
Distributional Reinforcement Learning via Sinkhorn Iterations

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Distributional reinforcement learning (RL) is a class of state-of-the-art algorithms
2 that estimate the whole distribution of the total return rather than only its expect-
3 ation. The representation manner of each return distribution and the choice of
4 distribution divergence are pivotal for the empirical success of distributional RL.
5 In this paper, we propose a new class of *Sinkhorn distributional RL (Sinkhorn-*
6 *DRL)* algorithm that learns a finite set of statistics, i.e., deterministic samples,
7 from each return distribution and then leverages Sinkhorn iterations to evaluate
8 the Sinkhorn distance between the current and target Bellman distributions. Re-
9 markably, Sinkhorn divergence interpolates between the Wasserstein distance and
10 Maximum Mean Discrepancy (MMD). This allows our proposed SinkhornDRL
11 algorithm to find a sweet spot leveraging the geometry of optimal transport based
12 distance and the unbiased gradient estimates of MMD. Finally, experiments on
13 the suit of 55 Atari games reveal the competitive performance of SinkhornDRL
14 algorithm as opposed to existing state-of-the-art algorithms.

15 1 Introduction

16 Classical reinforcement learning (RL) algorithms are normally based on the expectation of discounted
17 cumulative rewards that an agent observes while interacting with the environment. Recently, a new
18 class of RL algorithms called *distributional RL* estimates the full distribution of total returns and has
19 exhibited the state-of-the-art performance in a wide range of environments [2, 8, 7, 24, 26, 17].

20 From the literature of distributional RL, it is easily recognized that algorithms based on either
21 Wasserstein distance or MMD have gained great attention due to their superior performance. As such,
22 their mutual connection from the perspective of mathematical properties intrigues us to explore further
23 in order to design new algorithms. Particularly, Wasserstein distance, long known to be a powerful
24 tool to compare probability distributions with non-overlapping supports, has recently emerged as an
25 appealing contender in various machine learning applications. It is known that Wasserstein distance
26 was long disregarded because of its computational burden in its original form to solve an expensive
27 network flow problem. However, recent works [21, 14] have shown that this cost can be largely
28 mitigated by settling for cheaper approximations through strongly convex regularizers. The benefit of
29 this regularization has opened the path to wider applications of the Wasserstein distance in relevant
30 learning problems, including the design of distributional RL algorithms.

31 The Sinkhorn divergence [21] introduces the entropic regularization on the Wasserstein distance,
32 allowing it tractable for the evaluation especially in high-dimensions. It has been successfully applied
33 in numerous crucial machine learning developments, including the Sinkhorn-GAN [14] and Sinkhorn-
34 based adversarial training [23]. More importantly, it has been shown that Sinkhorn divergence
35 interpolates Wasserstein distance and MMD, and their equivalence form can be well established in the
36 limit cases [11, 18, 17]. However, a Sinkhorn-based distributional RL algorithm has not yet been

37 formally proposed and its connection with algorithms based on Wasserstein distance and MMD is
 38 also less studied. Therefore, a natural question is *can we design a new class of distributional RL*
 39 *algorithms via Sinkhorn divergence, thus bridging the gap between existing two main branches of*
 40 *distributional RL algorithms?* Moreover, the dominant quantile-based algorithms, e.g., QR-DQN [8],
 41 aimed at approximating Wasserstein distance, suffers from the non-crossing issue in the quantile
 42 estimation [26], while sample-based Sinkhorn algorithm can naturally circumvent this problem.

43 In this paper, we propose a novel distributional RL algorithm based on *Sinkhorn divergence*. Firstly,
 44 we point out the key roles of distribution divergence and representation of value distribution in the
 45 design of distributional RL. After a detailed introduction of our proposed SinkhornDRL algorithm,
 46 we theoretically analyze its convergence guarantee and moment matching behavior of distributional
 47 Bellman operators under Sinkhorn divergence. Thus, a regularized MMD equivalence form of
 48 Sinkhorn divergence is derived, interpreting the empirical success of our algorithms in real applications.
 49 Finally, we compare the performance of our SinkhornRL algorithm with typical baselines on 55 Atari
 50 games, verifying the competitive performance of our proposal. Our approach inspires researchers
 51 to find a trade-off that simultaneously leverages the geometry of the Wasserstein distance and the
 52 favorable unbiased gradient estimate property of MMD while designing new distributional RL
 53 algorithms in the future.

54 2 Preliminary Knowledge

55 2.1 Distributional Reinforcement Learning

56 In the classical RL setting, an agent interacts with an environment via a Markov decision pro-
 57 cess (MDP), a 5-tuple $(\mathcal{S}, \mathcal{A}, R, P, \gamma)$, where \mathcal{S} and \mathcal{A} are the state and action spaces, respectively. P
 58 is the environment transition dynamics, R is the reward function and $\gamma \in (0, 1)$ is the discount factor.

59 **From Value function to Value distribution.** Given a policy π , the discounted sum of future
 60 rewards is a random variable $Z^\pi(s, a) = \sum_{t=0}^{\infty} \gamma^t R(s_t, a_t)$, where $s_0 = s$, $a_0 = a$, $s_{t+1} \sim$
 61 $P(\cdot|s_t, a_t)$, and $a_t \sim \pi(\cdot|s_t)$. In the control setting, expectation-based RL is based on the action-
 62 value function $Q^\pi(s, a)$, which is the expectation of $Z^\pi(s, a)$, i.e., $Q^\pi(s, a) = \mathbb{E}[Z^\pi(s, a)]$. By
 63 contrast, distributional RL focuses on the action-value distribution, the full distribution of $Z^\pi(s, a)$,
 64 and the incorporation of additional distributional knowledge intuitively interprets its empirical success.

65 **Distributional Bellman operators.** For the policy evaluation in expectation-based RL, the action-
 66 value function is updated via the Bellman operator $\mathcal{T}^\pi Q(s, a) = \mathbb{E}[R(s, a)] + \gamma \mathbb{E}_{s' \sim p, \pi} [Q(s', a')]$.
 67 In distributional RL, the action-value distribution of $Z^\pi(s, a)$ is updated via the distributional Bellman
 68 operator \mathfrak{T}^π

$$\mathfrak{T}^\pi Z(s, a) = R(s, a) + \gamma Z(s', a'), \quad (1)$$

69 where $s' \sim P(\cdot|s, a)$ and $a' \sim \pi(\cdot|s')$. The equality in Eq. 1 implies that random variables of
 70 both sides are equal in distribution. The distributional Bellman operator \mathfrak{T}^π is contractive under
 71 certain distribution divergence metrics, but the distributional Bellman optimality operator \mathfrak{T} can only
 72 converge to a set of optimal non-stationary value distributions in a weak sense [9].

73 2.2 Divergences between Measures

74 **Optimal Transport (OT) and Wasserstein Distance** The optimal transport (OT) metric between
 75 two probability measures (μ, ν) supported on two metric spaces is defined as the solution of the linear
 76 program:

$$\min_{\Pi \in \Pi(\mu, \nu)} \int c(x, y) d\Pi(x, y), \quad (2)$$

77 where c is the cost function and Π is the joint distribution with marginals (μ, ν) . Wasserstein distance
 78 (a.k.a. earth mover distance) is a special case of optimal transport with the Euclidean norm as the
 79 cost function. In particular, given two scalar random variables X and Y , p -Wasserstein metric W_p
 80 between the distributions of X and Y can be simplified as

$$W_p(X, Y) = \left(\int_0^1 |F_X^{-1}(\omega) - F_Y^{-1}(\omega)|^p d\omega \right)^{1/p}, \quad (3)$$

81 where F^{-1} is the inverse cumulative distribution function of a random variable. The desirable
 82 geometric property of Wasserstein distance allows it to recover full support of measures, but it suffers
 83 from the curse of dimension [13, 1].

84 **Maximum Mean Discrepancy** The squared Maximum Mean Discrepancy (MMD) MMD_k^2 with
 85 the kernel k is formulated as

$$\text{MMD}_k^2 = \mathbb{E}[k(X, X')] + \mathbb{E}[k(Y, Y')] - 2\mathbb{E}[k(X, Y)], \quad (4)$$

86 where $k(\cdot, \cdot)$ is a continuous kernel on \mathcal{X} . X' (resp. Y') is a random variable independent of X
 87 (resp. Y). If k is a trivial kernel, MMD degenerates to the energy distance. Mathematically, the “flat”
 88 geometry that MMD induces on the space of probability measures does not faithfully lift the ground
 89 distance [11], but MMD is cheaper to compute than OT and has a smaller sample complexity, i.e.,
 90 approximating the distance with samples of measures [13]. We provide the detailed introduction of
 91 more distribution divergences in Appendix A.

92 3 Roles of Distribution Divergence and Representation in distributional RL

93 3.1 Distributional RL: From Neural Q-Fitted Iteration to Neural Z-Fitted Iteration

94 **Neural Q-Fitted Iteration.** It is known that Deep Q Learning [16] can be simplified into *Neural*
 95 *Q-Fitted Iteration* [10] under tricks of experience replay and the target network Q_{θ^*} , where we update
 96 parameterized $Q_{\theta}(s, a)$ in each iteration k :

$$Q_{\theta}^{k+1} = \underset{Q_{\theta}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n [y_i - Q_{\theta}^k(s_i, a_i)]^2, \quad (5)$$

97 where the target $y_i = r(s_i, a_i) + \gamma \max_{a \in \mathcal{A}} Q_{\theta^*}^k(s'_i, a)$ is fixed within every T_{target} steps to update
 98 target network Q_{θ^*} by letting $\theta^* = \theta$ and the experience buffer induces independent samples
 99 $\{(s_i, a_i, r_i, s'_i)\}_{i \in [n]}$. In an ideal case that neglects the non-convexity and TD approximation errors,
 100 we have $Q_{\theta}^{k+1} = \mathcal{T}Q_{\theta}^k$, which is exactly equivalent to updating under Bellman optimality operator.

101 **Neural Z-Fitted Iteration.** Analogous to neural Q-fitted iteration, we can also simplify value-based
 102 distributional RL methods based on a parameterized Z_{θ} into a *Neural Z-fitted Iteration* as

$$Z_{\theta}^{k+1} = \underset{Z_{\theta}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n d_p(Y_i, Z_{\theta}^k(s_i, a_i)), \quad (6)$$

103 where the target $Y_i = R(s_i, a_i) + \gamma Z_{\theta^*}^k(s'_i, \pi_Z(s'))$ with $\pi_Z(s') = \operatorname{argmax}_{a'} \mathbb{E}[Z_{\theta^*}^k(s', a')]$ is
 104 fixed within every T_{target} steps to update target network Z_{θ^*} , and d_p is a divergence metric between
 105 two distributions.

106 3.2 Key Roles of d_p and Z_{θ}

107 Within the Neural Z-fitted Iteration framework proposed in Eq. 6, we observe that the choice of
 108 representation manner on Z_{θ} and the metric d_p are pivotal for the distributional RL algorithms. For
 109 instance, QR-DQN [8] approximates Wasserstein distance W_p , which leverages quantiles to represent

Algorithm	d_p Distribution Divergence	Representation Z_{θ}	Convergence Rate of \mathfrak{T}^{π}	Sample Complexity of d_p
C51 [2]	Cramér distance	Histogram	$\sqrt{\gamma}$	\searrow
QR-DQN [8]	Wasserstein distance	Quantiles	γ	$\mathcal{O}(n^{-\frac{1}{d}})$
MMDDRL [17]	MMD	Samples	$\gamma^{\alpha/2}$ with kernel k_{α}	$\mathcal{O}(1/n)$
SinkhornDRL (ours)	Sinkhorn divergence	Samples	$\gamma (\epsilon \rightarrow 0)$ $\gamma^{\alpha/2} (\epsilon \rightarrow \infty)$	$\mathcal{O}(n^{\frac{\epsilon}{\epsilon + d/2} \sqrt{\pi}})$ ($\epsilon \rightarrow 0$) $\mathcal{O}(n^{-\frac{1}{2}})$ ($\epsilon \rightarrow \infty$)

Table 1: Comparison between typical distributional RL algorithms under different distribution divergences and representation of Z_{θ} . $k_{\alpha} = -\|x - y\|^{\alpha}$ in MMDDRL, d is the sample dimension and $\kappa = 2\beta d + \|c\|_{\infty}$, where the cost function c is β -Lipschitz [13]. Sample complexity of MMD can be improved to $\mathcal{O}(1/n)$ using kernel herding technique [5].

110 the distribution of Z_θ . C51 [2] represents Z_θ via a categorical distribution under the convergence of
 111 Cramér distance [3, 19], while MMD distributional RL (MMDDL) [17] learns samples to represent
 112 the distribution of Z_θ based on MMD. We compare characteristics of these distribution divergence,
 113 including the convergence rate and sample complexity, in Table 1. Theoretical results regarding
 114 Sinkhorn divergence is based on [13] and the detailed convergence proof of other distances is also
 115 provided in Appendix A. In summary, we argue that d_p and Z_θ are two crucial factors in distributional
 116 RL design, based on which we introduce our Sinkhorn distributional RL.

117 4 Sinkhorn Distributional RL (SinkhornDRL)

118 In this section, we firstly introduce Sinkhorn divergence and apply it in distributional RL. Next, we
 119 conduct a theoretical analysis about the convergence speed and a new moment matching manner of
 120 our algorithm under the Sinkhorn divergence. Finally, a practical Sinkhorn iteration algorithm is
 121 introduced to evaluate the Sinkhorn divergence.

122 4.1 Sinkhorn Divergence and Genetic Algorithm

123 We design Sinkhorn distributional RL algorithm via Sinkhorn divergence. Sinkhorn divergence [21] is
 124 a tractable loss to approximate the optimal transport problem by leveraging an entropic regularization
 125 to turn the original Wasserstein distance into a differentiable and more robust quantity. The resulting
 126 loss can be computed using Sinkhorn fixed point iterations, which is naturally suitable for modern deep
 127 learning frameworks. In particular, the entropic smoothing generates a family of losses interpolating
 128 between Wasserstein distance and Maximum Mean Discrepancy (MMD). As such, it allows us to find
 129 a sweet trade-off that simultaneously leverages the geometry of Wasserstein distance on the one hand,
 130 and the favorable high-dimensional sample complexity and unbiased gradient estimates of MMD. We
 131 introduce the entropic regularized Wasserstein distance $\mathcal{W}_{c,\varepsilon}(\mu, \nu)$ as

$$\min_{\Pi \in \Pi(\mu, \nu)} \int c(x, y) d\Pi(x, y) + \varepsilon \text{KL}(\Pi | \mu \otimes \nu), \quad (7)$$

132 where $\text{KL}(\Pi | \mu \otimes \nu) = \int \log \left(\frac{\Pi(x, y)}{d\mu(x) d\nu(y)} \right) d\Pi(x, y)$ is a strongly convex regularization. The impact
 133 of this entropy regularization is similar to ℓ_2 ridge regularization in linear regression. Next, the
 134 Sinkhorn loss [11, 14] between two measures μ and ν is defined as

$$\overline{\mathcal{W}}_{c,\varepsilon}(\mu, \nu) = 2\mathcal{W}_{c,\varepsilon}(\mu, \nu) - \mathcal{W}_{c,\varepsilon}(\mu, \mu) - \mathcal{W}_{c,\varepsilon}(\nu, \nu). \quad (8)$$

135 As demonstrated by [11], the Sinkhorn divergence $\overline{\mathcal{W}}_{c,\varepsilon}(\mu, \nu)$ is convex, smooth and positive definite
 136 that metrizes the convergence in law. In statistical physics, $\mathcal{W}_{c,\varepsilon}(\mu, \nu)$ can be re-factored as a
 137 projection problem:

$$\mathcal{W}_{c,\varepsilon}(\mu, \nu) := \min_{\Pi \in \Pi(\mu, \nu)} \text{KL}(\Pi | \mathcal{K}), \quad (9)$$

138 where \mathcal{K} is the Gibbs distribution with the density function satisfies $d\mathcal{K}(x, y) = e^{-\frac{c(x, y)}{\varepsilon}} d\mu(x) d\nu(y)$.
 139 This problem is often referred to as the “static Schrödinger problem” [15, 20] as it was initially
 140 considered in statistical physics.

141 **Distributional RL with Sinkhorn Divergence and Particle Representation.** The key of apply-
 142 ing Sinkhorn divergence in distributional RL is to simply leverage the Sinkhorn loss $\overline{\mathcal{W}}_{c,\varepsilon}$ to measure
 143 the distance between the current action-value distribution $Z_\theta(s, a)$ and the target distribution
 144 $\mathfrak{T}^\pi Z_\theta(s, a)$, yielding $\overline{\mathcal{W}}_{c,\varepsilon}(Z_\theta(s, a), \mathfrak{T}^\pi Z_\theta(s, a))$ for each s, a pairs. In terms of the representation
 145 for $Z_\theta(s, a)$, we employ the unrestricted statistics, i.e., deterministic samples, due to its superiority in
 146 MMDDL [17], instead of using predefined statistic functionals, e.g., quantiles in QR-DQN [8] or
 147 histogram partitions in C51 [2]. More concretely, we use neural networks to generate samples that
 148 approximate the value distribution. This can be expressed as $Z_\theta(s, a) := \{Z_\theta(s, a)_i\}_{i=1}^N$, where N
 149 is the number of generated samples. We refer to the samples $\{Z_\theta(s, a)_i\}_{i=1}^N$ as *particles*. Then we
 150 leverage the Dirac mixture $\frac{1}{N} \sum_{i=1}^N \delta_{Z_\theta(s, a)_i}$ to approximate the true density function of $Z^\pi(s, a)$,
 151 thus minimizing the Sinkhorn divergence between the approximate distribution and its distributional
 152 Bellman target. A detailed and generic distributional RL algorithm with Sinkhorn divergence and
 153 particle representation is provided in Algorithm 1.

Algorithm 1 Generic Sinkhorn distributional RL Update

Require: Number of generated samples N , the cost function c and hyperparameter ε .

Input: Sample transition (s, a, r', s')

- 1: **if** Policy evaluation **then**
- 2: $a^* \sim \pi(\cdot|s')$.
- 3: **else**
- 4: $a^* \leftarrow \arg \max_{a' \in \mathcal{A}} \frac{1}{N} \sum_{i=1}^N Z_\theta(s', a')_i$
- 5: **end if**
- 6: $\mathfrak{T}Z_i \leftarrow r + \gamma Z_{\theta^*}(s', a^*)_i, \forall 1 \leq i \leq N$

Output: $\overline{W}_{c,\varepsilon} \left(\{Z_\theta(s, a)_i\}_{i=1}^N, \{\mathfrak{T}Z_\theta(s, a)_j\}_{j=1}^N \right)$

154 **Remark.** By comparing the state-of-the-art MMDDRL algorithm [17], our Sinkhorn distributional
155 RL simply modifies the distribution divergence. Hence, we can also easily extend our generic
156 Sinkhorn algorithm to DQN-like architecture as well as IQN [7] and FQF [24]. A following question
157 is whether there is any theoretical connection between Sinkhorn distributional RL and algorithms
158 based on MMD and Wasserstein distance. We provide this crucial analysis in Section 4.2

159 4.2 Theoretical Analysis under Sinkhorn Divergence

160 **Convergence Analysis.** Firstly, we denote the supreme form of Sinkhorn divergence as $\overline{W}_{c,\varepsilon}^\infty(\mu, \nu)$:

161

$$\overline{W}_{c,\varepsilon}^\infty(\mu, \nu) = \sup_{(x,a) \in \mathcal{S} \times \mathcal{A}} \overline{W}_{c,\varepsilon}(\mu(x, a), \nu(x, a)). \quad (10)$$

162 We will use $\overline{W}_{c,\varepsilon}^\infty(\mu, \nu)$ to establish the convergence of \mathfrak{T}^π in Theorem 1.

163 **Theorem 1.** *If we leverage Sinkhorn loss $\overline{W}_{c,\varepsilon}(\mu, \nu)$ in Eq. 8 as the distribution divergence in*
164 *distributional RL, and **choose the unrectified kernel** $k_\alpha := -\|x - y\|^\alpha$ as $-c$ ($\alpha > 0$), it holds that*

165 (1) *As $\varepsilon \rightarrow 0$, $\overline{W}_{c,\varepsilon}(\mu, \nu) \rightarrow 2W_\alpha(\mu, \nu)$. When $\varepsilon = 0$, \mathfrak{T}^π is a γ -contraction under $\overline{W}_{c,\varepsilon}^\infty$.*

166 (2) *As $\varepsilon \rightarrow +\infty$, $\overline{W}_{c,\varepsilon}(\mu, \nu) \rightarrow \text{MMD}_{k_\alpha}^2(\mu, \nu)$. When $\varepsilon = +\infty$, \mathfrak{T}^π is $\gamma^{\alpha/2}$ -contractive under $\overline{W}_{c,\varepsilon}^\infty$.*

167 (3) *For any $\varepsilon \in (0, +\infty)$, \mathfrak{T}^π is a **closely non-expansive operator** under $\overline{W}_{c,\varepsilon}^\infty$, and the difference*
168 *term $\Delta(\gamma) \rightarrow 0$ as $\gamma \rightarrow 1$.*

169 Proof is provided in Appendix B. Theorem 1 (1) and (2) are follow-up conclusions in terms of the
170 convergence behavior of \mathfrak{T}^π based on the interpolation relationship between Sinkhorn divergence with
171 Wasserstein distance and MMD [14]. Our key theoretical contribution is for the general $\varepsilon \in (0, \infty)$,
172 the convergence behavior is determined by the “joint” KL divergence in Eq. 9 between the optimal
173 joint distribution Π^* and the Gibbs distribution associated with the cost function c . We conclude
174 that \mathfrak{T}^π is a **close** non-expansive operator and the different term $\Delta(\gamma) \rightarrow 0$ as $\gamma \rightarrow 1$. Note that γ is
175 normally very close to 1 in practice, and this is beneficial for the convergence of \mathfrak{T}^π under $\overline{W}_{c,\varepsilon}^\infty$.

176 **Remark on Theorem 1 (3).** If we consider to use Gaussian kernel, we can not guarantee \mathfrak{T}^π is
177 closely non-expansive for any $\varepsilon \in (0, \infty)$. This conclusion is consistent with those discussed
178 in MMDDRL [17], where \mathfrak{T}^π is generally not a contraction operator under MMD equipped with
179 Gaussian kernels as a counterexample has been pointed out in MMDDRL (when $\varepsilon \rightarrow +\infty$). When
180 $\varepsilon \rightarrow 0$, the γ -contractive \mathfrak{T}^π under Wasserstein distance is also not contradictory to Theorem 1
181 (3). Moreover, although we can only obtain that \mathfrak{T}^π is closely non-expansive, the expectation of
182 Z^π remains a γ -contraction (see Appendix B). In experiments, we thereby use k_α and we can also
183 demonstrate the appealing empirical performance of our SinkhornDRL algorithm in Section 5.

184 **Regularized Moment Matching under Sinkhorn Divergence.** We further examine the potential
185 reason behind the empirical success for SinkhornDRL, although only a non-expansive contraction
186 can be guaranteed for the general case when $\varepsilon \in (0, +\infty)$ as shown in Theorem 1. Inspired by the
187 similar manner in MMDDRL [17], we find that the Sinkhorn divergence with the Gaussian kernel
188 can also promote to match all moments between two distributions. More specifically, the Sinkhorn
189 divergence can be rewritten as a regularized moment matching form in Proposition 1.

190 **Proposition 1.** For $\varepsilon \in (0, +\infty)$, Sinkhorn divergence $\overline{W}_{c,\varepsilon}(\mu, \nu)$ associated with Gaussian kernels
 191 $k(x, y) = \exp(-(x - y)^2/(2\sigma^2))$ as $-c$, can be equivalent to

$$\overline{W}_{c,\varepsilon}(\mu, \nu) := \sum_{n=0}^{\infty} \frac{1}{\sigma^{2n} n!} \left(\tilde{M}_n(\mu) - \tilde{M}_n(\nu) \right)^2 + \varepsilon \mathbb{E} \left[\log \frac{(\Pi_{\varepsilon}^*(X, Y))^2}{\Pi_{\varepsilon}^*(X, X') \Pi_{\varepsilon}^*(Y, Y')} \right], \quad (11)$$

192 where Π_{ε}^* denotes the optimal Π determined by ε by evaluating the Sinkhorn divergence via
 193 $\min_{\Pi \in \Pi(\mu, \nu)} \overline{W}_{c,\varepsilon}(\mu, \nu)$. $\tilde{M}_n(\mu) = \mathbb{E}_{x \sim \mu} \left[e^{-x^2/(2\sigma^2)} x^n \right]$, and similarly for $\tilde{M}_n(\nu)$.

194 We provide the proof of Proposition 1 in Appendix C. Similar to MMDDRL associated with a
 195 Gaussian kernel [17], Sinkhorn divergence approximately performs a regularized moment matching
 196 scaled by $e^{-x^2/(2\sigma^2)}$. This similar moment matching impact intuitively explains the empirical success
 197 of SinkhornDRL as MMDDRL, although the contraction of both MMD with Gaussian kernel [17]
 198 and Sinkhorn divergence for general $\varepsilon \in (0, +\infty)$ may not be guaranteed.

199 **Equivalence to Regularized MMD distributional RL.** Based on Proposition 1, we can immedi-
 200 ately establish the connection between Sinkhorn divergence and MMD in Corollary 1, indicating that
 201 minimizing Sinkhorn divergence between two distributions is equivalent to minimizing a regularized
 202 squared MMD.

203 **Corollary 1.** For $\varepsilon \in (0, +\infty)$ and denote Π_{ε}^* as the optimal Π by evaluating the Sinkhorn divergence,
 204 it holds that

$$\overline{W}_{c,\varepsilon} := \text{MMD}_{-c}^2(\mu, \nu) + \varepsilon \mathbb{E} \left[\log \frac{(\Pi_{\varepsilon}^*(X, Y))^2}{\Pi_{\varepsilon}^*(X, X') \Pi_{\varepsilon}^*(Y, Y')} \right], \quad (12)$$

205 where we use $\overline{W}_{c,\varepsilon}$ to replace $\overline{W}_{c,\varepsilon}(\mu, \nu)$ for short.

206 Proof of Corollary 1 is provided in Appendix C. It is worthy of noting that this equivalence is
 207 established for the general case when $\varepsilon \in (0, +\infty)$, and it does not hold in the limit cases when
 208 $\varepsilon \rightarrow 0$ or $+\infty$. For example, when $\varepsilon \rightarrow +\infty$, the second part including ε in Eq. 12 is not expected to
 209 dominate. This is owing to the fact that the regularization term would be 0 as $\Pi_{\varepsilon}^* \rightarrow \mu \otimes \nu$ when
 210 $\varepsilon \rightarrow +\infty$. In summary, even though the Sinkhorn divergence was initially proposed to serve as an
 211 entropy regularized Wasserstein distance, it turns out that it is equivalent to a regularized MMD, as
 212 revealed in Corollary 1. This connection provides strong evidence for our empirical results, in which
 213 SinkhornDRL achieves competitive performance as opposed to MMDDRL.

214 4.3 Distributional RL via Sinkhorn Iterations

215 The theoretical analysis in Section 4.2 sheds light on the behavior of distributional RL with Sinkhorn
 216 divergence, but another crucial issue we need to address is how to evaluate the Sinkhorn loss
 217 effectively. Due to the advantages of Sinkhorn divergence that both enjoys geometry property of
 218 optimal transport and the computational effectiveness of MMD, we can utilize Sinkhorn's algorithm,
 219 i.e., Sinkhorn Iterations [21, 14], to evaluate the Sinkhorn loss. Notably, Sinkhorn iteration with
 220 L steps yields a differentiable and solvable efficiently loss function as the main burden involved
 221 in it is the matrix-vector multiplication, which streams well on the GPU with simply adding extra
 222 differentiable layers on the typical deep neural network, such as a DQN architecture.

223 Specifically, given two sample sequences $\{Z_i\}_{i=1}^N, \{\mathfrak{Z}Z_j\}_{j=1}^N$ in the distributional RL algorithm, the
 224 optimal transport distance is equivalent to the form:

$$\min_{P \in \mathbb{R}_+^{N \times N}} \left\{ \langle P, \hat{c} \rangle; P \mathbf{1}_N = \mathbf{1}_N, P^{\top} \mathbf{1}_N = \mathbf{1}_N \right\}, \quad (13)$$

225 where the empirical cost function $\hat{c}_{i,j} = c(Z_i, \mathfrak{Z}Z_j)$. By adding entropic regularization on optimal
 226 transport distance, Sinkhorn divergence can be viewed to restrict the search space of P in the
 227 following scaling form:

$$P_{i,j} = a_i \mathcal{K}_{i,j} b_j, \quad (14)$$

228 where $\mathcal{K}_{i,j} = e^{-\hat{c}_{i,j}/\varepsilon}$ is the Gibbs kernel defined in Eq. 9. This allows us to leverage iterations
 229 regarding the vectors a and b . More specifically, we initialize $b_0 = \mathbf{1}_N$, and then the Sinkhorn
 230 iterations are expressed as

$$a_{l+1} \leftarrow \frac{\mathbf{1}_N}{\mathcal{K} b_l} \quad \text{and} \quad b_{l+1} \leftarrow \frac{\mathbf{1}_N}{\mathcal{K}^{\top} a_{l+1}}, \quad (15)$$

Algorithm 2 Sinkhorn Iterations to Approximate $\overline{\mathcal{W}}_{c,\varepsilon} \left(\{Z_i\}_{i=1}^N, \{\mathfrak{T}Z_j\}_{j=1}^N \right)$

Input: Two samples sequences $\{Z_i\}_{i=1}^N, \{\mathfrak{T}Z_j\}_{j=1}^N$, number of Sinkhorn iterations L and hyperparameter ε .

- 1: $\hat{c}_{i,j} = c(Z_i, \mathfrak{T}Z_j)$ for $\forall i = 1, \dots, N, j = 1, \dots, N$
- 2: $\mathcal{K}_{i,j} = \exp(-\hat{c}_{i,j}/\varepsilon)$
- 3: $b_0 \leftarrow \mathbf{1}_N$
- 4: **for** $l = 1, 2, \dots, L$ **do**
- 5: $a_l \leftarrow \frac{\mathbf{1}_N}{\mathcal{K}b_{l-1}}, b_l \leftarrow \frac{\mathbf{1}_N}{\mathcal{K}a_l}$
- 6: **end for**
- 7: $\widehat{\mathcal{W}}_{c,\varepsilon} \left(\{Z_i\}_{i=1}^N, \{\mathfrak{T}Z_j\}_{j=1}^N \right) = \langle (K \odot \hat{c})b, a \rangle$

Return: $\widehat{\mathcal{W}}_{c,\varepsilon} \left(\{Z_i\}_{i=1}^N, \{\mathfrak{T}Z_j\}_{j=1}^N \right)$

231 where $\dot{\cdot}$ indicates an entry-wise division. It has been proven that Sinkhorn iteration asymptotically
232 converges to the true loss in a linear rate [14, 12, 6]. We provide a detailed algorithm description of
233 Sinkhorn iterations in Algorithm 2. With the efficient and differential Sinkhorn iterations, we can
234 easily evaluate the Sinkhorn divergence and thus let our algorithm enjoy its theoretical advantages. In
235 practice, we need to choose L and ε , and we conduct a rigorous sensitivity analysis in Section 5.

236 5 Experiments

237 We demonstrate the effectiveness of SinkhornDRL as described in Algorithm 1 on the full 55 Atari
238 2600 games. Specifically, we leverage the same architecture as QR-DQN [8], and replace the quantiles
239 output with N particles, i.e., samples. In contrast to MMDDRL, SinkhornDRL only changes the
240 distribution divergence from MMD to Sinkhorn divergence, and therefore the potential superiority in
241 the performance can be attributed to the advantages of Sinkhorn divergence. In Section 5.1, we make
242 a rigorous comparison between SinkhornDRL with other typical distributional RL algorithms from
243 the perspectives of learning curves and final ratio improvement of returns. An extensive sensitivity
244 analysis in terms of multiple hyperparameters in SinkhornDRL is provided in Section 5.2.

245 **Baselines.** Due to the interpolation characteristic of Sinkhorn divergence between Wassertein
246 distance and MMDDRL, we choose three typical distributional RL algorithms as classic baselines,
247 including QR-DQN [8] that approximates the Wasserstein distance, C51 [2] and MMDDRL [17], as
248 well as DQN [16]. MMDDRL algorithm is implemented with the same architecture as QRDQN, and
249 leverages Gaussian kernels $k_h(x, y) = \exp(-(x - y)^2/h)$ with the kernel mixture trick covering a
250 range of bandwidths h , which is same as the basic setting in the original MMDDQN paper [17]. We
251 deploy all algorithms on 55 Atari 2600 games, and reported results are averaged over 3 seeds with
252 the shade indicating the standard deviation.

253 **Hyperparameter settings.** For a fair comparison with QR-DQN, C51 and MMDDRL, we used
254 the same hyperparameters: the number of generated samples $N = 200$, Adam optimizer with
255 $\text{lr} = 0.00005, \epsilon_{\text{Adam}} = 0.01/32$. We used a target network to compute the distributional Bellman
256 target, which fits well in the neural Z-fitted iteration framework. In addition, we choose number of
257 Sinkhorn iterations $L = 10$ and smoothing hyperparameter $\varepsilon = 10.0$ in Section 5.1 as they are not
258 sensitive within a proper interval as demonstrated in Section 5.2. We choose the unrectified kernel as
259 the cost function, i.e., $-c = k_\alpha$, and select $\alpha = 2$ in k_α in our SinkhornDRL algorithm.

260 5.1 Performance of SinkhornDRL

261 Figure 1 illustrates that SinkhornDRL can achieve the competitive performance across 55 Atari games
262 compared with various baseline algorithms with different metrics d_p and representation manners on
263 Z_θ . On a large number of games, e.g., Tennis, Seaquest and Atlantis, SinkhornDRL can significantly
264 outperform other baselines, especially on Tennis where other algorithms even fail to converge. The
265 improvement of SinkhornDRL over MMDDRL empirically verifies the regularization advantage of
266 the Sinkhorn as analyzed in Corollary 1. On some games, e.g., Breakout, Pong and SpaceInvaders,

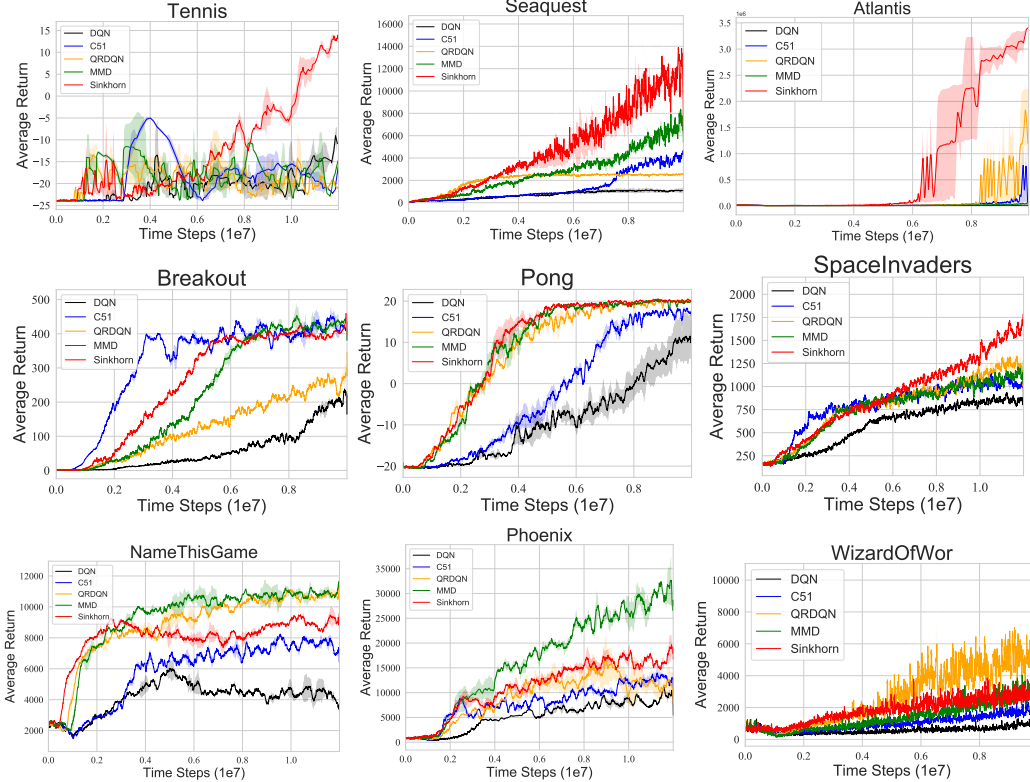


Figure 1: Learning curves of SinkhornDRL algorithm compared with DQN, C51, QR-DQN and MMD, on nine typical Atari games over 3 seeds.

267 SinkhornDRL is on par with MMDDRL and other baselines, while on the last row in Figure 1,
 268 SinkhornDRL is slightly inferior to the state-of-the-art algorithm. We provide learning curves of all
 269 typical distributional RL algorithms on all 55 Atari games in Appendix E, where SinkhornDRL still
 270 achieves the competitive performance in general.

271 To further demonstrate theoretical properties of SinkhornDRL in Theorem 1, we conduct a ratio im-
 272 provement comparison across 55 Atari games between SinkhornDRL with QR-DQN and MMDDRL,
 273 respectively. Figure 2 showcases that by comparing with QR-DQN (left), SinkhornDRL achieves
 274 better performance across more than half of considered games. More importantly, the superiority
 275 of SinkhornDRL is significant across a large amount of games, including Venture, Seaquest, Tennis
 276 and Phoenix. This empirical outperformance verifies the effectiveness and potential of smoothing
 277 Wasserstein distance in distributional RL, e.g., Sinkhorn divergence. In contrast with MMDDRL, the
 278 superiority of SinkhornDRL is reduced with the performance improvement only on a small proportion
 279 of games, while a remarkable boost of performance for SinkhornDRL on a large amount of games
 280 can be easily observed. We also report mean and median of best human-normalized scores in Table 2
 281 of Appendix D, where SinkhornDRL achieves almost state-of-the-art performance as MMDDRL on
 282 average.

283 Therefore, we conclude that SinkhornDRL is competitive with the state-of-the-art distributional
 284 RL algorithms, e.g., MMDDRL, and can be extremely superior over existing algorithms on a large
 285 proportion of games. This empirical success can be owing to theoretical advantage of Sinkhorn
 286 divergence that simultaneously makes full use of the data geometry from Wasserstein distance and
 287 the unbiased gradient estimate property from MMD, which coincides with results in Theorem 1.

288 5.2 Sensitivity Analysis and Computational Cost

289 Figure 3 (a) suggests the performance of our algorithm is robust to ε in a certain range, e.g., [1, 500],
 290 facilitating its deployment in practice. If we increase ε , SinkhornDRL’s performance tends to MMD,

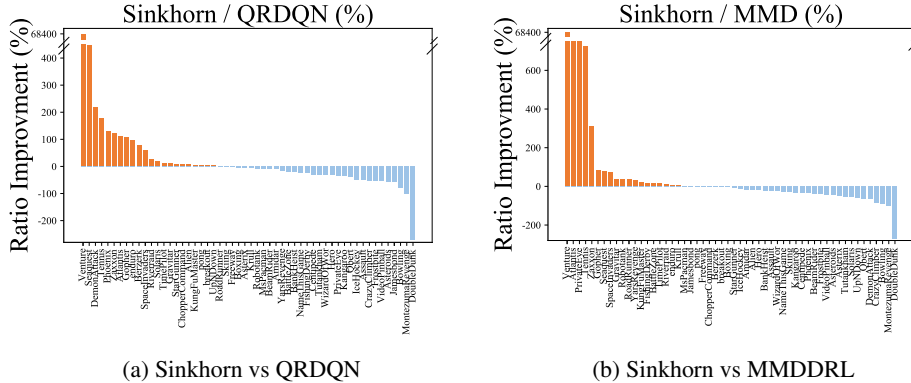


Figure 2: Ratio improvement of return for Sinkhorn distributional RL algorithm over QRDQN (left) and MMDDRL (right) over 3 seeds. For example, the ratio improvement is calculated by $(\text{Sinkhorn} - \text{QRDQN}) / \text{QRDQN}$ in the left.

291 while if we gradually decline ε , SinkhornDRL’s performance tends to QR-DQN. It is also noted that
 292 Sinkhorn iterations in Algorithm 2 will suffer from the numerical instability issue under an overly
 293 small or large ε . More results with the discussion are provided in Appendix F. It is also illustrated
 294 that our algorithm is insensitive to the number of iterations L and samples N as well, but an overly
 295 large N can slightly worsen the performance of SinkhornDRL, and at the same time increases the
 296 computational burden. Therefore, a proper number of samples, e.g., 200, is sufficient to attain an
 297 appealing performance with the computational effectiveness.

298 For the computation cost, SinkhornDRL indeed increases around 50% computation cost compared
 299 with QR-DQN and C51, but only slightly increases the cost (by around 20%) in contrast to MMDDRL.
 300 Detailed comparison is given in Appendix F.

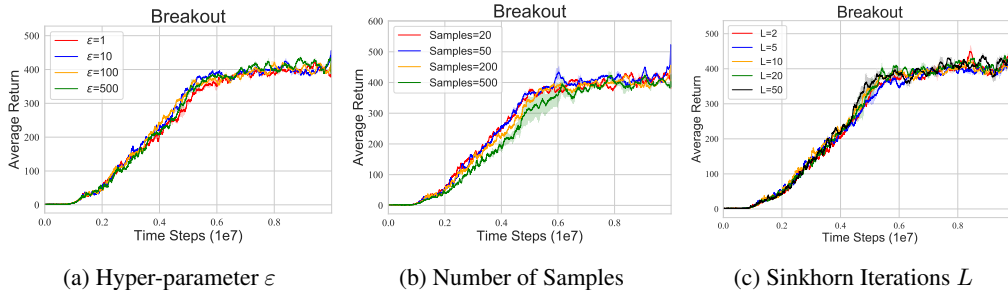


Figure 3: Sensitivity analysis of SinkhornDRL on Breakout regarding ε , number of samples, and number of iteration L . Learning curves are reported over 3 seeds.

301 6 Discussions and Conclusion

302 The main limitation of our proposal is that the superiority over existing state-of-the-art algorithms may
 303 not be sufficiently significant. To extend our algorithm for better performance, implicit generative
 304 models, including parameterizing the cost function in Sinkhorn loss, can be further incorporated. We
 305 leave it as the future work. Moreover, other divergences, e.g., those that can also smooth Wasserstein
 306 distance, can also be applied into the design of distributional RL algorithms in the future.

307 In this paper, a novel family of distributional RL algorithms based on Sinkhorn Divergence is proposed
 308 that accomplishes a competitive performance compared with the-state-of-the-art distributional RL
 309 algorithms on 55 Atari games. Theoretical analysis about the convergence and moment matching
 310 behavior is provided along with a rigorous empirical verification. Albeit being associated with MMD
 311 algorithm, distributional RL with Sinkhorn divergence is complementary to previous algorithms,
 312 leading to an important contribution among the research community.

313 **References**

- 314 [1] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial
315 networks. In *International conference on machine learning*, pages 214–223. PMLR, 2017.
- 316 [2] Marc G Bellemare, Will Dabney, and Rémi Munos. A distributional perspective on reinforce-
317 ment learning. *International Conference on Machine Learning (ICML)*, 2017.
- 318 [3] Marc G Bellemare, Ivo Danihelka, Will Dabney, Shakir Mohamed, Balaji Lakshminarayanan,
319 Stephan Hoyer, and Rémi Munos. The cramer distance as a solution to biased wasserstein
320 gradients. *arXiv preprint arXiv:1705.10743*, 2017.
- 321 [4] Pablo Samuel Castro, Subhodeep Moitra, Carles Gelada, Saurabh Kumar, and Marc G Bellemare.
322 Dopamine: A research framework for deep reinforcement learning. *CoRR abs/1812.06110*,
323 2018.
- 324 [5] Yutian Chen, Max Welling, and Alex Smola. Super-samples from kernel herding. *UAI, 109–116*.
325 *AUAI Press*, 2012.
- 326 [6] Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. *Advances in*
327 *neural information processing systems*, 26, 2013.
- 328 [7] Will Dabney, Georg Ostrovski, David Silver, and Rémi Munos. Implicit quantile networks for
329 distributional reinforcement learning. *International Conference on Machine Learning (ICML)*,
330 2018.
- 331 [8] Will Dabney, Mark Rowland, Marc G Bellemare, and Rémi Munos. Distributional reinforcement
332 learning with quantile regression. *Association for the Advancement of Artificial Intelligence*
333 *(AAAI)*, 2018.
- 334 [9] Odin Elie and Charpentier Arthur. *Dynamic Programming in Distributional Reinforcement*
335 *Learning*. PhD thesis, Université du Québec à Montréal, 2020.
- 336 [10] Jianqing Fan, Zhaoran Wang, Yuchen Xie, and Zhuoran Yang. A theoretical analysis of deep
337 q-learning. In *Learning for Dynamics and Control*, pages 486–489. PMLR, 2020.
- 338 [11] Jean Feydy, Thibault Séjourné, François-Xavier Vialard, Shun-ichi Amari, Alain Trouvé, and
339 Gabriel Peyré. Interpolating between optimal transport and mmd using sinkhorn divergences.
340 In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 2681–2690.
341 PMLR, 2019.
- 342 [12] Joel Franklin and Jens Lorenz. On the scaling of multidimensional matrices. *Linear Algebra*
343 *and its applications*, 114:717–735, 1989.
- 344 [13] Aude Genevay, Lénaïc Chizat, Francis Bach, Marco Cuturi, and Gabriel Peyré. Sample
345 complexity of sinkhorn divergences. In *The 22nd International Conference on Artificial*
346 *Intelligence and Statistics*, pages 1574–1583. PMLR, 2019.
- 347 [14] Aude Genevay, Gabriel Peyré, and Marco Cuturi. Learning generative models with sinkhorn
348 divergences. In *International Conference on Artificial Intelligence and Statistics*, pages 1608–
349 1617. PMLR, 2018.
- 350 [15] Christian Léonard. A survey of the schrödinger problem and some of its connections with
351 optimal transport. *arXiv preprint arXiv:1308.0215*, 2013.
- 352 [16] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G
353 Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al.
354 Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- 355 [17] Thanh Tang Nguyen, Sunil Gupta, and Svetha Venkatesh. Distributional reinforcement learning
356 with maximum mean discrepancy. *Association for the Advancement of Artificial Intelligence*
357 *(AAAI)*, 2020.
- 358 [18] Aaditya Ramdas, Nicolás García Trillos, and Marco Cuturi. On wasserstein two-sample testing
359 and related families of nonparametric tests. *Entropy*, 19(2):47, 2017.

- 360 [19] Mark Rowland, Marc Bellemare, Will Dabney, Rémi Munos, and Yee Whye Teh. An analysis
361 of categorical distributional reinforcement learning. In *International Conference on Artificial*
362 *Intelligence and Statistics*, pages 29–37. PMLR, 2018.
- 363 [20] Ludger Rüschendorf and Wolfgang Thomsen. Closedness of sum spaces and the generalized
364 schrödinger problem. *Theory of Probability & Its Applications*, 42(3):483–494, 1998.
- 365 [21] Richard Sinkhorn. Diagonal equivalence to matrices with prescribed row and column sums.
366 *The American Mathematical Monthly*, 74(4):402–405, 1967.
- 367 [22] Gábor J Székely. E-statistics: The energy of statistical samples. *Bowling Green State University,*
368 *Department of Mathematics and Statistics Technical Report*, 3(05):1–18, 2003.
- 369 [23] Eric Wong, Frank Schmidt, and Zico Kolter. Wasserstein adversarial examples via projected
370 sinkhorn iterations. In *International Conference on Machine Learning*, pages 6808–6817.
371 PMLR, 2019.
- 372 [24] Derek Yang, Li Zhao, Zichuan Lin, Tao Qin, Jiang Bian, and Tie-Yan Liu. Fully parameterized
373 quantile function for distributional reinforcement learning. *Advances in neural information*
374 *processing systems*, 32:6193–6202, 2019.
- 375 [25] Shangtong Zhang. Modularized implementation of deep rl algorithms in pytorch. <https://github.com/ShangtongZhang/DeepRL>, 2018.
376
- 377 [26] Fan Zhou, Jianing Wang, and Xingdong Feng. Non-crossing quantile regression for distri-
378 butional reinforcement learning. *Advances in Neural Information Processing Systems*, 33,
379 2020.
- 380 [27] Florian Ziel. The energy distance for ensemble and scenario reduction. *arXiv preprint*
381 *arXiv:2005.14670*, 2020.

382 Checklist

- 383 1. For all authors...
- 384 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
385 contributions and scope? [Yes]
- 386 (b) Did you describe the limitations of your work? [Yes] We provide the discussion about
387 the limitation of our proposal in Section 6.
- 388 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 389 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
390 them? [Yes]
- 391 2. If you are including theoretical results...
- 392 (a) Did you state the full set of assumptions of all theoretical results? [Yes] Please refer to
393 Appendix B and C.
- 394 (b) Did you include complete proofs of all theoretical results? [Yes] Please refer to
395 Appendix B and C.
- 396 3. If you ran experiments...
- 397 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
398 mental results (either in the supplemental material or as a URL)? [Yes]
- 399 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
400 were chosen)? [Yes] Our implementation is adapted from Pytorch distributional RL
401 modules [25].
- 402 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
403 ments multiple times)? [Yes]
- 404 (d) Did you include the total amount of compute and the type of resources used (e.g., type
405 of GPUs, internal cluster, or cloud provider)? [Yes] We provide the comparison of
406 computational cost in Figure 12 of Appendix F.

- 407 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 408 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 409 (b) Did you mention the license of the assets? [N/A]
- 410 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 411
- 412 (d) Did you discuss whether and how consent was obtained from people whose data you're
- 413 using/curating? [N/A]
- 414 (e) Did you discuss whether the data you are using/curating contains personally identifiable
- 415 information or offensive content? [N/A]
- 416 5. If you used crowdsourcing or conducted research with human subjects...
- 417 (a) Did you include the full text of instructions given to participants and screenshots, if
- 418 applicable? [N/A]
- 419 (b) Did you describe any potential participant risks, with links to Institutional Review
- 420 Board (IRB) approvals, if applicable? [N/A]
- 421 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 422 spent on participant compensation? [N/A]