

# Multivariate Distributional Reinforcement Learning Using Sliced Divergences

KU LEUVEN



Baptiste Debes, Tinne Tuytelaars

KU Leuven



## Multivariate distributional RL

- Reward function can be multivariate  $R \in \mathbb{R}^d$
- Multivariate Temporal Difference
  - Sobolev TD
  - Successor features

$$S' \sim P(\cdot | s, a), \quad A' \sim \pi(\cdot | S'), \quad Z(S', A') \sim \mu(S', A')$$

$$(\mathcal{T}^\pi \mu)(s, a) = \text{Law}[R(s, a) + \Gamma(s, a)Z(S', A')]$$

$$\Gamma(s, a) = \gamma I_d = \begin{bmatrix} \gamma & 0 & \dots & 0 \\ 0 & \gamma & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \gamma \end{bmatrix}$$

$$\Gamma(s, a) = \begin{bmatrix} \gamma_{11}(s, a) & \gamma_{12}(s, a) & \dots & \gamma_{1d}(s, a) \\ \gamma_{21}(s, a) & \gamma_{22}(s, a) & \dots & \gamma_{2d}(s, a) \\ \vdots & \vdots & \ddots & \vdots \\ \gamma_{d1}(s, a) & \gamma_{d2}(s, a) & \dots & \gamma_{dd}(s, a) \end{bmatrix}$$

Regular multivariate RL

General case

## Contraction property

- Divergence between return laws

$$D : \mathcal{P}(\mathbb{R}^d) \times \mathcal{P}(\mathbb{R}^d) \rightarrow \mathbb{R}_+$$

- Contraction criterion

$$\bar{D}(\eta_1, \eta_2) := \sup_{(s, a)} D(\eta_1(s, a), \eta_2(s, a))$$

$$\bar{D}(\mathcal{T}^\pi \eta_1, \mathcal{T}^\pi \eta_2) \leq \kappa \bar{D}(\eta_1, \eta_2), \quad \kappa < 1$$

- Baseline is Wasserstein

$$\bar{W}_p(\mathcal{T}^\pi \eta_1, \mathcal{T}^\pi \eta_2) \leq \bar{L} \bar{W}_p(\eta_1, \eta_2), \quad \bar{L} := \sup_{(s, a)} \|\Gamma(s, a)\|_{\text{op}} < 1$$

Exact empirical estimator is **intractable**  $\mathcal{O}(n^3 \log n)$

## Slicing

- Setup

- Base 1D divergence  $\Delta : \mathcal{P}(\mathbb{R}) \times \mathcal{P}(\mathbb{R}) \rightarrow \mathbb{R}_+$
- Two distributions  $\mu, \nu \in \mathcal{P}(\mathbb{R}^d)$
- Random direction  $\theta \in \mathbb{S}^{d-1}$

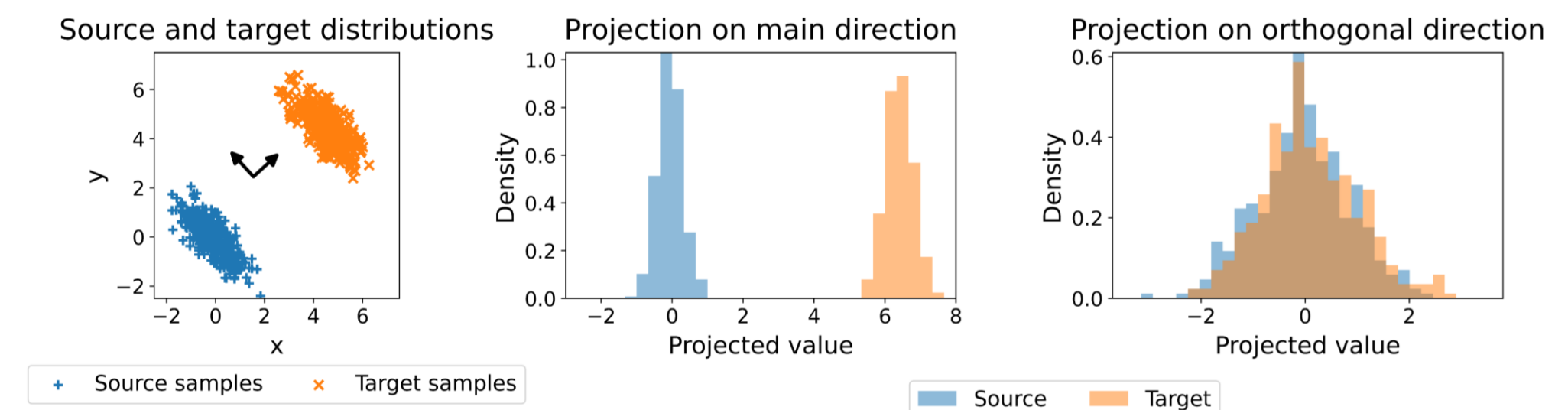
- Core idea : compare projections

$$\Delta((P_\theta)_\# \mu, (P_\theta)_\# \nu) \quad P_\theta(x) = \langle \theta, x \rangle$$

- Aggregate

$$\theta_1, \dots, \theta_K \stackrel{\text{i.i.d.}}{\sim} \sigma, \quad \sigma = \text{Unif}(\mathbb{S}^{d-1})$$

$$\hat{\mathbf{S}}\Delta_p^p(\mu, \nu) = \frac{1}{K} \sum_{k=1}^K \Delta^p((P_{\theta_k})_\# \mu, (P_{\theta_k})_\# \nu)$$



## Base 1D divergences

- Cramér distance

- Distance between CDFs

$$C_2(\mu, \nu) = \int_{\mathbb{R}} |F_\mu(u) - F_\nu(u)|^2 du$$

- Maximum Mean Discrepancy

- Kernel-based  $k : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$

$$\text{MMD}_k^2(\mu, \nu) = \mathbb{E}_{x, x' \sim \mu} [k(x, x')] + \mathbb{E}_{y, y' \sim \nu} [k(y, y')] - 2\mathbb{E}_{x \sim \mu, y \sim \nu} [k(x, y)]$$

Divergence	Time complexity
Wasserstein $W_p$	$\mathcal{O}(n \log n)$
Cramér $C_2$	$\mathcal{O}(n \log n)$
MMD $_k$	$\mathcal{O}(n^2)$

## Contraction

- Result #1: standard case  $\Gamma = \gamma I$

- $\mathbf{S}\Delta_p$  inherits contraction from  $\Delta$

- Result #2: dense case  $\Gamma(s, a), \bar{L} < 1$

- $\mathbf{S}\Delta_p$  fails norm-based contraction

- Result #3: max-slicing

- Search for most discriminative direction

$$\text{MS}\Delta(\mu, \nu) = \sup_{\theta \in \mathbb{S}^{d-1}} \Delta((P_\theta)_\# \mu, (P_\theta)_\# \nu)$$

- MS $\Delta$  recovers norm-based contraction

## Stochastic training

- True Bellman target

$$\mu = \mathbb{E}_{S', A'} [\text{Law}(R + \Gamma Z^-(S', A'))]$$

- Sampled Bellman target

$$\hat{\mu} = \text{Law}(R + \Gamma Z^-(s', a'))$$

- Unbiased gradient

$$(\mathbf{U}) \quad \mathbb{E}_{S', A'} [\nabla_\phi D(\hat{\mu}, \nu_\phi)] = \nabla_\phi D(\mu, \nu_\phi)$$

Method	Contraction $\bar{L} < 1$	(U)
$\mathbf{S}\Delta_p$	×	if base $\Delta$ has (U)
MS $\Delta$	✓	×
$\mathbf{W}_p$	✓	×

Open question: can we get both ?

## Experiments

