \times \mathcal{A}} \int_{\mathbb{R}} \mathcal{R}(\mathrm{d}r|x,a)p(x'|x,a)\pi(a'|x') \sum_{m=0}^k \binom{k}{m} \gamma^{k-m}s_{k-m}(\eta_\pi(x',a'))r^m \\
&= \mathbb{E}\left[\sum_{m=0}^k \binom{k}{m}\gamma^{k-m}s_{k-m}(\eta_\pi(X_1,A_1))R_0^m\Bigg| X_0 = x, A_0 = a\right].
\end{aligned}$$

Thus, $s_k(\eta_\pi(x,a))$ can be expressed in terms of $R_0$ and $s_{1:K}(\eta_\pi(X_1,A_1))$, as required. $\qquad\square$

**Theorem 4.3.** *The only finite sets of statistics of the form $s(\mu) = \mathbb{E}_{Z \sim \mu}[h(Z)]$ that are Bellman closed are given by collections of statistics $s_1, \ldots, s_K : \mathscr{P}(\mathbb{R}) \to \mathbb{R}$ with the property that the linear span $\{\sum_{k=0}^K \alpha_k s_k | \alpha_k \in \mathbb{R} \ \forall k\}$ is equal to the linear span of the set of moment functionals $\{\mu \mapsto \mathbb{E}_{Z \sim \mu}\left[Z^l\right] | l = 0, \ldots, L\}$, for some $L \leq K$, where $s_0$ is the constant functional equal to 1.*

*Proof.* Suppose $s_1, \ldots, s_K : \mathscr{P}(\mathbb{R}) \to \mathbb{R}$ form a Bellman closed set of statistical functionals of the form $s_k(\mu) = \mathbb{E}_{Z \sim \mu}\left[h_k(Z)\right]$ for some measurable $h_k : \mathbb{R} \to \mathbb{R}$, for each $k = 1, \ldots, K$. Now note that for any MDP $(\mathcal{X}, \mathcal{A}, p, \gamma, \mathcal{R})$, we have the following equation:

$$s_k(\eta_\pi(x,a)) = \sum_{(x',a') \in \mathcal{X} \times \mathcal{A}} \int_{\mathbb{R}} \mathcal{R}(\mathrm{d}r|x,a)p(x'|x,a)\pi(a'|x')s_k((f_{r,\gamma})_{\#}\eta_\pi(x',a')),$$

for all $(x,a) \in \mathcal{X} \times \mathcal{A}$, and for each $k = 1, \ldots, K$. By assumption of Bellman closedness, the right-hand side of this equation may be written as a function of $\mathcal{R}(x,a)$, $\gamma$, and the collection of statistics $(s_{1:K}(\eta_\pi(x',a'))|(x',a') \in \mathcal{X} \times \mathcal{A})$. Since this must hold across all valid sets of return distributions, it must the case that each $s_k((f_{r,\gamma})_{\#}\eta_\pi(x',a'))$ may be

written as a function of $r$, $\gamma$ and $s_{1:K}(\eta_\pi(x', a'))$; we will write $s_k((f_{r,\gamma})_{\#}\eta_\pi(x', a')) = g(r, \gamma, s_{1:K}(\eta_\pi(x', a')))$ for some $g$.

We next claim that $g(r, \gamma, s_{1:K}(\eta_\pi(x', a')))$ is affine in $s_{1:K}(\eta_\pi(x', a'))$. To see this, note that both $s_k((f_{r,\gamma})_{\#}\eta_\pi(x', a'))$ and $s_{1:K}(\eta_\pi(x', a'))$ are affine as functions of the distribution $\eta_\pi(x', a')$, by assumption on the form of the statistics $s_{1:K}$. Therefore $g(r, \gamma, \cdot)$ too is affine on the (convex) codomain of $s_{1:K}$.

Thus, we have

$$\mathbb{E}_{Z \sim \eta_\pi(x', a')}\left[h_k(r + \gamma Z)\right] = a_0(r, \gamma) + \sum_{k'=1}^{K} a_{k'}(r, \gamma)\mathbb{E}_{Z \sim \eta_\pi(x', a')}\left[h_{k'}(Z)\right], \tag{13}$$

for some functions $a_{0:K} : \mathbb{R} \times [0, 1) \to \mathbb{R}$. By taking $\eta_\pi(x', a')$ to be a Dirac delta at an arbitrary real number, we obtain

$$h_k(r + \gamma x) = a_0(r, \gamma) + \sum_{k'=1}^{K} a_{k'}(r, \gamma)h_{k'}(x) \quad \text{for all } x \in \mathbb{R}. \tag{14}$$

In particular, the function $h_k(\gamma x)$ lies in the span of the functions $h_1, \ldots, h_K, \mathbb{1}$, where $\mathbb{1}$ is the constant function at $1$. Further, $h_k(r + \gamma x)$ lies in this span for all $r \in \mathbb{R}$, and so the collection of functions $\{x \mapsto h_k(r + \gamma x)|r \in \mathbb{R}\}$ lies in a finite-dimensional subspace of functions. We may now appeal to Theorem 1 of Engert (1970), which states that any finite-dimensional space of functions which is closed under translation is spanned by a set of functions of the form

$$\bigcup_{j=1}^{J}\{x \mapsto x^\ell \exp(\lambda_j x) \mid 0 \le \ell \le L_j\}, \tag{15}$$

for some finite subset $\{\lambda_1, \ldots, \lambda_J\}$ of $\mathbb{C}$. From this, we deduce that each function $x \mapsto h_k(x)$ may be expressed as a linear combination of functions of the form appearing in the set in expression (15). Further, enforcing the condition that the linear span must be closed under composition with $f_{r,\gamma}$ with $\gamma \in [0, 1)$ rules out any values of $\lambda_j$ above which are not zero. Therefore, the linear span of the functions $h_1, \ldots, h_K, \mathbb{1}$ must be equal to the span of some set of monomials $x \mapsto x^\ell$, $0 \le \ell \le L$, for some $L \in \mathbb{N}$, and hence the statement of the theorem follows. $\square$

**Lemma 4.4.** *The sets of statistics learnt under (i) CDRL, and (ii) QDRL, are not Bellman closed.*

*Proof.* (i) This follows as a special case of Theorem 4.3, since the statistics learnt by CDRL are expectations, as shown in Lemma 3.2.

(ii) Quantiles cannot be expressed as expectations, and so we cannot appeal to Theorem 4.3. We instead proceed by describing a concrete counterexample to Bellman closedness. Fix a number $K \in \mathbb{N}$ of quantiles. Consider an MDP with a single action, and an initial state $x_0$ which transitions to one of two terminal states $x_1$, $x_2$ with equal probability. Suppose there is no immediate reward at state $x_0$. We consider two different possibilities for reward distributions at states $x_1$, $x_2$, and show that these two possibilities yield the same quantiles for the return distributions at states $x_1$ and $x_2$, but different quantiles for the return distribution at state $x_0$; thus demonstrating that finite sets of quantiles are not Bellman closed.

Firstly, suppose rewards are drawn from $\mathrm{Unif}([0, 1])$ at state $x_1$ and $\mathrm{Unif}([1/K, 1 + 1/K])$ at $x_2$, so that the $\frac{2k-1}{2K}$-quantile of the return at states $x_1$ and $x_2$ are $\frac{2k-1}{2K}$ and $\frac{2k+1}{2K}$, for each $k = 1, \ldots, K$. Then the return distribution at state $x_0$ is the mixture $\frac{1}{2}\mathrm{Unif}([0, \gamma]) + \frac{1}{2}\mathrm{Unif}([\gamma/K, \gamma + \gamma/K])$, and hence the $\frac{1}{2K}$-quantile is $\frac{\gamma}{K}$. Now, suppose instead that the reward distribution at state $x_1$ is $\frac{1}{K}\sum_{k=1}^{K}\delta_{\frac{2k-1}{2K}}$ and the reward distribution at state $x_2$ is $\frac{1}{K}\sum_{k=1}^{K}\delta_{\frac{2k+1}{2K}}$. Then the $\frac{1}{2K}$-quantile of the return distribution at state $x_0$ is $\frac{3\gamma}{2K}$. $\square$

### B.3. Proofs of results from Section 4.2

In this section, we use operator notation reviewed in Section A. In both proofs, the supremum-Wasserstein distance will be of use, defined as $\overline{W}_1(\mu_1, \mu_2) = \sup_{(x,a) \in \mathcal{X} \times \mathcal{A}} W_1(\mu_1(x, a), \mu_2(x, a))$ for all $\mu_1, \mu_2 \in \mathscr{P}_1(\mathbb{R})^{\mathcal{X} \times \mathcal{A}}$. Before proving theorem 4.6, we state and prove an auxiliary lemma.

**Lemma B.1.** *Let $\Pi_\mathcal{C}$ be the Cramér projection for equally-spaced support points $z_1 < \cdots < z_K$, defined in Appendix Section A.2. (i) $\Pi_\mathcal{C}$ is a non-expansion in $W_1$. (ii) For any distribution $\mu \in \mathscr{P}(\mathbb{R})$ supported on $[z_1, z_K]$, we have $W_1(\Pi_\mathcal{C}\mu, \mu) \le \frac{z_K - z_1}{2(K-1)}$.*

*Proof.* In the proof of the first claim, we use the following characterisation of the Cramér projection (Rowland et al., 2018). For any distribution $\mu \in \mathscr{P}(\mathbb{R})$ with CDF $F_\mu$, the CDF of $\Pi_{\mathcal{C}}\mu$ is given by $F_{\Pi_{\mathcal{C}}\mu}(v) = \frac{1}{z_{k+1}-z_k} \int_{z_k}^{z_{k+1}} F_\mu(t)\mathrm{d}t$ for $v \in [z_k, z_{k+1})$, $k = 1, \ldots, K-1$, with $F_{\Pi_{\mathcal{C}}\mu}$ equal to 0 on $(\infty, z_1)$ and equal to 1 on $[z_K, \infty)$.

(i) Let $\mu_1, \mu_2 \in \mathscr{P}(\mathbb{R})$. We compute

$$W_1(\mu_1, \mu_2) \geq \sum_{k=1}^{K-1} \int_{z_k}^{z_{k+1}} |F_{\mu_1}(t) - F_{\mu_1}(t)|\mathrm{d}t \geq \sum_{k=1}^{K} (z_{k+1} - z_k)|F_{\Pi_{\mathcal{C}}\mu_1}(z_k) - F_{\Pi_{\mathcal{C}}\mu_1}(k)| = W_1(\Pi_{\mathcal{C}}\mu_1, \Pi_{\mathcal{C}}\mu_2),$$

as required. The first inequality comes from expressing the Wasserstein distance between two distributions as the $L^1$ distance between their CDFs, and truncating the corresponding integral at $z_1$ and $z_K$. The second inequality follows from Jensen's inequality.

(ii) We first introduce some notation. Let $l, u : [z_1, z_K] \rightarrow \{z_1, \ldots, z_K\}$ be functions such that $l(y)$ is the largest element of $\{z_1, \ldots, z_K\}$ which is less than or equal to $y$, and $u(y)$ is the smallest element of $\{z_1, \ldots, z_K\}$ which is greater than or equal to $y$, for all $y \in [z_1, z_K]$. A valid coupling between $\mu$ and $\Pi_{\mathcal{C}}$ is then given as follows. Let $Y \sim \mu$, and conditional on $Y$, let $p \sim \text{Bernoulli}\left(\frac{Y-l(Y)}{u(Y)-l(Y)}\right)$ if $Y \notin \{z_1, \ldots, z_K\}$, and $p = 1$ almost surely conditional on $Y \in \{z_1, \ldots, z_K\}$. Then define $Z = pl(Y) + (1-p)u(Y)$. It is straightforward to check that the marginal distribution of $Z$ is $\Pi_{\mathcal{C}}\mu$, and we can straightforwardly upper-bound the transport cost associated with this coupling, by observing that for each possible value $y$ of $Y$, the contribution to the transport cost is 0 if $y \in \{z_1, \ldots, z_K\}$, and $\frac{u(y)-y}{u(y)-l(y)}(y - l(y)) + \frac{y-l(y)}{u(y)-l(y)}(u(y) - y) \leq \frac{u(y)-l(y)}{2} = \frac{z_K-z_1}{2(K-1)}$. Therefore, integrating over the distribution of $Y$ gives a transport cost of at most $\frac{z_K-z_1}{2(K-1)}$, which gives the required bound on the Wasserstein distance. $\square$

**Theorem 4.6.** *Consider the class $\mathcal{M}$ of MDPs with a fixed discount factor $\gamma \in [0, 1)$, and immediate reward distributions supported on $[-R_{max}, R_{max}]$. The set of statistics and imputation strategy corresponding to CDRL with evenly spaced bin locations at $-R_{max}/(1-\gamma) = z_1 < \cdots < z_K = R_{max}/(1-\gamma)$ is $\varepsilon$-approximately Bellman closed for $\mathcal{M}$, where $\varepsilon = \frac{\gamma}{2(1-\gamma)(K-1)}$.*

*Proof.* For the CDRL statistics, we have $s_{z_k, z_{k+1}}(\eta_\pi(x, a)) = s_{z_k, z_{k+1}}(\Pi_{\mathcal{C}}\eta_\pi(x, a))$ for $k = 1, \ldots, K$ and all $(x, a) \in \mathcal{X} \times \mathcal{A}$. Further, since $\Pi_{\mathcal{C}}\eta_\pi(x, a)$ is supported on $\{z_1, \ldots, z_K\}$ for all $(x, a) \in \mathcal{X} \times \mathcal{A}$, we have that $s_{z_k, z_{k+1}}(\Pi_{\mathcal{C}}\eta_\pi(x, a)) = F_{\Pi_{\mathcal{C}}\eta_\pi(x,a)}^{-1}(z_k)$. Let $(\eta(x, a)|(x, a) \in \mathcal{X} \times \mathcal{A})$ be the set of approximate distributions learnt by CDRL. As noted in Appendix Section A.2, $\eta$ is the fixed point of the projected Bellman operator $\Pi_{\mathcal{C}}\mathcal{T}^\pi$, and $\eta_\pi$ is the fixed point of the Bellman operator

$\mathcal{T}^\pi$. We now compute:

$$\frac{1}{K-1}\sum_{k=1}^{K-1}\left|s_{z_k,z_{k+1}}(\eta(x,a)) - s_{z_k,z_{k+1}}(\eta_\pi(x,a))\right|$$

$$=\frac{1}{K-1}\sum_{k=1}^{K-1}\left|s_{z_k,z_{k+1}}(\eta(x,a)) - s_{z_k,z_{k+1}}(\Pi_\mathcal{C}\eta_\pi(x,a))\right|$$

$$=\frac{1}{K-1}\sum_{k=1}^{K-1}\left|F^{-1}_{\eta(x,a)}(z_k) - F^{-1}_{\Pi_\mathcal{C}\eta_\pi(x,a)}(z_k)\right|$$

$$=\frac{1}{2R_{\max}/(1-\gamma)}\frac{2R_{\max}/(1-\gamma)}{K-1}\sum_{k=1}^{K-1}\left|F^{-1}_{\eta(x,a)}(z_k) - F^{-1}_{\Pi_\mathcal{C}\eta_\pi(x,a)}(z_k)\right|$$

$$=\frac{1}{2R_{\max}/(1-\gamma)}W_1(\eta, \Pi_\mathcal{C}\eta_\pi(x,a))$$

$$\overset{(a)}{=}\frac{1}{2R_{\max}/(1-\gamma)}W_1(\Pi_\mathcal{C}\mathcal{T}^\pi\eta, \Pi_\mathcal{C}\mathcal{T}^\pi\eta_\pi(x,a))$$

$$\overset{(b)}{\leq}\frac{1}{2R_{\max}/(1-\gamma)}\gamma\overline{W}_1(\eta, \eta_\pi)$$

$$\overset{(c)}{\leq}\frac{1}{2R_{\max}/(1-\gamma)}\gamma\frac{1}{1-\gamma}\overline{W}_1(\Pi_\mathcal{C}\eta_\pi, \eta_\pi)$$

$$\overset{(d)}{\leq}\frac{1}{2R_{\max}/(1-\gamma)}\gamma\frac{1}{1-\gamma}\frac{R_{\max}}{(1-\gamma)(K-1)}$$

$$=\frac{\gamma}{2(1-\gamma)(K-1)}\,,$$

as required. Here, (a) follows since $\eta$ is the fixed point of $\Pi_\mathcal{C}\mathcal{T}^\pi$ and $\eta_\pi$ is the fixed point of $\mathcal{T}^\pi$. (b) follows since $\Pi_\mathcal{C}$ is a non-expansion in $\overline{W}_1$, by Lemma B.1.(i), and $\mathcal{T}^\pi$ is a $\gamma$-contraction in $\overline{W}_1$. (c) follows from the following argument:

$$\overline{W}_1(\eta, \eta_\pi) \leq \overline{W}_1(\eta, \Pi_\mathcal{C}\eta_\pi) + \overline{W}_1(\Pi_\mathcal{C}\eta_\pi, \eta_\pi)$$

$$=\overline{W}_1(\Pi_\mathcal{C}\mathcal{T}^\pi\eta, \Pi_\mathcal{C}\mathcal{T}^\pi\eta_\pi) + \overline{W}_1(\Pi_\mathcal{C}\eta_\pi, \eta_\pi)$$

$$\leq\gamma\overline{W}_1(\eta, \eta_\pi) + \overline{W}_1(\Pi_\mathcal{C}\eta_\pi, \eta_\pi)$$

$$\implies \overline{W}_1(\eta, \eta_\pi) \leq \frac{1}{1-\gamma}\overline{W}_1(\Pi_\mathcal{C}\eta_\pi, \eta_\pi)\,.$$

Finally, (d) follows from Lemma B.1.(ii). $\qquad\square$

Before giving a proof of Theorem 4.7, we first state and prove a lemma that will be useful.

**Lemma B.2.** *Let* $\tau_k = \frac{2k-1}{2K}$ *for* $k = 1,\ldots,K$, *and consider the corresponding Wasserstein-1 projection operator* $\Pi_{W_1} : \mathscr{P}(\mathbb{R}) \to \mathscr{P}(\mathbb{R})$, *defined by*

$$\Pi_{W_1}(\mu) = \frac{1}{K}\sum_{k=1}^{K}\delta_{F^{-1}_\mu(\tau_k)}\,,$$

*for all* $\mu \in \mathscr{P}(\mathbb{R})$, *where* $F^{-1}_\mu$ *is the inverse c.d.f. of* $\mu$. *Let* $\eta_1, \eta_2 \in \mathscr{P}(\mathbb{R})$, *such that* $\sup(supp(\eta_i)) - \inf(supp(\eta_i)) \leq I$ *for* $i = 1,2$. *Then we have:*

$$(i)\; W_1(\Pi_{W_1}\eta_1, \eta_1) \leq \frac{I}{K}\,;$$

$$(ii)\; W_1(\Pi_{W_1}\eta_1, \Pi_{W_1}\eta_2) \leq W_1(\eta_1, \eta_2) + \frac{2I}{K}\,.$$

*Proof.* We start by proving (i). Let $F_{\eta_1}^{-1}$ be the inverse c.d.f of $\eta_1$. We have

$$
\begin{aligned}
W_1(\mu, \Pi_{W_1}\mu) &= \sum_{i=0}^{K-1} \frac{1}{K} \mathbb{E}_{X \sim \mu}\left[\left|X - F_{\eta_1}^{-1}\left(\frac{2i+1}{2K}\right)\right| \middle| F_{\eta_1}^{-1}\left(\frac{i}{K}\right) \le X \le F_{\eta_1}^{-1}\left(\frac{i+1}{K}\right)\right] \\
&\le \frac{1}{K}\left(F_{\eta_1}^{-1}(1) - F_{\eta_1}^{-1}(0)\right) \\
&= \frac{I}{K}
\end{aligned}
$$

We can now prove (ii), using the triangle inequality and (i):

$$
W_1(\Pi_{W_1}\eta_1, \Pi_{W_1}\eta_2) \le W_1(\Pi_{W_1}\eta_1, \eta_1) + W_1(\eta_1, \eta_2) + W_1(\eta_2, \Pi_{W_1}\eta_2)
$$

$$
\le W_1(\eta_1, \eta_2) + \frac{2I}{K}.
$$

$\square$

**Theorem 4.7.** *Consider the class of MDPs $\mathcal{M}$ with a fixed discount factor $\gamma \in [0, 1)$, and immediate reward distributions supported on $[-R_{max}, R_{max}]$. Then the collection of quantile statistics $s_k(\mu) = F_\mu^{-1}(\frac{2k-1}{2K})$ for $k = 1, \ldots, K$, together with the standard QDRL imputation strategy, is $\varepsilon$-approximately Bellman closed for $\mathcal{M}$, where $\varepsilon = \frac{2R_{max}(5-2\gamma)}{(1-\gamma)^2 K}$.*

*Proof.* Let $(\hat{s}_{1:K}(x, a)|(x, a) \in \mathcal{X} \times \mathcal{A})$ be the collection of statistics learnt under QDRL. We denote by $\eta(x, a)$ the distribution imputed from the statistics $\hat{s}_{1:K}(x, a)$, for each $(x, a) \in \mathcal{X} \times \mathcal{A}$. As noted in Appendix Section A.3, $\eta$ is the fixed point of the projected Bellman operator $\Pi_{W_1}\mathcal{T}^\pi$, and $\eta_\pi$ is the fixed point of $\mathcal{T}^\pi$. We begin by noting that if all immediate reward distributions have support contained within $[-R_{max}, R_{max}]$, then the true and learnt reward distributions are supported on $[-R_{max}/(1-\gamma), R_{max}/(1-\gamma)]$, and further, so are the distributions $\mathcal{T}^\pi\eta(x, a)$ for each $(x, a) \in \mathcal{X} \times \mathcal{A}$. We thus compute

$$
\begin{aligned}
&\sup_{(x,a)\in\mathcal{X}\times\mathcal{A}} \frac{1}{K}\sum_{k=1}^{K} |s_k(\eta_\pi(x, a)) - \hat{s}_k(x, a)| \\
&= \sup_{(x,a)\in\mathcal{X}\times\mathcal{A}} W_1(\Pi_{W_1}\eta(x, a), \Pi_{W_1}\eta_\pi(x, a)) \\
&\le \sup_{(x,a)\in\mathcal{X}\times\mathcal{A}} W_1(\eta(x, a), \eta_\pi(x, a)) + \frac{4R_{max}}{K(1-\gamma)},
\end{aligned}
$$

with the inequality following from Lemma B.2(ii). From here, we note that

$$
\begin{aligned}
\sup_{(x,a)\in\mathcal{X}\times\mathcal{A}} W_1(\eta(x, a), \eta_\pi(x, a)) &\overset{(a)}{\le} \sup_{(x,a)\in\mathcal{X}\times\mathcal{A}} [W_1(\eta(x, a), \Pi_{W_1}\eta_\pi(x, a)) + W_1(\Pi_{W_1}\eta_\pi(x, a), \eta_\pi(x, a))] \\
&\overset{(b)}{\le} \sup_{(x,a)\in\mathcal{X}\times\mathcal{A}} W_1(\eta(x, a), \Pi_{W_1}\eta_\pi(x, a)) + \frac{2R_{max}}{K(1-\gamma)} \\
&\overset{(c)}{=} \sup_{(x,a)\in\mathcal{X}\times\mathcal{A}} W_1(\Pi_{W_1}\mathcal{T}^\pi\eta(x, a), \Pi_{W_1}\mathcal{T}^\pi\eta_\pi(x, a)) + \frac{2R_{max}}{K(1-\gamma)} \\
&\overset{(d)}{\le} \sup_{(x,a)\in\mathcal{X}\times\mathcal{A}} W_1(\mathcal{T}^\pi\eta(x, a), \mathcal{T}^\pi\eta_\pi(x, a)) + \frac{4R_{max}}{K(1-\gamma)} + \frac{2R_{max}}{K(1-\gamma)} \\
&\overset{(e)}{\le} \sup_{(x,a)\in\mathcal{X}\times\mathcal{A}} \gamma W_1(\eta(x, a), \eta_\pi(x, a)) + \frac{6R_{max}}{K(1-\gamma)} \\
\implies \sup_{(x,a)\in\mathcal{X}\times\mathcal{A}} W_1(\eta(x, a), \eta_\pi(x, a)) &\le \frac{6R_{max}}{K(1-\gamma)^2}.
\end{aligned}
$$

Here, (a) follows from the triangle inequality, (b) follows from Lemma B.2(i). (c) follows since $\eta$ is the fixed point of $\Pi_{W_1}\mathcal{T}^\pi$ and $\eta_\pi$ is the fixed point of $\mathcal{T}^\pi$. (d) follows from Lemma B.2(ii), where we use the fact that the support of the distributions

constituting the fixed points of $\Pi_{W_1}\mathcal{T}^\pi$ and $\mathcal{T}^\pi$ necessarily are supported on $[-R_{\max}/(1-\gamma), R_{\max}/(1-\gamma)]$. (e) follows from the $\gamma$-contractivity of the Bellman operator $\mathcal{T}^\pi$ with respect to the metric $\sup_{(x,a)\in\mathcal{X}\times\mathcal{A}} W_1(\mu_1(x,a), \mu_2(x,a))$, for $\mu_1, \mu_2 \in \mathscr{P}(\mathbb{R})^{\mathcal{X}\times\mathcal{A}}$ (Bellemare et al., 2017). Hence, we obtain

$$\sup_{(x,a)\in\mathcal{X}\times\mathcal{A}} \frac{1}{K}\sum_{k=1}^K |s_k(\eta_\pi(x,a)) - \hat{s}_k(x,a)| \leq \frac{6R_{\max}}{K(1-\gamma)^2} + \frac{4R_{\max}}{K(1-\gamma)}$$

$$= \frac{2R_{\max}(5-2\gamma)}{K(1-\gamma)^2}.$$

$\square$

### B.4. Proofs of results from Section 4.3

**Lemma 4.8.** *(i) Under CDRL updates using support locations $z_1 < \cdots < z_K$, if all approximate reward distributions have support bounded in $[z_1, z_K]$, expected returns are exactly learnt. (ii) Under QDRL updates, expected returns are not exactly learnt.*

*Proof.* (i) The statistics learnt by CDRL are of the form $s_k(\mu) = \mathbb{E}_{Z\sim\mu}\left[h_{z_k, z_{k+1}}(Z)\right]$, for $k = 1, \ldots, K-1$. We observe that the mean functional $m(\mu) = \mathbb{E}_{Z\sim\mu}[Z]$ is contained in the linear span of $s_{0:K-1}$, where $s_0(\mu) = 1$ for all $\mu$. Indeed,

$$m = R_{\max}s_0 - \left(\frac{R_{\max} - R_{\min}}{K}\right)\sum_{k=1}^{K-1} s_k,$$

since

$$x = R_{\max} - \left(\frac{R_{\max} - R_{\min}}{K}\right)\sum_{k=1}^{K-1} h_{z_k, z_{k+1}}(x)$$

for all $x \in [-R_{\min}, R_{\max}]$. Since the singleton set consisting of the mean functional is Bellman closed, it follows that whatever distribution is imputed, the effective update to the mean of the distribution at the current state is the same as updating according to the classical Bellman update for the mean.

(ii) We note that the mean is not encoded by a finite set of quantiles, and hence it is impossible for expected returns to be correctly in general. To make this concrete, fix a number $K$ of quantiles to be learnt, and consider a single state, two action MDP, with reward distribution $\frac{4K-1}{4K}\delta_0 + \frac{1}{4K}\delta_1$ for the first action, and reward distribution $\delta_{1/8K}$ for the second action. Fitting quantiles at $\tau \in \{\frac{2k-1}{2K} | k = 1, \ldots, K\}$ results in all quantiles for the first distribution being equal to 0, and thus the imputed distribution is $\delta_0$, resulting in a imputed mean of 0. By contrast, for the second distribution, all quantiles are fitted at $1/8K$, resulting in an imputed distribution of $\delta_{1/8K}$ and an imputed mean of $1/8K$. Thus, a QDRL control algorithm will act greedily with respect to these imputed means and select the second action, which is sub-optimal as the first action has higher expected reward. $\square$

## C. Additional theoretical results

In this section, we provide several examples to illustrate the point made in Section 4.2 that in general, it is not possible to simultaneously achieve low approximation error on all statistics in a non-Bellman closed collection.

**Lemma C.1.** *For a fixed $K \in \mathbb{N}$, let $s_{1:K-1}$ be the statistics corresponding to CDRL (with fixed discount factor $\gamma \in [0, 1)$) with equally spaced support $R_{min} = z_1 < \ldots < z_K = R_{max}$. As earlier in the paper, we denote by $\hat{s}_k(x,a)$ the relevant learnt value of the statistic concerned. Then we have:*

$$\sup_{\substack{\mathcal{M} \, MDP \\ \pi \, policy}} \sup_{\substack{x\in\mathcal{X} \\ a\in\mathcal{A}}} \sup_{k=1,\ldots,K-1} |\hat{s}_k(x,a) - s_k(\eta_\pi(x,a))| \not\to 0$$

*as $K \to \infty$.*

*Proof.* We work with a particular family of MDPs with two states $x_1, x_2$, one action in each state, with $x_1$ transitioning to $x_2$ with probability 1, and $x_2$ terminal. In such MDPs, there is only one policy, which we denote by $\pi$; and we drop notational dependence on actions for clarity. No rewards are received at state $x_1$; we specify the rewards received at state $x_2$ below. We take a discount factor $\gamma = \frac{2^m}{2^m+1}$ for some $k \in \mathbb{N}$. Fix $L \in \mathbb{N}$, and consider CDRL updates with bin locations at $z_k = \frac{k}{2^L}$ for $k = 0, \ldots, 2^L$. Specifically, consider learning the statistic

$$\mathbb{E}_{Z \sim \eta_\pi(x_1)} \left[ h_{\frac{1}{2}, \frac{1}{2} + \frac{1}{2^L}} (Z) \right] .$$

Since there are no rewards recieved at state $x_1$, at convergence the estimate of this statistic (which we denote by $\hat{s}(x_1)$) is equal to

$$\mathbb{E}_{Z \sim \hat{\eta}(x_2)} \left[ h_{\frac{1}{2}, \frac{1}{2} + \frac{1}{2^L}} (\gamma Z) \right] = \mathbb{E}_{Z \sim \hat{\eta}(x_2)} \left[ h_{\frac{\gamma^{-1}}{2}, \frac{\gamma^{-1}}{2} + \frac{\gamma^{-1}}{2^L}} (Z) \right] = \mathbb{E}_{Z \sim \hat{\eta}(x_2)} \left[ h_{\frac{1}{2} + \frac{1}{2^{m+1}}, \frac{1}{2} + \frac{1}{2^{m+1}} + \frac{1}{2^L} + \frac{1}{2^{L+m}}} (Z) \right]$$

where $\hat{\eta}(x_2)$ is the approximate return distribution learnt at state $x_2$. Now, consider two possible reward distributions at state $x_2$:

$$\rho_A = \delta_{\frac{1}{2} + \frac{1}{2^{m+1}} + \frac{3}{2^{L+1}}} , \text{ and } \rho_B = \frac{1}{2} \left( \delta_{\frac{1}{2} + \frac{1}{2^{m+1}} + \frac{1}{2^L}} + \delta_{\frac{1}{2} + \frac{1}{2^{m+1}} + \frac{2}{2^L}} \right) .$$

Under these two reward distributions, the fitted distribution $\eta(x_2)$ is the same, namely $\rho_B$, and thus the estimate $\hat{s}(x_1)$ is the same. Our aim is to show that for these two different reward distributions, the difference of the true values of the statistic $\hat{s}(x_1)$ is independent of $L$, and hence the value of $\hat{s}(x_1)$ cannot converge to the true statistic as $L \to \infty$. To achieve this, and finish the proof, we calculate directly. In the case where the reward distribution at state $x_2$ is $\rho_A$, we have (assuming $L > m + 1$)

$$s(\eta_\pi(x_1)) = \mathbb{E}_{Z \sim \rho_A} \left[ h_{\frac{1}{2} + \frac{1}{2^{m+1}}, \frac{1}{2} + \frac{1}{2^{m+1}} + \frac{1}{2^L} + \frac{1}{2^{L+m}}} (Z) \right] = 0 .$$

In the case where the reward distribution at state $x_2$ is $\rho_B$, we have

$$s(\eta_\pi(x_1)) = \mathbb{E}_{Z \sim \rho_B} \left[ h_{\frac{1}{2} + \frac{1}{2^{m+1}}, \frac{1}{2} + \frac{1}{2^{m+1}} + \frac{1}{2^L} + \frac{1}{2^{L+m}}} (Z) \right] =$$
$$= \frac{1}{2} \left( \frac{(\frac{1}{2} + \frac{1}{2^{m+1}} + \frac{1}{2^L}) - (\frac{1}{2} + \frac{1}{2^{m+1}} + \frac{1}{2^L} + \frac{1}{2^{L+m}})}{(\frac{1}{2} + \frac{1}{2^{m+1}}) - (\frac{1}{2} + \frac{1}{2^{m+1}} + \frac{1}{2^L} + \frac{1}{2^{L+m}})} \right)$$
$$= \frac{1}{2} \left( \frac{1}{2^m + 1} \right) .$$

$\square$

**Lemma C.2.** *For a fixed $K \in \mathbb{N}$, let $s_{1:K-1}$ be the statistical functionals corresponding to by QDRL (with fixed discount factor $\gamma \in [0, 1)$). As earlier in the paper, we denote by $\hat{s}_k(x, a)$ the relevant* learnt *value of the statistic concerned. Then we have:*

$$\sup_{\substack{\mathcal{M} \text{ MDP} \\ \pi \text{ policy}}} \sup_{\substack{x \in \mathcal{X} \\ a \in \mathcal{A}}} \sup_{k=1,\ldots,K} |\hat{s}_k(x, a) - s_k(\eta_\pi(x, a))| \not\to 0$$

*as $K \to \infty$.*

*Proof.* We work with a particular family of MDPs with three states $x_0, x_1, x_2$, one action in each state, with $x_0$ transitioning to $x_1$ with probability $\frac{1}{2} - \varepsilon$ and with $x_0$ transitioning to $x_2$ with probability $\frac{1}{2} + \varepsilon$ with $\varepsilon \ll 1$. We take $x_1$ and $x_2$ to be terminal, no rewards are received at state $x_0$; we specify the rewards received at state $x_1$ and $x_2$ below. We suppose in the following that $K$ is odd.

The reward distributions at state $x_1$ and $x_2$ are given by

$$\rho_1 = \left( \frac{1}{2K} - \varepsilon \right) \delta_0 + \left( \frac{2K-1}{2K} + \varepsilon \right) \delta_1 , \text{ and } \rho_2 = \left( \frac{1}{2K} - \varepsilon \right) \delta_0 + \left( \frac{2K-1}{2K} + \varepsilon \right) \delta_{-1} .$$

Under these reward distributions the fitted return distributions are:

$$\eta(x_1) = \delta_1 \,, \text{ and } \eta(x_2) = \delta_{-1} \,.$$

Therefore, we have

$$s_{\frac{K+1}{2}}(\eta_\pi(x_0)) = 0 \,, \text{ and } \hat{s}_{\frac{K+1}{2}}(x_0) = -\gamma \,.$$

$\square$

## D. ER-DQN experimental details

### D.1. ER-DQN architecture

As discussed in Section 5.3, the ER-DQN architecture matches the exact architecture of QR-DQN (Dabney et al., 2017). The Q-network, for a given input $x$, outputs expectiles $e_{\tau_{1:K}}(x, a)$ for each $a \in \mathcal{A}$. In our experiments with 11 statistics, we take $\tau_{1:11}$ to be linearly spaced with $\tau_1 = 0.01$, $\tau_{11} = 0.99$. Note that we have $\tau_6 = 0.5$, and thus this expectile statistic is in fact the mean. For the purposes of control, greedy actions at a state $x \in \mathcal{X}$ are thus selected according to $\arg\max_{a \in \mathcal{A}} e_{\tau_6}(x, a)$, rather than averaging over statistics as in QR-DQN. For the imputation strategy, we take the root-finding problem in Expression (7), and use a call to the SciPy `root` routine with default parameters.

### D.2. Training details

We use the Adam optimiser with a learning rate of 0.00005, after testing learning rates 0.00001, 0.00003, 0.00005, 0.00007, and 0.0001 on a subset of 6 Atari games. All other hyperparameters in training correspond to those used in (Dabney et al., 2017). In particular, the target distribution is computed from a target network. Note that each training pass requires a call to the SciPy optimiser to compute the imputed samples, and thus in general will be more computationally expensive than other deep distributional Q-learning-style agents, such as C51 and QR-DQN. However, by parallelising the optimiser calls for a minibatch of transitions across several CPUs, we found that training times when using 11 expectiles to be comparable to training times of QR-DQN.

For ER-DQN Naive, we found that results were slightly improved by using 201 expectiles compared to 11, so include results with this larger number of statistics in the main paper. We take $\tau_{1:201}$ according to the same prescription as for QR-DQN: linearly spaced, with $\tau_1 = 1/(2 \times 201)$ and $\tau_{201} = 1 - 1/(2 \times 201)$.

### D.3. Environment details

We use the Arcade Learning Environment (Bellemare et al., 2013) to train and evaluate ER-DQN on a selection of 57 Atari games. The precise parameter settings of the environment are exactly the same as in the experiments performed on QR-DQN, to allow for direct comparison.

### D.4. Detailed results

In addition to the human normalised mean/median results presented in the main paper, we include training curves for all 4 evaluated agents on all 57 Atari games in Figure 11, and raw maximum scores attained in Table 1.

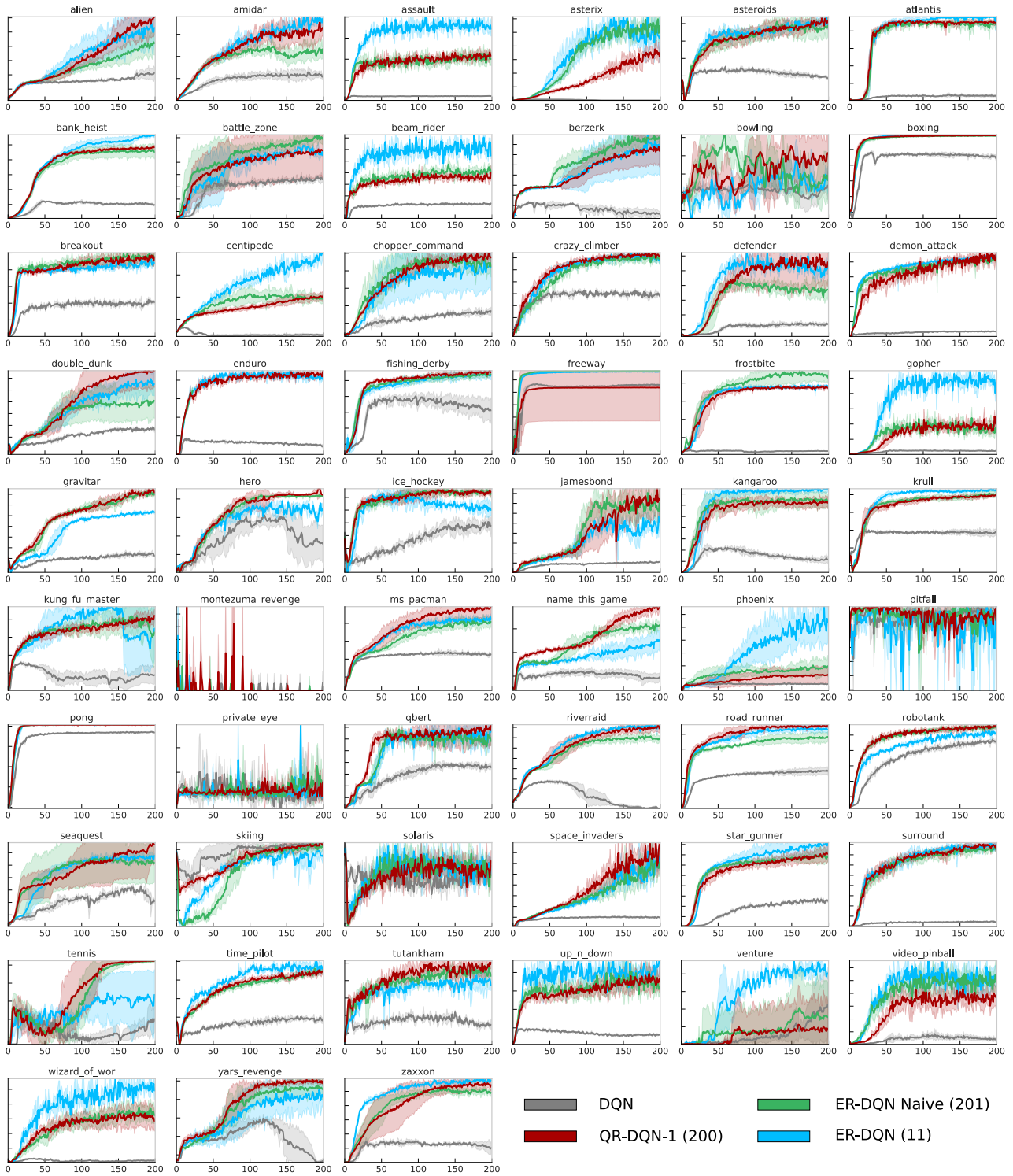*Figure 11.* Training curves for DQN, QR-DQN-1, ER-DQN Naive, and ER-DQN on all 57 Atari games.

| GAMES | QR-DQN-1 | ER-DQN Naive | ER-DQN |
|---|---|---|---|
| alien | **7279.5** | 5056.2 | 6212.0 |
| amidar | 2235.8 | 1528.8 | **2313.0** |
| assault | 17653.9 | 19156.2 | **25826.8** |
| asterix | 306055.9 | 366152.1 | **434743.6** |
| asteroids | 3484.4 | 3250.9 | **3793.2** |
| atlantis | 947995.0 | 939050.0 | **974408.3** |
| bank_heist | 1185.7 | 1132.5 | **1326.5** |
| battle_zone | 33987.2 | **40805.3** | 35098.5 |
| beam_rider | 25095.7 | 29542.5 | **48230.1** |
| berzerk | 2151.2 | 2626.6 | **2749.8** |
| bowling | 58.0 | **63.4** | 53.1 |
| boxing | 99.5 | 99.4 | **99.9** |
| breakout | 505.2 | **538.6** | 509.8 |
| centipede | 11465.1 | 12325.3 | **22505.9** |
| chopper_command | **12767.2** | 11765.8 | 11886.1 |
| crazy_climber | 159244.2 | 158369.9 | **161040.2** |
| defender | **41098.7** | 32225.2 | 36473.5 |
| demon_attack | **114530.2** | 108496.2 | 111921.2 |
| double_dunk | **16.5** | 4.0 | 16.3 |
| enduro | 2294.1 | 1923.9 | **2339.5** |
| fishing_derby | **21.6** | 18.4 | 20.2 |
| freeway | 27.2 | **34.0** | 33.9 |
| frostbite | 4068.1 | **5408.0** | 4233.7 |
| gopher | 82060.6 | 86874.1 | **115828.3** |
| gravitar | 937.0 | **942.8** | 680.9 |
| hero | **23799.1** | 21916.6 | 20374.5 |
| ice_hockey | **-1.7** | -1.9 | -2.7 |
| jamesbond | 5298.5 | **5440.4** | 4113.6 |
| kangaroo | 14827.6 | 15371.1 | **15954.4** |
| krull | 10591.2 | 10738.0 | **11318.5** |
| kung_fu_master | 49695.5 | 52080.6 | **58802.2** |
| montezuma_revenge | **0.1** | 0.0 | 0.0 |
| ms_pacman | **5860.4** | 4856.1 | 5048.5 |
| name_this_game | **20509.1** | 17064.9 | 13090.9 |
| phoenix | 15475.2 | 25177.3 | **91189.4** |
| pitfall | **0.0** | **0.0** | **0.0** |
| pong | 21.0 | **21.0** | 21.0 |
| private_eye | **531.3** | 388.3 | 176.3 |
| qbert | **17573.5** | 14536.0 | 17418.4 |
| riverraid | 18125.3 | 15726.4 | **18472.2** |
| road_runner | **67084.8** | 57168.0 | 64577.7 |
| robotank | **58.0** | 56.7 | 54.8 |
| seaquest | 16143.3 | 13501.0 | **19401.0** |
| skiing | -16869.1 | -15085.4 | **-10528.6** |
| solaris | 2615.3 | 2483.3 | **2810.6** |
| space_invaders | 11873.3 | 10099.6 | **14265.7** |
| star_gunner | 76556.3 | 75404.8 | **88900.3** |
| surround | 8.4 | 8.2 | **8.6** |
| tennis | **22.8** | 22.7 | 5.8 |
| time_pilot | 9902.0 | 10009.6 | **11675.5** |
| tutankham | **282.8** | 256.7 | 237.9 |
| up_n_down | **44893.6** | 35169.7 | 32083.3 |
| venture | 266.5 | 476.7 | **1107.0** |
| video_pinball | 570852.7 | 603852.1 | **727091.1** |
| wizard_of_wor | 21667.1 | 24397.5 | **36049.8** |
| yars_revenge | **27264.3** | 26056.7 | 24099.4 |
| zaxxon | 11707.1 | 11120.2 | **12264.4** |

*Table 1.* Raw max test scores across all 57 Atari games, starting with 30 no-op actions.

# E. Additional experimental results

In Section 5.1, we saw that the expectiles learned by EDRL-Naive on an $N$-Chain with length 15 collapsed, whereas the expectiles learnt by EDRL were reasonable approximations to the true expectiles of the return distribution. This resulted in lower average expectile estimation error with the latter expectiles, as described in Definition 4.5. In Figure 12, we supplement this by plotting Wasserstein distance between an imputed distribution for the learnt statistics and the true return distribution. This gives an alternate metric which additionally indicates how well the collection of learnt statistics summarises the full return distribution. Under this metric, we observe that increasing the number of expectiles always leads to improved performance under EDRL, whilst for EDRL-Naive, poor Wasserstein reconstruction error is observed for large numbers of expectiles and/or distance from the goal state.

We also include results for $N$-chain environments with different reward distributions, observing qualitatively similar phenomena as those noted for Figure 12. Specifically, we use two additional variants of the reward distribution at the goal state: uniform and Gaussian. We plot average expectile error in Figure 13, and Wasserstein distance between imputed and true return distributions in Figure 14.
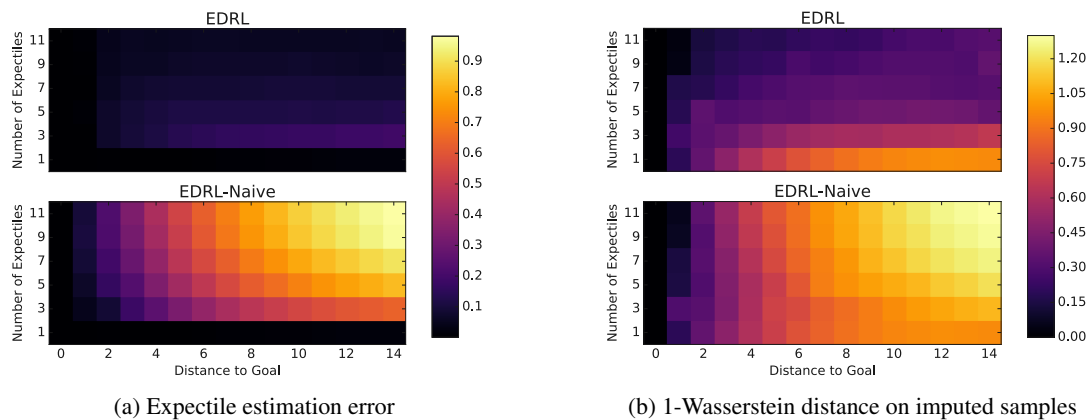


(a) Expectile estimation error

(b) 1-Wasserstein distance on imputed samples

*Figure 12.* Expectile estimation error and 1-Wasserstein distance between imputed samples and the true return distribution for varying numbers of learned expectiles and different $N$-Chain lengths.
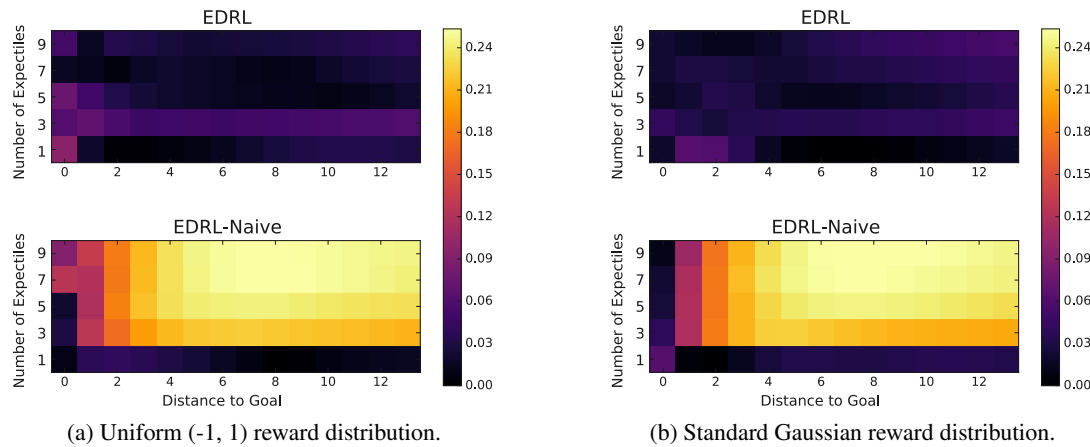


(a) Uniform (-1, 1) reward distribution.

(b) Standard Gaussian reward distribution.

*Figure 13.* Expectile estimation error for varying number of expectiles and different chain lengths. Different terminal reward distributions.

(a) Uniform (-1, 1) reward distribution.

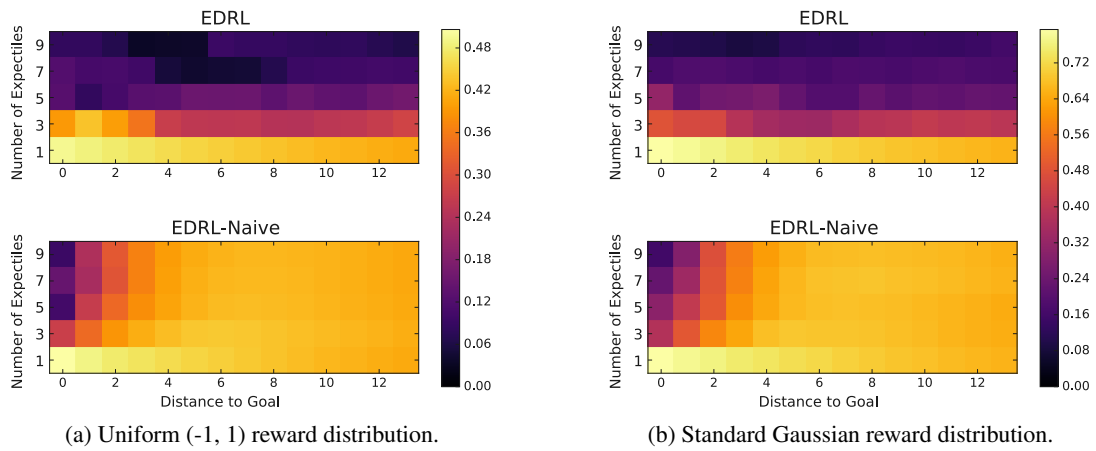(b) Standard Gaussian reward distribution.

*Figure 14.* 1-Wasserstein distance for varying number of expectiles and different chain lengths. Different terminal reward distributions.