

Division of Operations Research, Osaka Univ.

---

# Data-Driven Distributionally Robust Optimization

Zhicheng Zhang, D 1  
Fujisaki Lab Seminar Camp

September 27, 2021

## Robust Optimization

A robust optimization problem can be written as

$$\begin{aligned} \inf_{x \in \mathbb{X}} \quad & c^\top x \\ \text{s.t.} \quad & f(x, \xi) > 0, \quad \forall \xi \in \Xi \end{aligned} \tag{1}$$

where  $x \in \mathbb{R}^n$  is the decision variable,  $c$  is the given objective direction,  $f(x, \xi) : \mathbb{R}^n \times \Xi \mapsto \mathbb{R}$  defining the design constraints and parameterized by uncertainty instances  $\xi \in \Xi$ ;

- Semi-infinite optimization can be solvable under some regularity conditions on uncertainty set and constraints;
- In general, it is intractable; Solution is overly conservatism;
- Take probabilistic (chance) relaxation on constraints.

## Chance-constrained Optimization

The general chance-constrained optimization is given by,

$$\begin{aligned} \text{(CCO)} \quad & \inf_{x \in \mathbb{X}} \quad c^\top x \\ & \text{s.t.} \quad \mathbb{P}\{\xi \in \Xi \mid f(x, \xi) > 0\} \leq \epsilon \end{aligned} \quad (2)$$

where  $\mathbb{P}$  is a (known) probability distribution on  $\Xi$ , and  $\epsilon \in (0, 1)$  is the violation of tolerance level.

- Min-max form of CCO: minimize the max cost  $\ell(x, \xi)$  with max taken over a reduced uncertainty set  $\Xi_\epsilon \subset \Xi$  having probability  $\mathbb{P}\{\Xi_\epsilon\} = 1 - \epsilon$ , namely,

$$\inf_{x \in \mathbb{X}} \sup_{\xi \in \Xi_\epsilon} \ell(x, \xi) \quad (3)$$

- The probability of violation (**risk**) is defined as:  
 $V(x) = \mathbb{P}\{\xi \in \Xi \mid f(x, \xi) > 0\} \Rightarrow V(x) \leq \epsilon$  (Robustness)

## Stochastic Optimization

Stochastic optimization is shown as follows

$$\inf_{x \in \mathbb{X}} \mathbb{E}_{\mathbb{P}} [\ell(x, \xi)] \quad (4)$$

The objective is to find a data-driven solution  $\hat{x}_N$  of (4), constructed using the dataset  $\hat{\Xi} = \{\hat{\xi}_i\}_{i=1}^N \subset \Xi$ , that has a *finite-sample guarantee* given by

$$\mathbb{P}^N \left\{ \mathbb{E}_{\mathbb{P}} [\ell(x, \xi)] \leq \hat{J}_N \right\} \geq 1 - \beta \quad (5)$$

where  $\hat{J}_N$  might depend on  $\hat{\Xi}$  and  $\beta \in (0, 1)$  is the parameter governing  $\hat{x}_N$  and  $\hat{J}_N$ . In order to identify the  $\hat{x}_N$  with low  $\hat{J}_N$  and  $\beta \in (0, 1)$ , the strategy is to design an ambiguity set  $\mathcal{P}$ .

## Distributionally Robust Optimization

The distributionally robust optimization aims to minimize the *worst-case* expected cost according to the well-defined ambiguity set containing all distribution measures, that is

$$\begin{aligned} \hat{J}_N : \inf_{x \in \mathbb{X}} \sup_{\mathbb{P} \in \mathcal{P}} \mathbb{E}_{\mathbb{P}} [\ell(x, \xi)] \\ \iff \inf_{\lambda \geq 0, x \in \mathbb{X}} \left\{ \lambda \gamma^2 + \frac{1}{N} \sum_{i=1}^N \max_{\xi \in \Xi} \left( \ell(x, \xi) - \lambda \|\xi - \hat{\xi}_i\|^2 \right) \right\} \end{aligned} \quad (6)$$

- In practice, the “true” probability distribution of uncertain model parameters/data may not be known
- Considers a set of probability distributions (ambiguity set).
- Determines decisions that provide hedging against the worst-case distribution by solving a min-max problem.
- An intermediate approach between stochastic programming and traditional robust optimization

## Moment-based ambiguity set

- Typically do not contain the true distribution.
- Conservative solutions: very different distributions can have the same lower moments and the use of higher moments can be impractical
- May be not guaranteed to converge to the true probability distributions the number of uncertain data tends to infinity.

$$\mathcal{P} \doteq \left\{ \mathbb{P} \in \mathcal{M}(\Xi) \left| \begin{array}{l} \mathbb{P}[\xi \in \Xi] = 1 \\ \mathbb{E}_{\mathbb{P}}[\xi] = \mu \\ \mathbb{E}_{\mathbb{P}}[(\xi - \mu)(\xi - \mu)^T] = \Sigma \end{array} \right. \right\}$$

$$\iff \int_{\Xi} \begin{bmatrix} \xi \\ 1 \end{bmatrix} \begin{bmatrix} \xi \\ 1 \end{bmatrix}^T d\mathbb{P}(\xi) = \begin{bmatrix} \Sigma + \mu\mu^T & \mu \\ \mu^T & 1 \end{bmatrix}$$

Strong duality condition ( $\inf \sup = \sup \inf$ )  $\Rightarrow$  SDP (Solvable)

## Statistical or probabilistic metric-based ambiguity set

Given a radius  $\gamma \geq 0$ . Define

$$\mathcal{P} \doteq \{\mathbb{P} \in \mathcal{M}(\Xi) \mid \rho(\mathbb{P}, \mathbb{P}_0) \leq \gamma\} \quad (7)$$

where  $\mathbb{P}_0$  is the *reference distribution*,  $\rho(\cdot, \cdot)$  denotes the some probabilistic distance between two distributions.

E.g. Let  $\mathbb{Q}_1$  and  $\mathbb{Q}_2$  be two probability distributions over a space  $\mathcal{M}$  such that  $\mathbb{Q}_1$  is *absolutely continuous* with respect to  $\mathbb{Q}_2$ . Then, for a convex function  $\phi$  such that  $\phi(1) = 0$ , the  $\phi$ -divergence of  $\mathbb{Q}_1$  from  $\mathbb{Q}_2$  is defined as

$$D_\phi(\mathbb{Q}_1 \parallel \mathbb{Q}_2) = \int_{\mathcal{M}} \phi\left(\frac{d\mathbb{Q}_1}{d\mathbb{Q}_2}\right) d\mathbb{Q}_2, \Rightarrow \sum_i \xi_{2,i} \phi\left(\frac{\xi_{1,i}}{\xi_{2,i}}\right)$$

where  $\frac{d\mathbb{Q}_1}{d\mathbb{Q}_2}$  is the Radon–Nikodym derivative of  $\mathbb{Q}_1$  w.r.t.  $\mathbb{Q}_2$ .

- $\phi$ -divergence is asymmetric, i.e.,  $D_\phi(\mathbb{Q}_1 \parallel \mathbb{Q}_2) \neq D_\phi(\mathbb{Q}_2 \parallel \mathbb{Q}_1)$
- Kullback-Leibler divergence ( $\phi_{KL} = t \log t$ ), total variation ( $\phi_{TV} = \frac{1}{2}|t - 1|$ ), etc.

# Wasserstein Metric (Earth Mover's Distance)

The ( $p$ -) Wasserstein metric  $\rho_w : \mathcal{M}(\Xi) \times \mathcal{M}(\Xi) \mapsto \mathbb{R}_+$  :

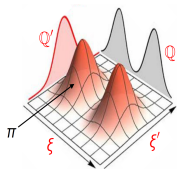
$$\rho_w(Q_1, Q_2) = \inf_{\Pi \in \mathcal{M}(\Xi^2)} \left\{ \int_{\Xi^2} \|\xi_1 - \xi_2\| \Pi(d\xi_1, d\xi_2) \mid \Pi \in \mathcal{H}(Q_1, Q_2) \right\}$$

where  $\mathcal{H}(Q_1, Q_2)$  is the set of all distributions on  $\Xi \times \Xi$  with marginals  $Q_1$  and  $Q_2$ , and  $\|\cdot\|$  represents an *arbitrary* norm.

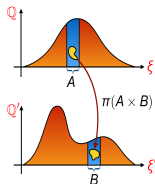
Given  $\gamma > 0$ , denote  $\mathcal{B}_\gamma(\hat{\mathbb{P}}_N) := \{Q \in \mathcal{M}(\Xi) \mid \rho_w(\hat{\mathbb{P}}_N, Q) \leq \gamma\}$ , which can be viewed as the Wasserstein ball of radius  $\gamma$  centered at the *empirical distribution*  $\hat{\mathbb{P}}_N := \frac{1}{N} \sum_{k=1}^N \delta_{\hat{\xi}_i}$  where  $\delta_{\hat{\xi}_i}$  is the unit point mass at  $\hat{\xi}_i$ .  $\Rightarrow \hat{\mathbb{P}}_N(\hat{\xi}_i) = \frac{1}{N}$

- How to guarantee  $\mathbb{P}^N \{\rho(\mathbb{P}, \mathbb{P}_0) > \gamma\} \leq \delta$  ?  
And  $\mathbb{P}_0$  can choose empirical distribution  $\hat{\mathbb{P}}_N$
- Convergence:  $\mathbb{P}^\infty \{\lim_{N \rightarrow \infty} \rho_w(\mathbb{P}, \hat{\mathbb{P}}_N) = 0\} = 1$ .



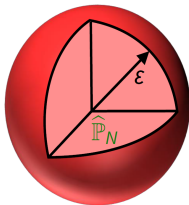


$\Pi(Q, Q') =$  set of couplings with marginals  $Q$  and  $Q'$



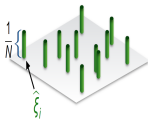
$\pi(A \times B) = \begin{cases} \text{mass moved from} \\ \text{source region } A \text{ to} \\ \text{target region } B \end{cases}$

$\|\xi - \xi'\|^p = \begin{cases} \text{price paid for moving} \\ \text{mass from } \xi \text{ to } \xi' \end{cases}$



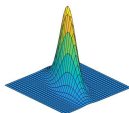
Contains every  $Q$  obtainable by re-shaping  $\hat{P}_N$  at a cost of at most  $\epsilon$

Non-parametric estimators:



Empirical distribution:  $\hat{P}_N = \frac{1}{N} \sum_{i=1}^N \delta_{\hat{\xi}_i}$

Parametric estimators:



Elliptical distribution:  $\hat{P}_N = \mathcal{E}_g(\hat{\mu}_N, \hat{\Sigma}_N)$

density generator

Density function:  $f(\xi) = C \det(\hat{\Sigma}_N)^{-1} g((\xi - \hat{\mu}_N) \hat{\Sigma}_N^{-1} (\xi - \hat{\mu}_N))$

## Scenario Program (SP)

One considers  $N$  i.i.d. random samples of the uncertainty  $\{\xi_i\}_{i=1}^N$ , and builds a *scenario program* (SP):

$$\begin{aligned}
 (SP) \quad J_N^* &= \inf_{x \in \mathbb{X}} \quad c^\top x \\
 \text{s.t.} \quad & f(x, \xi_i) \leq 0, \quad i = 1, \dots, N \\
 J_j^* \Leftarrow & \text{s.t.} \quad f(x, \xi_i) \leq 0, \quad i = \{1, \dots, N\} \setminus j
 \end{aligned}$$

An optimal solution  $x_N^*$  to this problem, if it exists, is a random variable which depends on the multiextraction  $\Xi^N$ .

- Min-max form of SP:  $\inf_{x \in \mathbb{X}} \sup_{i=1, \dots, N} \ell(x, \xi_i), \quad \xi_i \in \Xi$
- Support constraints:  $J_j^* < J_N^*$  (its removal changes solution)
- Goal: Find a data-driven solution  $x_N^* \in \mathbb{X}$ , such that hold a guarantee  $\mathbb{P}^N \{V(x_N^*) < \epsilon\} \geq 1 - \beta$

## A Priori Guarantee (*before* obtaining $x_N^*$ )

Given  $N \geq n$ , a prior guarantees for SP is

$$\mathbb{P}^N\{V(x_N^*) > \epsilon\} \leq \sum_{i=1}^{n-1} \binom{N}{i} \epsilon^i (1 - \epsilon)^{N-i} \leq \beta \quad (8)$$

- Fully-supported problems (the cardinality of support scenarios set  $\mathcal{S}$  is exactly  $n$ , i.e.,  $s_N^* = |\mathcal{S}| = n$ .)
- Helly's dimension:  $\text{ess sup}_{\xi \in \Xi^N} |\mathcal{S}(\xi)| \leq h < \bar{h}$ , the upper bounds  $\bar{h}$  on  $h$  is much easier to calculate.
- Sample size  $N$  is **sufficient** (If Large-size: sequential)

Conclude a procedure for a priori feasibility guarantee on SP:

- ① Exploring the problem structure of SP and obtain  $|\mathcal{S}| = s_N^*$ ;
- ② Select the sample complexity  $N(\epsilon, \beta, s_N^*)$  using (8);
- ③ Obtain optimal solution  $x_N^*$  and optimal objective value  $J_N^*$ .

If the resulting risk  $V(\beta, s_N^*, N) > \epsilon$ , we repeat this process with more scenarios until reaching  $\epsilon(\beta, s_N^*, N) \leq \epsilon$ . If the number of available scenarios is *limited*, then it might be *impossible* to obtain a solution  $x_N^*$  such that  $V(x_N^*) \leq \epsilon$ . The objective is

$$\mathbb{P}^N \{ V(x_N^*) < \epsilon(s_N^*) \} \geq 1 - \beta. \quad (9)$$

### A Posteriori Guarantee (*after* obtaining $x_N^*$ )

A *wait-and-judge* method to give a *a-posteriori* guarantee on sample complexity, which is based on a polynomial equation in variable  $t$ , for any  $k = 1, 2, \dots, n$ ,

$$\frac{\beta}{N+1} \sum_{i=1}^{n-1} \binom{i}{k} t^{i-k} \binom{N}{k} t^{N-i} = 0. \quad (10)$$

has exactly one solution  $\epsilon(k) \in (0, 1)$ .

- It is *not* fully-support thus difficult to calculate *a priori* bounds on number of support scenarios ( $s_N^* = |\mathcal{S}| < ?$ ).
- Sampling point is *insufficient*, it is difficult to meet the sample complexity from the a-priori guarantees.

SP is also related to the *sample average approximation* (SAA). A *sampling-and-discarding* method is that draw  $N$  scenarios and discard any  $k$  of them, then use the SP with remaining  $N - k$  samples, and the associated solution is denoted as  $x_{N,k}^*$ .

This removal method can be any algorithm that

$$\mathcal{A}\{\xi_1, \dots, \xi_N\} = \{i_1, \dots, i_k\}$$

of the  $k$  indexes of the  $k$  discarded constraints.

The objective is to guarantee that

$$\mathbb{P}^N\{V(x_{N,k}^*) < \epsilon\} \geq 1 - \beta. \quad (11)$$

## Sampling and Discarding (Non-degeneracy)

Given parameters  $N, \epsilon$  and  $\beta$ , and find the largest removal constraints  $k$  such that

$$\binom{k+n-1}{k} \sum_{i=1}^{k+n-1} \binom{N}{i} \epsilon^i (1-\epsilon)^{N-i} \leq \beta \quad (12)$$

holds, then the solution to SAA with  $\epsilon = N/k$  is feasible to chance constrained optimization with probability at least  $1 - \beta$ .

*W.p. 1, the SP has a unique irreducible support subsample<sup>1</sup>, consisting precisely of support constraints.*

## Generalization

For an algorithm  $\mathcal{A}_N$ , let  $\beta \in (0, 1)$  and  $\epsilon : \{0, 1, \dots, N\} \mapsto [0, 1]$  be a function such that

$$\sum_{k=0}^N \binom{N}{k} (1 - \epsilon(k))^{N-k} = \beta, \quad \epsilon(N) = 1. \quad (13)$$

Then, for any  $\mathcal{A}_N$ ,  $\mathcal{G}_N$ , and probability  $\mathbb{P}$ , it holds that

$$\mathbb{P}^N \{V(x_N^*) > \epsilon(s_N^*)\} \leq \beta. \quad (14)$$

<sup>1</sup>A support subsample  $S = (\xi_{i_1}, \dots, \xi_{i_k})$  with  $i_1 < i_2 < \dots < i_k$  for  $(\xi_1, \dots, \xi_N)$  is a  $k$ -tuple of elements extracted from  $(\xi_1, \dots, \xi_N)$ , which yields the same solution as the full sample, that is,  $\mathcal{A}_k(\xi_{i_1}, \dots, \xi_{i_k}) = \mathcal{A}_N(\xi_1, \dots, \xi_N)$ . A support subsample is said to be *irreducible* if no element can be further removed from  $S$  without changing the solution. Meanwhile, suppose that

This motivates us to approximate the distributionally CCO

$$\begin{aligned}
 & \inf_{x \in \mathbb{X}} c^\top x \\
 \text{s.t.} \quad & \sup_{\mathbb{P} \in \mathcal{P}} \mathbb{P}\{\xi \in \Xi \mid f(x, \xi) > 0\} \leq \epsilon
 \end{aligned} \tag{15}$$

through SP, we call this as DRO via SP, and it can be defined as

$$\begin{aligned}
 & \inf_{x \in \mathbb{X}} c^\top x \\
 \text{s.t.} \quad & f(x, q) > 0, \\
 & \forall q \quad \|q - \xi_i\| \leq \gamma, \quad \forall \xi_i \in \Xi, \quad i = 1, \dots, N
 \end{aligned} \tag{16}$$

where  $\xi_i, i = 1, \dots, N$  is i.i.d. samples drawn according to the central measure  $\mathbb{P}_0$  and the *arbitrary* norm on the space  $\mathcal{M}(\Xi)$  is related to the probability metric  $\rho(\cdot, \cdot)$ .

In min-max case, it is written as

$$\inf_{x \in \mathbb{X}} \sup_{i=1, \dots, N} \ell(x, \xi_i) \quad (17)$$

where  $\{\xi_i\}_i^N$  is an i.i.d. sequence of scenarios randomly sampled from a reference distribution  $\mathbb{P}_0$ , and the true distribution  $\mathbb{P} \in \mathcal{B}_\gamma(\mathbb{P}_0) := \{\mathbb{P} \in \mathcal{M}(\Xi) \mid \rho(\mathbb{P}, \mathbb{P}_0) \leq \gamma\}$ .

Based on finite-sample guarantee,

$$\mathbb{P}^N \{\rho_w(\mathbb{P}, \hat{\mathbb{P}}_N) \leq \gamma_N\} \geq 1 - \delta_N$$

consider the relation between

$$J_N^{SP} = \inf_{x \in \mathbb{X}} \sup_{i=1, \dots, N} \ell(x, \xi_i)$$

$$J^{CCO} = \inf_{x \in \mathbb{X}} \sup_{\xi \in \Xi_\epsilon} \ell(x, \xi)$$

$$J_N^{DRCCO} = \inf_{x \in \mathbb{X}} \sup_{\mathbb{P} \in \mathcal{P}} \ell(x, \xi), \quad \mathcal{P} = \mathcal{B}_{\gamma_N(\delta_N)}(\hat{\mathbb{P}}_N)$$