

So far we studied moment-based Uncertainty sets. In some applications, we may be more interested in distributions that are "close" to a reference distribution that is induced by the data.

e.g.: Consider the scenario where we have iid realizations  $\xi_1, \dots, \xi_N$  of a random vector  $\xi$ . Define the empirical measure by

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

Intuitively, as  $N \rightarrow \infty$ ,  $\hat{\mathbb{P}}_N$  should tend to the true distribution  $\mathbb{P}^*$  of  $\xi$ .

\* For fixed measurable  $A \subseteq \mathbb{R}^n$ ,  $\hat{\mathbb{P}}_N(A) \rightarrow \mathbb{P}^*(A)$  (SLLN)

\* Uniform convergence of  $\hat{\mathbb{P}}_N$  to  $\mathbb{P}^*$ : Vapnik-Chervonenkis.

\* How do we describe a "neighborhood" of the reference distribution?

1) Divergence-based measures

e.g.:  $f$ -divergence

Let  $f: \mathbb{R}_+ \rightarrow \mathbb{R}$  be convex with  $f(1) = 0$ . Let  $\mathbb{P}$  and  $\mathbb{Q}$  be probability distributions on a probability space  $(\Xi, \mathcal{B})$  with  $\mathbb{P}$  being the reference distribution. Suppose that  $\mathbb{Q}$  is absolutely continuous wrt  $\mathbb{P}$  (intuitively,  $\mathbb{P}, \mathbb{Q}$  have the same support).

The  $f$ -divergence of  $\mathbb{Q}$  from  $\mathbb{P}$  is defined as

$$D_f(\mathbb{Q} \parallel \mathbb{P}) = \int f\left(\frac{d\mathbb{Q}}{d\mathbb{P}}\right) d\mathbb{P}$$

e.g.: KL-divergence:  $f(t) = t \log t$

TV:  $f(t) = \frac{1}{2}|t-1|$

Then, we can define the "ball"

$$\mathcal{B}_{f, \varepsilon}(\mathbb{P}) = \{\mathbb{Q} : D_f(\mathbb{Q} \parallel \mathbb{P}) \leq \varepsilon\}$$

See, e.g.,

Ben-Tal et al.: Robust Solutions of Optimization Problems Affected by Uncertain Probabilities. *Manag. Sci.* 59(2): 341-357, 2013.

## 2) Probability metric-based measures

e.g.: Wasserstein distance

Let  $\mathbb{P}, \mathbb{Q}$  be as above. Let  $d(\cdot, \cdot)$  be a norm on  $\mathbb{R}^n$ . Define the d-Wasserstein distance between  $\mathbb{P}$  and  $\mathbb{Q}$  by

$$W(\mathbb{Q}, \mathbb{P}) = \inf_{\Pi \in \mathcal{M}(\Xi \times \Xi)} \left\{ \int_{\Xi \times \Xi} \underbrace{d(\xi, \xi')}_{\text{transport cost}} \Pi(d\xi, d\xi') : \underbrace{\Pi(\Xi, d\xi') = \mathbb{P}(d\xi'), \Pi(d\xi, \Xi) = \mathbb{Q}(d\xi)}_{\text{marginal distributions}} \right\}$$

Then, we can define the ball

$$B_\epsilon(\mathbb{P}) = \{ \mathbb{Q} : W(\mathbb{Q}, \mathbb{P}) \leq \epsilon \}$$

Wasserstein distance-based DRO has attracted much attention due to its connection to various problems in machine learning.

### Application: Logistic Regression

Setup. \*  $x \in \mathbb{R}^n$ : feature vector ;  $y \in \{+1, -1\}$  binary label

\* conditional distribution of  $y$  given  $x$ :

$$\Pr(y|x) = (1 + \exp(-y \cdot \beta^T x))^{-1},$$

where  $\beta$  is the regression parameter.

\* Suppose that  $N$  training samples  $\{\xi_i = (x_i, y_i)\}_{i=1}^N$  of the underlying random vector  $\xi = (x, y)$  are observed.

Then, the MLE of  $\beta$  is given by

$$\hat{\beta}_{MLE} = \arg\min_{\beta} \frac{1}{N} \sum_{i=1}^N \ell_{\beta}(\xi_i) = \arg\min_{\beta} \mathbb{E}_{\mathbb{P}_N} [\ell_{\beta}(\xi)],$$

where  $\ell_{\beta}(x, y) = \log(1 + \exp(-y \beta^T x))$  is the log-loss.

To avoid overfitting, a typical approach is regularization:

$$\hat{\beta} = \arg\min_{\beta} \frac{1}{N} \sum_{i=1}^N \ell_{\beta}(\xi_i) + \frac{\lambda}{2} \|\beta\|_2^2$$

To avoid overfitting, a typical approach is regularization:

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \quad \mathbb{E}_{\mathbb{P}_N}^{\wedge} [l_{\beta}(\xi)] + \varepsilon R(\beta)$$

e.g.:  $R(\beta) = \|\beta\|_1$  ,  $R(\beta) = \|\beta\|_2^2$  , ...

sparsity-inducing      ridge

One way of understanding the overfitting phenomenon is that  $\hat{\beta}_{MLE}$  above does not account for the **Unseen** data. This motivates the following formulation:

$$\min_{\beta} \sup_{\mathbb{B}_{\varepsilon}(\mathbb{P}_N^{\wedge})} \mathbb{E}_{\mathbb{P}} [l_{\beta}(\xi)]$$

The hope is that if  $\varepsilon$  is chosen appropriately, then the true distribution  $\mathbb{P}^* \in \mathbb{B}_{\varepsilon}(\mathbb{P}_N^{\wedge})$ , so that the solution  $\beta^*$  takes into account the effect of  $\mathbb{P}^*$ .

Here, for  $\xi = (x, y)$ ,  $\xi' = (x', y')$ , we take

$$d(\xi, \xi') = \|x - x'\| + \frac{\kappa}{2} |y - y'|$$

as the transport cost, where  $\|\cdot\|$  is an arbitrary norm and  $\kappa > 0$  is a parameter that specifies the relative emphasis between feature mismatch and label uncertainty.

\* In particular, if  $\kappa = +\infty$ , then the label needs to be a deterministic function of the feature, and label measurements are exact.

Recall our problem of interest:

$$\inf_{\beta \in \mathbb{R}^n} \sup_{Q \in B_\varepsilon(\hat{P}_N)} \mathbb{E}_{(x,y) \sim Q} [\ell_\beta(x,y)],$$

where

$$\ell_\beta(x,y) = \log(1 + \exp(-y \beta^T x)),$$

$$B_\varepsilon(\hat{P}_N) = \{Q \in \mathcal{M}(\Xi) : W(Q, \hat{P}_N) \leq \varepsilon\}, \quad \hat{P}_N(\xi) = \frac{1}{N} \sum_{i=1}^N \delta_{(\hat{x}_i, \hat{y}_i)}(\xi),$$

$$W(Q, P) = \inf_{\pi \in \mathcal{M}(\Xi \times \Xi)} \left\{ \int_{\Xi \times \Xi} d(\xi, \xi') \pi(d\xi, d\xi') : \pi(\Xi, d\xi') = P(d\xi'), \pi(d\xi, \Xi) = Q(d\xi) \right\},$$

$$\Xi = \mathbb{R}^n \times \{-1, +1\}, \quad \xi = (x, y), \quad d(\xi, \xi') = \|x - x'\| + \frac{\kappa}{2} |y - y'|$$

As before, consider the inner sup problem:

$$\begin{aligned} \sup_{Q \in B_\varepsilon(\hat{P}_N)} \mathbb{E}_{(x,y) \sim Q} [\ell_\beta(x,y)] &= \sup_{\pi \in \mathcal{M}(\Xi \times \Xi)} \int_{\Xi} \ell_\beta(\xi) \underbrace{\pi(d\xi, \Xi)}_{Q(d\xi)} \\ \text{s.t. } \int_{\Xi \times \Xi} d(\xi, \xi') \pi(d\xi, d\xi') &\leq \varepsilon \\ \pi(\Xi, d\xi') &= \hat{P}_N(d\xi') \end{aligned}$$

\* Since  $\hat{P}_N$  is discrete, we have

$$\begin{aligned} Q(d\xi) &= \pi(d\xi, \Xi) = \int_{\xi' \in \Xi} \pi(d\xi, d\xi') \\ &= \sum_{i=1}^N \pi(d\xi | \xi' = (\hat{x}_i, \hat{y}_i)) \cdot \hat{P}_N(\hat{x}_i, \hat{y}_i) \\ &= \frac{1}{N} \sum_{i=1}^N Q^i(d\xi) \quad \uparrow \text{conditional distribution of } \xi \text{ given } \xi' = (\hat{x}_i, \hat{y}_i) \end{aligned}$$

Similarly,

$$\pi(d\xi, d\xi') = \pi(d\xi | \xi') \cdot \hat{P}_N(\xi') = \frac{1}{N} \sum_{i=1}^N \delta_{(\hat{x}_i, \hat{y}_i)}(\xi') Q^i(d\xi)$$

Hence,

$$\sup_{\pi \in \mathcal{M}(\Xi \times \Xi)} \int_{\Xi} \ell_\beta(\xi) \underbrace{\pi(d\xi, \Xi)}_{Q(d\xi)} \quad \sup_{Q^i \geq 0} \frac{1}{N} \sum_{i=1}^N \int_{\Xi} \ell_\beta(\xi) Q^i(d\xi)$$

$$\begin{aligned}
& \sup_{\pi \in \mathcal{M}(\Xi \times \Xi)} \int_{\Xi} \ell_p(\xi) \underbrace{\pi(d\xi, \Xi)}_{Q(d\xi)} \\
& \text{s.t. } \int_{\Xi \times \Xi} d(\xi, \xi') \pi(d\xi, d\xi') \leq \varepsilon \\
& \quad \pi(\Xi, d\xi') = \hat{\mathbb{P}}_n(d\xi')
\end{aligned}
=
\begin{aligned}
& \sup_{Q^i \geq 0} \frac{1}{N} \sum_{i=1}^N \int_{\Xi} \ell_p(\xi) Q^i(d\xi) \\
& \text{s.t. } \frac{1}{N} \sum_{i=1}^N \int_{\Xi} d(\xi, (\hat{x}_i, \hat{y}_i)) Q^i(d\xi) \leq \varepsilon, \\
& \quad \int_{\Xi} Q^i(d\xi) = 1.
\end{aligned}$$

Now, recall  $\xi \in \mathbb{R}^n \times \{-1, +1\}$ . We decompose

$$Q^i(d\xi) = Q_{+1}^i(dx) + Q_{-1}^i(dx) \quad \text{where } Q_{+1}^i(dx) = Q(dx, y=1)$$

Then, we rewrite

$$\begin{aligned}
& \sup_{Q^i \geq 0} \frac{1}{N} \sum_{i=1}^N \int_{\Xi} \ell_p(\xi) Q^i(d\xi) \\
& \text{s.t. } \frac{1}{N} \sum_{i=1}^N \int_{\Xi} d(\xi, (\hat{x}_i, \hat{y}_i)) Q^i(d\xi) \leq \varepsilon, \\
& \quad \int_{\Xi} Q^i(d\xi) = 1.
\end{aligned}
=
\begin{aligned}
& \sup_{Q^i \geq 0} \frac{1}{N} \sum_{i=1}^N \int_{\mathbb{R}^n} \ell_p(x, +1) Q_{+1}^i(dx) + \ell_p(x, -1) Q_{-1}^i(dx) \\
& \text{s.t. } \frac{1}{N} \sum_{i=1}^N \left[ \int_{\mathbb{R}^n} d((x, +1), (\hat{x}_i, \hat{y}_i)) Q_{+1}^i(dx) \right. \\
& \quad \left. + \int_{\mathbb{R}^n} d((x, -1), (\hat{x}_i, \hat{y}_i)) Q_{-1}^i(dx) \right] \leq \varepsilon, \\
& \quad \int_{\Xi} Q_{+1}^i(dx) + Q_{-1}^i(dx) = 1
\end{aligned}$$

$$\begin{aligned}
& \sup_{Q^i \geq 0} \frac{1}{N} \sum_{i=1}^N \int_{\mathbb{R}^n} \ell_p(x, +1) Q_{+1}^i(dx) + \ell_p(x, -1) Q_{-1}^i(dx) \\
& \text{s.t. } \frac{1}{N} \left[ \sum_{i: \hat{y}_i = +1} \int_{\mathbb{R}^n} (\|x - \hat{x}_i\| Q_{+1}^i(dx) + (\|x - \hat{x}_i\| + x) Q_{-1}^i(dx) \right. \\
& \quad \left. + \sum_{i: \hat{y}_i = -1} \int_{\mathbb{R}^n} (\|x - \hat{x}_i\| + x) Q_{+1}^i(dx) + \|x - \hat{x}_i\| Q_{-1}^i(dx) \right] \leq \varepsilon, \\
& \quad \int_{\Xi} Q_{+1}^i(dx) + Q_{-1}^i(dx) = 1
\end{aligned}$$

$$\begin{aligned}
& \sup_{Q^i \geq 0} \frac{1}{N} \sum_{i=1}^N \int_{\mathbb{R}^n} \ell_p(x, +1) Q_{+1}^i(dx) + \ell_p(x, -1) Q_{-1}^i(dx) \\
& \text{s.t. } \frac{1}{N} \int_{\mathbb{R}^n} k \sum_{i: \hat{y}_i = +1} Q_{-1}^i(dx) + k \sum_{i: \hat{y}_i = -1} Q_{+1}^i(dx) \\
& \quad + \sum_{i=1}^N \|x - \hat{x}_i\| (Q_{+1}^i(dx) + Q_{-1}^i(dx)) \leq \varepsilon, \quad (2)
\end{aligned}$$

$$\int_{\mathbb{R}^n} Q_{+1}^i(dx) + Q_{-1}^i(dx) = 1. \quad (s_i)$$

Now, take the dual of (A) to get (Exercise)

$$\begin{aligned} \inf_{\lambda, s_i} \quad & \lambda \varepsilon + \frac{1}{N} \sum_{i=1}^N s_i \\ \text{s.t.} \quad & \sup_{x \in \mathbb{R}^n} l_{\beta}(x, +1) - \lambda \|x - \hat{x}_i\| - \frac{1}{2} \lambda \kappa (1 - \hat{y}_i) \leq s_i, \quad \text{--- (D)} \\ & \sup_{x \in \mathbb{R}^n} l_{\beta}(x, -1) - \lambda \|x - \hat{x}_i\| - \frac{1}{2} \lambda \kappa (1 + \hat{y}_i) \leq s_i, \\ & \lambda \geq 0. \end{aligned}$$

Exercise: Verify that strong duality holds for any  $\varepsilon > 0$

Claim: For every  $\lambda > 0$ ,

$$\sup_{x \in \mathbb{R}^n} l_{\beta}(x, \pm 1) - \lambda \|x - x'\| = \begin{cases} l_{\beta}(x', \pm 1) & \text{if } \|\beta\|_* \leq \lambda, \\ -\infty & \text{o/w} \end{cases}$$

(Exercise, use conjugate functions)

It then follows that (D) is equivalent to

$$\begin{aligned} \inf_{\lambda, s_i} \quad & \lambda \varepsilon + \frac{1}{N} \sum_{i=1}^N s_i & \inf_{\lambda, s_i} \quad & \lambda \varepsilon + \frac{1}{N} \sum_{i=1}^N s_i \\ \text{s.t.} \quad & l_{\beta}(\hat{x}_i, +1) - \frac{1}{2} \lambda \kappa (1 - \hat{y}_i) \leq s_i, & = \text{s.t.} \quad & l_{\beta}(\hat{x}_i, \hat{y}_i) \leq s_i, \\ & l_{\beta}(\hat{x}_i, -1) - \frac{1}{2} \lambda \kappa (1 + \hat{y}_i) \leq s_i, & & l_{\beta}(\hat{x}_i, \hat{y}_i) - \lambda \kappa \leq s_i, \\ & \|\beta\|_* \leq \lambda. & & \|\beta\|_* \leq \lambda. \end{aligned}$$

\* Putting back the original  $\inf_{\beta \in \mathbb{R}^n}$ , we obtain a convex optimization problem in the variables  $(\beta, \lambda, s_i)$ .

\* Suppose that we take  $\kappa = +\infty$  (all labels are deterministic). Then, we get

$$\inf_{\beta \in \mathbb{R}^n} \frac{1}{N} \sum_{i=1}^N l_{\beta}(\hat{x}_i, \hat{y}_i) + \varepsilon \|\beta\|_*$$

This provides a DRO interpretation of regularization

### References:

- 1) Shafieezadeh-Abadeh, Mohajerin Esfahani, Kuhn. Distributionally Robust Logistic Regression. NIPS 2015
- 2) Shafieezadeh-Abadeh, Kuhn, Mohajerin Esfahani: Regularization via Mass Transportation. JMLR 2019.

### Further Reading

- 1) Kuhn et al.: Wasserstein Distributionally Robust Optimization: Theory and Applications in Machine Learning. INFORMS Tutorials in Operations Research, 2019
- 2) Rahimian, Mehrotra: Distributionally Robust Optimization: A Review. arXiv 2019.