## **Distributionally Robust Optimization using Integral Probability Metrics** and Reproducing Kernel Hilbert Spaces

Weierstrass Institute for Applied Analysis and Stochastics Berlin, Germany & Max Planck Institute for Intelligent Systems Tübingen, Germany

Based on joint work with

Wittawat Jitkrittum (Google Research), Moritz Diehl (Uni. Freiburg), Bernhard Schölkopf (MPI Tübingen)

SIAM Conference on Optimization (OP21) July 21, 2021

## **Jia-Jie Zhu**

jj-zhu.github.io



andte Analysis und Stochastik



Code: https://github.com/jj-zhu/kdro



# Robust optimization

# Empirical risk minimization (ERM) (sample average approximation (SAA))



- Strength: high-performance (optimal)
- Weakness: fragile adversarial attacks, sim2real transfer, safety/off-policy in RL

# Robust optimization (RO) (robust control, games)



### Combine the strengths of ERM and RO: distributionally robust optimization (DRO) (ERM) min $\mathbb{E}$ $l(\theta, \xi)$ min sup $\mathbb{E}_{pl}(\theta, \xi)$ (dro) $\theta \quad \xi \sim \hat{P}$ $\theta$ [Delage and Ye 2010, Scarf 1958] (RO) min sup $l(\theta, \xi)$

Robustifies against a set of probability measures  $\mathscr{K}$  (*ambiguity set*), e.g.,

E∈U

- $\mathscr{K}$  can be a metric-ball centered at  $\hat{P}$ , e.g., using the popular Wasserstein metric, sets in RKHSs [this talk].
  - One way of constructing ambiguity region: one can quantify the empirical mean convergence rate  $\gamma(\hat{P}, P_{\text{true}}) \leq \epsilon$ .
    - Active research area: choosing better ambiguity regions
  - This talk provides a functional analysis and optimization perspective instead of statistics

Find the worst-case distribution! Problem of Moments [Stieltjes, Hausdorff, Hamburger, ...]





# Learning with kernels

- A kernel is a symmetric function  $k: \mathscr{X} \times \mathscr{X} \to \mathbb{R}$ , e.g., Gaussian kernel  $k(x, x') = \exp\left(-\|x - x'\|_{2}^{2} / 2\sigma^{2}\right).$
- A p.d. k corresponds to a Hilbert space  $\mathcal{H}$  (RKHS), which satisfies the **reproducing property**  $f(x) = \langle f, \phi(x) \rangle_{\mathcal{H}}, \forall f \in \mathcal{H}, x \in \mathcal{X},$  $\phi(x) := k(x, \cdot)$  is the **canonical feature** of  $\mathcal{H}$ .
- If  $\mathscr{H}$  is a large (dense in *C*),  $\gamma_{\mathscr{H}}$  is a metric on  $\mathscr{P}$ .
- We can generalize to the more general **integral probability**  $\bullet$ metric (IPM)

$$\mathsf{IPM}(\mathscr{F}; P, Q) := \sup_{f \in \mathscr{F}} \int f d(P - Q).$$

Special cases:

 $\mathscr{F} = \{f : ||f||_{\mathscr{H}} \leq 1\} \rightarrow \mathsf{Maximum Mean Discrepancy (MMD)}$  $\mathscr{F} = \{f : \|f\|_{\text{lip}} \le 1\} \longrightarrow \text{Wasserstein (type-1)}$ 



## Smooth is robust: Kernel DRO

(DRO) min sup 
$$\mathbb{E}_{P}l(\theta,\xi) \stackrel{\text{left}}{\sim} \sim P - \frac{1}{2} \frac{1$$

(P) min sup 
$$\left\{ \mathbb{E}_{P} l(\theta, \xi) \colon \int \phi \ dP = \mu, \mu \in \mathscr{C} \right\}$$

**Theorem (Kernel DRO duality, Zhu et al. '20)**. DRO (P) is equivalent to solving

(D) min  $\delta_{\mathscr{C}}^*(f)$  subject to  $l(\theta, \cdot) \leq f$ ,  $\theta, f \in \mathscr{H}$ 

 $\delta_{\mathscr{C}}^*(f)$  is the support function, e.g.,  $\mathbb{E}_{\hat{P}}f + \epsilon ||f||_{\mathscr{H}}$ . (Note: no need to estimate  $||l(\theta, \cdot)||_{\mathscr{H}}!$ )

Geometric intuition

Smoothness of  $f \leftrightarrow$  Distributional robustness ( $\leftrightarrow$  Size of  $\mathscr{H}$ ) Intuition: flatten the curve, smooth is robust



 $\mathcal{H}$ : special case SDP/SOS; generalization to IPM, e.g., W-1

 $l(\theta, \cdot)$ 



## Distributionally robust nonlinear optimization for machine learning and control

(DRNO) min sup  $\mathbb{E}_{P}l(\theta,\xi)$ 

[ZKNS '21] P∈.%  $\theta$ *l*: general nonlinear function, i.e., loss with DNN,  $l \notin \mathcal{H}$ . Kernel DRO handles this by finding a majorant  $f \in \mathcal{H}$ , with  $l(\theta, \cdot)$ no need to estimate  $||l(\theta, \cdot)||_{\mathscr{W}}$ (D)  $\min_{\theta, f \in \mathcal{H}} \delta_{\mathscr{C}}^{*}(f)$  subject to  $l(\theta, \cdot) \leq f$ DRO for stochastic model predictive control (MPC) with nonlinear constraints [NSZ '21] Test error on Fashion-MNIST, 5-layer CNN Test error on CIFAR-10, 20-layer ResNet 100 100 80 80 error (%) ERM error (%) 0.0--- WRM y=1.0 60 60 ARKS  $\sigma = 0.5$  $x_2$ IDRO WRM y=0.5 40 40 test test -2.5 -ARKS  $\sigma$ =0.9 20 20 0 10 50.2 0.2 0.0 0.1 0.3 0.0 0.1perturbation  $\Delta$ perturbation  $\Delta$  $x_1$ 

Code: jj-zhu.github.io/research

### **Adversarially Robust Kernel Smoothing**











# Conclusions

- A generalized dual program for solving DRO with general ambiguity sets and IPM-balls, with weak assumptions on the loss function (no need to estimate  $\|l(\theta, \cdot)\|_{\mathscr{H}}$ 
  - Kernel DRO: Maximizing w.r.t. a distribution  $\rightarrow$  finding a smooth surrogate function. For example, (D) min  $\mathbb{E}_{\hat{P}}f + \epsilon \|f\|_{\mathcal{H}}$ s.t.  $l(\theta, \cdot)$
- Takeaway  $\bullet$ 
  - Large (universal) RKHSs as dual spaces for DRO
  - Flatten the curve, smooth is robust

## **Future directions**

 $) \leq f$ 

 $l(\theta, \cdot)$ 

- Generalization and statistical bounds of Kernel DRO
  - Lam-Zeng 2021, Zhu in prep
- Kernel SIP, chance constraints...
  - Marteau-Ferey-Bach-Rudi 2020, Zhu et al. 2021, in prep (related: Lasserre moment-SOS)
- Applications to high-dim. data, deep models, adversarial learning, fairness, control...
  - Kernel DRO offers unique benefits but is not nearly as popular as the Wasserstein distance.







1. Zhu, J.-J., Kouridi, C., Nemmour, Y. & Schölkopf, B. Adversarially Robust Kernel Smoothing. arXiv:2102.08474 [cs, math, stat] (2021). https://arxiv.org/abs/2102.08474

2. Zhu, J.-J., Jitkrittum, W., Diehl, M. & Schölkopf, B. Kernel Distributionally Robust Optimization. Proceedings of the 24thInternational Conference on Artificial Intelligence and Statistics (AISTATS) 2021, San Diego, California, USA. PMLR: Volume 130. https://arxiv.org/abs/2006.06981

3. Nemmour, Y., Schölkopf, B. & Zhu, J.-J. Approximate Distributionally Robust Nonlinear Optimization with Application to Model Predictive Control: A Functional Approach. in Learning for Dynamics and Control 1255–1269 (PMLR, 2021). http://proceedings.mlr.press/v144/nemmour21a.html

4. Marteau-Ferey, U., Bach, F. & Rudi, A. Non-parametric Models for Non-negative Functions. arXiv:2007.03926 [cs, math, stat] (2020).

5. Lasserre, J.-B. The Moment-SOS hierarchy and the Christoffel-Darboux kernel. arXiv:2011.08566 [math, stat] (2020).

6. Lam, H. & Zeng, Y. Complexity-Free Generalization via Distributionally Robust Optimization. arXiv:2106.11180 [cs, math, stat] (2021).

## Related references

Code: jj-zhu.github.io/research

### Jia-Jie Zhu <u>jj-zhu.github.io</u>

Weierstrass Institute, Berlin & Max Planck Institute, Tübingen Germany

Ph.D. positions available in Berlin, Germany Robust machine learning and data-driven optimization & control



## **Co-authors**



Wittawat Jitkrittum (Google Research) Moritz Diehl (Uni. Freiburg) Bernhard Schölkopf (MPI Tübingen)

SIAM OP21