

# Sinkhorn Distributionally Robust Optimization

Jie Wang<sup>†</sup>, Rui Gao<sup>‡</sup>, Yao Xie<sup>†</sup>

<sup>†</sup> H. Milton Stewart School of Industrial and Systems Engineering, Georgia Institute of Technology

<sup>‡</sup> Department of Information, Risk, and Operations Management, McCombs School of Business, University of Texas at Austin

## Contributions

- Distributionally robust optimization with entropic regularized Wasserstein distance (Sinkhorn distance).
- Ambiguity set contains only absolutely continuous distributions.
- Computationally efficient first-order optimization algorithm.

## Decision-Making Under Uncertainty

- Objective: Find decision  $\theta$  to minimize the risk

$$\mathcal{R}(\theta; \mathbb{P}) = \mathbb{E}_{\mathbb{P}}[f_{\theta}(z)].$$

- Available Information:

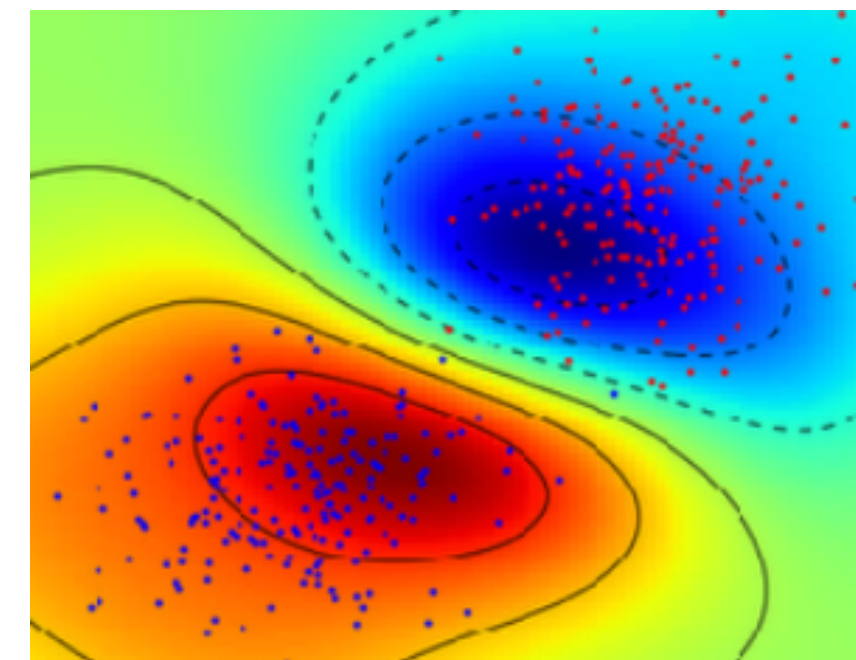
Structural :  $\mathbb{P}$  is supported on  $\Omega \subseteq \mathbb{R}^d$   
Statistical :  $\hat{x}_1, \dots, \hat{x}_n \sim \mathbb{P}$



Supply Chain Mgmt.



Portfolio Mgmt.



Machine Learning

- Sample Average Approximation (SAA):

$$\inf_{\theta \in \Theta} \left\{ \mathcal{R}(\theta; \hat{\mathbb{P}}_n) \triangleq \mathbb{E}_{\hat{\mathbb{P}}_n}[f_{\theta}(z)] \right\}, \quad \text{where } \hat{\mathbb{P}}_n = \frac{1}{n} \sum_{i=1}^n \delta_{\hat{x}_i}.$$

- Wasserstein Distributionally Robust Optimization (DRO):

$$\inf_{\theta \in \Theta} \left\{ \sup_{\mathbb{P} \in \mathcal{P}} \mathbb{E}_{\mathbb{P}}[f_{\theta}(z)] \right\}, \quad \text{where } \mathcal{P} = \{\mathbb{P} : W(\hat{\mathbb{P}}, \mathbb{P}) \leq \rho\}.$$

- Facts about Wasserstein DRO:

- For WDRO with  $n$ -point nominal distribution, the worst-case distribution is supported on  $n+1$  points.
- Finite-dimensional convex reformulation is available if the objective is a pointwise maximum of finitely many concave functions.
- Some cases the same performance as SAA.

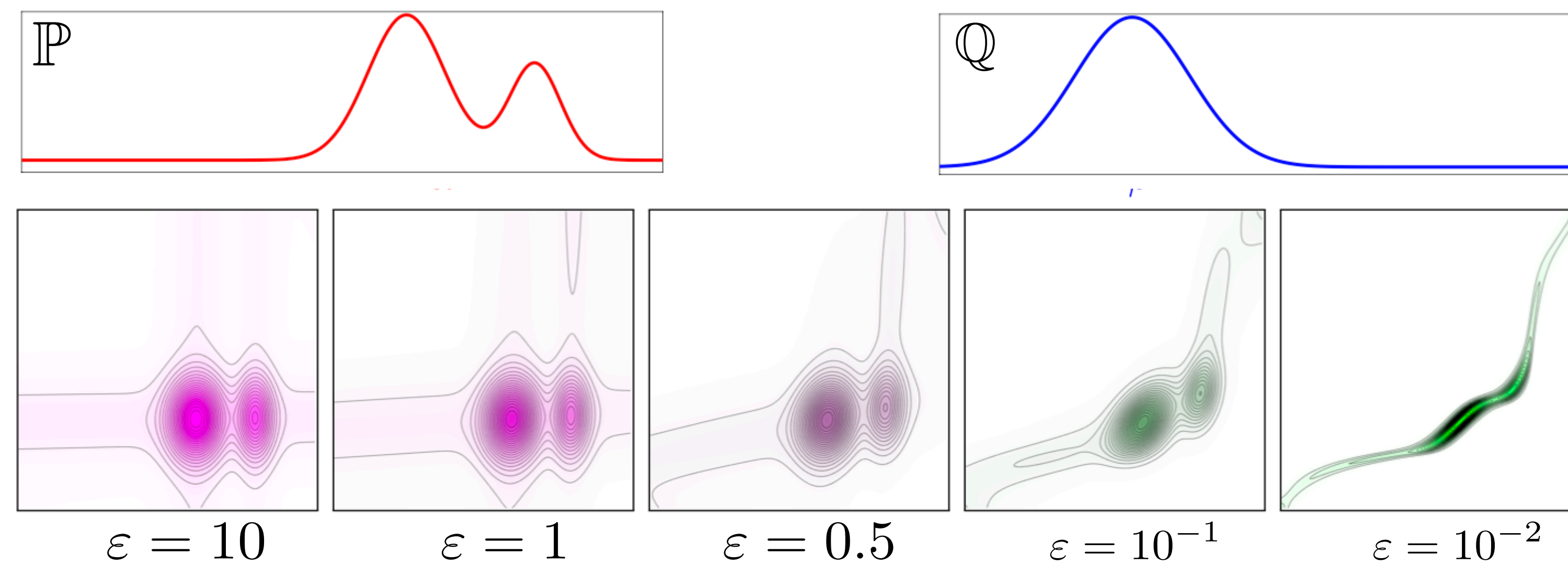
## Sinkhorn Distance and Robust Formulation

- Sinkhorn Distance [Cuturi 2013]:

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

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

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

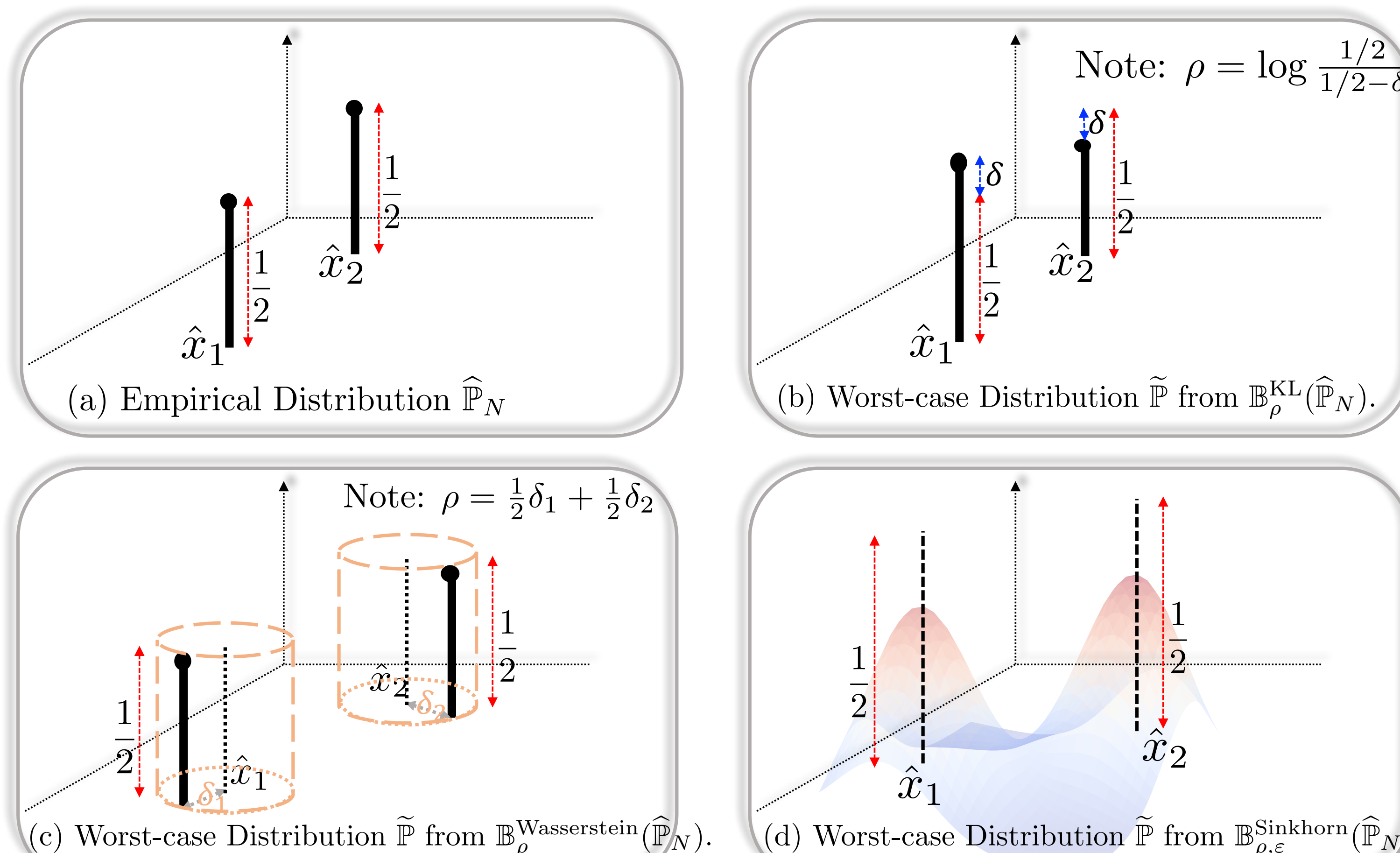


- Sinkhorn DRO:

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

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

- Visualization of Worst-Case Distributions:



## Theorem: Strong Dual Reformulation

Assume that

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

Then  $V_{\mathbb{P}} = V_{\mathbb{D}}$ :

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

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

where

$$\bar{\rho} = \rho + \varepsilon \int_{\Omega} \log \left( \int_{\Omega} e^{-c(x,z)/\varepsilon} d\nu(z) \right) d\hat{\mathbb{P}}(x),$$

$$d\mathbb{Q}_x(z) = \frac{e^{-c(x,z)/\varepsilon}}{\int_{\Omega} e^{-c(x,u)/\varepsilon} d\nu(u)} d\nu(z).$$

## Proof Sketch of Strong Duality

1. First show the **weak duality** result  $V_{\mathbb{P}} \leq V_{\mathbb{D}}$ .
2. Construct **primal feasible solution**  $\tilde{\mathbb{P}}$  with  $V_{\mathbb{P}} \geq \mathbb{E}_{z \sim \tilde{\mathbb{P}}}[f(z)] = V_{\mathbb{D}}$ .

## Geometry of Worst-Case Distribution:

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

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

- Worst-case distribution  $\tilde{\mathbb{P}} = \int \gamma_x d\hat{\mathbb{P}}(x)$ .

## Algorithm for Sinkhorn Robust Learning

$$V^* = \min_{\lambda \geq 0} \left\{ \lambda \bar{\rho} + \underbrace{\min_{\theta \in \Theta} \mathbb{E}_{x \sim \hat{\mathbb{P}}} \left[ \lambda \varepsilon \log \left( \mathbb{E}_{z \sim \mathbb{Q}_x} \left[ e^{f_{\theta}(z)/(\lambda \varepsilon)} \right] \right) \right]}_{V(\lambda)} \right\}$$

**Bisection Search on  $\lambda$ :** Estimating  $V(\lambda)$  up to accuracy  $O(\delta)$  for  $O(\text{Poly}(\log \frac{1}{\delta}))$  times to find  $\delta$ -optimal solution of  $V^*$ .

## Stochastic Approximation for Solving $V(\lambda)$

- 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 Stochastic Mirror Descent (BSMD): for  $t = 1, \dots, T$ ,

$$\begin{cases} v(\theta_t) \leftarrow (\text{biased}) \text{ gradient/subgradient estimate of } F(\theta_t) \\ \theta_{t+1} \leftarrow \text{Prox}_{\theta_t}(\gamma v(\theta_t)) \end{cases}$$

**Remark:** Gradient estimators should **optimally** balance the **bias-variance** trade-off.

- Complexity of finding  $\delta$ -optimal solution or  $\delta$ -critical point:

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	$\tilde{O}(\delta^{-2})$	$\tilde{O}(\delta^{-4})$

## Example: Mean-Risk Portfolio Optimization

