

Day 2: Data-Driven Modeling and Optimization of Dynamical Systems under Uncertainty

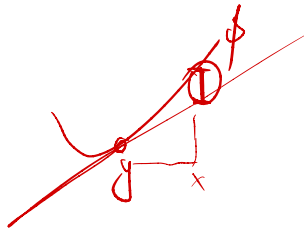
Jia-Jie Zhu

WIAS-Berlin

July 14, 2022

Mirror descent

$$\begin{aligned} \mathcal{P} &\rightarrow \mathcal{P}^* \\ \phi &\rightarrow \psi \end{aligned}$$



We also define the Bregman divergence associated with ϕ as

$$B_\phi(x, y) = \phi(x) - \phi(y) - \nabla \phi(y)^\top (x - y).$$

Mirror descent update:

$$x_{t+1} = \operatorname{argmin}_x \left\{ f(x_t) + g_t^\top (x - x_t) + \frac{1}{\eta_k} B_\phi(x, x_t) \right\}$$

$$\min_x f(x)$$

$$g_t = \nabla f(x_t) \in \partial f(x_t)$$

$$\|x_t - x\|_2^2 \leq R^2$$

Example of MD: Euclidean norm

Mirror map $\Phi(x) = \frac{1}{2}\|x\|_2^2$.

Example of MD: Euclidean norm

Mirror map $\Phi(x) = \frac{1}{2}\|x\|_2^2$.

MD:

$$x_{k+1} = \operatorname{argmin}_x \left\{ f(x_k) + g_k^\top (x - x_k) + \frac{1}{2\eta_k} \|x - x_k\|_2^2 \right\}$$

This is equivalent to (sub-)GD with step η_k .

Example of MD: negative entropy

Mirror map $\Phi(x)$ = $\sum_{i=1}^d \underbrace{x(i)} \log \underbrace{x(i)}$.

$$\begin{array}{ccccccc} | & & | & & & & | \\ x_1 & & x_2 & & \dots & & x_d \end{array}$$

Example of MD: negative entropy

$$(\mathbf{x} \succ 0, \mathbf{1})$$

Mirror map $\Phi(x) = \sum_{i=1}^d x(i) \log x(i)$. The resulting Bregman divergence is the KL-divergence $D_\Phi(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^d x(i) \log \frac{x(i)}{y(i)}$.

Example of MD: negative entropy

Mirror map $\Phi(x) = \sum_{i=1}^d x(i) \log x(i)$. The resulting Bregman divergence is the KL-divergence $D_\Phi(x, y) = \sum_{i=1}^d x(i) \log \frac{x(i)}{y(i)}$.

If we restrict x to a (discrete probability) simplex $1^\top x = 1$, then the MD update

$$x_{t+1} = \operatorname{argmin}_{1^\top x = 1} \left\{ f(x_t) + \underline{g_t^\top (x - x_t)} + \frac{1}{\eta} B_\Phi(x, x_t) \right\}$$

has the closed-form

$$x_{i,k+1} = \frac{x_{i,k} e^{-\eta g_k}}{\sum_{j=1}^n x_{j,k} e^{-\eta g_k}}, \quad i = 1, 2, \dots, n$$

MD convergence

Let Φ be ρ -strongly convex in $\|\cdot\|$. Let $R^2 = \sup_x \Phi(x) - \Phi(x_1)$, and f be convex and L -Lipschitz w.r.t. $\|\cdot\|$. Then mirror descent with $\eta = \frac{R}{L} \sqrt{\frac{2\rho}{t}}$ satisfies

$$f\left(\frac{1}{t} \sum_{s=1}^t x_s\right) - f(x^*) \leq RL \sqrt{\frac{2}{\rho t}}.$$

MD convergence

Let Φ be ρ -strongly convex in $\|\cdot\|$. Let $R^2 = \sup_x \Phi(x) - \Phi(x_1)$, and f be convex and L -Lipschitz w.r.t. $\|\cdot\|$. Then mirror descent with $\eta = \frac{R}{L} \sqrt{\frac{2\rho}{t}}$ satisfies

$$f\left(\frac{1}{t} \sum_{s=1}^t x_s\right) - f(x^*) \leq RL \sqrt{\frac{2}{\rho t}}.$$

For example, for negative entropy $R \leq \sqrt{\log(d)}$.

Stochastic (sub-)gradient descent

Stochastic oracle:

$$\mathbb{E} \tilde{g}(x) \in \partial f(x).$$

Then, SGD update rule:

$$x_{t+1} = x_t - \eta_t \tilde{g}(x_t).$$

$$g \in \partial f(x) \quad (= \nabla f(x)) \quad (\text{ERM})$$

$$\nabla_x \left\{ \frac{1}{n} \sum_{i=1}^n f(x, \xi_i) \right\}$$

$$\nabla_x f(x, \xi_i)$$

$$\xi_i \sim P_0$$

$$\nabla_x \mathbb{E}_{P_0} f(x, \xi)$$

Stochastic (sub-)gradient descent

Stochastic oracle:

$$\mathbb{E} \tilde{g}(x) \in \partial f(x).$$



Then, SGD update rule:

$$x_{t+1} = x_t - \eta_t \tilde{g}(x_t).$$

Convergence. Let f be α -strongly convex, and assume that the stochastic oracle is such that $\mathbb{E} \|\tilde{g}(x)\|^2 \leq B^2$. Then SGD with $\eta_k = \frac{2}{\alpha(k+1)}$ satisfies

$$\mathbb{E} f \left(\sum_{k=1}^t \frac{2k}{t(t+1)} x_k \right) - f(x^*) \leq \frac{2B^2}{\alpha(t+1)}.$$

$$f(x_t)$$

(Skip) stochastic mirror descent (S-MD)

Let $x_1 \in \operatorname{argmin}_{\mathcal{X}} \Phi(x)$, and

$$x_{t+1} = \operatorname{argmin}_x f(x_t) + \tilde{g}(x_t)^\top (x - x_t) + \frac{1}{\eta} B_\Phi(x, x_t).$$

Convergence. Let Φ be a mirror map 1-strongly convex. Let $R^2 = \sup_x \Phi(x) - \Phi(x_1)$. Let f be convex and β -smooth w.r.t. $\|\cdot\|$. Assume that the stochastic oracle is such that $\mathbb{E}\|\nabla f(x) - \tilde{g}(x)\|_*^2 \leq \sigma^2$. Then S-MD with stepsize $\frac{1}{\beta+1/\eta}$ and $\eta = \frac{R}{\sigma} \sqrt{\frac{2}{t}}$ satisfies

$$\mathbb{E}f\left(\frac{1}{t} \sum_{s=1}^t x_{s+1}\right) - f(x^*) \leq R\sigma \sqrt{\frac{2}{t}} + \frac{\beta R^2}{t}.$$

Compare the basic gradient descent with SGD for average loss functions

GD:

$$f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x).$$

$$x_{t+1} = x_t - \frac{\eta}{n} \sum_{i=1}^n \nabla f_i(x),$$

$$x_{t+1} = x_t - \eta \nabla f_{i_t}(x),$$

where i_t is drawn uniformly random.

$$f(x, \underline{\xi}_i) \quad \{\underline{\xi}_i\} \quad \{x_i, y_i\} \quad \mathcal{O}(e^{-t})$$

$E p_0$

$$\mathcal{O}(\underline{\underline{\theta}}_t^{-\alpha})$$

$$\underline{\xi}_i \sim p_2$$

$$\mathcal{N}(\underline{\xi}_i)$$

$$x_t = (X_f^{T+1} X_f^T)^{-1} X_f^T y$$

Optimization under distributional uncertainty

Empirical risk minimization

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N l(f_{\theta}; \xi_i), \quad \xi_i \sim P_0$$

Handwritten notes: $f_{\theta} \in H$, $\mathbb{E}_{P_0} l(f, \xi)$

- Robust under statistical fluctuation, e.g., we can bound $\mathbb{E}_{P_0} l(f_{\hat{\theta}}, \xi)$

Optimization under distributional uncertainty

Empirical risk minimization

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N l(f_{\theta}, \xi_i), \quad \xi_i \sim P_0$$

- ▶ Robust under statistical fluctuation, e.g., we can bound $\mathbb{E}_{P_0} l(f_{\hat{\theta}}, \xi)$
- ▶ Not robust under data distribution shifts, when $Q (\neq P_0)$

Optimization under distributional uncertainty

Empirical risk minimization

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N l(f_{\theta}, \xi_i), \quad \xi_i \sim P_0$$

- ▶ Robust under statistical fluctuation, e.g., we can bound $\mathbb{E}_{P_0} l(f_{\hat{\theta}}, \xi)$
- ▶ Not robust under data distribution shifts, when $Q (\neq P_0)$

Distributionally robust optimization and learning

$$\min_{f \in \mathcal{H}} \sup_{Q \in \mathcal{M}} \mathbb{E}_Q l(f, \xi)$$

- ▶ Minimize risk under a local worst-case distribution Q

Optimization under distributional uncertainty

Empirical risk minimization

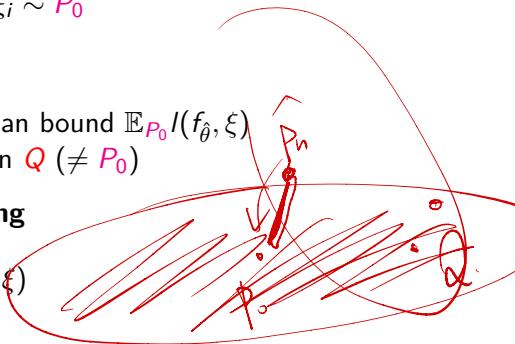
$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N l(f_{\theta}, \xi_i), \quad \xi_i \sim P_0$$

- ▶ Robust under statistical fluctuation, e.g., we can bound $\mathbb{E}_{P_0} l(f_{\hat{\theta}}, \xi)$
- ▶ Not robust under data distribution shifts, when $Q (\neq P_0)$

Distributionally robust optimization and learning

$$\min_{f \in \mathcal{H}} \sup_{Q \in \mathcal{M}} \mathbb{E}_Q L(f, \xi)$$

- ▶ Minimize risk under a local worst-case distribution Q
- ▶ Distribution shift described by an ambiguity set \mathcal{M}



Optimization under distributional uncertainty

Empirical risk minimization

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N l(f_{\theta}, \xi_i), \quad \xi_i \sim P_0$$

- ▶ Robust under statistical fluctuation, e.g., we can bound $\mathbb{E}_{P_0} l(f_{\hat{\theta}}, \xi)$
- ▶ Not robust under data distribution shifts, when $Q (\neq P_0)$

Distributionally robust optimization and learning

$$\min_{f \in \mathcal{H}} \sup_{Q \in \mathcal{M}} \mathbb{E}_Q L(f, \xi)$$

- ▶ Minimize risk under a local worst-case distribution Q
- ▶ Distribution shift described by an ambiguity set \mathcal{M}
- ▶ We can bound performance beyond statistical fluctuation (classical learning theory)

DRO with divergence or metric ball constraints

Let \mathcal{D} denote a divergence measure or metric on the \mathcal{P} , we consider the data-driven DRO problem

$$\min_{\theta} \max_{\substack{\mathcal{P} \\ \mathcal{D}(\mathcal{P}, \hat{\mathcal{P}}) \leq \epsilon}} \mathbb{E}_{\mathcal{P}} l(\theta, \xi)$$

$$\hat{\mathcal{P}} = \frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{z}_i}$$

$$\mathcal{P}^+$$

$$\hat{\mathcal{P}}$$

Φ -divergence (skip in class)

Tsybakov - -

Wasserstein distance

p -Wasserstein distance between probability measures μ, ν on \mathbb{R}^d (with p finite moments) is defined through the following Kantorovich problem

$$W_p(P, Q)^p := \min \left\{ \int_{\mathbb{R}^d \times \mathbb{R}^d} |x_0 - x_1|^p d\Pi \mid \pi_{\#}^{(1)} \Pi = P, \pi_{\#}^{(2)} \Pi = Q \right\}.$$

Wasserstein distance

p -Wasserstein distance between probability measures μ_0, μ_1 on \mathbb{R}^d (with p finite moments) is defined through the following Kantorovich problem

$$W_p(P, Q)^p := \min \left\{ \int_{\mathbb{R}^d \times \mathbb{R}^d} |x_0 - x_1|^p d\Pi \mid \pi_{\#}^{(1)} \Pi = P, \pi_{\#}^{(2)} \Pi = Q \right\}.$$

The dual Kantorovich problem

$$\max \left\{ \int_X \phi dP + \int_X \psi dQ : \phi, \psi \in C(X), \phi(x) + \psi(y) \leq c(x, y) \forall x, y \right\}.$$

∞ -dim
dual

Wasserstein distance

p -Wasserstein distance between probability measures μ_0, μ_1 on \mathbb{R}^d (with p finite moments) is defined through the following Kantorovich problem

$$W_p(P, Q)^p := \min \left\{ \int_{\mathbb{R}^d \times \mathbb{R}^d} |x_0 - x_1|^p d\Pi \mid \pi_{\#}^{(1)} \Pi = P, \pi_{\#}^{(2)} \Pi = Q \right\}.$$

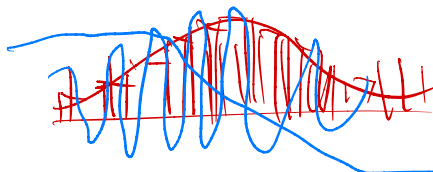
The dual Kantorovich problem

$$\max \left\{ \int_X \phi dP + \int_X \psi dQ : \phi, \psi \in C(X), \phi(x) + \psi(y) \leq c(x, y), \forall x, y \right\}.$$

Convergence depends on dimensions, e.g., (See also [Weed and Bach 2017])

$$W_1(\hat{P}, P_0) \geq \mathcal{O}(n^{-\frac{1}{d}}).$$

$$\xi_i \sim P_0$$



(Skip) The dynamic formulation of the Wasserstein distance



$$W_2(P, Q)^2 = \min \left\{ \int_0^1 \int_{\mathbb{R}^d} \underbrace{|v_t|^2}_{\text{velocity squared}} d\mu_t dt \mid \mu_0 = P, \mu_1 = Q, \underbrace{\frac{d}{dt} \mu_t + \operatorname{div}(v_t \mu_t)}_{\text{continuity equation}} = 0 \right\}$$

c-transform

Given a function $f : X \rightarrow \overline{\mathbb{R}}$ we define its $\{c\text{-transform}\}$ (or $c\text{-conjugate function}$) by

$$f^c(y) = \inf_{x \in X} \underbrace{c(x, y) - f(x)}_{\langle x, y \rangle}.$$

- Moreover, we say that a function ψ is $c\text{-concave}$ if there exists ϕ such that $\psi = \phi^c$ and we denote by $\Psi_c(X)$ the set of $c\text{-concave}$ functions.

c -transform

Given a function $f : X \rightarrow \overline{\mathbb{R}}$ we define its $\{c\text{-transform}\}$ (or c -conjugate function) by

$$f^c(y) = \inf_{x \in X} c(x, y) - f(x).$$

- ▶ Moreover, we say that a function ψ is c -concave if there exists ϕ such that $\psi = \phi^c$ and we denote by $\Psi_c(X)$ the set of c -concave functions.
- ▶ Using this transform, we can write down the so-called semi-dual formulation

$$\max \left\{ \int_X \phi \, d\mu + \int_X \phi^c \, d\nu : \phi \in \Psi_c(X) \right\}.$$

c-transform

Given a function $f : X \rightarrow \overline{\mathbb{R}}$ we define its $\{c\text{-transform}\}$ (or c -conjugate function) by

$$\underline{f^c(y) = \inf_{x \in X} c(x, y) - f(x).}$$

- ▶ Moreover, we say that a function ψ is c -concave if there exists ϕ such that $\psi = \phi^c$ and we denote by $\Psi_c(X)$ the set of c -concave functions.
- ▶ Using this transform, we can write down the so-called semi-dual formulation

$$\max \left\{ \int_X \phi \, d\mu + \int_X \phi^c \, d\nu : \phi \in \Psi_c(X) \right\}.$$

- ▶ Exercise. Can you see the relation with the α -strong convexity and β -smoothness we talked about?

Dual reformulation for 2-Wasserstein DRO

The primal DRO problem is intractable

$$\min_{\theta} \sup_{W_2(P, \hat{P}) \leq \epsilon} \mathbb{E}_P l(\theta, \xi).$$



$$\mathcal{D}(P, Q) = 0$$

when $P \neq Q$

$$\mathcal{D}(P, P) = 0$$



$$\max_{C \in \mathcal{R}} C$$

$$\text{s.t. } \mathbb{E}_P l(\theta, \xi) \leq C, \quad \forall P \in \mathcal{P}_2(\hat{P})$$

$$\max_{\xi} l(\theta, \xi)$$

$$\text{Moment Constr. } \mathbb{E}_P \xi^p = \mathbb{E}_Q \xi^p$$



Dual reformulation for 2-Wasserstein DRO

The primal DRO problem is intractable

$$\hat{P}_N = \frac{1}{N} \sum_{i=1}^N \delta_{\xi_i}$$

$$\min_{\theta} \sup_{W_2(P, \hat{P}) \leq \epsilon} \mathbb{E}_P l(\theta, \xi).$$

Fortunately, it has a strong dual as follows, i.e., the two problems have the same optimal value

$$\min_{\theta \in \mathcal{H}} \frac{1}{N} \sum_{i=1}^N (l_{\theta})^{\lambda \|\cdot\|^2}(\xi_i) + \lambda \epsilon^2$$

$$(P) \leq (D)$$

C-transform of $l(f_0, \cdot)$

$$\sup_{\lambda} \{ l(f_0, \xi) - \frac{\lambda}{2} \|\xi - \hat{\xi}\|^2 \}$$

Dual reformulation for 2-Wasserstein DRO

The primal DRO problem is intractable

$$\min_{\theta} \sup_{W_2(P, \hat{P}) \leq \epsilon} \mathbb{E}_P l(\theta, \xi).$$

Fortunately, it has a strong dual as follows, i.e., the two problems have the same optimal value

$$\min_{\theta, \lambda > 0} \frac{1}{N} \sum_{i=1}^N (l_{\theta})^{\lambda \|\cdot\|^2}(\xi_i) + \lambda \epsilon^2$$

This can be shown to motivate a stochastic gradient algorithm for DRO.

PA
① Sup.
② inf

Kernel maximum mean discrepancy (MMD)

$$\text{MMD}(P, Q) := \left\| \int \underline{k(x, \cdot)} dP - \int \underline{k(x, \cdot)} dQ \right\|_{\mathcal{H}}$$

$$\underline{\int k(x, \cdot) d(P-Q)}$$

Kernel maximum mean discrepancy (MMD)

$$\text{MMD}(P, Q) := \left\| \int k(x, \cdot) dP - \int k(x, \cdot) dQ \right\|_{\mathcal{H}}$$

Given two samples from the distribution of interest

$$x_i \sim P, i = 1 \dots M; y_j \sim Q, j = 1 \dots N,$$

$$\begin{aligned} \text{MMD}(P, Q)^2 &= \mathbb{E}_{x, x' \sim P} k(x, x') + \mathbb{E}_{y, y' \sim Q} k(y, y') - 2\mathbb{E}_{x \sim P, y \sim Q} k(x, y) \\ &\approx \frac{1}{MN} \sum_{i=1}^M \sum_{j=1}^N k(x_i, x_j) + \frac{1}{MN} \sum_{i=1}^M \sum_{j=1}^N k(y_i, y_j) - 2 \frac{1}{MN} \sum_{i=1}^M \sum_{j=1}^N k(x_i, y_j) \end{aligned}$$

This is particularly handy in, e.g., training deep generative models.

$\hat{P} = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$
 bedroom
 \downarrow
 $\min_Q (P, Q)$
 MMD-GAN
 W-GAN
 $\frac{1}{m} \sum_{i=1}^m \mathbb{E} [z_i]$



Kernel maximum mean discrepancy (MMD)

$$\text{MMD}(P, Q) := \left\| \int k(x, \cdot) dP - \int k(x, \cdot) dQ \right\|_{\mathcal{H}}.$$

Given two samples from the distribution of interest

$$x_i \sim P, i = 1 \dots M; y_j \sim Q, j = 1 \dots N,$$

$$\begin{aligned} \text{MMD}(P, Q)^2 &= \mathbb{E}_{x, x' \sim P} k(x, x') + \mathbb{E}_{y, y' \sim Q} k(y, y') - 2\mathbb{E}_{x \sim P, y \sim Q} k(x, y) \\ &\approx \frac{1}{MN} \sum_{i=1}^M \sum_{j=1}^N k(x_i, x'_j) + \frac{1}{MN} \sum_{i=1}^M \sum_{j=1}^N k(y_i, y'_j) - 2 \frac{1}{MN} \sum_{i=1}^M \sum_{j=1}^N k(x_i, y_j) \end{aligned}$$

This is particularly handy in, e.g., training deep generative models.

Just like W_p , the MMD has a dual formulation

$$\text{MMD}(P, Q) = \sup_{\|f\|_{\mathcal{H}} \leq 1} \int f d(P - Q).$$

Kernel maximum mean discrepancy (MMD)

$$\text{MMD}(P, Q) := \left\| \int \underline{k}(x, \cdot) dP - \int \underline{k}(x, \cdot) dQ \right\|_{\mathcal{H}}.$$

Given two samples from the distribution of interest

$$x_i \sim P, i = 1 \dots M; y_j \sim Q, j = 1 \dots N,$$

$$\begin{aligned} \text{MMD}(P, Q)^2 &= \mathbb{E}_{x, x' \sim P} k(x, x') + \mathbb{E}_{y, y' \sim Q} k(y, y') - 2\mathbb{E}_{x \sim P, y \sim Q} k(x, y) \\ &\approx \frac{1}{MN} \sum_{i=1}^M \sum_{j=1}^N k(x_i, x'_j) + \frac{1}{MN} \sum_{i=1}^M \sum_{j=1}^N k(y_i, y'_j) - 2\frac{1}{MN} \sum_{i=1}^M \sum_{j=1}^N k(x_i, y_j) \end{aligned}$$

This is particularly handy in, e.g., training deep generative models.

Just like W_p , the MMD has a dual formulation

$$\text{MMD}(P, Q) = \sup_{\|f\|_{\mathcal{H}} \leq 1} \int f d(P - Q).$$

Convergence

$$\text{MMD}(\hat{P}, P_0) \leq \mathcal{O}(n^{-\frac{1}{2}}).$$

†

Dual reformulation for Kernel DRO

$\ell \in \mathcal{H}$

The DRO problem with MMD constraint

$$\min_{\theta} \sup_{\text{MMD}(P, \hat{P}) \leq \epsilon} \mathbb{E}_P \ell(\theta, \xi)$$

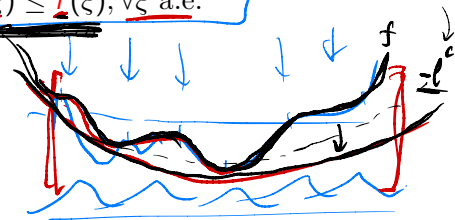
$$\phi(x) + \psi(y) \leq c(x, y)$$

can be reformulated using a Kantorovich-type duality as

$$\min_{\theta, f \in \mathcal{H}} \frac{1}{N} \sum_{i=1}^N f(\xi_i) + \epsilon \|f\|_{\mathcal{H}} \quad \text{s.t. } \ell(\theta, \xi) \leq f(\xi), \forall \xi \text{ a.e.}$$

$$\sum \alpha_i \delta(x_i, \cdot)$$

"Moreau"



$f \in \mathcal{H}$
 ~~$f \in \mathcal{H}$~~

Dual reformulation for Kernel DRO

The DRO problem with MMD constraint

$$\min_{\theta} \sup_{\text{MMD}(P, \hat{P}) \leq \epsilon} \mathbb{E}_P l(\theta, \xi)$$

can be reformulated using a Kantorovich-type duality as

$$\min_{\theta, \underline{f} \in \mathcal{H}} \frac{1}{N} \sum_{i=1}^N \underline{f}(\xi_i) + \epsilon \|\underline{f}\|_{\mathcal{H}} \quad \text{s.t. } \underline{l(\theta, \xi) \leq \underline{f}(\xi), \forall \xi \text{ a.e.}}$$

Optionally, we may consider to solve the problem with a relaxed (albeit with statistical guarantee) constraint

$$\underline{l(\theta, \xi_i) \leq \underline{f}(\xi_i), \quad i = 1, \dots, N.$$

Dual reformulation for Kernel DRO

The DRO problem with MMD constraint

$$\min_{\theta} \sup_{\text{MMD}(P, \hat{P}) \leq \epsilon} \mathbb{E}_P l(\theta, \xi)$$

can be reformulated using a Kantorovich-type duality as

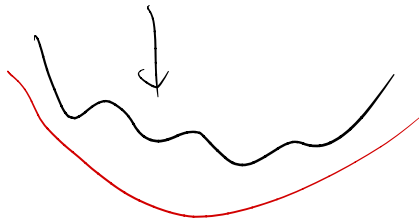
$$\min_{\theta, \underline{f} \in \mathcal{H}} \frac{1}{N} \sum_{i=1}^N \underline{f}(\xi_i) + \epsilon \|\underline{f}\|_{\mathcal{H}} \quad \text{s.t. } l(\theta, \xi) \leq \underline{f}(\xi), \forall \xi \text{ a.e.}$$

Optionally, we may consider to solve the problem with a relaxed (albeit with statistical guarantee) constraint

$$l(\theta, \xi_i) \leq \underline{f}(\xi_i), \quad i = 1, \dots, N.$$

The solution has the nice property of being

- ▶ a kernel interpolant of the loss $l(\theta, \xi)$ at the data points ξ_i 's.



Dual reformulation for Kernel DRO

The DRO problem with MMD constraint

$$\hat{P}_N$$

$$\underbrace{\underbrace{\xi}_P \underbrace{\nu}_P}$$

$$\min_{\theta} \sup_{\text{MMD}(P, \hat{P}) \leq \epsilon} \mathbb{E}_P l(\theta, \xi)$$

$$\int f dP - \hat{P}$$

can be reformulated using a Kantorovich-type duality as

$$\min_{\theta, f \in \mathcal{H}} \frac{1}{N} \sum_{i=1}^N f(\xi_i) + \epsilon \|f\|_{\mathcal{H}} \quad \text{s.t.} \quad l(\theta, \xi) \leq f(\xi), \forall \xi \text{ a.e.}$$

Optionally, we may consider to solve the problem with a relaxed (albeit with statistical guarantee) constraint

$$l(\theta, \xi_i) \leq f(\xi_i), \quad i = 1, \dots, N.$$

The solution has the nice property of being

- ▶ a kernel interpolant of the loss $l(\theta, \xi)$ at the data points ξ_i 's.
- ▶ a witness (optimal test) function between \hat{P} and the underlying worst-case distribution P .

Going beyond DRO

Going beyond DRO

- ▶ Big gap between theory and practice in large-scale learning models

Going beyond DRO

- ▶ Big gap between theory and practice in large-scale learning models
- ▶ Continuous optimization and gradient flows of nonlinear functionals of measures and distributions

$$\min_{\mu \in \mathcal{M}} F(\mu).$$

Handwritten red annotations: A circle around the minimization symbol and the set \mathcal{M} , a circle around the functional $F(\mu)$, and a box around the expression $F(\mu)$ with a red arrow pointing from the circle around $F(\mu)$ to the box.

End! I