

## Motivation of Study

### Why Distributionally Robust Optimization(DRO)

DRO improves model robustness against distribution shift, which has important applications many several ML fields, including

- **Adversarial attack:** Gradient Attack will cause distribution shift over training data
- **Self-Supervised Learning:** Wrong selection of negative samples will cause distribution shift in embedded image-text pairs (i.e., CLIP model).
- **Reinforcement Learning:** Environment is subject to change, need to force policy shift for safety issue in real applications.

In this work, we study the information-divergence regularized DRO problem

$$\min_{x \in \mathbf{R}^d} \sup_{\mathbb{Q}} \left\{ \mathbb{E}_{\xi \sim \mathbb{Q}} [\ell(x; \xi)] - \lambda W_{\varepsilon}(\mathbb{P}, \mathbb{Q}) \right\}, \quad (1)$$

$\ell(x; \xi)$  represents loss function under shifted distribution  $\mathbb{Q}$ ,  $W_{\varepsilon}(\mathbb{P}, \mathbb{Q})$  represents information divergence among nominal distribution  $\mathbb{P}$  and shifted distribution  $\mathbb{Q}$ .

**Challenges:**  $\sup_{\mathbb{Q}}$  is maximized over distribution  $\rightarrow$  Hard to find explicit  $\mathbb{Q}^*$  in practice.

### Choice of $W_{\varepsilon}(\mathbb{P}, \mathbb{Q})$ : Generalized Sinkhorn Distance

Denote  $\Gamma(\mathbb{P}, \mathbb{Q})$  as the set of joint distributions that have marginal distributions  $\mathbb{P}, \mathbb{Q}$ . For a fixed regularization parameter  $\varepsilon > 0$  and a cost metric  $c: \Omega \times \Omega \rightarrow \mathbf{R}$ , the generalized Sinkhorn distance is defined as

$$W_{\varepsilon}(\mathbb{P}, \mathbb{Q}) = \inf_{\gamma \in \Gamma(\mathbb{P}, \mathbb{Q})} \left\{ \mathbb{E}_{(\xi, \xi') \sim \gamma} [c(\xi, \xi')] + \varepsilon D_f(\gamma \mid \mathbb{P} \otimes \nu) \right\},$$

where  $D_f$  corresponds to the  $f$ -divergence, that is,

$$D_f(\gamma \mid \mathbb{P} \otimes \nu) = \int f\left(\frac{d\gamma(\xi, \xi')}{d\mathbb{P}(\xi)d\nu(\xi')}\right) d\nu(\xi)d\mathbb{P}(\xi').$$

And  $\frac{d\gamma(\xi, \xi')}{d\mathbb{P}(\xi)d\nu(\xi')}$  represents density ratio of  $\gamma$  with respect to  $\mathbb{P} \otimes \nu$  evaluated at  $(\xi, \xi')$ .

### Why Generalized Sinkhorn Distance?

- **vs. KL:** 1. Symmetric; 2. Allows sample to have different probability support.
- **vs. Wasserstein Distance** Convex Programming  $\rightarrow$  easier to solve.
- **vs. Original Sinkhorn Distance**  $f$ -divergence is more general than  $KL$ -divergence.

## Our Contributions

### TL; DR

- **Generalize** Sinkhorn distance based on the class of  $f$ -divergence measures, which allows to use a broader range of divergences to model the ambiguity set.
- **Derive** an equivalent dual formulation with strong duality guarantee. The dual formulation shares novel structures, but it can be solved efficiently using nested stochastic programming.
- **Design** a Nested-SGD algorithm with convergence guarantee, which enables to solve large-scale problems.

Ghadimi, S. and Lan, G. (2013). Stochastic first- and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368.

Jin, J., Zhang, B., Wang, H., and Wang, L. (2021). Non-convex distributionally robust optimization: Non-asymptotic analysis. In *Advances in Neural Information Processing Systems*, pages 2771–2782. Curran Associates, Inc.

## Dual Problem Formulation, Assumptions and Structures

Denote  $\gamma_{\xi}(\xi)$  as conditional probability over  $\xi$ . We decompose the joint distribution as  $\gamma(\xi, \xi) = \gamma_{\xi}(\xi)\mathbb{P}(\xi)$  From principle of interchangeability, the primal problem in (1) can be rewritten as

$$\min_{x \in \mathbf{R}^d} \mathbb{E}_{\xi \sim \mathbb{P}} \left[ \sup_{\gamma_{\xi}} \left( \mathbb{E}_{\xi \sim \gamma_{\xi}} [\ell(x; \xi) - \lambda c(\xi, \xi)] - \lambda \varepsilon D_f(\gamma_{\xi} \mid \nu) \right) \right]. \quad (2)$$

By inverse C.D.F sampling, the inner supremum term  $\sup_{\gamma_{\xi}}(\cdot)$  has the following equivalent dual formulation

$$\min_{\eta \in \mathbf{R}} \underbrace{\left\{ L_{\xi}(x, \eta) := \lambda \varepsilon \mathbb{E}_{\xi \sim \nu} \left[ f^* \left( \frac{\ell(x; \xi) - \lambda c(\xi, \xi) - \eta}{\lambda \varepsilon} \right) \right] + \eta \right\}}_{L_{\xi}(x, \eta_{\xi}^*(\xi))}, \quad (3)$$

where  $\nu$  is the reference measure of  $\xi \sim \mathbb{Q}$ ;  $\eta$  is dual variable,  $f^*$  denotes the conjugate function of  $f$  and  $\eta_{\xi}^*(\xi) \in \arg \min_{\eta} L_{\xi}(x, \eta)$ .

For simplicity, denote

$$\Psi_{\xi}(x) := L_{\xi}(x, \eta_{\xi}^*(\xi)), \quad L_{\xi, \xi}(x, \eta) := \lambda \varepsilon f^* \left( \frac{\ell(x; \xi) - \lambda c(\xi, \xi) - \eta}{\lambda \varepsilon} \right) + \eta.$$

Then, the **dual problem** of (1) can be written as the following problem

$$\min_{x \in \mathbf{R}^d} \mathbb{E}_{\xi \sim \mathbb{P}} [\Psi_{\xi}(x)], \quad \text{where } \Psi_{\xi}(x) = L_{\xi}(x, \eta^*(\xi)). \quad (4)$$

### Challenges:

- **Double Expectation over different probability measure**,  $\xi \sim \nu$  and  $\xi \sim \mathbb{P} \rightarrow$  **Nested Structure!**
- **In-context inner minimizer  $\eta^*(\xi)$  subject to change with each  $\xi$ .**  $\rightarrow$  **Sample Inefficiency!**

In this work, we adapt following assumptions

- **Lipschitz Continuous and Smooth**  $\ell(\cdot; \xi)$  For every  $\xi$ ,  $\ell(\cdot; \xi)$  is  $G$ -Lipschitz continuous, and  $\ell(\cdot; \xi)$  is differentiable and  $L$ -smooth.
- **Smoothness of  $f^*$**  Function  $f^*(\cdot)$  is differentiable and  $M$ -smooth.
- **Bounded Variance of  $\ell$**  For every  $x$ , the variance of  $\ell(x; \cdot)$  is bounded by  $\sigma^2$ .
- **Bounded Variance of  $c$**  For every  $\xi$ , the variance of  $c(\xi, \cdot)$  is bounded by  $\delta^2$ . And for every  $\xi$ , the variance of  $c(\cdot, \xi)$  is bounded by  $\delta^2$ .

### Why Dual Formulation can be Solved by Nested Stochastic Programming? Two Fundamental Conclusions!

- **Gradient Equivalence between  $\nabla \Psi_{\xi}(x)$  and  $\nabla_1 L_{\xi}(x, \eta_{\xi}^*(\xi))$**  (Jin et al., 2021) Let Assumptions hold and consider any fixed  $x$  and  $\xi$ . Then, the function  $\Psi_{\xi}(x)$  is differentiable and satisfies  $\nabla \Psi_{\xi}(x) = \nabla_1 L_{\xi}(x, \eta_{\xi}^*(\xi))$ , where  $\eta_{\xi}^*(\xi) \in \arg \min_{\eta} L_{\xi}(x, \eta)$ .
- **Approximation Error Relationship** Suppose we obtain  $x$  and  $\eta_{\xi}(\xi)$  such that the gradient taken over second argument satisfies

$$\|\nabla_2 L_{\xi}(x, \eta_{\xi}(\xi))\| \leq \varepsilon_1. \quad (5)$$

Then, for any  $\xi$ , the gradient taken over first argument satisfies

$$\|\nabla \Psi_{\xi}(x) - \nabla_1 L_{\xi}(x, \eta_{\xi}(\xi))\| \leq G \varepsilon_1. \quad (6)$$

**Conclusion: As long as  $\eta_{\xi}^*(\xi)$  is near-optimal, we can guarantee  $\nabla_1 L_{\xi}(x, \eta_{\xi}(\xi))$  approximate  $\nabla \Psi_{\xi}(x)$  with controllable error!**

## Proposed Algorithms, Properties and Convergence

### Algorithm 1 Nested-SGD for solving $\mathbb{E}_{\xi \sim \mathbb{P}} [\Psi_{\xi}(x)]$

- 1: **Input:**  $T \in \mathbb{N}$ , initialization  $x_0, \eta_0$ , learning rate  $\gamma$
- 2: **for**  $t = 0$  **to**  $T - 1$  **do**
- 3:   Sample  $\{\xi\}$  and  $\{\xi\}_{B_1}$  with batch size  $B_1$
- 4:   Construct estimator  $\eta_{x_t}(\xi)$  via Algorithm 2
- 5:   Compute gradient estimator  $\hat{g}_t^B$  for  $\nabla \Psi_{\xi}(x)$
- 6:   Update  $x_{t+1} = x_t - \gamma \hat{g}_t^B$
- 7: **end for**
- 8: **Output:**  $\bar{x}_T$ , where  $\bar{t}$  is sampled from  $\{0, \dots, T - 1\}$  uniformly at random

### Algorithm 2 Construct Estimator $\eta_{\xi}(\xi)$

- 1: **Input:**  $D \in \mathbb{N}$ , learning rate  $\alpha_d$
- 2: **for**  $d = 0$  **to**  $D - 1$  **do**
- 3:   Utilize the  $\xi$  sampled in Algorithm 1
- 4:   Sample  $\{\xi\}_{B_2}$  with batch size  $B_2$
- 5:   Compute gradient estimator  $v_d^B$  for  $\nabla_2 L_{\xi, \xi}(x, \eta)$
- 6:   Update  $\eta_{x_t}^{d+1}(\xi) = \eta_{x_t}^d(\xi) - \alpha_d v_d^B$
- 7: **end for**
- 8: **Output:**  $\eta_{x_t}^{\bar{d}}(\xi)$ , where  $\bar{d} \in \{0, \dots, D - 1\}$  corresponds to the index with minimal gradient norm

- **Directional Smoothness:** For variable  $x$  and  $\eta$ , the following smoothness conditions hold. For any  $x, x'$ , it holds that

$$\mathbb{E}_{\xi \sim \mathbb{P}} \|\nabla \Psi_{\xi}(x) - \nabla_1 L_{\xi}(x', \eta_{\xi}^*(\xi))\|^2 \leq K^2 \|x - x'\|^2, \quad (7)$$

where  $K = G^2(\lambda \varepsilon)^{-1} M + L$ .

For any  $x$  and any  $\eta, \eta'$ , it holds that

$$\mathbb{E}_{\xi \sim \nu} \|\nabla_2 L_{\xi, \xi}(x, \eta) - \nabla_2 L_{\xi, \xi}(x, \eta')\|^2 \leq K'^2 \|\eta - \eta'\|^2, \quad (8)$$

where  $K' = M(\lambda \varepsilon)^{-1}$ .

- **Affine Bounded Variance:** For mini-batch gradient estimator  $\hat{g}_t^B$  used in Algorithm 1, it satisfies

$$\mathbb{E}_{\xi \sim \mathbb{P}, \xi_B \sim \nu} \|\hat{g}_t^B\|^2 \leq R_{B_1} + \frac{8G^2 \varepsilon_1^2}{B_1} + \|\nabla_1 \mathbb{E}_{\xi \sim \mathbb{P}} [L_{\xi}(x_t, \eta_{\xi}(\xi))]\|^2, \quad (9)$$

where  $R_{B_1} = O(\frac{G^2 + G^2 M^2 (\lambda \varepsilon)^{-2} \sigma^2}{B_1} + G^2 M^2 \varepsilon^{-2} \delta^{-2})$ .

For mini-batch gradient estimator  $v_d^B$  used in Algorithm 2, it satisfies

$$\mathbb{E}_{\xi_B \sim \nu} \|v_d^B\|^2 \leq \frac{R_2}{B_2} + \|\nabla_2 L_{\xi}(x_t, \eta_{x_t}^d(\xi))\|^2, \quad (10)$$

where  $R_2 = 2M^2(\lambda \varepsilon)^{-2}(\sigma^2 + \lambda^2 \delta^2)$ .

### Convergence of Main Algorithm

Let Assumptions hold. Denote  $\Delta = \mathbb{E}_{\xi \sim \mathbb{P}} [\Psi_{\xi}(x_0)] - \inf_x \mathbb{E}_{\xi \sim \mathbb{P}} [\Psi_{\xi}(x)]$ . Run Nested-SGD Algorithm 1 for  $T$  iterations with learning rate  $\gamma = \min \left\{ \frac{1}{3K}, \sqrt{\frac{2\Delta}{KR_{B_1}T}} \right\}$  and error threshold  $\varepsilon_1(t) = \Theta(G^{-1}T^{-\frac{1}{2}})$  for all  $t$ . Then, the convergence result is

$$\mathbb{E} \|\nabla \mathbb{E}_{\xi \sim \mathbb{P}} [\Psi_{\xi}(x_t)]\|^2 \leq O\left(\sqrt{\frac{\Delta K R_{B_1}}{T}}\right) + O\left(\frac{\Delta K}{T}\right) + O\left(\frac{B_1^{-1} \sqrt{\Delta K / R_{B_1}}}{T^{3/2}}\right). \quad (11)$$

Moreover, to achieve  $\mathbb{E} \|\nabla \mathbb{E}_{\xi \sim \mathbb{P}} [\Psi_{\xi}(x_t)]\| \leq \delta_1$ , choose  $B_1 = \Theta(1)$ , then the sample complexity of Algorithm 1 is  $\Omega(\Delta K R_{B_1} \delta_1^{-4})$ .

For Algorithm 2(1-dimension stochastic programming), the convergence analysis follows the standard analysis (Ghadimi and Lan, 2013). The difference is we use  $B_2 = \Theta(\varepsilon^{-2})$  mini-batch size to ensure convergence.