

# Sinkhorn Distributionally Robust Optimization

Jie Wang<sup>1</sup>, Rui Gao<sup>2</sup>, and Yao Xie<sup>1</sup>

<sup>1</sup>Georgia Institute of Technology

<sup>2</sup>The University of Texas at Austin

2022 INFORMS Annual Meeting

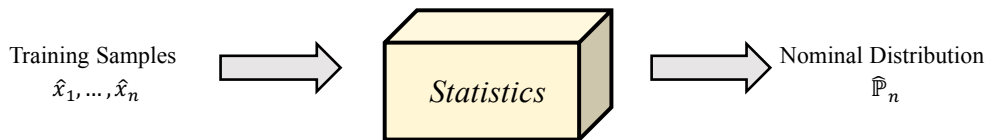
# Decision-Making Under Uncertainty

$$\begin{aligned}\text{Risk :} & \quad \mathcal{R}(\theta; \mathbb{P}) = \mathbb{E}_{\mathbb{P}}[f_{\theta}(z)] \\ \text{Optimal Risk :} & \quad \mathcal{R}(\Theta; \mathbb{P}) = \inf_{\theta \in \Theta} \mathbb{E}_{\mathbb{P}}[f_{\theta}(z)]\end{aligned}$$

► Available Information:

$$\begin{aligned}\text{Structual :} & \quad \mathbb{P} \text{ is supported on } \Omega \subseteq \mathbb{R}^d \\ \text{Statistical :} & \quad \hat{x}_1, \dots, \hat{x}_n \sim \mathbb{P}\end{aligned}$$

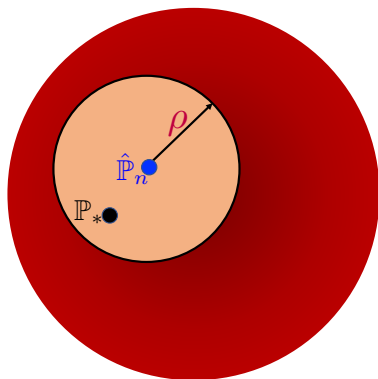
► Nominal Problem:



- Non-parametric estimators:  $\hat{\mathbb{P}}_n = \frac{1}{n} \sum_{i=1}^n \delta_{\hat{x}_i}$ .
- Kernel density estimators:  $\hat{\mathbb{P}}_n = \frac{1}{n} \sum_{i=1}^n K(\hat{x}_i)$ .

## Wasserstein DRO

**Definition:**  $\mathcal{P} = \{\mathbb{P} : W(\mathbb{P}, \hat{\mathbb{P}}_n) \leq \rho\}$ .



Contain each  $\mathbb{P}$  such  
that  $W(\mathbb{P}, \hat{\mathbb{P}}_n) \leq \rho$

Worst-case risk :

$$\sup_{\mathbb{P} \in \mathcal{P}} \mathbb{E}_{\mathbb{P}}[f_{\theta}(z)]$$

Robust Optimal Risk :

$$\inf_{\theta \in \Theta} \sup_{\mathbb{P} \in \mathcal{P}} \mathbb{E}_{\mathbb{P}}[f_{\theta}(z)]$$

# Limitations of Wasserstein DRO

- ▶ Worst-case distribution is **discrete**:

*For WDRO with  $n$ -point nominal distribution, the worst-case distribution is supported on  $n + 1$  points<sup>1</sup>.*

- ▶ Tractability for **limited** scenarios:

*Finite-dimensional convex reformulation is available if the objective is a pointwise maximum of finitely many concave functions<sup>2</sup>.*

- ▶ Some cases the **same performance** as SAA<sup>2</sup>.

---

<sup>1</sup>Rui Gao and Anton J. Kleywegt. “Distributionally Robust Stochastic Optimization with Wasserstein Distance”. In: *arXiv preprint arXiv:1604.02199* (Apr. 2016).

<sup>2</sup>Peyman Mohajerin Esfahani and Daniel Kuhn. “Data-driven distributionally robust optimization using the Wasserstein metric: performance guarantees and tractable reformulations”. In: *Mathematical Programming* 171.1 (July 2017), pp. 115–166.

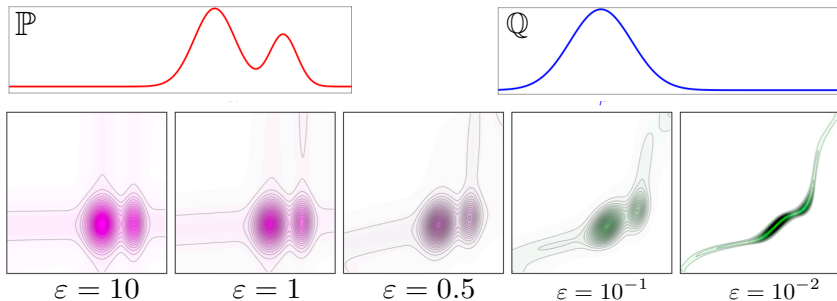
# Sinkhorn Distance

- Sinkhorn Distance [Cuturi 2013]:

$$W_{\varepsilon}(\mathbb{P}, \mathbb{Q}) = \inf_{\gamma \in \Gamma(\mathbb{P}, \mathbb{Q})} \left\{ \mathbb{E}_{(X,Y) \sim \gamma} [c(X,Y)] + \varepsilon H(\gamma | \mathbb{P} \otimes \nu) \right\}.$$

- Relative Entropy between  $\gamma$  and  $\mathbb{P} \otimes \nu$ :

$$H(\gamma | \mathbb{P} \otimes \nu) = \int \log \left( \frac{d\gamma(x,y)}{d\mathbb{P}(x) d\nu(y)} \right) d\gamma(x,y).$$



## Highlights of Sinkhorn Distance

Probability distance between distributions in  $\mathbb{R}^d$  using  $n$  samples:

	MMD	Wasserstein	Sinkhorn
<b>Computation</b>	$O(n)$	$\tilde{O}(n^3)$	$\tilde{O}(n^2)$ [Altschuler, Niles-Weed, and Rigollet 2017]
<b>Sample Complexity</b>	$O(n^{-1/2})$	$O(n^{-1/d})$	$O(e^{\kappa/\epsilon} n^{-1/2} \epsilon^{-\lfloor d/2 \rfloor})$ [Genevay et al. 2019]

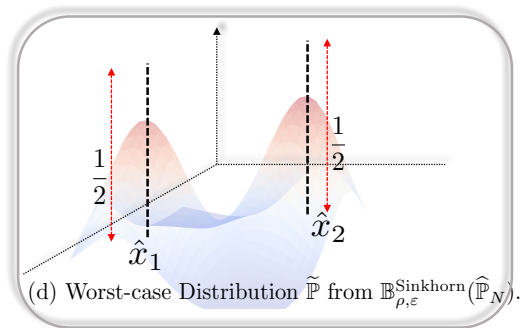
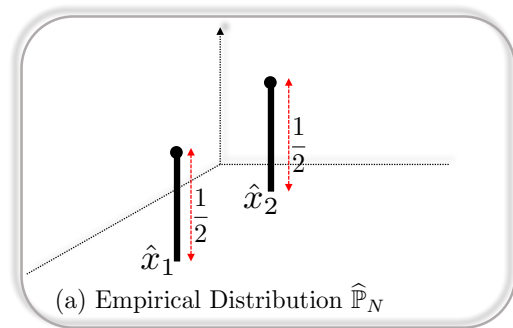
- ▶ Fast algorithms for implementation;
- ▶ Sharp sample complexity rate;
- ▶ Encourage stochastic optimal transport (helpful in some applications, e.g., domain adaptation [Courty, Flamary, and Tuia 2014]).

# Main Framework

► Sinkhorn DRO:

$$\inf_{\theta} \sup_{\mathbb{P} \in \mathbb{B}_{\rho, \varepsilon}(\hat{\mathbb{P}})} \mathbb{E}_{z \sim \mathbb{P}}[f_{\theta}(z)],$$

$$\text{where } \mathbb{B}_{\rho, \varepsilon}(\hat{\mathbb{P}}) = \{\mathbb{P} : W_{\varepsilon}(\hat{\mathbb{P}}, \mathbb{P}) \leq \rho\}.$$

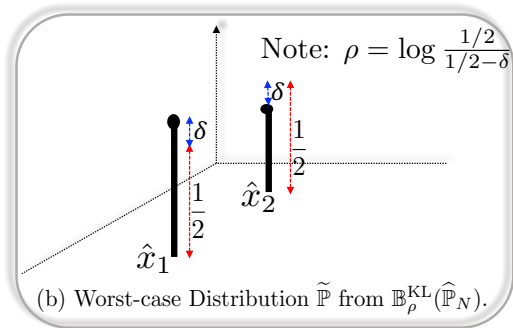
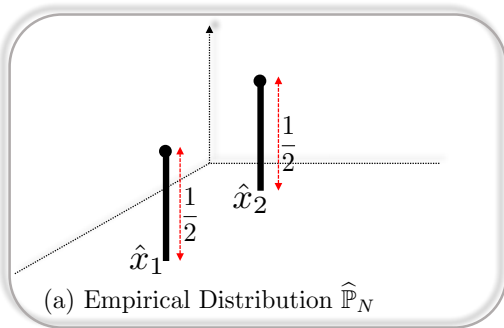


# General DRO Models

► KL-DRO:

$$\inf_{\theta} \sup_{\mathbb{P} \in \mathbb{B}_{\rho}^{\text{KL}}(\hat{\mathbb{P}})} \mathbb{E}_{z \sim \mathbb{P}}[f_{\theta}(z)],$$

$$\text{where } \mathbb{B}_{\rho}^{\text{KL}}(\hat{\mathbb{P}}) = \{\mathbb{P} : D_{\text{KL}}(\hat{\mathbb{P}}, \mathbb{P}) \leq \rho\}.$$



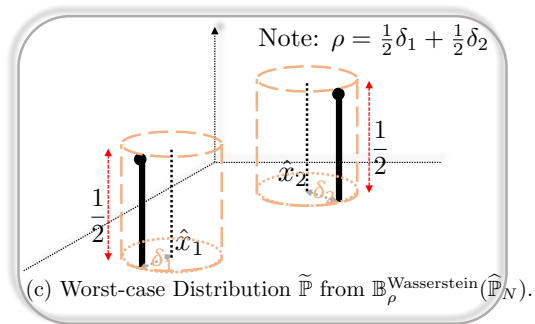
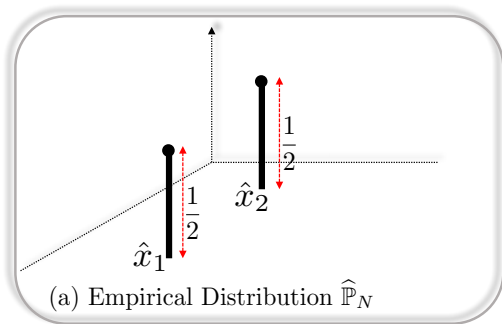


# General DRO Models

## ► Wasserstein-DRO:

$$\inf_{\theta} \sup_{\mathbb{P} \in \mathbb{B}_{\rho}^{\text{Wasserstein}}(\hat{\mathbb{P}})} \mathbb{E}_{z \sim \mathbb{P}}[f_{\theta}(z)],$$

$$\text{where } \mathbb{B}_{\rho}^{\text{Wasserstein}}(\hat{\mathbb{P}}) = \{\mathbb{P} : W(\hat{\mathbb{P}}, \mathbb{P}) \leq \rho\}.$$



## Ongoing Outline

- ▶ Sinkhorn DRO:

$$\inf_{\theta} \sup_{\mathbb{P} \in \mathbb{B}_{\rho, \varepsilon}(\hat{\mathbb{P}})} \mathbb{E}_{z \sim \mathbb{P}}[f_{\theta}(z)],$$

where  $\mathbb{B}_{\rho, \varepsilon}(\hat{\mathbb{P}}) = \{\mathbb{P} : W_{\varepsilon}(\hat{\mathbb{P}}, \mathbb{P}) \leq \rho\}$ .

- ▶ Duality Formulation for Sinkhorn DRO
- ▶ First-order Optimization Algorithm
- ▶ Properties and Numerical Results

# Tractable Formulation

Assume that

- (I)  $\nu\{z : 0 \leq c(x, z) < \infty\} = 1$  for  $\widehat{\mathbb{P}}$ -almost every  $x$ ;
- (II) The integral  $\int e^{-c(x, z)/\varepsilon} d\nu(z) < \infty$  for  $\widehat{\mathbb{P}}$ -almost every  $x$ ;
- (III)  $\Omega$  is a measurable space, and the function  $f : \Omega \rightarrow \mathbb{R} \cup \{\infty\}$  is measurable.

Consider the primal

$$V_P = \sup_{\mathbb{P} \in \mathbb{B}_{\rho, \varepsilon}(\widehat{\mathbb{P}})} \mathbb{E}_{z \sim \mathbb{P}}[f(z)], \quad \text{where } \mathbb{B}_{\rho, \varepsilon}(\widehat{\mathbb{P}}) = \{\mathbb{P} : W_\varepsilon(\widehat{\mathbb{P}}, \mathbb{P}) \leq \rho\}. \quad (\text{Sinkhorn DR0})$$

It admits the **strong dual reformulation**:

$$V_D = \inf_{\lambda > 0} \lambda \bar{\rho} + \lambda \varepsilon \int_{\Omega} \log \left( \mathbb{E}_{Q_x} \left[ e^{f(z)/(\lambda \varepsilon)} \right] \right) d\widehat{\mathbb{P}}(x),$$

where

$$\begin{aligned} \bar{\rho} &= \rho + \varepsilon \int_{\Omega} \log \left( \int_{\Omega} e^{-c(x, z)/\varepsilon} d\nu(z) \right) d\widehat{\mathbb{P}}(x), \\ dQ_x(z) &= \frac{e^{-c(x, z)/\varepsilon}}{\int_{\Omega} e^{-c(x, u)/\varepsilon} d\nu(u)} d\nu(z). \end{aligned}$$

## Interpretation of Worst-case Distribution

$$\tilde{\mathbb{P}} = \arg \max_{\mathbb{P}} \left\{ \mathbb{E}_{z \sim \mathbb{P}}[f(z)] : W_{\varepsilon}(\hat{\mathbb{P}}, \mathbb{P}) \leq \rho \right\}$$

- ▶ For each  $x \in \text{supp}(\hat{\mathbb{P}})$ , optimal transport maps it to a (conditional) distribution  $\gamma_x$  such that

$$\frac{d\gamma_x(z)}{d\nu(z)} = \alpha_x \cdot \exp \left( (f(z) - \lambda^* c(x, z)) / (\lambda^* \varepsilon) \right),$$

where  $\alpha_x$  is the normalizing constant.

- ▶ Closed-form expression on  $\tilde{\mathbb{P}}$ :

$$\frac{d\tilde{\mathbb{P}}(z)}{d\nu(z)} = \int \alpha_x \cdot \exp \left( (f(z) - \lambda^* c(x, z)) / (\lambda^* \varepsilon) \right) d\hat{\mathbb{P}}(x).$$

**Worst-case distribution  $\tilde{\mathbb{P}}$  support on whole space, while W-DRO is discrete.**

## Toy Example: Newsvendor

Newsvendor problem: ( $\beta$ : Demand); ( $u \min\{\beta, z\}$ : Earning); ( $k\beta$ : Loss).

$$\min_{\beta} \mathbb{E}_{\mathbb{P}_*} [k\beta - u \min\{\beta, z\}], \quad k = 5, u = 7.$$

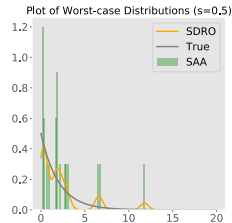
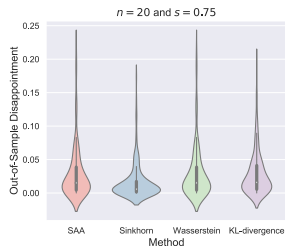
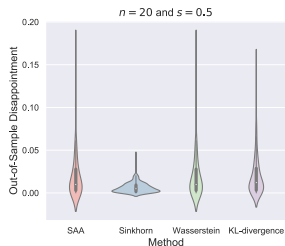
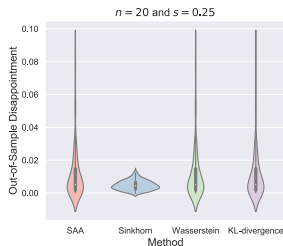


# Performance and Visualization

News vendor problem:

$$\min_{\beta} \mathbb{E}_{\mathbb{P}_*} [k\beta - u \min\{\beta, \zeta\}], \quad k = 5, u = 7.$$

$\mathbb{P}_* \sim \exp(1/s)$  with  $s \in \{0.25, 0.5, 0.75\}$ . Access to  $n = 20$  samples.

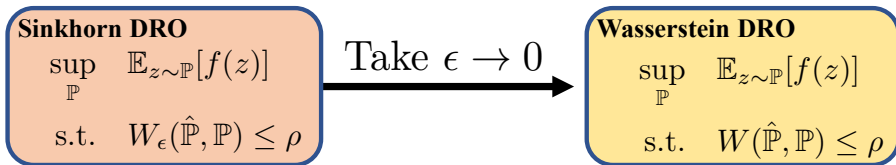


## Connection of Sinkhorn DRO with Wasserstein DRO

When  $\varepsilon \rightarrow 0$ , the dual objective of Sinkhorn DRO converges into

$$\lambda \rho + \int \text{ess-sup}_v (f(\cdot) - \lambda c(x, \cdot)) d\hat{\mathbb{P}}(x).$$

When  $\text{supp}(v) = \Omega$ ,



# Optimization Algorithm for Sinkhorn DRO

- Based on strong duality,

$$\begin{aligned} & \min_{\theta \in \Theta} \sup_{\mathbb{P}} \left\{ \mathbb{E}_{z \sim \mathbb{P}} [f_{\theta}(z)] : W_{\varepsilon}(\widehat{\mathbb{P}}, \mathbb{P}) \leq \rho \right\} \\ &= \min_{\lambda \geq 0} \left\{ \lambda \bar{\rho} + \underbrace{\min_{\theta \in \Theta} \mathbb{E}_{x \sim \widehat{\mathbb{P}}} \left[ \lambda \varepsilon \log \left( \mathbb{E}_{z \sim Q_x} \left[ e^{f_{\theta}(z)/(\lambda \varepsilon)} \right] \right) \right]}_{V(\lambda)} \right\} \end{aligned}$$

- Solve the Monte-Carlo approximated formulation<sup>3</sup>:

$$V(\lambda) \approx \min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \lambda \varepsilon \log \left( \frac{1}{m} \sum_{j=1}^m e^{f_{\theta}(z_{i,j})/(\lambda \varepsilon)} \right),$$

where  $\{\hat{x}_i\}_{i=1}^n \sim \widehat{\mathbb{P}}$  and  $\{z_{i,j}\}_{j=1}^m$  are i.i.d. samples generated from  $Q_{\hat{x}_i}$ .

- **Cons:** It requires  $\tilde{O}(\delta^{-3})$  samples to obtain  $\delta$ -optimal solution.

---

<sup>3</sup>Alexander Shapiro, Darinka Dentcheva, and Andrzej Ruszczyński. *Lectures on stochastic programming: modeling and theory*. SIAM, 2014.



# Optimization Algorithm for Sinkhorn DRO

- Based on strong duality,

$$\begin{aligned} & \min_{\theta \in \Theta} \sup_{\mathbb{P}} \left\{ \mathbb{E}_{z \sim \mathbb{P}} [f_{\theta}(z)] : W_{\varepsilon}(\widehat{\mathbb{P}}, \mathbb{P}) \leq \rho \right\} \\ &= \min_{\lambda \geq 0} \left\{ \lambda \bar{\rho} + \underbrace{\min_{\theta \in \Theta} \mathbb{E}_{x \sim \widehat{\mathbb{P}}} \left[ \lambda \varepsilon \log \left( \mathbb{E}_{z \sim \mathbb{Q}_x} \left[ e^{f_{\theta}(z)/(\lambda \varepsilon)} \right] \right) \right]}_{V(\lambda)} \right\} \end{aligned}$$

- Solve the Monte-Carlo approximated formulation<sup>3</sup>:

$$V(\lambda) \approx \min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \lambda \varepsilon \log \left( \frac{1}{m} \sum_{j=1}^m e^{f_{\theta}(z_{i,j})/(\lambda \varepsilon)} \right),$$

where  $\{\hat{x}_i\}_{i=1}^n \sim \widehat{\mathbb{P}}$  and  $\{z_{i,j}\}_{j=1}^m$  are i.i.d. samples generated from  $\mathbb{Q}_{\hat{x}_i}$ .

- **Cons:** It requires  $\tilde{O}(\delta^{-3})$  samples to obtain  $\delta$ -optimal solution.

---

<sup>3</sup>Alexander Shapiro, Darinka Dentcheva, and Andrzej Ruszczyński. *Lectures on stochastic programming: modeling and theory*. SIAM, 2014.

# Optimization Algorithm for Sinkhorn DRO: Biased Gradient Update

- Goal: to solve the optimization

$$\min_{\theta \in \Theta} \left\{ F(\theta) := \mathbb{E}_{x \sim \hat{\mathbb{P}}} \left[ \lambda \varepsilon \log \left( \mathbb{E}_{z \sim \mathbb{Q}_x} \left[ e^{f_{\theta}(z)/(\lambda \varepsilon)} \right] \right) \right] \right\}.$$

- Biased gradient update: for each iteration  $t$ ,
  - Construct a gradient estimate<sup>4</sup> of  $F(\theta_t)$ , denoted as  $v(\theta_t)$ ;
  - Update  $\theta_{t+1} = \mathbf{Proximal}_{\theta_t}(\gamma_t v(\theta_t))$ .

**Estimator of solution:** randomly selected from (or average over)  $\{\theta_t\}_{t=1}^T$

---

<sup>4</sup>Yifan Hu, Xin Chen, and Niao He. “On the Bias-Variance-Cost Tradeoff of Stochastic Optimization”. In: *Advances in Neural Information Processing Systems*. Dec. 2021.

# Optimization Algorithm for Sinkhorn DRO: Biased Gradient Update

- ▶ Goal: to solve the optimization

$$\min_{\theta \in \Theta} \left\{ F(\theta) := \mathbb{E}_{x \sim \hat{\mathbb{P}}} \left[ \lambda \varepsilon \log \left( \mathbb{E}_{z \sim \mathbb{Q}_x} \left[ e^{f_{\theta}(z)/(\lambda \varepsilon)} \right] \right) \right] \right\}.$$

- ▶ Biased gradient update: for each iteration  $t$ ,
  - ▶ Construct a gradient estimate<sup>4</sup> of  $F(\theta_t)$ , denoted as  $v(\theta_t)$ ;
  - ▶ Update  $\theta_{t+1} = \mathbf{Proximal}_{\theta_t}(\gamma_t v(\theta_t))$ .

**Estimator of solution:** randomly selected from (or average over)  $\{\theta_t\}_{t=1}^T$

---

<sup>4</sup>Yifan Hu, Xin Chen, and Niao He. “On the Bias-Variance-Cost Tradeoff of Stochastic Optimization”. In: *Advances in Neural Information Processing Systems*. Dec. 2021.

# Optimization Algorithm for Sinkhorn DRO: Biased Gradient Update

- Goal: to solve the optimization

$$\min_{\theta \in \Theta} \left\{ F(\theta) := \mathbb{E}_{x \sim \hat{\mathbb{P}}} \left[ \lambda \varepsilon \log \left( \mathbb{E}_{z \sim \mathbb{Q}_x} \left[ e^{f_{\theta}(z)/(\lambda \varepsilon)} \right] \right) \right] \right\}.$$

- Biased gradient update: for each iteration  $t$ ,
  - Construct a gradient estimate<sup>4</sup> of  $F(\theta_t)$ , denoted as  $v(\theta_t)$ ;
  - Update  $\theta_{t+1} = \mathbf{Proximal}_{\theta_t}(\gamma_t v(\theta_t))$ .

**Estimator of solution:** randomly selected from (or average over)  $\{\theta_t\}_{t=1}^T$

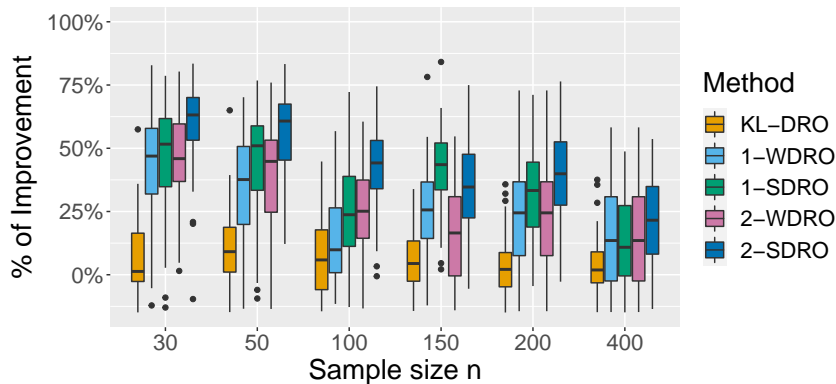
Estimators	Convex Nonsmooth	Convex Smooth	Nonconvex Smooth
Vanilla SGD	$O(\delta^{-3})$	$O(\delta^{-3})$	$O(\delta^{-6})$
V-MLMC	N/A	$\tilde{O}(\delta^{-2})$	$\tilde{O}(\delta^{-4})$
RT-MLMC	N/A	$\tilde{O}(\delta^{-2})$	$\tilde{O}(\delta^{-4})$

<sup>4</sup>Yifan Hu, Xin Chen, and Niao He. “On the Bias-Variance-Cost Tradeoff of Stochastic Optimization”. In: *Advances in Neural Information Processing Systems*. Dec. 2021.

# Numerical Results

Portfolio Optimization:

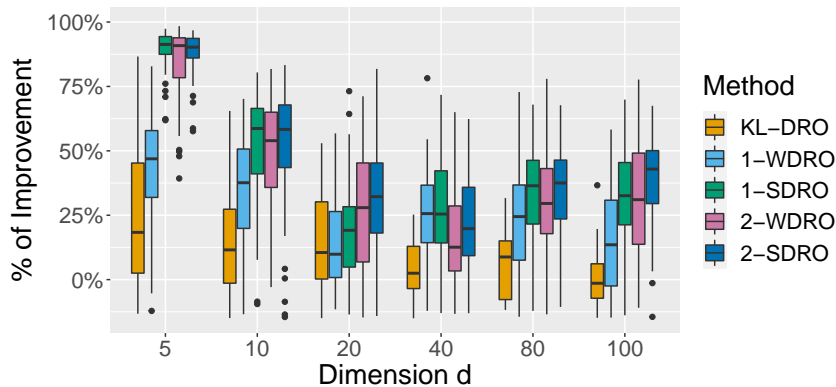
$$\begin{aligned} \inf_x \quad & \mathbb{E}_{\mathbb{P}_*} [-\langle x, \zeta \rangle] + \rho \cdot \mathbb{P}_* \text{-CVaR}_\alpha(-\langle x, \zeta \rangle) \\ \text{s.t.} \quad & x \in \mathcal{X} = \{x \in \mathbb{R}_+^D : x^T \mathbf{1} = 1\}. \end{aligned}$$



# Numerical Results

Portfolio Optimization:

$$\begin{aligned} \inf_x \quad & \mathbb{E}_{\mathbb{P}_*} [-\langle x, \zeta \rangle] + \rho \cdot \mathbb{P}_* \text{-CVaR}_\alpha(-\langle x, \zeta \rangle) \\ \text{s.t.} \quad & x \in \mathcal{X} = \{x \in \mathbb{R}_+^D : x^T \mathbf{1} = 1\}. \end{aligned}$$



## Numerical Simulation Results

### Semi-supervised Learning:

- ▶ Train classifiers based on data with labels and without labels;
- ▶ Two performance measures:
  - ▶ Training error for data without labels;
  - ▶ Testing error.

	SAA	Sinkhorn	Wasserstein	KL-divergence
Breast Cancer	$.20 \pm .068$	<b><math>.12 \pm .068</math></b>	$.17 \pm .073$	$.19 \pm .038$
	$.19 \pm .073$	<b><math>.11 \pm .067</math></b>	$.17 \pm .075$	$.19 \pm .073$
Magic	$.28 \pm .082$	<b><math>.25 \pm .091</math></b>	$.27 \pm .077$	$.26 \pm .078$
	$.28 \pm .064$	<b><math>.25 \pm .074</math></b>	$.27 \pm .058$	$.27 \pm .066$
QSAR Bio	$.25 \pm .057$	<b><math>.22 \pm .063</math></b>	$.23 \pm .073$	$.25 \pm .037$
	$.25 \pm .062$	<b><math>.22 \pm .065</math></b>	$.23 \pm .079$	$.25 \pm .042$
Spambase	$.19 \pm .038$	<b><math>.14 \pm .046</math></b>	$.16 \pm .036$	$.18 \pm .034$
	$.19 \pm .032$	<b><math>.14 \pm .036</math></b>	$.16 \pm .028$	$.18 \pm .042$

# Take Home Message

Sinkhorn DRO is a great notion of DRO models:

- ▶ Inherit **geometric properties** from optimal transport;
- ▶ **Absolutely continuous** worst-case distribution thanks to **entropic regularization**;
- ▶ **Improve the out-of-sample performance** of Wasserstein DRO;
- ▶ Optimization by **Monte Carlo approximation** and **first order method**;
- ▶ **More applications in operations research** with Sinkhorn DRO can be explored!





# Sinkhorn Distributionally Robust Optimization

To be Submitted to Operations Research – INFORMS PUBs

Online Available: [arxiv.org/abs/2109.11926](https://arxiv.org/abs/2109.11926)



**SCAN ME**