
A Cramér Distance perspective on Quantile Regression based Distributional Reinforcement Learning

Alix Lhéritier

Amadeus SAS, F-06902 Sophia Antipolis, France

Nicolas Bondoux

Abstract

Distributional reinforcement learning (DRL) extends the value-based approach by approximating the full distribution over future returns instead of the mean only, providing a richer signal that leads to improved performances. Quantile Regression (QR) based methods like QR-DQN project arbitrary distributions into a parametric subset of staircase distributions by minimizing the 1-Wasserstein distance. However, due to biases in the gradients, the quantile regression loss is used instead for training, guaranteeing the same minimizer and enjoying unbiased gradients. Non-crossing constraints on the quantiles have been shown to improve the performance of QR-DQN for uncertainty-based exploration strategies. The contribution of this work is in the setting of fixed quantile levels and is twofold. First, we prove that the Cramér distance yields a projection that coincides with the 1-Wasserstein one and that, under non-crossing constraints, the squared Cramér and the quantile regression losses yield collinear gradients, shedding light on the connection between these important elements of DRL. Second, we propose a low complexity algorithm to compute the Cramér distance.

1 INTRODUCTION

Distributional Reinforcement Learning (DRL) extends the value-based approach of DQN (Mnih et al., 2015) by considering the full distribution of returns as a learning signal allowing to take into account all the complexity of the randomness coming from the rewards,

the transitions and the policy, which is hidden when considering the mean only. Even when a policy aims at maximizing the expected return, considering the full distribution provides an advantage in the presence of approximations, allowing to learn better representations and helping to reduce state aliasing (Bellemare et al., 2017a). With this new approach comes a generalization of the *Bellman operator*—the *distributional Bellman operator*—, whose contraction properties are key for guaranteeing the stability of DRL algorithms.

How distributions are represented and learned is also a key point, since some choices can break the contraction property (Rowland et al., 2018, Lemma 2). Some approaches use staircase parametric representations whose steps correspond to fixed quantile values like in C51 (Bellemare et al., 2017a) or to fixed quantile levels like in QR-DQN (Dabney et al., 2018b). Alternatively, FQN (Yang et al., 2019) fully parameterize the staircase distributions. IQN (Dabney et al., 2018a) follows a different approach by approximating the quantile function with a neural network that takes the quantile level as input and must therefore be sampled during training.

DRL methods resort to different notions of distance or divergence between distributions in order to practically learn them but also to analyze the effect on the contraction property of the distributional Bellman operator. In Rowland et al. (2018), a Hilbert space endowed with the ℓ_2 norm on cumulative distribution functions has been shown to be a natural framework to analyze the effect of the fixed quantile value representation of C51. In Bellemare et al. (2017b), the squared ℓ_2 distance, called *Cramér distance* in that work,¹ has been proposed for Generative Adversarial Networks but also for machine learning in general due its unbiased gradients. In Dabney et al. (2018b), the Wasserstein distance has been used for defining how a general distribution should be represented with fixed quantile levels and also to analyze the effect on the contraction property of the distributional Bellman operator. However, due

Proceedings of the 25th International Conference on Artificial Intelligence and Statistics (AISTATS) 2022, Valencia, Spain. PMLR: Volume 151. Copyright 2022 by the author(s).

¹In this work, we follow Rowland et al. (2018) and use the term Cramér distance for the ℓ_2 distance.

to the biased gradients of the Wasserstein distance, the quantile regression loss is used to train the network, guaranteeing the same minimizer as the 1-Wasserstein distance and enjoying unbiased gradients.

When estimating multiple quantiles, one faces the issue of crossing quantiles, i.e., a violation of the monotonicity of the quantile function. In QR-DQN, crossing quantiles make the learning signal noisy, affecting disambiguation of states as shown in Zhou et al. (2020). This issue has been addressed in the statistical literature of quantile regression (see, e.g. Koenker et al. (1994); He (1997); Liu and Wu (2009); Hall et al. (1999); Dette and Volgushev (2008); Bondell et al. (2010)) but also, more generally, in the machine learning literature on how to represent and learn monotonic functions (see, e.g., (Gupta et al., 2016, Table 1)), with different approaches like including penalties in the loss function or enforcing monotonicity by design. Methods that take sampled quantile levels as input during training like Tagasovska and Lopez-Paz (2019) or Dabney et al. (2018a), have been shown to alleviate the problem. In the DRL literature, Zhou et al. (2020, 2021) enforce monotonicity with special neural network designs obtaining improved results with respect to QR-DQN, in the setting of uncertainty-based exploration.

In this work, we analyze QR-based methods from a Cramér distance perspective and propose its square as an alternative loss function. In Section 2, we expose the necessary background. In Section 3, we show that the Cramér distance projection coincides with the 1-Wasserstein one, yielding a contraction guarantee. In Section 4, we propose an alternative expression of the Cramér distance allowing to show that the QR and the Cramér losses are essentially equivalent for gradient based optimization under monotonicity constraints. In Section 5, we propose another alternative expression of the Cramér distance based on quantile sorting, leading to an $O(N \log N)$ algorithm in contrast to the $O(N^2)$ complexity of the QR loss. In Section 7, we experimentally compare the different losses, illustrating the theory and the algorithm but also hinting at future research directions discussed in Section 8.

2 BACKGROUND

We consider the classical model of agent-environment interactions (Puterman, 2014), i.e., a Markov Decision Process (MDP) $(\mathcal{S}, \mathcal{A}, R, P, \gamma)$, with \mathcal{S} and \mathcal{A} being the state and action space, $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ being the reward function, $P(s'|s, a) : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ being the probability of transitioning from state s to state s' after taking action a and $\gamma \in [0, 1)$ the discount factor. A stochastic policy $\pi(\cdot|s) : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ maps a state s to a distribution over \mathcal{A} .

2.1 Q-Learning

For a fixed policy π , the *return* $Z^\pi(s, a)$ is a random variable (RV) representing the discounted cumulative rewards the agent gains from a state s by taking the action a and then following the policy π , i.e., $Z^\pi(s, a) \equiv \sum_{t=0}^{\infty} \gamma^t R(s_t, a_t)$ with $s_0 = s, a_0 = a$ and $s_{t+1} \sim P(\cdot | s_t, a_t), a_t \sim \pi(\cdot | s_t)$. The usual goal in reinforcement learning (RL) is to find an optimal policy π^* maximizing the *state-action value function* $Q^\pi(s, a) \equiv \mathbb{E} Z^\pi(s, a)$, i.e., $Q^{\pi^*}(s, a) = \max_{\pi} Q^\pi(s, a) \equiv Q^*(s, a) \forall s, a$. *Q-Learning* (Watkins and Dayan, 1992) is an off-policy reinforcement learning algorithm that directly learns the optimal state-action value function using the *Bellman optimality operator*

$$(\mathcal{T}Q)(s, a) \equiv \mathbb{E}R(s, a) + \gamma \mathbb{E}_P \max_{a' \in \mathcal{A}} Q(s', a'). \quad (1)$$

In the evaluation case, the *Bellman operator* \mathcal{T}^π (Bellman, 1957; Watkins and Dayan, 1992) is defined as

$$(\mathcal{T}^\pi Q)(s, a) \equiv \mathbb{E}R(s, a) + \gamma \mathbb{E}_{P, \pi} Q(s', a'). \quad (2)$$

These operators are contractions and their repeated application to some initial value function Q_0 converges exponentially to Q^* or Q^π , respectively (Bertsekas and Tsitsiklis, 1996). However, when Q is represented by a neural network that is trained on batches of sampled transitions (s, a, r, s') as in most deep learning studies, a gradient update is preferred since it allows for the dissipation of noise introduced in the target by stochastic approximation (Bertsekas and Tsitsiklis, 1996; Kushner and Yin, 2003). DQN (Mnih et al., 2015) iteratively trains the network by minimizing the squared *temporal difference (TD)* error $\frac{1}{2} [r + \gamma \max_{a'} Q_{\omega^-}(s', a') - Q_{\omega}(s, a)]^2$ over samples (s, a, r, s') , where ω^- is the target network, which is a copy of ω , synchronized with it periodically. When using an ε -greedy policy, the samples are obtained while the agent interacts with the environment choosing actions uniformly at random with probability ε and otherwise according to $\arg \max_a Q_{\omega}(s, a)$.

2.2 Distributional reinforcement learning

In order to extend the previous concepts to DRL, the distributional Bellman operator and *optimality operator* (Bellemare et al., 2017a) are defined as

$$\begin{aligned} (\mathcal{T}^\pi Z)(s, a) &\stackrel{D}{=} R(s, a) + \gamma Z(s', a'), a' \sim \pi(\cdot | s') \quad (3) \\ (\mathcal{T}Z)(s, a) &\stackrel{D}{=} R(s, a) + \gamma Z\left(s', \arg \max_{a' \in \mathcal{A}} \mathbb{E}_p Z(s', a')\right) \\ &\text{with } s' \sim p(\cdot | s, a), \end{aligned}$$

where $Y \stackrel{D}{=} U$ denotes equality of probability laws, i.e., the RV Y is distributed according to the same law as

U . In order to characterize the contraction properties of these operators, some notion of distance between indexed collections of distributions is necessary. The p -Wasserstein distance between two RV U and Y is the ℓ_p metric between their inverse cumulative distribution functions (inverse CDFs) (Müller, 1997), i.e.,

$$d_p(U, Y) \equiv \left(\int_0^1 |F_Y^{-1}(\omega) - F_U^{-1}(\omega)|^p d\omega \right)^{1/p}$$

where, for a RV Y , the *inverse CDF* $F_Y^{-1}(\omega) \equiv \inf \{y \in \mathbb{R} : \omega \leq F_Y(y)\}$ where $F_Y(y) \equiv \Pr(Y \leq y)$ is the CDF of Y .² Then, the maximal Wasserstein metric between two indexed collections of distributions Z_1 and Z_2 is defined as $\bar{d}_p(Z_1, Z_2) \equiv \sup_{s,a} d_p(Z_1(s, a), Z_2(s, a))$. (Bellemare et al., 2017a, Lemma 3) shows that \mathcal{T}^π is a contraction in \bar{d}_p , i.e.,

$$\bar{d}_p(\mathcal{T}^\pi Z_1, \mathcal{T}^\pi Z_2) \leq \gamma \bar{d}_p(Z_1, Z_2). \quad (4)$$

The case of the distributional optimality operator \mathcal{T} is more involved. In general, it is not a contraction (Bellemare et al., 2017a). However, based on the fact that \mathcal{T}^π is a contraction, Bellemare et al. (2017a) proves that, if the optimal policy is unique, then the iterates $Z_{k+1} \leftarrow \mathcal{T}Z_k$ converge to Z^{π^*} (in p -Wasserstein metric, $\forall s, a$) and, under some conditions, \mathcal{T} has a unique fixed point corresponding to an optimal value distribution.

2.3 Finite support projection

Previous approaches of DRL project return distributions $Z(s, a)$ onto a space of distributions of finite support, modeled by a mixture of Diracs over N support points $\theta_i(s, a), i = 1..N$, i.e.,

$$Z_\theta(s, a) \equiv \sum_{i=1}^N p_i(s, a) \delta_{\theta_i(s, a)} \quad (5)$$

which yields a staircase CDF $\sum_{i=1}^N p_i(s, a) \mathbb{1}_{z \geq \theta_i(s, a)}$. Different approaches have been followed to parameterize these distributions depending on whether p_i and θ_i are learned or fixed. In this work, we consider p_i fixed and θ_i a learned parameter.

In order to analyze how arbitrary distributions are mapped into these finite representations, different projection operators are defined as minimizers of some distance between distributions. For instance, in Dabney et al. (2018b), the 1-Wasserstein projection Π_{W_1} is used and it is shown that the resulting projected Bellman operator remains a contraction, i.e.,

$$\bar{d}_\infty(\Pi_{W_1} \mathcal{T}^\pi Z_1, \Pi_{W_1} \mathcal{T}^\pi Z_2) \leq \gamma \bar{d}_\infty(Z_1, Z_2). \quad (6)$$

²For $p = \infty$, $d_\infty(Y, U) \equiv \sup_{\omega \in [0,1]} |F_Y^{-1}(\omega) - F_U^{-1}(\omega)|$.

However, since Wasserstein distances suffer from biased gradients (Bellemare et al., 2017b,a), the *quantile regression (QR) loss* is used in practice, guaranteeing the same minimizer and enjoying unbiased gradients (Dabney et al., 2018b). Given a target distribution \bar{F} , the QR loss, which allows to learn the parameters $\{\theta_1, \dots, \theta_N\}$ of $F(z) \equiv \frac{1}{N} \sum_{i=1}^N \mathbb{1}_{z \geq \theta_i}$, is defined as

$$\mathcal{L}_{\text{QR}}(F, \bar{F}) \equiv \sum_{i=1}^N \mathbb{E}_{Z \sim \bar{F}} [\rho_{\hat{\tau}_i}(Z - \theta_i)] \quad (7)$$

$$\text{with } \rho_\tau(u) \equiv u(\tau - \mathbb{1}_{u < 0}) \quad (8)$$

where $\hat{\tau}_i$ are the midpoints of a uniform grid of N quantile levels, i.e., $\hat{\tau}_i \equiv \frac{2i-1}{2N}$. Note that this definition makes θ_i an estimate of the $\hat{\tau}_i$ -quantile. As we shall see in the next section (cf. Remark 1), this correspondence is not enforced by the Cramér projection. Improved empirical results have been reported in Dabney et al. (2018b) by Huberizing the QR loss, i.e., by replacing $\rho_\tau(u)$ by $\rho_\tau^\kappa(u) = |\tau - \mathbb{1}_{u < 0}| \mathcal{L}_\kappa(u)$ where $\mathcal{L}_\kappa(u)$ is the *Huber loss* (Huber, 1964)

$$\mathcal{L}_\kappa(u) \equiv \begin{cases} \frac{1}{2}u^2, & \text{if } |u| \leq \kappa \\ \kappa(|u| - \frac{1}{2}\kappa), & \text{otherwise} \end{cases}. \quad (9)$$

3 CRAMÉR AND 1-WASSERSTEIN PROJECTION EQUIVALENCE

The ℓ_p distance between two RV U and Y is the ℓ_p metric between their CDFs, i.e.,

$$\ell_p(U, Y) \equiv \left(\int_{-\infty}^{\infty} |F_Y(z) - F_U(z)|^p dz \right)^{1/p}.$$

The *Cramér distance* corresponds to the ℓ_p distance for $p = 2$. We now show that, given an arbitrary distribution and a grid of quantile levels, there is a staircase representation that minimizes the ℓ_p distance, which puts the quantile values at the inverse of the quantile level midpoints. We first introduce an auxiliary Lemma.

Lemma 1. *For any $\tau, \tau' \in [0, 1]$ with $\tau < \tau'$ and CDF F with inverse F^{-1} , let $t \equiv F^{-1}(\tau)$ and $t' \equiv F^{-1}(\tau')$ and consider the scaled and vertically shifted Heaviside step function*

$$H_\theta^{\tau, \tau'}(z) \equiv \tau + (\tau' - \tau) \mathbb{1}_{z \geq \theta}.$$

Then, for any $p \in \mathbb{R}, p > 1$, the set of $\theta \in [t, t']$ minimizing

$$\int_t^{t'} |F(z) - H_\theta^{\tau, \tau'}|^p dz \quad (10)$$

is given by

$$\left\{ \theta \in [t, t'] \mid F(\theta) = \left(\frac{\tau + \tau'}{2} \right) \right\}. \quad (11)$$

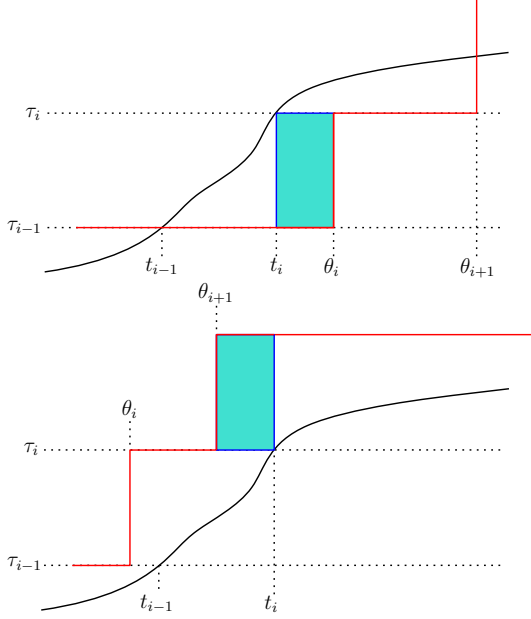


Figure 1: **Intuition for proving** $t_{i-1} \leq \theta_i^* \leq t_i$. The ℓ_p distance can be decreased by moving θ_i to t_i , in the first situation, and θ_{i+1} to t_i , in the second one. The shaded area represents the decrease for $p = 1$.

If F^{-1} is the inverse CDF, then $F^{-1}((\tau + \tau')/2)$ is always a valid minimizer, and if F^{-1} is continuous at $(\tau + \tau')/2$, then $F^{-1}((\tau + \tau')/2)$ is the unique minimizer.

Proof. A visual intuition of the proof is shown in Fig. 2. See Appendix A for details. \square

Theorem 1. Given $p_i \geq 0, i = 1..N$ such that $\sum_i p_i = 1$, the ℓ_p distance between F and a mixture of Heaviside step functions $F_N(z) = \sum_{i=1}^N p_i \mathbb{1}_{z \geq \theta_i}$ is minimized with $\theta_i = F^{-1}((\tau_i + \tau_{i-1})/2)$ where τ_i are the quantile levels $\tau_i = \sum_{j=1}^i p_j$ and F^{-1} is the inverse CDF.

Proof. Let $t_i \equiv F^{-1}(\tau_i)$. We first prove that an optimal θ^* satisfies $t_{i-1} \leq \theta_i^* \leq t_i$. See Fig. 1 for an intuition.

Without loss of generality, we assume that $\theta_1^* \leq \dots \leq \theta_N^*$. Let us suppose that there is an optimal F_N with $\theta_1 \geq t_1$. We can write the p -th power of the ℓ_p distance as

$$\begin{aligned} \ell_p^p(F, F_N) &= \int_{-\infty}^{t_1} |F(z) - F_N(z)|^p dz \\ &+ \int_{t_1}^{\theta_2} |F(z) - F_N(z)|^p dz + \int_{\theta_2}^{\infty} |F(z) - F_N(z)|^p dz \end{aligned} \quad (12)$$

The value of the middle term strictly decreases when θ_1 decreases toward t_1 (while the other terms are unaf-

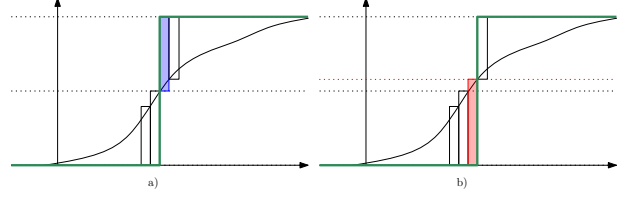


Figure 2: **Midpoint minimizer.** a) The black curve is approximated by one step function (in green) located at the inverse of the mid-point. The rectangles represent an approximation of the ℓ_p distance. b) If we move the step function to the right, the blue rectangle will be replaced by the larger red one.

fected) since

$$\begin{aligned} \int_{t_1}^{\theta_2} |F(z) - F_N(z)|^p dz &= \int_{t_1}^{\theta_2} |F(z) - H_{\theta_1}^{0, \tau_1}(z)|^p dz \\ &= \int_{t_1}^{\theta_1} F(z)^p dz + \int_{\theta_1}^{\theta_2} (F(z) - \tau_1)^p dz \end{aligned} \quad (13)$$

and $F(z)^p > (F(z) - \tau_1)^p$. In consequence $\theta_1 = t_1$; It proves that no optimal exist for $\theta_1 > t_1$, and thus that we have $\theta_1 \leq t_1$.

By induction, we assume that $\theta_{n-1}^* \leq t_{n-1}$. As before, we suppose, that there is an optimal F_N with $\theta_n \geq t_n$ and we observe that the value of the term

$$\int_{t_n}^{\theta_{n+1}} |F(z) - F_N(z)|^p dz \quad (14)$$

$$= \int_{t_n}^{\theta_{n+1}} |F(z) - H_{\theta_n}^{\tau_{n-1}, \tau_n}(z)|^p dz \quad (15)$$

$$= \int_{t_n}^{\theta_n} (F(z) - \tau_{n-1})^p dz + \int_{\theta_n}^{\theta_{n+1}} (F(z) - \tau_n)^p dz$$

strictly decreases when θ_n decreases toward t_n since $(F(z) - \tau_{n-1})^p > (F(z) - \tau_n)^p$. In consequence $\theta_n = t_n$; it proves that no optimal exist for $\theta_n > t_n$, and thus that we have $\theta_n \leq t_n \forall n \in \{1..N\}$. Analogously, starting by θ_N and going backwards, we can prove that $\theta_n \geq t_{n-1} \forall n \in \{1..N\}$. This allows us to show that the optimization problem has an optimal substructure and thus it amounts to solving independent minimization problems of the form (10), i.e.,

$$\begin{aligned} \min_{\theta_1, \dots, \theta_N} \ell_p^p(F, F_N) &= \min_{\theta_1, \dots, \theta_N} \sum_{i=1}^N \int_{t_{i-1}}^{t_i} |F(z) - F_N(z)|^p dz \\ &= \sum_{i=1}^N \min_{\theta_i} \int_{t_{i-1}}^{t_i} |F(z) - H_{\theta_i}^{\tau_{i-1}, \tau_i}(z)|^p dz \end{aligned} \quad (16)$$

with $t_0 \equiv -\infty$. \square

Remark 1. For simplicity, we chose $\theta_i = F^{-1}((\tau_i + \tau_{i-1})/2)$, however any permutation σ in the symmetric group of size N makes $\hat{\theta}_i \equiv \theta_{\sigma(i)}$ a minimizer too.

We define the ℓ_p projection of an arbitrary CDF F with inverse CDF F^{-1} onto a grid of quantile levels as

$$\Pi_{\ell_p} F \equiv F_N^*(z) = \sum_{i=1}^N p_i \mathbb{1}_{z \geq \theta_i^*} \quad (17)$$

with $\theta_i^* = F^{-1}((\tau_i + \tau_{i-1})/2)$. Therefore, it is equivalent to the 1-Wasserstein projection and to QR loss minimization (Dabney et al., 2018b, Lemma 2), which implies the following corollary.

Corollary 1. The Cramér projected distributional Bellman operator is a contraction in \bar{d}_∞ i.e.

$$\bar{d}_\infty(\Pi_{\ell_p} \mathcal{T}^\pi Z_1, \Pi_{\ell_p} \mathcal{T}^\pi Z_2) \leq \gamma \bar{d}_\infty(Z_1, Z_2). \quad (18)$$

Proof. It follows directly from Eq. (6) (Bellemare et al., 2017a, Lemma 3) and Theorem 1. \square

4 CRAMÉR AND QR LOSS OPTIMIZATION EQUIVALENCE

In order to put in evidence the relationship between the gradients of the QR loss and the squared Cramér distance—which we refer to as *Cramér loss*—, we first present an alternative formula for the latter.

Lemma 2. Given two staircase distributions $F(z) = \frac{1}{N} \sum_{i=1}^N \mathbb{1}_{z \geq \theta_i}$ and $\bar{F}(z) = \frac{1}{N} \sum_{i=1}^N \mathbb{1}_{z \geq \bar{\theta}_i}$ such that $\theta_1 \leq \dots \leq \theta_N$ and $\bar{\theta}_1 \leq \dots \leq \bar{\theta}_N$. Let $u_{ij} \equiv \bar{\theta}_j - \theta_i$ and $\delta_{ij} \equiv \mathbb{1}_{u_{ij} < 0}$. The squared Cramér distance between the distributions can be expressed as

$$\int_{-\infty}^{\infty} (F(z) - \bar{F}(z))^2 dz = \frac{1}{N^2} \sum_{i=1}^N \left(|u_{ii}| + \sum_{j=i+1}^N \delta_{ij} 2|u_{ij}| + \sum_{j=1}^{i-1} (1 - \delta_{ij}) 2|u_{ij}| \right). \quad (19)$$

Proof. We compute the squared Cramér distance in a constructive way. The idea is to cover the area between the two curves with rectangular tiles as in Fig. 3 to compute the integral by pieces. A tile of height i/N and width u corresponds to the term $u(i/N)^2$. We start from a) and replace parts of tiles to arrive to b).

First, we formally define our tiling operator T . Second, we show that it is well built: the sum of the tiles given by T is equal to the squared Cramér distance between the two curves. Third, we derive Eq. (19) by using that tiling operator.

(Tiling operator) First consider an interval $u^+ \equiv [t_1, t_2]$ such that $\bar{F}(t_1) = F(t_1)$, $\bar{F}(t_2) = F(t_2)$ and

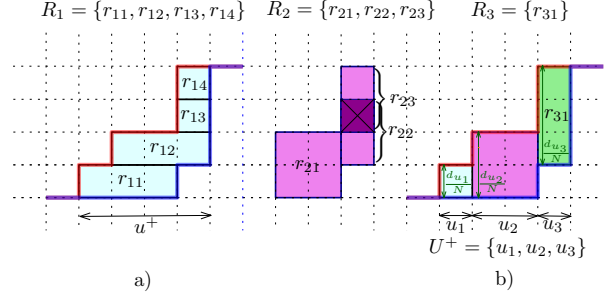


Figure 3: **Computing the Cramér distance between \bar{F} (red) and F (blue) with the tiling operator T .** a) starting point represents $T_1 = \frac{1}{N^2} \sum_{r \in R_1} u_r$. b) ending point represents the squared Cramér distance $\frac{1}{N^2} (u_1^2 + u_2^2 + u_3^2)$, where u_i is the width of each rectangles in b). Only the leftmost part of r_{11} remains in b), the rest has been replaced by taller rectangles occupying the whole height $\frac{d_{u_i}}{N}$. The middle diagram illustrates the effect of the tiling operator T_2 yielding the final rectangle in the middle and, on the right, two overlapping rectangles—that need to be replaced by a taller one—and an oversubstrated rectangle (with a cross). The result of $T_1 + T_2 + T_3$ is shown in b), a rectangle of height $3/N$ has been added, the two overlapping rectangles have been removed and the oversubstrated rectangle has been added back.

$\bar{F}(z) > F(z) \forall z \in (t_1, t_2)$. Let us define the tiling operator T_h for $h \geq 1$

$$\begin{aligned} T_h(F, \bar{F}, u^+) &\equiv \frac{1}{N^2} \sum_{r \in R_h} u_r (h^2 - 2(h-1)^2 + \mathbb{1}_{h>1}(h-2)^2) \quad (20) \\ &= \frac{1}{N^2} \sum_{r \in R_h} u_r (2 - \mathbb{1}_{h=1}) \quad (21) \end{aligned}$$

where u_r is the width of a rectangle r and R_h is the set of rectangles of height h/N whose upper left and lower right angles are aligned with quantiles of, respectively, \bar{F} and F lying in u^+ ; more formally, $R_h \equiv \{r : r \text{ is an axis-parallel rectangle with opposite corners coordinates } (\theta_i, i/N) \text{ and } (\bar{\theta}_j, j/N) \forall i, j \text{ s.t. } \theta_i, \bar{\theta}_j \in u^+, j - i = h \text{ and } \theta_i > \bar{\theta}_j\}$. Note that these rectangles lie completely within the difference area since F and \bar{F} are monotonically increasing. Note that T_1 corresponds to the initial step depicted in Fig. 3 a). Intuitively, for $h > 1$, Eq. (20) represents the fact that the operator T_h replaces parts of width u_r of two tiles of height $(h-1)/N$ by a tile of height h/N and width u_r and fixes oversubstrated tiles of the step $h-2$.

(Soundness) Let $T^h(F, \bar{F}, u^+) \equiv \sum_{d=1}^h T_d(F, \bar{F}, u^+)$. We are going to express T^h as a sum over a set U^+ of left-closed right-open intervals constituting a partition of u^+ , s.t. for any $u \equiv [a, b) \in U^+$ the difference

between the CDFs is constant, i.e.,

$$\bar{F}(z) - F(z) = \frac{d_u}{N} > 0 \quad \forall z \in u, \quad (22)$$

and no quantile lies strictly within u , i.e., $\nexists k$ s.t. $\theta_k \in (a, b) \vee \bar{\theta}_k \in (a, b)$. See Fig. 3 b). We prove by induction the following property.

$$T^h(F, \bar{F}, u^+) = \frac{1}{N^2} \sum_{u \in U^+} |u| (\mathbb{1}_{d_u \leq h} d_u^2 + \mathbb{1}_{d_u > h} g_{u,h}) \quad (23)$$

with $g_{u,h} \equiv (d_u - h + 1)(2h - 1) + (h - 1)^2$. We first express T_h as a sum over U^+ . We can rearrange the sum in Eq. (21), by decomposing each width u_r as a sum of lengths of intervals in U^+ and by noting that for each $u \in U^+$ there are $\mathbb{1}_{d_u \geq h}(d_u - h + 1)$ rectangles in R_h with non-empty projection on u , as follows

$$T_h(F, \bar{F}, u^+) = \frac{1}{N^2} \sum_{u \in U^+} |u| \mathbb{1}_{d_u \geq h} (d_u - h + 1)(2 - \mathbb{1}_{h=1}) \quad (24)$$

In particular, for $h = 1$, we have

$$T_1(F, \bar{F}, u^+) = \frac{1}{N^2} \sum_{u \in U^+} |u| d_u, \quad (25)$$

which validates the base case since $T^1(F, \bar{F}, u^+) = T_1(F, \bar{F}, u^+)$ and $\mathbb{1}_{d_u \leq h} d_u^2 + \mathbb{1}_{d_u > h} g_{u,h} = d_u$. We now assume that the property (23) holds for $h - 1$ and note that $g_{u,h-1} + 2(d_u - h + 1) = g_{u,h}$. Then, for $h > 1$,

$$T^h(F, \bar{F}, u^+) = T^{h-1}(F, \bar{F}, u^+) + T_h(F, \bar{F}, u^+) \quad (26)$$

$$= \frac{1}{N^2} \sum_{u \in U^+} |u| (\mathbb{1}_{d_u \leq h-1} d_u^2 + \mathbb{1}_{d_u > h-1} g_{u,h-1} + \mathbb{1}_{d_u \geq h} 2(d_u - h + 1)) \quad (27)$$

$$= \frac{1}{N^2} \sum_{u \in U^+} |u| (\mathbb{1}_{d_u \leq h-1} d_u^2 + \mathbb{1}_{d_u \geq h} g_{u,h}) \quad (28)$$

$$= \frac{1}{N^2} \sum_{u \in U^+} |u| (\mathbb{1}_{d_u \leq h} d_u^2 + \mathbb{1}_{d_u > h} g_{u,h}) \quad (29)$$

since $\mathbb{1}_{d_u > h-1} = \mathbb{1}_{d_u \geq h}$ and $\mathbb{1}_{d_u = h} g_{u,h} = \mathbb{1}_{d_u = h} d_u^2$.

Since $\mathbb{1}_{d_u \leq N} = 1 - \mathbb{1}_{d_u > N} = 1$, the final tiling $T^N(F, \bar{F}, u^+)$ corresponds to the squared Cramér distance on the interval u^+ , i.e.,

$$T^N(F, \bar{F}, u^+) = \frac{1}{N^2} \sum_{u \in U^+} |u| d_u^2. \quad (30)$$

(Final derivation) Now, we are going to use (21) to get to the claimed expression. First note that for a rectangle $r \in R_h$ with upper leftmost and lower

rightmost angles corresponding, respectively, to $\bar{\theta}_j$ and θ_i , its width is $u_r = |u_{ij}|$. Since $\theta_1 \leq \dots \leq \theta_N$ and $\bar{\theta}_1 \leq \dots \leq \bar{\theta}_N$, when $\bar{F}(z) > F(z)$, each rectangle in R_h corresponds to exactly one pair $(\bar{\theta}_j, \theta_i)$ such that $(\delta_{ij} = 1) \wedge (i \leq j)$. By symmetry, the condition $(\delta_{ij} = 0) \wedge (j \leq i)$ allows us to consider intervals such that $\bar{F}(z) < F(z)$. This allows to express the sum (21) as sums over indices i, j . We consider the case $i = j$ separately to avoid double counting and also because it corresponds to $h = 1$. Therefore, from (21), we have

$$T^N(F, \bar{F}, \mathbb{R}) = \frac{1}{N^2} \left(\sum_{r \in R_1} u_r + \sum_{h=2}^N \sum_{r \in R_h} 2u_r \right) \quad (31)$$

$$= \frac{1}{N^2} \left(\sum_{i=1}^N |u_{ii}| + \sum_{i=1}^{N-1} \sum_{j=i+1}^N \delta_{ij} 2|u_{ij}| \right) \quad (32)$$

$$+ \sum_{j=1}^{N-1} \sum_{i=j+1}^N (1 - \delta_{ij}) 2|u_{ij}| \Big). \quad (33)$$

By rearranging the sums, we get Equation (19). \square

Corollary 2. *Under the conditions of Lemma 2,*

$$\frac{\partial \mathcal{L}_{\text{QR}}(F, \bar{F})}{\partial \theta_i} = \frac{1}{N} \left(\frac{1 - 2i}{2} + \sum_{j=1}^N \delta_{ij} \right) \quad (34)$$

$$\frac{\partial \ell_2^2(F, \bar{F})}{\partial \theta_i} = \frac{1}{N^2} \left(1 - 2i + 2 \sum_{j=1}^N \delta_{ij} \right). \quad (35)$$

Therefore, their gradients are collinear, i.e.

$$\nabla_{\theta} \mathcal{L}_{\text{QR}} = \frac{N}{2} \nabla_{\theta} \ell_2^2. \quad (36)$$

Proof. For a target distribution $\bar{F}(z) = \frac{1}{N} \sum_{i=1}^N \mathbb{1}_{z \geq \bar{\theta}_i}$, the quantile regression loss can be expressed as

$$\mathcal{L}_{\text{QR}}(F, \bar{F}) = \sum_{i=1}^N \frac{1}{N} \sum_{j=1}^N \rho_{\hat{\tau}_i}(\bar{\theta}_j - \theta_i) \quad (37)$$

$$= \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^N (\bar{\theta}_j - \theta_i)(\hat{\tau}_i - \delta_{ij}) \quad (38)$$

and thus

$$\frac{\partial \mathcal{L}_{\text{QR}}(F, \bar{F})}{\partial \theta_i} = \frac{1}{N} \sum_{j=1}^N (\delta_{ij} - \hat{\tau}_i) \quad (39)$$

$$= \frac{1}{N} \left(\frac{1 - 2i}{2} + \sum_{j=1}^N \delta_{ij} \right). \quad (40)$$

In order to obtain the partial derivative of the squared Cramér distance, first note that $\delta_{ij} |u_{ij}| = \delta_{ij} (\theta_i - \bar{\theta}_j)$,

$(1 - \delta_{ij})|u_{ij}| = (1 - \delta_{ij})(\bar{\theta}_j - \theta_i)$ and $|u_{ii}| = \delta_{ii}(\theta_i - \bar{\theta}_i) + (1 - \delta_{ii})(\bar{\theta}_i - \theta_i)$. By replacing these quantities in (19) and taking the derivative with respect to θ_i we obtain

$$\frac{\partial \ell_2^2(F, \bar{F})}{\partial \theta_i} \quad (41)$$

$$= \frac{1}{N^2} \left[2\delta_{ii} - 1 + 2 \left(\sum_{j=i+1}^N \delta_{ij} + \sum_{j=1}^{i-1} (\delta_{ij} - 1) \right) \right]$$

$$= \frac{1}{N^2} \left(2 \sum_{j=1}^N \delta_{ij} - 1 + 2 \sum_{j=1}^{i-1} (-1) \right) \quad (42)$$

$$= \frac{1}{N^2} \left(1 - 2i + 2 \sum_{j=1}^N \delta_{ij} \right). \quad (43)$$

□

Remark 2. *Therefore, gradient descent methods whose parameter updates are invariant to rescaling of the gradient like ADAM Kingma and Ba (2015), yield the same optimization path with both losses.*

Remark 3. *Huberization of the QR loss breaks the equivalence with the Cramér loss.*

5 ALGORITHM

Formula (19) allows to compute the squared Cramér distance between two staircase distributions $F(z) = \frac{1}{N} \sum_{i=1}^N \mathbb{1}_{z \geq \theta_i}$ and $\bar{F}(z) = \frac{1}{N} \sum_{i=1}^N \mathbb{1}_{z \geq \bar{\theta}_i}$ assuming the quantiles are ordered, i.e., $\theta_1 \leq \dots \leq \theta_N$ and $\bar{\theta}_1 \leq \dots \leq \bar{\theta}_N$. That formula involves two nested sums making it of quadratic complexity in N as the quantile regression loss. Alternatively, if we consider the sorted sequence of merged quantiles $\theta' \equiv \text{sort}(\{\theta_i\}_{i=1..N} \cup \{\bar{\theta}_i\}_{i=1..N})$, we have that $F(z) - \bar{F}(z)$ is constant between any two consecutive quantile values of θ' and the difference can be obtained by accumulating the increments from F and the decrements from \bar{F} , see Appendix B for an illustration and a formal proof. Therefore, we can express the Cramér loss between two staircase distributions as follows

$$\int_{-\infty}^{\infty} (F(z) - \bar{F}(z))^2 dz = \quad (44)$$

$$\sum_{i=1}^{2N-1} (\theta'_{i+1} - \theta'_i) \left(\sum_{j \text{ s.t. } \theta_j \leq \theta'_i} \frac{1}{N} - \sum_{j \text{ s.t. } \bar{\theta}_j \leq \theta'_i} \frac{1}{N} \right)^2$$

where θ'_i is the i -the element of θ' . Algorithm 1 implements this formula based on sorting the merged quantiles of both distributions, yielding $O(N \log N)$ complexity. Note that this algorithm does not require the input vectors θ and $\bar{\theta}$ to be ordered. This has an

important consequence on the network that outputs θ , since it is not required to be in a particular order as for the QR loss. This permutation equivalence creates symmetries in the loss landscape (see Fig. 7 in Appendix 1, for an illustration). Non-crossing architectures like Zhou et al. (2020, 2021) eliminate these symmetries by enforcing monotonicity on the output.

Algorithm 1: Cramér loss. The operators $[1 :]$ and $[-1]$ remove, respectively, the first and the last elements of the array. $\mathbf{1}_N$ denotes an array of N ones and $*$ denotes elementwise multiplication.

Input: $\theta \equiv [\theta_1, \dots, \theta_N]$, $\bar{\theta} \equiv [\bar{\theta}_1, \dots, \bar{\theta}_N]$: array
Output: $\int_{-\infty}^{\infty} (F(z) - \bar{F}(z))^2 dz$

```

 $\theta' \leftarrow \text{concat}(\theta, \bar{\theta})$ 
 $i_1, \dots, i_{2N} \leftarrow \text{argsort}(\theta')$ 
 $\theta' \leftarrow \theta'[i_1, \dots, i_{2N}]$ 
 $\Delta_z \leftarrow \theta'[1 :] - \theta'[: -1]$ 
 $\Delta_\tau \leftarrow \text{concat}(-\frac{1}{N} \mathbf{1}_N, \frac{1}{N} \mathbf{1}_N)$ 
 $\Delta_\tau \leftarrow \Delta_\tau[i_1, \dots, i_{2N}]$ 
 $\Delta_\tau \leftarrow \text{cumsum}(\Delta_\tau)[: -1]$ 
 $I \leftarrow \Delta_\tau * \Delta_\tau * \Delta_z$ 
return  $\text{sum}(I)$ 
    
```

6 CRAMÉR TD-LEARNING ON SAMPLED TRANSITIONS

In order to train a DRL agent using the Cramér loss, we extend temporal-difference (TD) learning to distributions. For this, we express distributional Bellman's equations in the language of distributions as in Rowland et al. (2018). Given a probability distribution $\nu \in \mathcal{P}(\mathbb{R})$ and a measurable function $f : \mathbb{R} \rightarrow \mathbb{R}$, the *push-forward measure* $f_{\#}\nu \in \mathcal{P}(\mathbb{R})$ is defined by $f_{\#}\nu(A) \equiv \nu(f^{-1}(A))$, for all Borel sets $A \subseteq \mathbb{R}$. Let $f_{r,\gamma}(x) \equiv r + \gamma x$ and η_π be the collection of return distributions for each state and action, associated with a policy π . The basis of DRL is given by the fixed point equation

$$\eta_\pi(s, a) = (\mathcal{T}^\pi \eta_\pi)(s, a) \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A}$$

where $\mathcal{T}^\pi : \mathcal{P}(\mathbb{R})^{\mathcal{S} \times \mathcal{A}} \rightarrow \mathcal{P}(\mathbb{R})^{\mathcal{S} \times \mathcal{A}}$ is the distributional Bellman operator on distributions³ defined as

$$(\mathcal{T}^\pi \eta)(s, a) \equiv \mathbb{E}_{r, s', a' | s, a} (f_{r,\gamma})_{\#} \eta(s', a')$$

for all $\eta \in \mathcal{P}(\mathbb{R})^{\mathcal{X} \times \mathcal{A}}$. For Cramér-based TD-learning, we should use a parametric distribution η_θ and a frozen version of it that we call η' and do stochastic gradient descent by approximating $\mathbb{E}_{s,a} \nabla_{\theta} \ell_2^2(\eta_\theta, \mathbb{E}_{r, s', a' | s, a} (f_{r,\gamma})_{\#} \eta'(s', a'))$. Let F_θ and

³Eq. (3) is expressed in the language of random variables.

$F_{r,s',a'}$ denote the CDFs of η_θ and $(f_{r,\gamma})_\# \eta'(s',a')$, respectively. Following the steps of the proof of (Belle-mare et al., 2017b, Theorem 2) (*unbiased gradients*):

$$\begin{aligned} & \mathbb{E}_{s,a} \nabla_\theta \ell_2^2(F_\theta, \mathbb{E}_{r,s',a'|s,a} F_{r,s',a'}) \\ & \stackrel{(a)}{=} \mathbb{E}_{s,a} \int_{-\infty}^{\infty} \nabla_\theta (F_\theta(x) - \mathbb{E}_{r,s',a'|s,a} F_{r,s',a'}(x))^2 dx \\ & \stackrel{(b)}{=} \mathbb{E}_{s,a} \mathbb{E}_{r,s',a'|s,a} \int_{-\infty}^{\infty} 2(F_\theta(x) - F_{r,s',a'}(x)) \nabla_\theta F_\theta(x) dx \\ & = \mathbb{E}_{s,a,r,s',a'} \nabla_\theta \ell_2^2(F_\theta, F_{r,s',a'}) \end{aligned} \quad (45)$$

where (a) and (b) hold assuming that F_θ and $F_{r,s',a'}$ have light enough tails (which is our case since they are mixtures of N Heaviside functions) to avoid infinite squared Cramér distances and expected gradients. In the control case, the expectation over a' is not needed anymore, since a' is deterministically given by the policy. Practically, Eq. (45) allows us to use the average gradient of $\ell_2^2(F_\theta, F_{r,s',a'})$ over batches of sample transitions for Cramér TD-learning.

7 EXPERIMENTS

In light of the previous results, we investigate how the differences between the Cramér and the QR losses affect the results in synthetic and Atari 2600 experiments, considering the presence or not of non-crossing constraints and Huberization. The code and the full output of the experiments are available at <https://github.com/alherit/cr-dqn>.

7.1 Synthetic experiment

We first propose an experiment that is simple but representative of the challenges that DRL faces. We consider an MDP with only one possible action at one state s that can transition to two possible states s_1 and s_2 with probabilities $2/3$ and $1/3$, respectively, each with a different return distribution—a Dirac located at -1 and 1 respectively. The goal is to learn the return distribution at s , i.e., the mixture distribution shown in red in Fig. 4. The figure shows the estimated distributions obtained after 1000 training iterations with the different losses and two architectures: a fully connected (FC) one as in QR-DQN and the non-crossing (NC) one of NC-QR-DQN with a comparable number of parameters (2712 and 2702, respectively). The networks output $N = 12$ quantiles, allowing to represent the mixture exactly. We repeat the experiment 100 times. We show the average 1-Wasserstein distance d_1 and the standard deviation to quantify how close are the learned distributions with respect to the true target. See Appendix C for details.

We see in Fig. 4g that, due to the biased gradients of

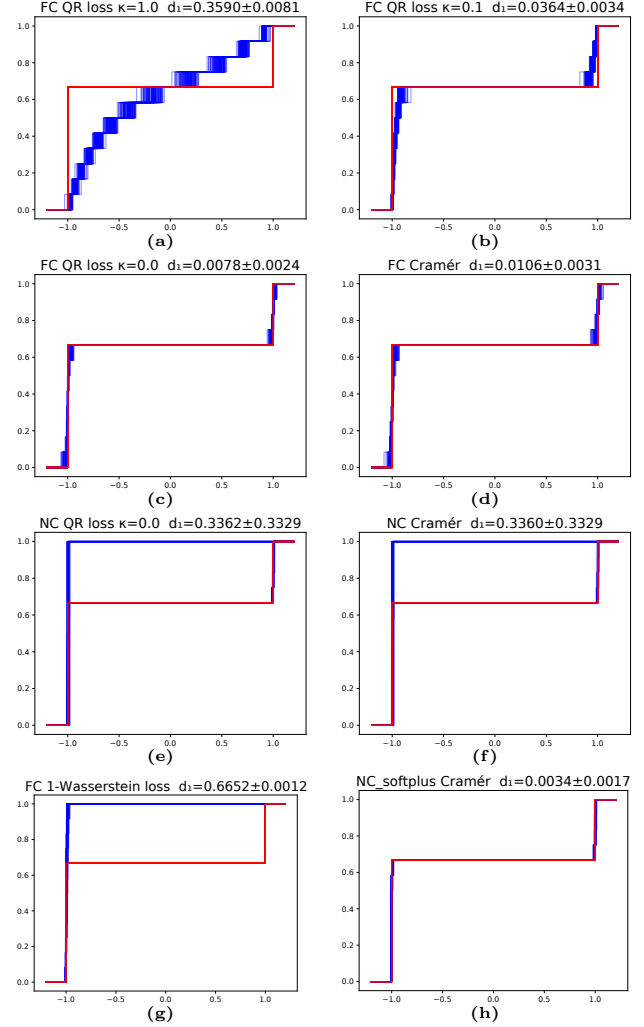


Figure 4: **Synthetic experiments.** The learned CDF for each trial is shown in blue. The average d_1 with the mixture of targets (in red) is shown for each case.

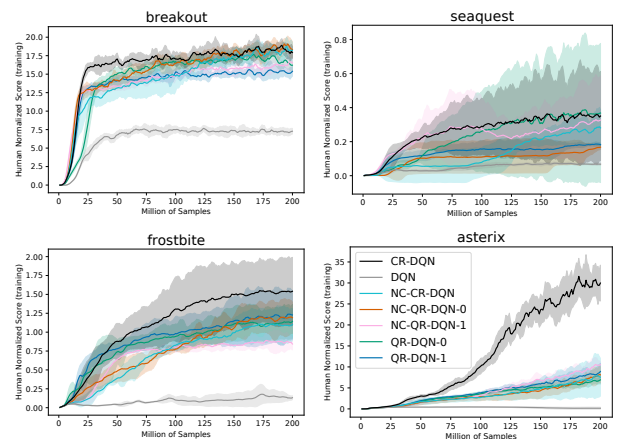


Figure 5: **Atari games.** CR: Cramér loss. The suffix is the value of κ . NC: non-crossing. Curves show mean and std. dev. of human-normalized online performance over 3 seeds, smoothed over a sliding window of 5 iter.

the 1-Wasserstein loss, the learned distribution converges to one of the Diracs instead of converging to the mixture. We also see that the Huberized QR loss yields a shrunken distribution (Fig. 4b), the effect being larger with $\kappa = 1$ (Fig. 4a). The standard QR loss (Fig. 4c) and the Cramér one (Fig. 4d) do not exhibit this effect but we see an oscillation around the actual step locations due to the lack of smoothness. The Cramér loss exhibits a slightly larger oscillation effect as suggested by the larger 1-Wasserstein distance. The non-crossing constraints make the QR loss (Fig. 4e) equivalent to the Cramér one (Fig. 4f) and reduce the oscillation effect but introduce another effect due to the specific architecture. More precisely, the ReLU activation that outputs the scale factor α (Zhou et al., 2020, Eq. (19)) is prone to the dying ReLU problem in this simple setting. This makes the NC architecture converge to one of Diracs in some of the trials. Replacing the ReLU activation by a SoftPlus solves the problem (Fig. 4h). Note that this problem is less likely to happen in more complex scenarios—with more states and actions—as in the Atari games considered next (see Glorot et al. (2011)).

7.2 Atari games

We consider four Atari games exhibiting different learning behaviours. Fig. 5 shows the online training performance (Machado et al., 2018) given by different combinations of networks and losses. The NC network (Zhou et al., 2020) and Algorithm 1 (denoted CR and used in CR-DQN and NC-CR-DQN) were implemented on top of the DQN_Z00 framework (Quan and Ostrovski, 2020) which also provides pre-computed results for the two reference algorithms QR-DQN (aka QR-DQN-1) (Dabney et al., 2018b) and DQN (Mnih et al., 2015). Equivalent hyperparameter values were used for all the methods, see Appendix C for details.

Although equivalent in theory, NC-QR-DQN-0 and NC-CR-DQN do not exactly match empirically because of GPU non-determinism and differences in numerical errors. See Appendix C for more details.

The permutation invariance of our sort-based algorithm makes the crossing quantile problem vanish, removing the need of non-crossing architectures that are prone to undesired effects as the dying ReLU problem. In these four games, the increased freedom of CR-DQN provides a significant advantage over the other methods with, in particular, a remarkable performance on Asterix.

To provide comparable results with existing work, we report, in Table 1, evaluation results over the full Atari 57 benchmark under the best agent protocol (see, e.g., Dabney et al. (2018b)) obtained with the pre-computed results provided in Quan and Ostrovski (2020) for the

	Seeds	Median
DQN	5	85%
C51	5	183%
QR-DQN-1	5	182%
IQN	5	220%
CR-DQN	3	201%

Table 1: Median of best scores across 57 Atari 2600 games, measured as percentages of human baseline (Nair et al., 2015) using reference values from DQN_Z00.

contenders. We observe that CR-DQN outperforms C51 (Bellemare et al., 2017a) and standard QR-DQN (Dabney et al., 2018b).

8 DISCUSSION

Our results shed light on QR-based algorithms by showing the equivalence of the Cramér projection with the 1-Wasserstein one, and that learning distributions with the QR loss under non-crossing constraints is essentially equivalent to learning with the Cramér loss. On the practical viewpoint, we proposed a low complexity algorithm that we tested on synthetic examples and Atari games using an unconstrained architecture and another one with non-crossing constraints.

In the unconstrained setting, symmetries creates a factorial number of optimal solutions (due to the permutations): in a stochastic optimization perspective, this could facilitate (since there are more places to find optimal solutions) and give more freedom to the deep network but it can also make the learning process unstable by jumping from one region to a symmetric one. In a constrained setting, Algorithm 1 computes an output that is equivalent to that of the QR-loss and, thus, it is subject to the same lack of smoothness that has been reported to hurt the performance in comparison to Huberized QR-loss. However, Huberization breaks the equivalence with the Cramér distance and introduces biases whose magnitude depends on the chosen κ and the scale of distributions, which can vary from one state to another. Another important point is that the architectures introducing monotonicity constraints can also introduce new effects depending on the design choices. As future work, we foresee investigating alternative approaches to smoothen the Cramér loss.

Acknowledgements

Thanks to Mourad Boudia, Eoin Thomas and Rodrigo Acuña-Agost for their insightful comments and to the anonymous reviewers whose suggestions have greatly improved this manuscript.

References

- Marc G Bellemare, Will Dabney, and Rémi Munos. A distributional perspective on reinforcement learning. In *International Conference on Machine Learning*, pages 449–458. PMLR, 2017a.
- Marc G Bellemare, Ivo Danihelka, Will Dabney, Shakir Mohamed, Balaji Lakshminarayanan, Stephan Hoyer, and Rémi Munos. The cramer distance as a solution to biased wasserstein gradients. *arXiv preprint arXiv:1705.10743*, 2017b.
- Richard Bellman. *Dynamic Programming*. Princeton University Press, 1957. ISBN 069107951X.
- Dimitri P Bertsekas and John N Tsitsiklis. *Neuro-dynamic programming*. Athena Scientific, 1996.
- Howard D Bondell, Brian J Reich, and Huixia Wang. Non-crossing quantile regression curve estimation. *Biometrika*, 97(4):825–838, 2010.
- Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym, 2016.
- Will Dabney, Georg Ostrovski, David Silver, and Rémi Munos. Implicit quantile networks for distributional reinforcement learning. In *International conference on machine learning*, pages 1096–1105. PMLR, 2018a.
- Will Dabney, Mark Rowland, Marc Bellemare, and Rémi Munos. Distributional reinforcement learning with quantile regression. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018b.
- Holger Dette and Stanislav Volgushev. Non-crossing non-parametric estimates of quantile curves. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 70(3):609–627, 2008.
- Xavier Glorot, Antoine Bordes, and Yoshua Bengio. Deep sparse rectifier neural networks. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 315–323. JMLR Workshop and Conference Proceedings, 2011.
- Maya Gupta, Andrew Cotter, Jan Pfeifer, Konstantin Vovodski, Kevin Canini, Alexander Mangylov, Wojciech Moczydlowski, and Alexander Van Esbroeck. Monotonic calibrated interpolated look-up tables. *The Journal of Machine Learning Research*, 17(1):3790–3836, 2016.
- Peter Hall, Rodney CL Wolff, and Qiwei Yao. Methods for estimating a conditional distribution function. *Journal of the American Statistical association*, 94(445):154–163, 1999.
- Xuming He. Quantile curves without crossing. *The American Statistician*, 51(2):186–192, 1997.
- Peter J. Huber. Robust Estimation of a Location Parameter. *The Annals of Mathematical Statistics*, 35(1):73 – 101, 1964.
- Diederick P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015.
- Roger Koenker, Pin Ng, and Stephen Portnoy. Quantile smoothing splines. *Biometrika*, 81(4):673–680, 1994.
- Harold Kushner and G George Yin. *Stochastic approximation and recursive algorithms and applications*, volume 35. Springer Science & Business Media, 2003.
- Yufeng Liu and Yichao Wu. Stepwise multiple quantile regression estimation using non-crossing constraints. *Statistics and its Interface*, 2(3):299–310, 2009.
- Marlos C. Machado, Marc G. Bellemare, Erik Talvitie, Joel Veness, Matthew J. Hausknecht, and Michael Bowling. Revisiting the arcade learning environment: Evaluation protocols and open problems for general agents. *Journal of Artificial Intelligence Research*, 61:523–562, 2018.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fiedjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- Alfred Müller. Integral probability metrics and their generating classes of functions. *Advances in Applied Probability*, pages 429–443, 1997.
- Arun Nair, Praveen Srinivasan, Sam Blackwell, Cagdas Alceice, Rory Fearon, Alessandro De Maria, Vedavyas Panneershelvam, Mustafa Suleyman, Charles Beattie, Stig Petersen, et al. Massively parallel methods for deep reinforcement learning. In *ICML Workshop on Deep Learning*, 2015.
- Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 2014.
- John Quan and Georg Ostrovski. DQN Zoo: Reference implementations of DQN-based agents, 2020. URL http://github.com/deepmind/dqn_zoo.
- Mark Rowland, Marc Bellemare, Will Dabney, Rémi Munos, and Yee Whye Teh. An analysis of categorical distributional reinforcement learning. In *International Conference on Artificial Intelligence and Statistics*, pages 29–37. PMLR, 2018.
- Natasa Tagasovska and David Lopez-Paz. Single-model uncertainties for deep learning. In *Advances in Neural Information Processing Systems*, pages 6417–6428, 2019.
- Hado van Hasselt, Arthur Guez, and David Silver. Deep reinforcement learning with double q-learning. *CoRR*, abs/1509.06461, 2015.
- Christopher JCH Watkins and Peter Dayan. Q-learning. *Machine learning*, 8(3-4):279–292, 1992.
- Derek Yang, Li Zhao, Zichuan Lin, Tao Qin, Jiang Bian, and Tie-Yan Liu. Fully parameterized quantile function for distributional reinforcement learning. *Advances in Neural Information Processing Systems*, 32:6193–6202, 2019.
- Fan Zhou, Jianing Wang, and Xingdong Feng. Non-crossing quantile regression for distributional reinforcement learning. *Advances in Neural Information Processing Systems*, 33:15909–15919, 2020.
- Fan Zhou, Zhoufan Zhu, Qi Kuang, and Liwen Zhang. Non-decreasing quantile function network with efficient exploration for distributional reinforcement learning. In Zhi-Hua Zhou, editor, *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21*, pages 3455–3461. International Joint Conferences on Artificial Intelligence Organization, 8 2021. Main Track.

Supplementary Material: A Cramér Distance perspective on Quantile Regression based Distributional Reinforcement Learning

A ADDITIONAL PROOFS

Lemma 1. For any $\tau, \tau' \in [0, 1]$ with $\tau < \tau'$ and CDF F with inverse F^{-1} , let $t \equiv F^{-1}(\tau)$ and $t' \equiv F^{-1}(\tau')$ and consider the scaled and vertically shifted Heaviside step function

$$H_{\theta}^{\tau, \tau'}(z) \equiv \tau + (\tau' - \tau)\mathbf{1}_{z \geq \theta}.$$

Then, for any $p \in \mathbb{R}, p > 1$, the set of $\theta \in [t, t']$ minimizing

$$\int_t^{t'} |F(z) - H_{\theta}^{\tau, \tau'}|^p dz \tag{10}$$

is given by

$$\left\{ \theta \in [t, t'] \mid F(\theta) = \left(\frac{\tau + \tau'}{2} \right) \right\}. \tag{11}$$

If F^{-1} is the inverse CDF, then $F^{-1}((\tau + \tau')/2)$ is always a valid minimizer, and if F^{-1} is continuous at $(\tau + \tau')/2$, then $F^{-1}((\tau + \tau')/2)$ is the unique minimizer.

Proof. A visual intuition of the proof is shown in Figure 2. We decompose the integral as follows

$$\int_t^{t'} |F(z) - H_{\theta}^{\tau, \tau'}(z)|^p dz = \int_t^{\theta} (F(z) - \tau)^p dz + \int_{\theta}^{t'} (\tau' - F(z))^p dz \tag{46}$$

$$= \lim_{a \rightarrow t} \int (F(z) - \tau)^p dz|_a^{\theta} + \lim_{b \rightarrow t'} \int (\tau' - F(z))^p dz|_{\theta}^b \tag{47}$$

where the limits are taken to cover the particular cases of $t = -\infty$ and $t' = \infty$. The last equation stems from the second fundamental theorem of calculus, which holds since the integrated functions are bounded and the set of points of discontinuity has measure zero (since F is a CDF). Since we are minimizing with respect to θ we can drop the constant terms and consider

$$\frac{d}{d\theta} \int (F(z) - \tau)^p dz|_{\theta} - \int (\tau' - F(z))^p dz|_{\theta} = (F(\theta) - \tau)^p - (\tau' - F(\theta))^p. \tag{48}$$

First note that for $\theta \in [t, t']$, we have $F(\theta) - \tau > 0$ and $\tau' - F(\theta) > 0$. Then, equating the derivative to zero yields

$$(F(\theta) - \tau)^p - (\tau' - F(\theta))^p = 0 \Leftrightarrow F(\theta) - \tau = \tau' - F(\theta) \Leftrightarrow F(\theta) = \frac{\tau + \tau'}{2}. \tag{49}$$

By replacing $=$ by $<$ (resp., $>$) in the previous equations, we see that the derivative is strictly negative (resp., strictly positive) if $F(\theta) < (\frac{\tau + \tau'}{2})$ (resp., $F(\theta) > (\frac{\tau + \tau'}{2})$), which proves the claim. If there is a jump in F making F^{-1} undefined at $\frac{\tau + \tau'}{2}$, the set defined in Eq. (11) becomes empty. However, the previous inequalities determining the sign of derivative still hold and the quantity to be minimized is a continuous function of θ (see Eq. (46)). Therefore, if we redefine F^{-1} to be the inverse CDF, it makes $F^{-1}((\tau + \tau')/2)$ always a valid minimizer. NB: if the standard inverse F^{-1} is undefined at τ or τ' , the whole derivation still holds if F^{-1} is redefined as the inverse CDF. \square

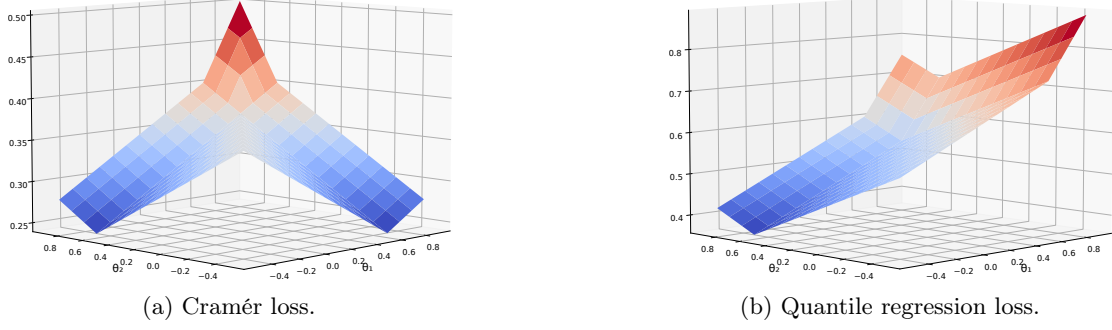


Figure 7: **Symmetry in the Cramér loss landscape (a) in comparison to the QR loss (b).** The loss landscape correspond to estimating the return distribution of a state s_0 with transitions to states s_1 and s_2 with probability $1/3$ and $2/3$, respectively, whose return distributions are Diracs located at -0.5 and 0.6 respectively, with $N = 3$. The plots are for a fixed $\theta_0 = -0.5$. Notice that when $\theta_0 \leq \theta_1 \leq \theta_2$, the two losses have collinear gradients as shown in Corollary 2.

B CORRECTNESS OF ALGORITHM 1

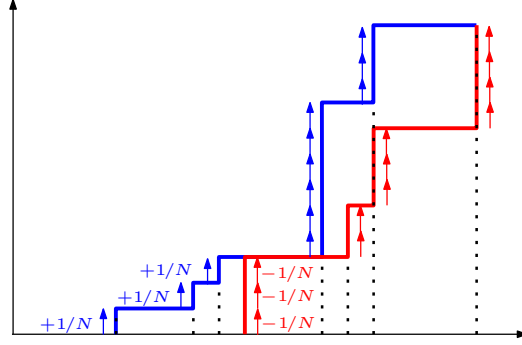


Figure 6: **Cramér loss algorithm.** Illustration of Δ_τ computation by accumulating increments/decrements.

Proposition 1. Given two distributions $F(z) = \frac{1}{N} \sum_{i=1}^N \mathbf{1}_{z \geq \theta_i}$, and $\bar{F}(z) = \frac{1}{N} \sum_{i=1}^N \mathbf{1}_{z \geq \bar{\theta}_i}$, Algorithm 1 computes

$$\int_{-\infty}^{\infty} (F(z) - \bar{F}(z))^2 dz = \sum_{i=1}^{2N-1} (\theta'_{i+1} - \theta'_i) \left(\sum_{j \text{ s.t. } \theta_j \leq \theta'_i} \frac{1}{N} - \sum_{j \text{ s.t. } \bar{\theta}_j \leq \theta'_i} \frac{1}{N} \right)^2. \quad (50)$$

Proof. Consider the sorted sequence of merged quantiles

$$\theta' \equiv \theta'_1, \dots, \theta'_{2N} \equiv \text{sort} \left(\{\theta_i\}_{i=1..N} \cup \{\bar{\theta}_i\}_{i=1..N} \right). \quad (51)$$

We have that $F(z) - \bar{F}(z) \equiv \Delta_i$ is constant in $[\theta'_i, \theta'_{i+1})$, $\forall i \in 1..2N - 1$ and is zero elsewhere. Therefore,

$$\int_{-\infty}^{\infty} (F(z) - \bar{F}(z))^2 dz = \sum_{i=1}^{2N-1} \int_{\theta'_i}^{\theta'_{i+1}} (F(z) - \bar{F}(z))^2 dz = \sum_{i=1}^{2N-1} \Delta_i^2 (\theta'_{i+1} - \theta'_i) \quad (52)$$

If $\theta'_i \leq z < \theta'_{i+1}$, then

$$F(z) = \frac{1}{N} \sum_{j=1}^N \mathbf{1}_{z \geq \theta_j} = \frac{1}{N} \sum_{j \text{ s.t. } \theta_j \leq \theta'_i} 1 \quad (53)$$

$$\bar{F}(z) = \frac{1}{N} \sum_{j=1}^N \mathbf{1}_{z \geq \bar{\theta}_j} = \frac{1}{N} \sum_{j \text{ s.t. } \bar{\theta}_j \leq \theta'_i} 1 \quad (54)$$

and thus

$$\Delta_i = \sum_{j \text{ s.t. } \theta_j \leq \theta'_i} \frac{1}{N} - \sum_{j \text{ s.t. } \bar{\theta}_j \leq \theta'_i} \frac{1}{N}, \quad (55)$$

which proves (50).

The algorithm computes the differences $(\theta'_{i+1} - \theta'_i)$ and stores them in Δ_z . After the steps

$$\Delta_\tau \leftarrow \text{concat} \left(-\frac{1}{N} \mathbf{1}_N, \frac{1}{N} \mathbf{1}_N \right) \quad (56)$$

$$\Delta_\tau \leftarrow \Delta_\tau [i_1, \dots, i_{2N}], \quad (57)$$

in words, the i -th element of the vector Δ_τ is $-\frac{1}{N}$ if θ'_i comes from $\bar{\theta}$ or $\frac{1}{N}$ otherwise, i.e.

$$\Delta_\tau [i] = \frac{1}{N} (-1)^{\mathbb{1}_{\exists j \theta'_i \equiv \bar{\theta}_j}} \quad (58)$$

where \equiv denotes symbol equality. See Fig. 6 for an illustration. After the final step

$$\Delta_\tau \leftarrow \text{cumsum}(\Delta_\tau)[:, -1], \quad (59)$$

the i -th element of the vector Δ_τ can be expressed as

$$\Delta_\tau [i] = \frac{1}{N} \sum_{k=1}^i (-1)^{\mathbb{1}_{\exists j \theta'_k \equiv \bar{\theta}_j}}. \quad (60)$$

If $\theta'_i \neq \theta'_{i+1}$, then $\Delta_\tau [i] = \Delta_i$. Otherwise, $\Delta_\tau [i] \neq \Delta_i$, but, since $\theta'_{i+1} - \theta'_i = 0$, the corresponding term in (50) is zero too. Therefore, the algorithm produces the claimed output. \square

C EXPERIMENTAL DETAILS

C.1 The networks

We describe here the two types of architecture used in the experiments. See Figure 8 for an illustration. QR-DQN (Dabney et al., 2018b) uses a series of convolutional layers each one followed by a ReLU activation in order to extract features from the input frames to obtain an embedded state $e(s) \in \mathbb{R}^{d'}$ (Fig. 8a). They are followed by a fully connected network with λ layers of η nodes each and an output layer of size $|\mathcal{A}| \times N$ (Fig. 8b).

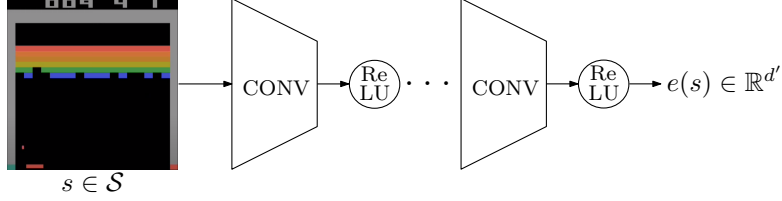
Following Zhou et al. (2020), the NC network used in the experiments replaces the fully connected network of QR-DQN by a *Non-Crossing Quantile Logit* (NCQL) network and a *Scale Factor* (SF) network (Fig. 8c). The NCQL network maps the embedded state $e(s)$ to $|\mathcal{A}| \times N$ -dimensional logits by using a fully connected network of λ layers with η nodes each, which is followed by a softmax transformation. Then a cumulated sum operator produces a non-decreasing sequence of normalized quantile values $\psi(s)[a, 1], \dots, \psi(s)[a, N]$ for each action a . The SF network produces an output in $|\mathcal{A}| \times 2$ representing the scale $\alpha(s)[a]$ and the location $\beta(s)[a]$ of the CDF, by mapping the embedded state $e(s)$ through a fully-connected network of λ layers and η nodes. A ReLU function is applied to the output corresponding to the scale $\alpha(s)[a]$ to ensure its non-negativity. The final quantile estimates are obtained by combining the outputs of the two networks as follows

$$q(s)[a, i] := \alpha(s)[a] \times \psi(s)[a, i] + \beta(s)[a]; i = 1, \dots, N, a = 1, \dots, |\mathcal{A}|. \quad (61)$$

Since in the synthetic experiment there is only one state, the feature extraction layers are removed and therefore QR-DQN turns into a standard fully-connected (FC) architecture. The NC architecture in this case boils down to the combined NCQL and SF networks.

C.2 Synthetic experiment

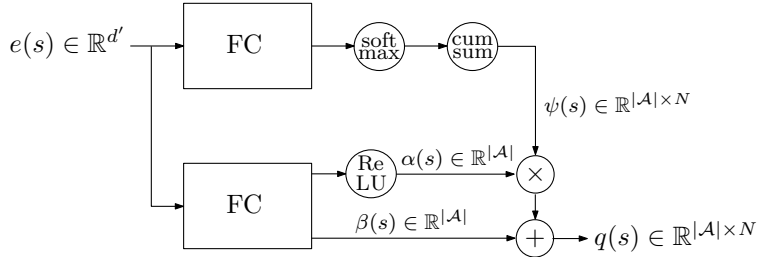
This experiment considers an MDP with only one possible action in one state s that can transition to two possible states s_1 and s_2 with probabilities $2/3$ and $1/3$, respectively, each with a different return distribution—a Dirac located at -1 and 1 respectively. The goal is to learn the return distribution at s .



(a) DQN backbone: feature extraction by a series of convolutional layers with ReLU activations.



(b) QR-DQN head: a fully-connected network.



(c) Non-Crossing (NC) head: combination of NCQL (upper part) and SF (lower part) networks.

Figure 8: Architectures used in the experiments.

Since we aim at learning the return distribution of only one state, the two networks (FC and NC) take a constant scalar input 1. The FC and NC networks have $\lambda = 2$ hidden layers of $\eta = 45$ and $\eta = 32$ nodes, respectively, and an output of $N = 12$ quantiles allowing to represent the mixture exactly.

We use the Adam optimizer (Kingma and Ba, 2015) with a learning rate of 1×10^{-3} and a batch size of 32.

C.3 Atari games

We implemented our algorithm on top of the DQN_Z00 (Quan and Ostrovski, 2020) framework, which integrates reference implementations of RL algorithms with the gym/atari-py RL environment (Brockman et al., 2016). DQN_Z00 provides pre-computed simulation results for each of these algorithms, each of them being run on 5 seeds and on the full set of 57 Atari 2600 games.

In order to implement the NC architecture, we replaced the fully connected network in the DQN_Z00 implementation of QR-DQN by the combination of the NCQL and SF networks.

Hyperparameters For model training, we set our hyperparameters with the values used in Dabney et al. (2018b) for the epsilon decay and experience replay settings. Notice that ADAM’s invariance (cf. Remark 2) is broken with the parameter ϵ used in the update step to avoid divisions by zero (Kingma and Ba, 2015): $\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon)$, where \hat{m}_t and \hat{v}_t are the first and second moment estimates at timestep t , which are scaled by a factor of c and c^2 respectively when the gradient is scaled by c . Since the gradient of the Cramér loss is $c = 2/N$ times the one of the QR loss (cf. Corollary 2), we use the adjusted $\epsilon' \equiv (2/N)\epsilon$ to have equivalent update steps. Each experiment consists in 200 iterations. Each iteration is made of a learning phase (1 million frames), followed by an evaluation phase, on 500 thousands frames. We thus use the same experiment procedure, and the same epsilon hyperparameter than the one used for the experiments provided with DQN_Z00; also, our neural network architecture uses the same three convolutional layers as the other algorithms implemented within DQN_Z00. The experiment settings being the same, our experiment performance can therefore be compared to the experiment data provided with DQN_Z00 for the other algorithms. Finally, the neural networks are defined by $\lambda = 1$, $\eta = 512$ and $N = 201$. Table 2 summarizes the hyperparameters and their values.

Table 2: Hyperparameters for $^*\text{-}\{C|Q\}\text{R-DQN}$ methods.

Hyperparameter	Value	Comment
replay_capacity	1e6	
min_replay_capacity_fraction	0.05	Min replay set size for learning
batch_size	32	
max_frames_per_episode	108000	= 30 min
num_action_repeats	4	In frames
num_stacked_frames	4	
exploration_epsilon_begin_value	1	
exploration_epsilon_end_value	0.01	
exploration_epsilon_decay_frame_fraction	0.02	
eval_exploration_epsilon	0.001	
target_network_update_period	4e4	
learning_rate	5e-5	
optimizer_epsilon (for $^*\text{-CR-DQN}$ and NC-QR-DQN-0)	0.01 / 32 * 2/N	ADAM’s parameter
optimizer_epsilon (otherwise)	0.01 / 32	ADAM’s parameter
additional_discount	0.99	Discount_rate multiplier
max_abs_reward	1	
max_global_grad_norm	10	Gradient clipping
num_iterations	200	
num_train_frames	1e6	Per iteration
num_eval_frames	5e5	Per iteration
learn_period	16	One learning step each 16 frames
num_quantiles	201	N
Convolutional layer 1	32, (8, 8), (4, 4)	num_features, kernel_shape, stride
Convolutional layer 2	64, (4, 4), (2, 2)	
Convolutional layer 3	64, (3, 3), (1, 1)	
n_layers	1	Number of hidden layers λ
n_nodes	512	Number of nodes η per hidden layer

Online training performance Performance during training protocol: this protocol, described in Machado et al. (2018), puts the emphasis on the learning quality. It consists in using normalized training scores to evaluate the algorithms. Human-normalization of score is given by van Hasselt et al. (2015): $\text{normalized_score} = \frac{\text{agent_score} - \text{random}}{\text{human} - \text{random}}$ where random and human are baseline scores, given for each game.

Detailed results Figure 9 shows the online training performance of CR-DQN in comparison to the pure distributional contenders C51 , QR-DQN (aka QR-DQN-1) and IQN , on the full Atari-57 benchmark. For C51 , QR-DQN and IQN , 5 seeds were used (provided by DQN_Z00 Quan and Ostrovski (2020)). For CR-DQN , 3 seeds were used.

On the empirical matching of NC-QR-DQN-0 and NC-CR-DQN In order to make NC-QR-DQN-0 and NC-CR-DQN as practically equivalent as possible for the experiments of Figure 5, the gradient of NC-QR-DQN-0 was scaled by a factor of $2/N$ to make the effect of gradient clipping by $\text{max_global_grad_norm}$ equivalent and the same optimizer_epsilon was used (see Table 2). Despite this, numerical errors and GPU non-determinism still produce different results.

A Cramér Distance perspective on Quantile Regression based Distributional Reinforcement Learning

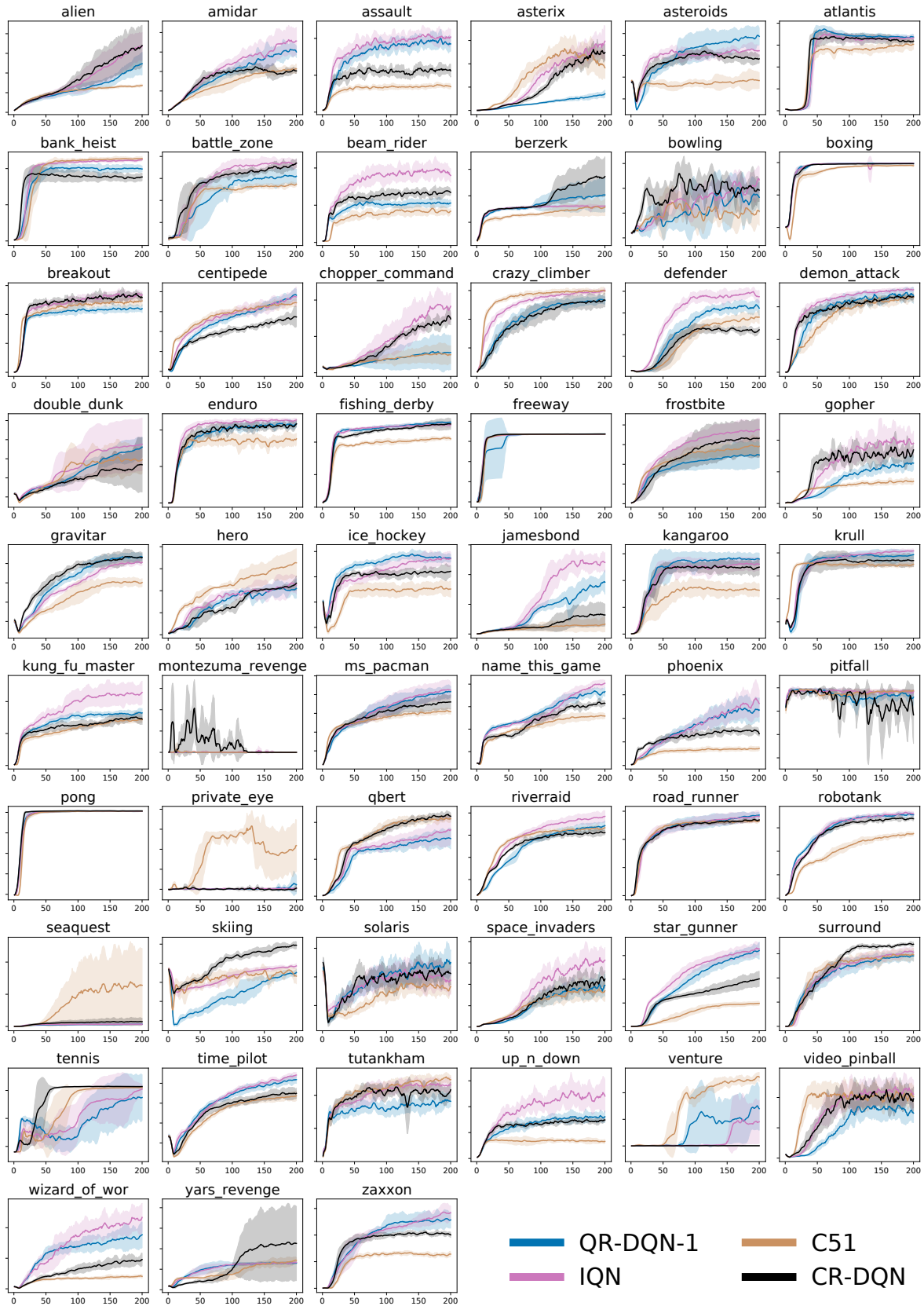


Figure 9: **Training performance on the Atari-57 benchmark.** Curves are averages over a number of seeds, smoothed over a sliding window of 5 iterations, and error bands give standard deviations.