

Erick's notes

I wish to warmly thank
Daniel Kuhn for sharing
this keynote
presentation

Data-Driven Distributionally Robust Optimization using the Wasserstein Metric

Daniel Kuhn

Risk Analytics and Optimization Chair
École Polytechnique Fédérale de Lausanne
rao.epfl.ch



SWISS NATIONAL SCIENCE FOUNDATION



ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

Stochastic Programming

$$\text{SP} : \min_{x \in \mathcal{X}} \mathbb{E}^{\mathbb{P}} [\ell(x, \xi)]$$

$$J^* = \min \text{SP}$$

$$x^* = \operatorname{argmin} \text{SP}$$

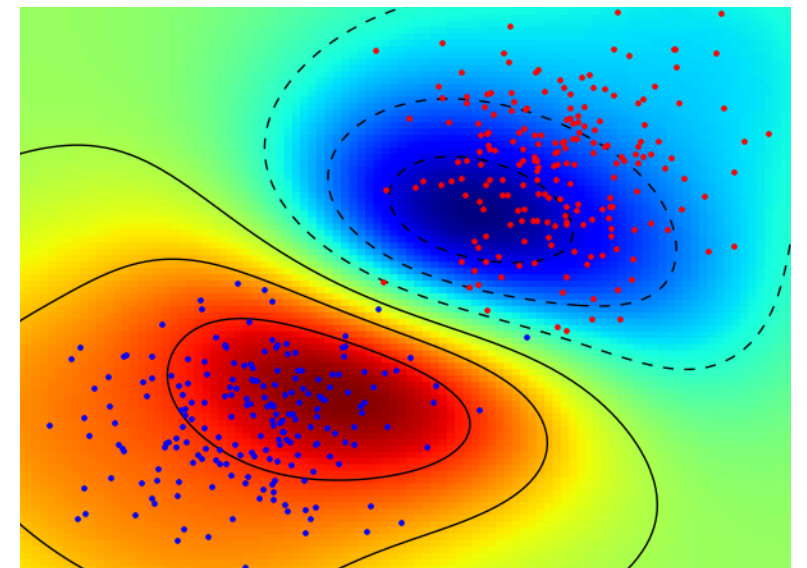
Applications:



Supply Chain Mgmt.

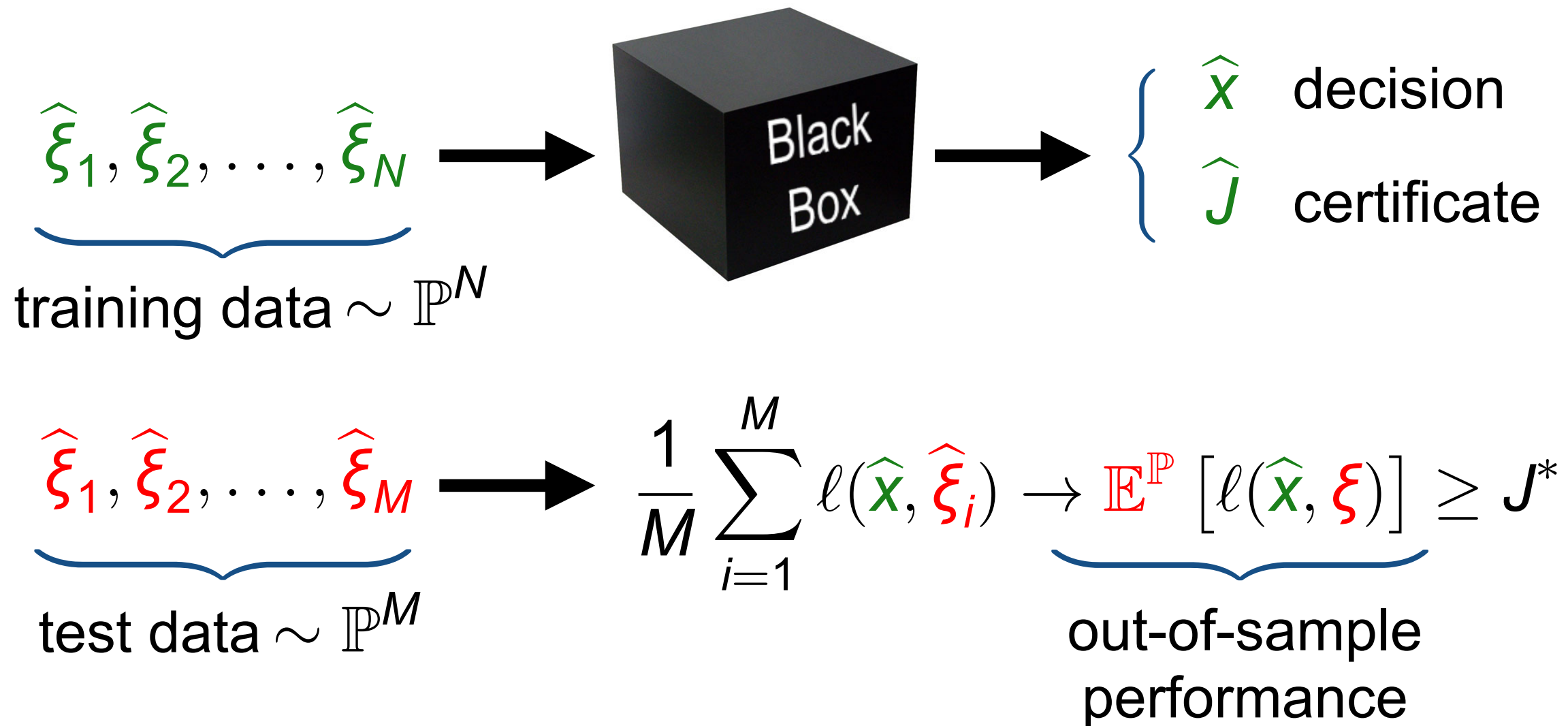


Portfolio Mgmt.



Machine Learning

Data-Driven Stochastic Programming

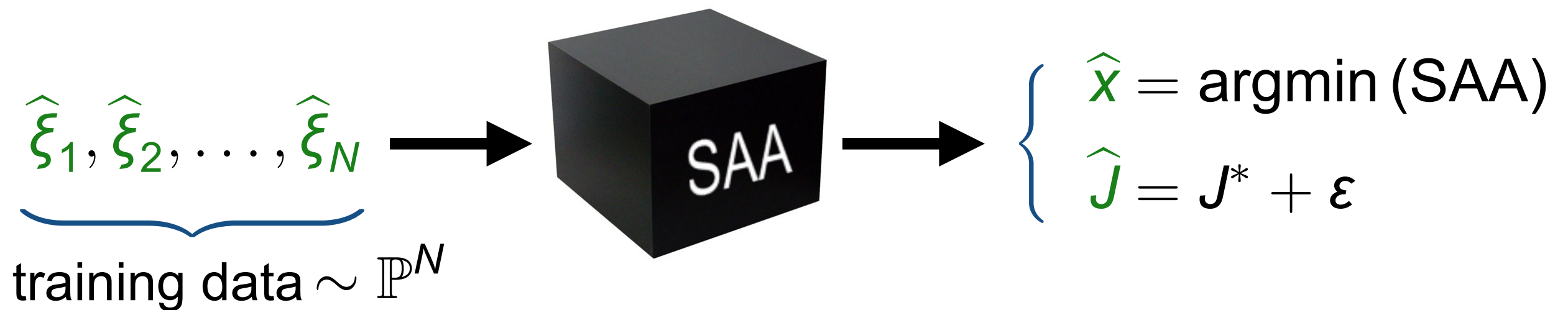


Aim: establish a finite sample guarantee

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

Sample Average Approximation (SAA)

$$\text{SAA : } \min_{\mathbf{x} \in \mathcal{X}} \frac{1}{N} \sum_{i=1}^N \ell(\mathbf{x}, \hat{\xi}_i)$$



Finite sample guarantee:²⁾

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

$$N \geq C \left(\frac{\text{diam}(\mathcal{X}) \text{lip}(\ell)}{\varepsilon} \right)^2 \left[\text{dim}(\mathbf{x}) \log \left(\frac{\text{diam}(\mathcal{X}) \text{lip}(\ell)}{\varepsilon} \right) + \log \left(\frac{C}{\beta} \right) \right]$$

²⁾ Shapiro & Nemirovski, *Springer*, 2005.

SAA with Scarce Data

Mean-risk portfolio problem

$$\min_{x \in \mathcal{X}} \left\{ \mathbb{E}^{\mathbb{P}} \left[-x^{\top} \xi \right] + \rho \mathbb{P}\text{-CVaR}_{\alpha}(-x^{\top} \xi) \right\}$$

▶ 10 assets

▶ $\rho = 10$

▶ $\alpha = 20\%$

▶ $\xi_i = \psi + \zeta_i$ where $\psi \sim \mathcal{N}(0, 2\%)$
and $\zeta_i \sim \mathcal{N}(i \times 3\%, i \times 2.5\%)$

Erick's notes

In our notation

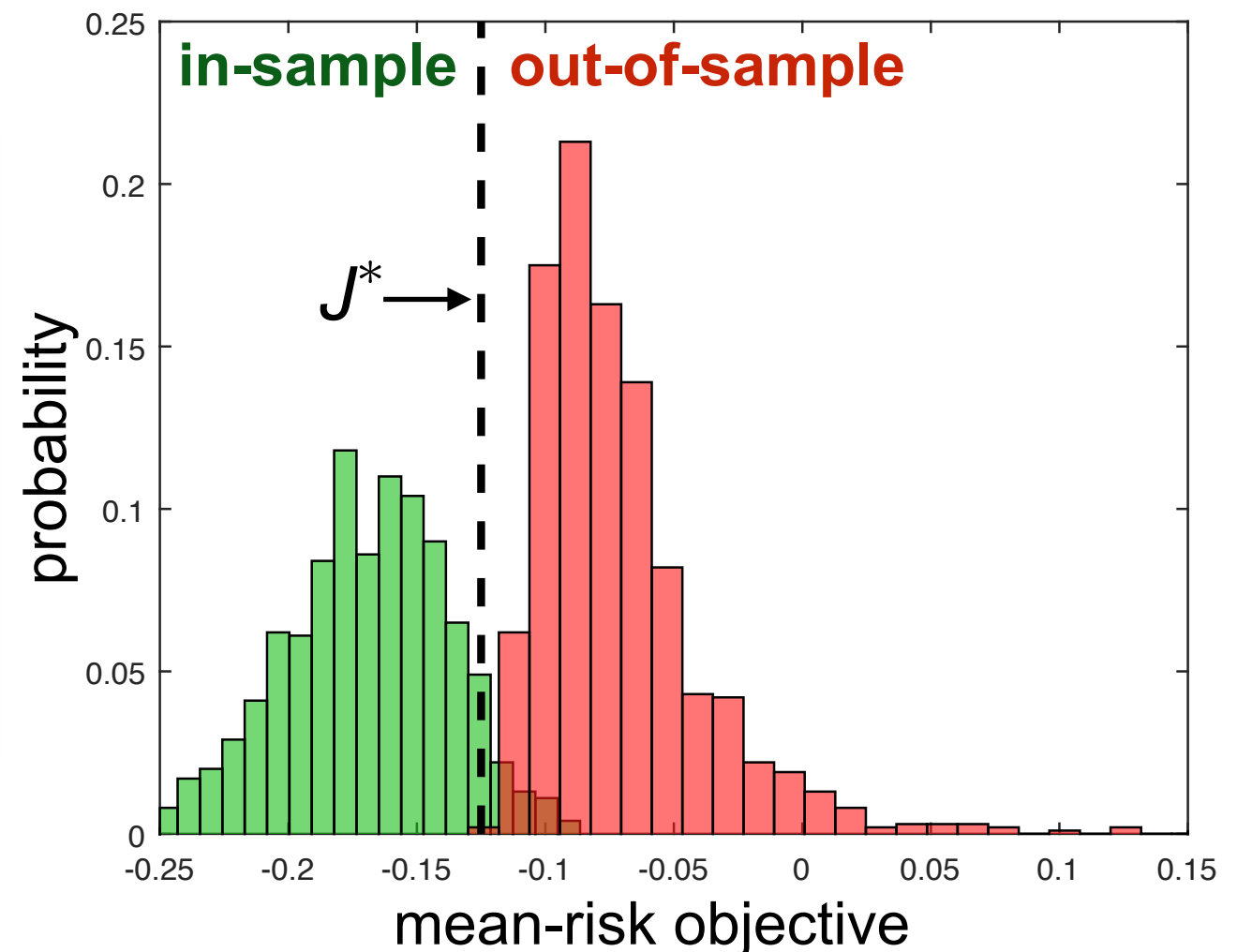
1-epsilon=80%

▶ 30 training samples

▶ in-sample: optimistic bias

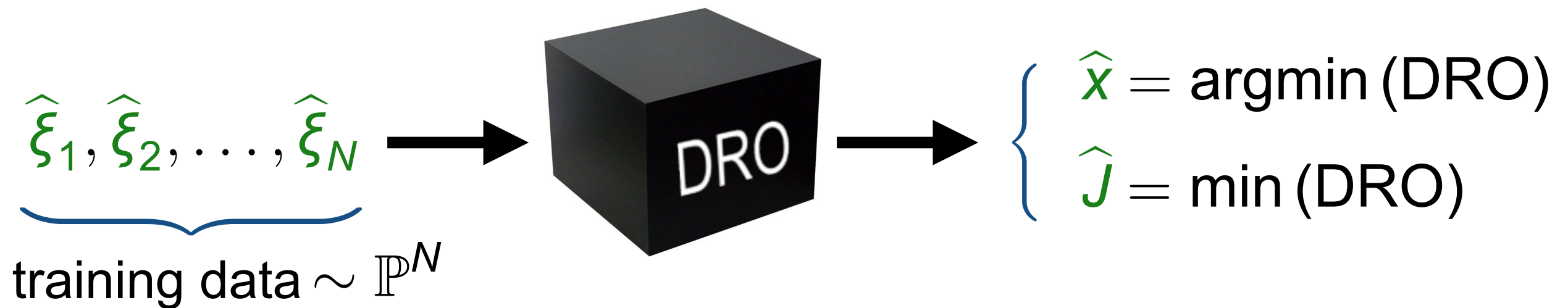
▶ out-of-sample: pessimistic bias

Performance of SAA solution



Distributionally Robust Optimization (DRO)

$$\text{DRO} : \min_{x \in \mathcal{X}} \sup_{Q \in \hat{\mathcal{P}}_N} \mathbb{E}^Q [\ell(x, \xi)]$$



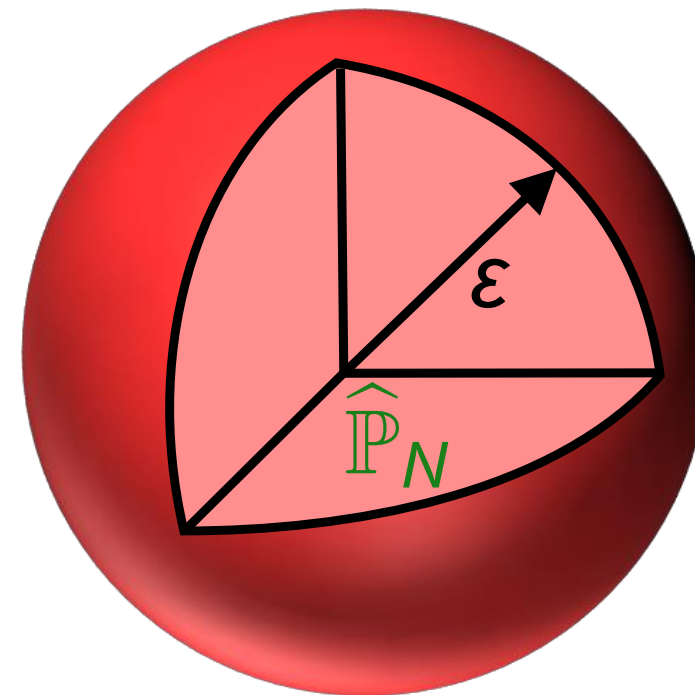
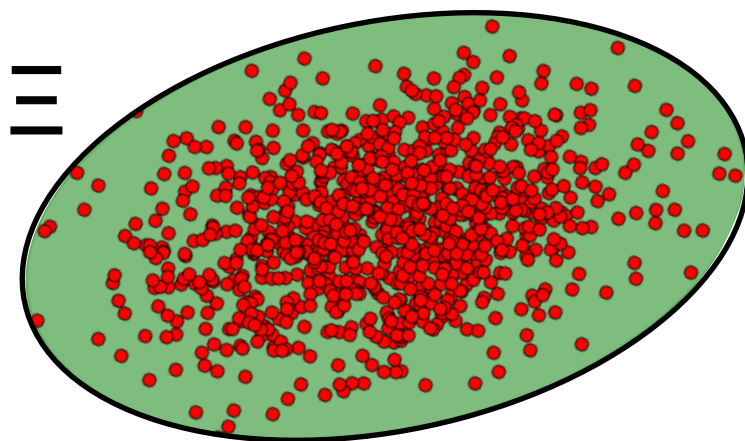
Desired properties:³⁾

- ▶ Finite sample guarantee: $\mathbb{P}^N \left\{ \mathbb{E}^{\mathbb{P}} [\ell(\hat{x}, \xi)] \leq \hat{J} \right\} \geq 1 - \beta$
- ▶ Asymptotic guarantee: $\mathbb{P}^\infty \left\{ \lim_{N \rightarrow \infty} \hat{x} = x^* \right\} = 1$
- ▶ Tractability: DRO is in the same complexity class as SAA

³⁾ Bertsimas, Gupta & Kallus, <http://arxiv.org>, 2014.

Wasserstein Ambiguity Set⁴⁾

$$\mathbb{B}_\varepsilon(\hat{\mathbb{P}}_N) = \left\{ \mathbb{Q} : \mathbb{Q}(\xi \in \Xi) = 1 \right\} \cap \left\{ \mathbb{Q} : d_W(\hat{\mathbb{P}}_N, \mathbb{Q}) \leq \varepsilon \right\}$$



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

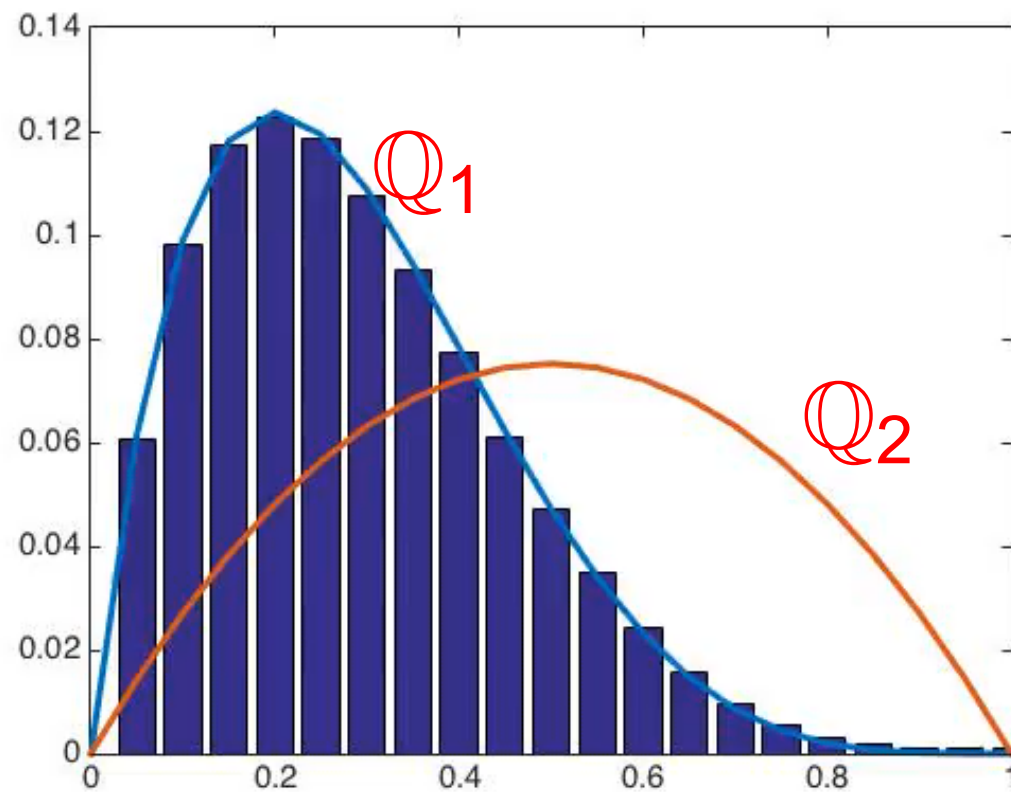
Wasserstein metric:
$$d_W(\mathbb{Q}_1, \mathbb{Q}_2) = \sup_{\text{lip}(f) \leq 1} \mathbb{E}^{\mathbb{Q}_1} [f(\xi)] - \mathbb{E}^{\mathbb{Q}_2} [f(\xi)]$$

⁴⁾ Pflug & Wozabal, *Quant. Finance*, 2007.

Kantorovich-Rubinstein Theorem⁶⁾

LP duality implies:

$$d_W(Q_1, Q_2) = \inf \left\{ \mathbb{E}^Q [\|\xi_1 - \xi_2\|] : \begin{array}{l} Q \text{ distribution of } \xi_1 \text{ and } \xi_2 \\ \text{with marginals } Q_1 \text{ and } Q_2 \end{array} \right\}$$



$d_W(Q_1, Q_2) =$ minimum cost of moving Q_1 to Q_2

⁶⁾ Kantorovich & Rubinstein, *Vestn. Lenin. U.*, 1958.

Finite-Sample Guarantee

Measure concentration theorem:⁵⁾ If the tail of \mathbb{P} decays exponentially at rate a , then

$$\mathbb{P}^N \left\{ d_W(\mathbb{P}, \hat{\mathbb{P}}_N) > \varepsilon \right\} \leq \begin{cases} C e^{-cN\varepsilon^{\dim(\xi)}} & \text{if } \varepsilon \leq 1, \\ C e^{-cN\varepsilon^a} & \text{if } \varepsilon > 1. \end{cases}$$

$$\varepsilon_N(\beta) = \begin{cases} (\log(C\beta^{-1}) / (cN))^{\frac{1}{a}} & \text{if } N < \log(C\beta^{-1}) / c \\ (\log(C\beta^{-1}) / (cN))^{\frac{1}{\dim(\xi)}} & \text{if } N \geq \log(C\beta^{-1}) / c \end{cases}$$

$$\implies \mathbb{P}^N \left\{ \mathbb{P} \in \mathbb{B}_{\varepsilon_N(\beta)}(\hat{\mathbb{P}}_N) \right\} \geq 1 - \beta$$

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

⁵⁾ Fournier & Guillin, *Probab. Theory Rel.*, 2014.

Asymptotic Guarantee

Convergence theorem: If the tail of \mathbb{P} decays exponentially and the Wasserstein radius is $\varepsilon_N(\beta_N)$ with $\beta_N \propto \exp(-N^\delta)$ for some $\delta > 0$, then we have:

- ▶ Convergence of optimal values:

$$\mathbb{P}^\infty \left\{ \lim_{N \rightarrow \infty} \hat{J}_N = J^* \right\} = 1$$

- ▶ Convergence of optimal solutions:

$$\mathbb{P}^\infty \left\{ \limsup_{N \rightarrow \infty} \{\hat{X}_N\} \subseteq \arg \min(\text{SP}) \right\} = 1$$

Convexity assumption:

- ▶ $-\ell$ is proper, convex and lsc;
- ▶ Ξ is convex and closed.

Worst-case expectation problem:

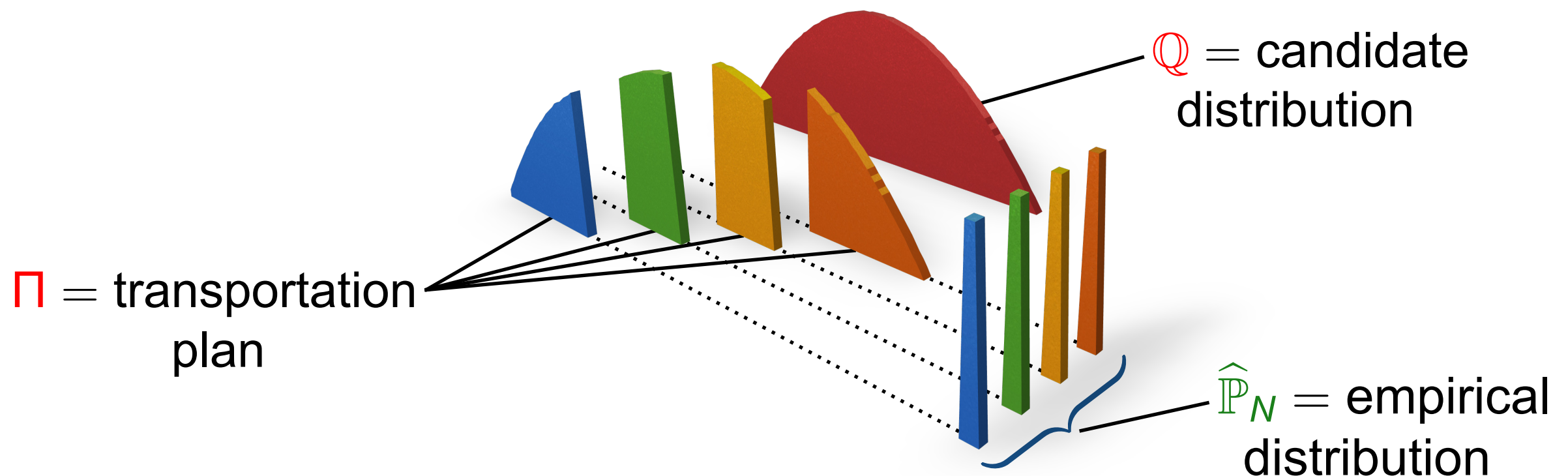
$$\sup_{Q \in \hat{\mathcal{P}}_N} \mathbb{E}^Q[l(\xi)]$$

Use the Kantorovich-Rubinstein theorem:

$$\sup_{\Pi, Q} \int_{\Xi} \ell(\xi) Q(d\xi)$$

$$\text{s.t.} \quad \int_{\Xi^2} \|\xi - \xi'\| \Pi(d\xi, d\xi') \leq \varepsilon$$

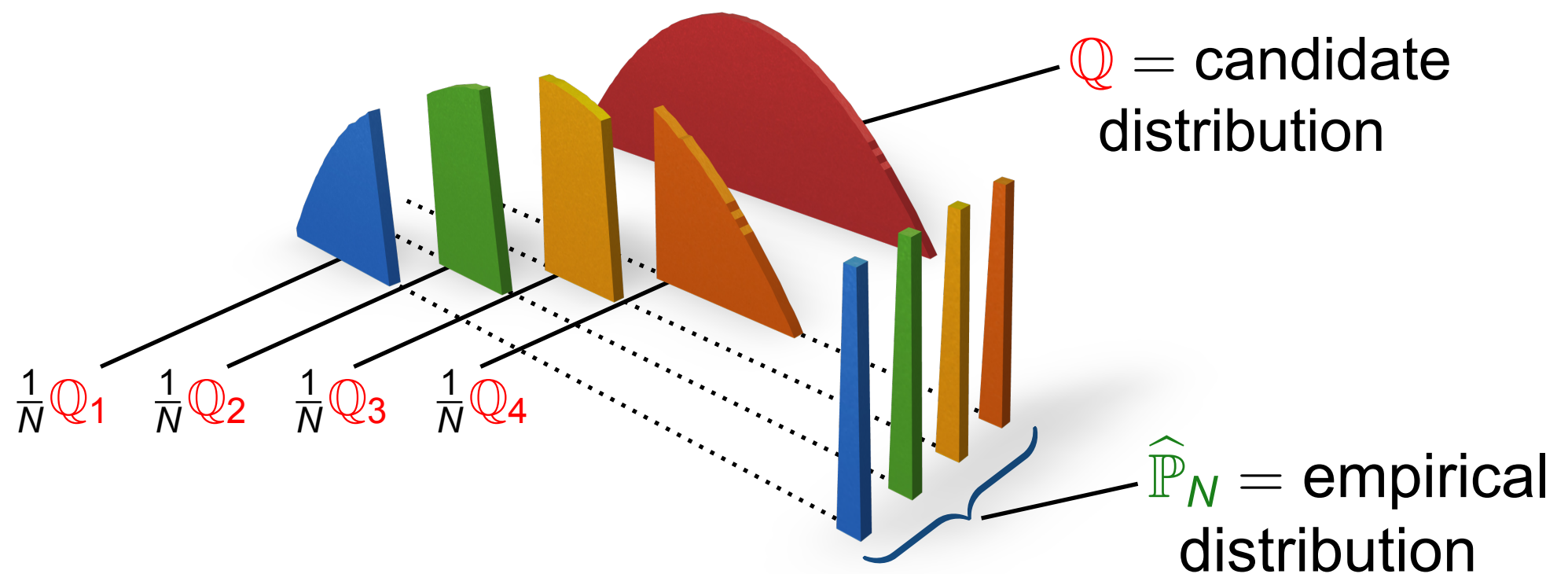
$\left\{ \begin{array}{l} \Pi \text{ is a joint distribution of } \xi \text{ and } \xi' \\ \text{with marginals } Q \text{ and } \hat{\mathbb{P}}_N, \text{ respectively} \end{array} \right.$



Tractability

Decompose Π into Q_1, \dots, Q_N :

$$\begin{aligned} \sup_{Q_i} \quad & \frac{1}{N} \sum_{i=1}^N \int_{\Xi} \ell(\xi) Q_i(d\xi) \\ \text{s.t.} \quad & \frac{1}{N} \sum_{i=1}^N \int_{\Xi} \|\xi - \hat{\xi}_i\| Q_i(d\xi) \leq \varepsilon \end{aligned}$$



Dual of the moment problem is a robust program:

$$\begin{aligned} \inf_{\lambda, s_i} \quad & \lambda \varepsilon + \frac{1}{N} \sum_{i=1}^N s_i \\ \text{s.t.} \quad & \ell(\xi) - \lambda \|\xi - \hat{\xi}_i\| \leq s_i \quad \forall \xi \in \Xi, \quad \forall i \leq N \\ & \lambda \geq 0 \end{aligned}$$

Introduce the indicator function of Ξ :

$$\begin{aligned} \inf_{\lambda, s_j} \quad & \lambda \varepsilon + \frac{1}{N} \sum_{i=1}^N s_i \\ \text{s.t.} \quad & \ell(\xi) - \lambda \|\xi - \hat{\xi}_i\| - \delta_{\Xi}(\xi) \leq s_i \quad \forall \xi \in \mathbb{R}^m, \quad \forall i \leq N \\ & \lambda \geq 0 \end{aligned}$$

Reformulate robust program as bilevel program:

$$\begin{aligned} \inf_{\lambda, s_j} \quad & \lambda \varepsilon + \frac{1}{N} \sum_{i=1}^N s_i \\ \text{s.t.} \quad & \sup_{\xi \in \mathbb{R}^m} \left(\ell(\xi) - \lambda \|\xi - \hat{\xi}_i\| - \delta_{\Xi}(\xi) \right) \leq s_i \quad \forall i \leq N \\ & \lambda \geq 0 \end{aligned}$$

Take the Fenchel dual of the lower-level problem:⁷⁾

$$\begin{aligned}
 & \inf_{\lambda, s_i, v_i, z_i} \quad \lambda \varepsilon + \frac{1}{N} \sum_{i=1}^N s_i \\
 & \text{s.t.} \quad \boxed{[-\ell]^*} (z_i - v_i) + \boxed{\sigma_{\Xi}}(v_i) - z_i^\top \hat{\xi}_i \leq s_i \quad \forall i \leq N \\
 & \quad \quad \boxed{\|z_i\|_*} \leq \lambda
 \end{aligned}$$

dual norm of z_i
convex conjugate of $-\ell$
support function of Ξ

⁷⁾ Ben-Tal, den Hertog & Vial, *Math. Program.*, 2015.

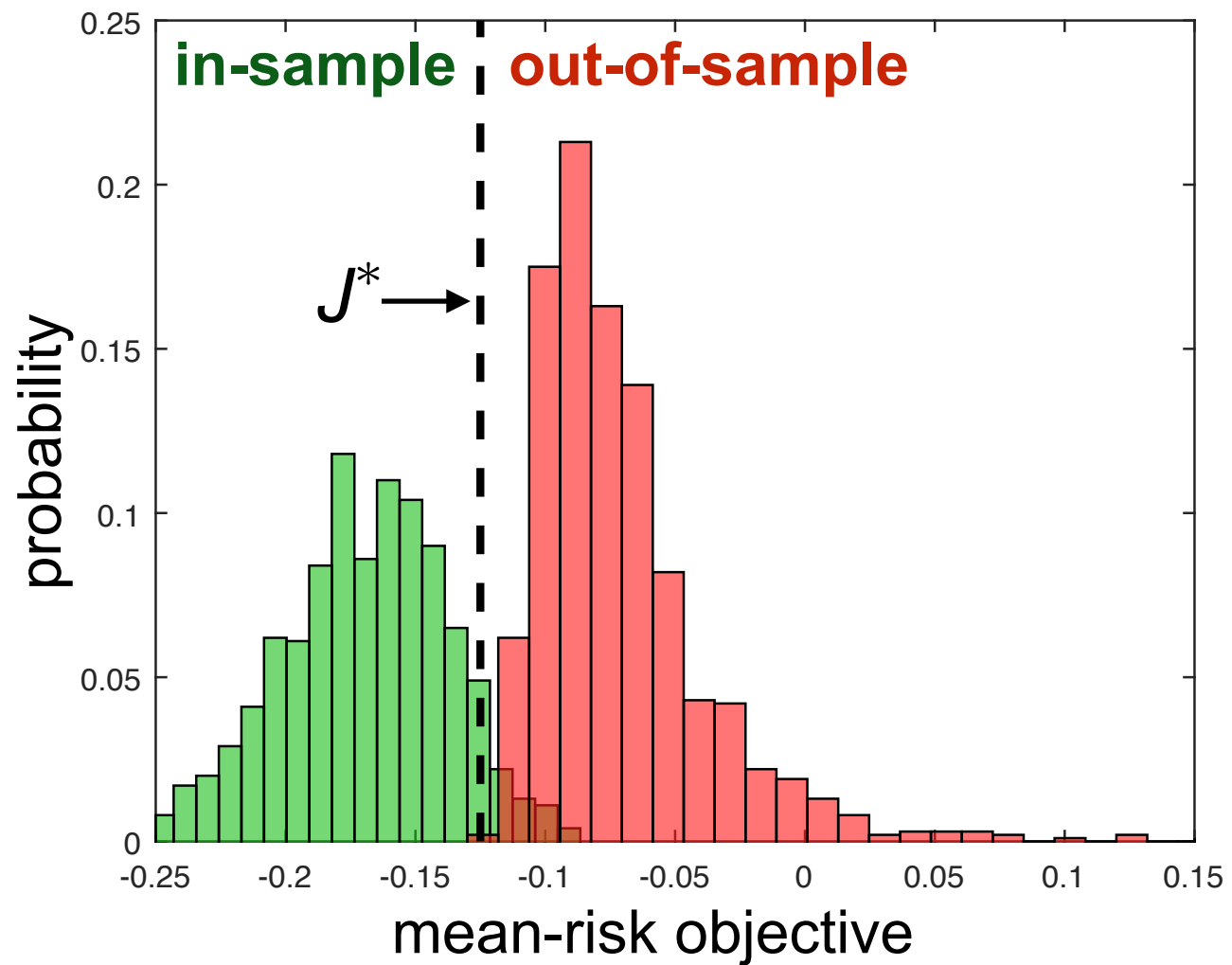
The worst-case expectation equals:

$$\begin{aligned} \inf_{\lambda, s_i, v_i, z_i} \quad & \lambda \varepsilon + \frac{1}{N} \sum_{i=1}^N s_i \\ \text{s.t.} \quad & [-\ell]^* (z_i - v_i) + \sigma_{\Xi}(v_i) - z_i^{\top} \hat{\xi}_i \leq s_i \quad \forall i \leq N \\ & \|z_i\|_* \leq \lambda \end{aligned}$$

- ▶ **Finite convex program**
- ▶ Problem size grows **polynomially** in input data
- ▶ Can be combined with minimization over $x \in \mathcal{X}$: the resulting problem is in the **same complexity class as SAA**

DRO with Scarce Data

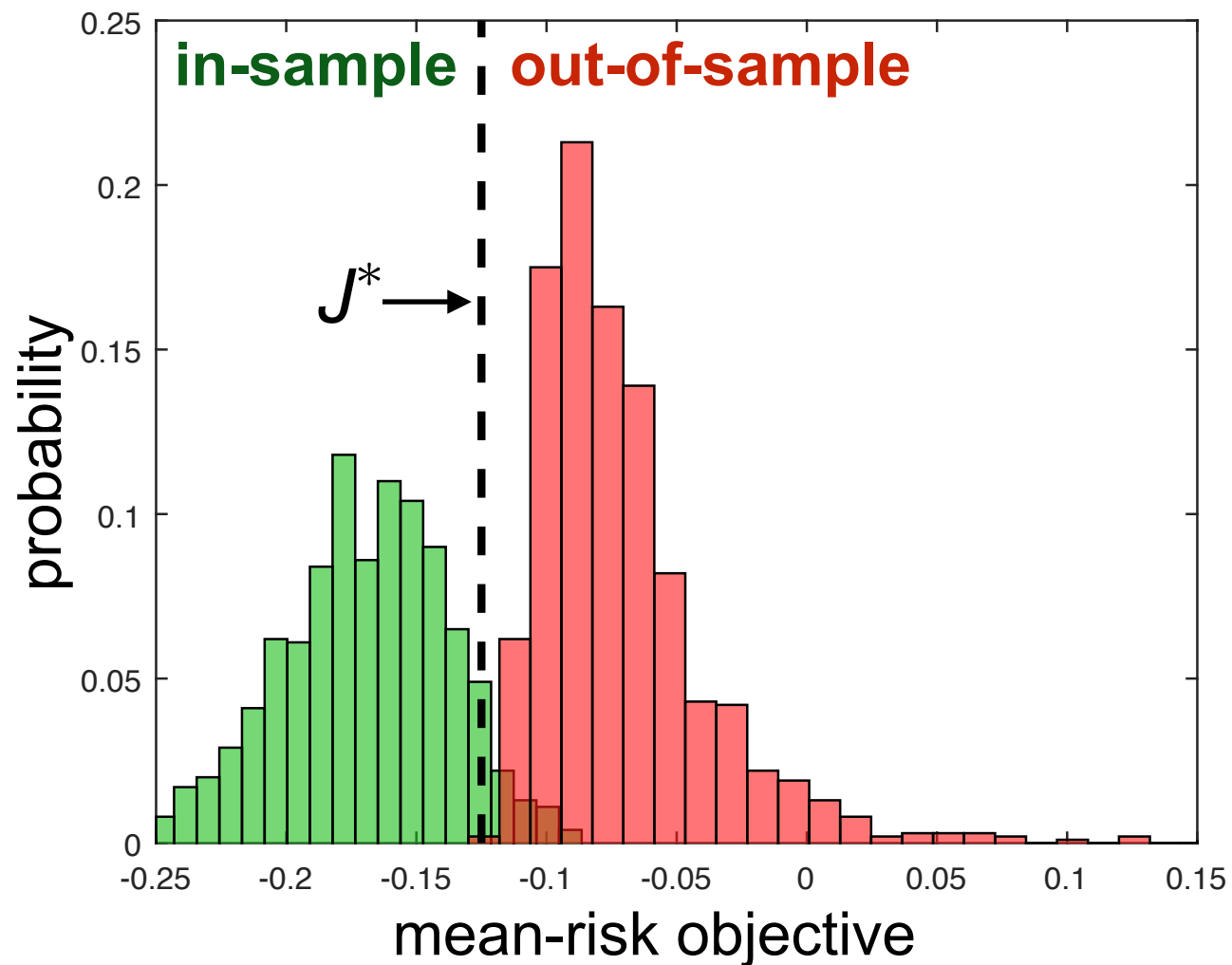
Performance of SAA solution



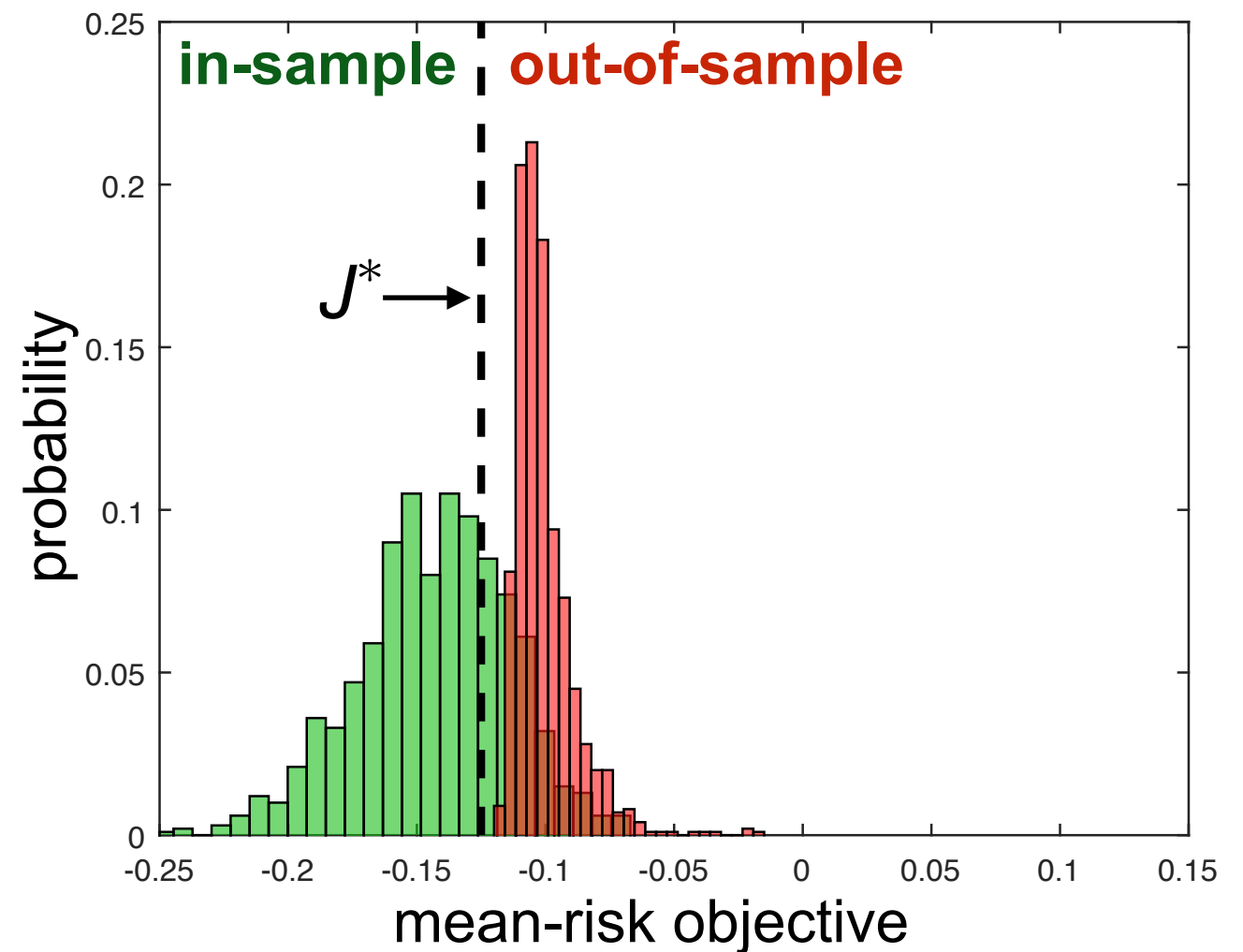
- ▶ in-sample: optimistic bias
- ▶ out-of-sample: pessimistic bias

DRO with Scarce Data

Performance of SAA solution



Performance of DRO solution



- ▶ in-sample: optimistic bias
- ▶ out-of-sample: pessimistic bias
- ▶ **DRO reduces bias & post-decision disappointment**

Worst-Case Distributions

Dualize finite convex program:

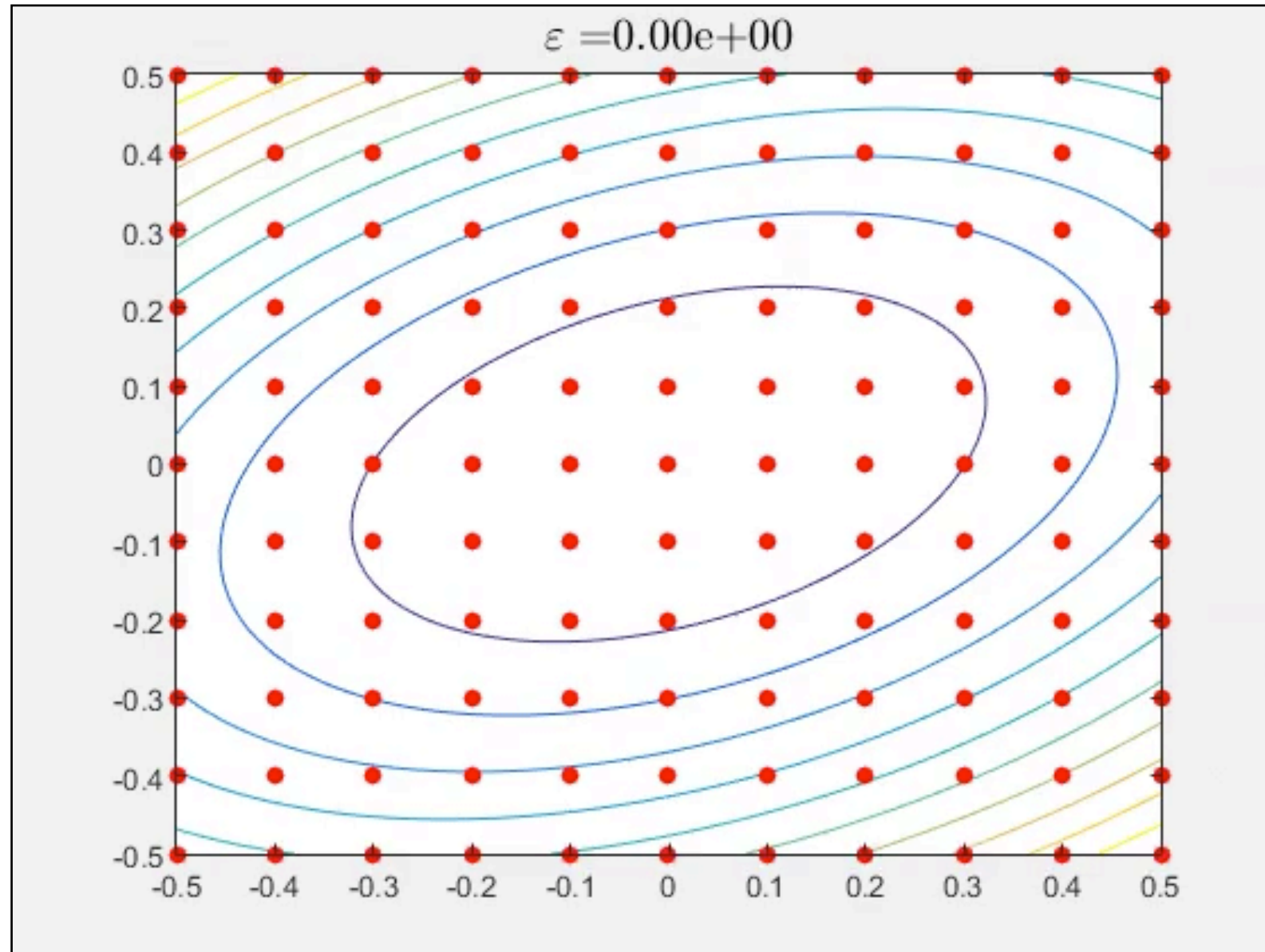
$$\begin{aligned} \sup_{Q \in \hat{\mathcal{P}}_N} \mathbb{E}^Q[\ell(\xi)] &= \sup_{q_i} \frac{1}{N} \sum_{i=1}^N \ell(\hat{\xi}_i - q_i) \\ \text{s.t.} \quad &\frac{1}{N} \sum_{i=1}^N \|q_i\| \leq \varepsilon \\ &\hat{\xi}_i - q_i \in \Xi \quad \forall i \leq N \end{aligned}$$

If q_1^*, \dots, q_N^* solves the dual problem, then

$$Q^* = \frac{1}{N} \sum_{i=1}^N \delta_{\hat{\xi}_i - q_i^*}$$

is a worst-case distribution.

Concave Quadratic Loss Function



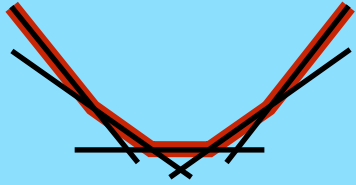
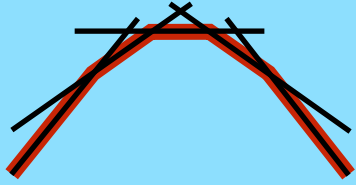
Generalized Loss Functions

Generalized convexity assumption:

► $l = \max\{l_1, \dots, l_J\}$ where each $-l_j$ is proper, convex and lsc.

Maxima & Minima of Affine Functions

Set $\Xi = \{\xi : C\xi \leq d\}$ and $a_j(\xi) = a_j^\top \xi - b_j$ for all $j \leq J$.

$\ell(\xi)$	$\sup_{Q \in \hat{\mathcal{P}}_N} \mathbb{E}^Q[\ell(\xi)]$
$\max_{j \leq J} a_j(\xi)$ 	$\inf_{\lambda, s_i, Y_{ij}} \lambda \varepsilon + \frac{1}{N} \sum_{i=1}^N s_i$ $\text{s.t. } b_j + a_j^\top \hat{\xi}_i + Y_{ij}^\top (d - C\hat{\xi}_i) \leq s_i \quad \forall i \leq N, \forall j \leq J$ $\ C^\top Y_{ij} - a_j\ _* \leq \lambda \quad \forall i \leq N, \forall j \leq J$ $Y_{ij} \geq 0 \quad \forall i \leq N, \forall j \leq J$
$\min_{j \leq J} a_j(\xi)$ 	$\inf_{\lambda, s_i, Y_i, \theta_i} \lambda \varepsilon + \frac{1}{N} \sum_{i=1}^N s_i$ $\text{s.t. } \theta_i^\top (b - A\hat{\xi}_i) + Y_i^\top (d - C\hat{\xi}_i) \leq s_i \quad \forall i \leq N$ $\ C^\top Y_i - A^\top \theta_i\ _* \leq \lambda \quad \forall i \leq N$ $e^\top \theta_i = 1, \quad Y_i \geq 0, \quad \theta_i \geq 0 \quad \forall i \leq N$

Two-Stage Stochastic Programming

Objective uncertainty:

$$\ell(\xi) = \inf_y \{ \underbrace{y^\top Q \xi : Wy \geq h}_{\text{concave piecewise affine}} \}$$

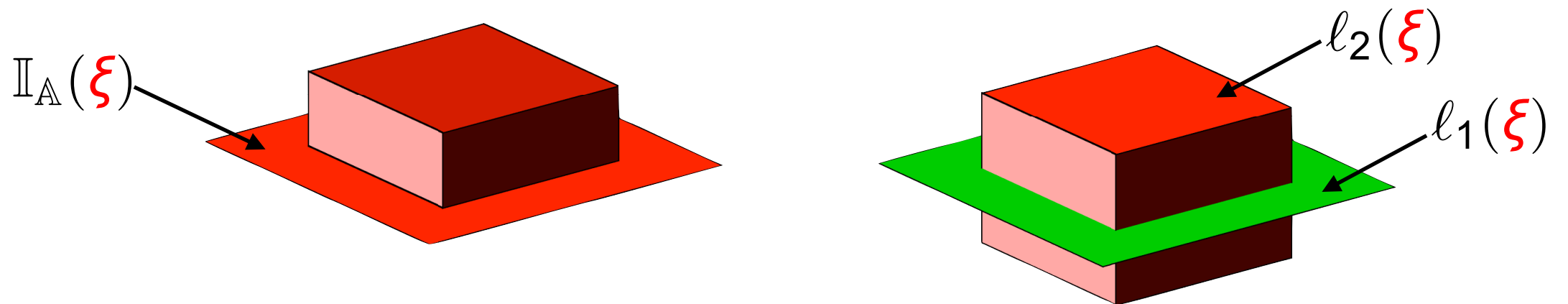
Constraint uncertainty:

$$\ell(\xi) = \inf_y \{ \underbrace{q^\top y : Wy \geq H\xi + h}_{\text{convex piecewise affine}} \}$$

Uncertainty Quantification

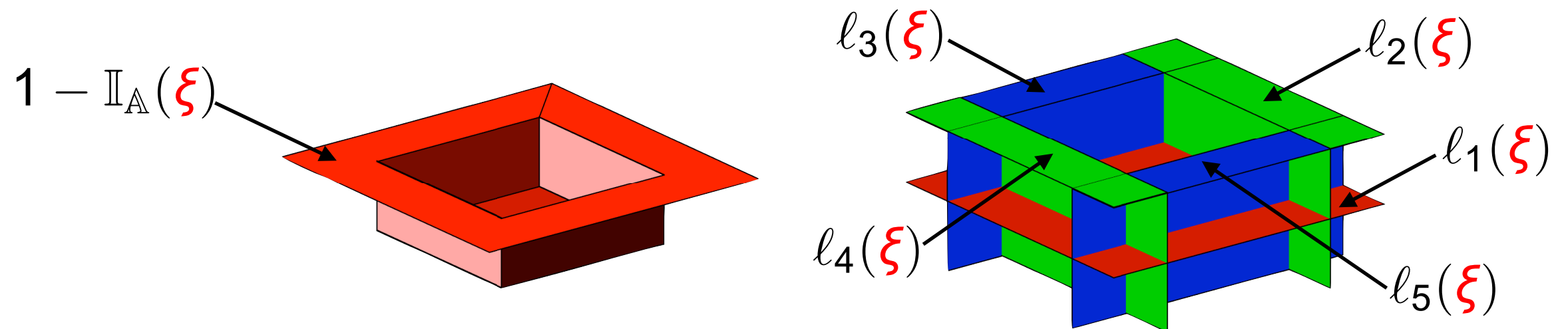
Probability of being inside a polytope:

$$\sup_{Q \in \hat{\mathcal{P}}_N} Q[\xi \in \mathbb{A}] = \sup_{Q \in \hat{\mathcal{P}}_N} \mathbb{E}^Q[\mathbb{I}_{\mathbb{A}}(\xi)] = \sup_{Q \in \hat{\mathcal{P}}_N} \mathbb{E}^Q[\max\{l_1(\xi), l_2(\xi)\}]$$



Probability of being outside of a polytope:

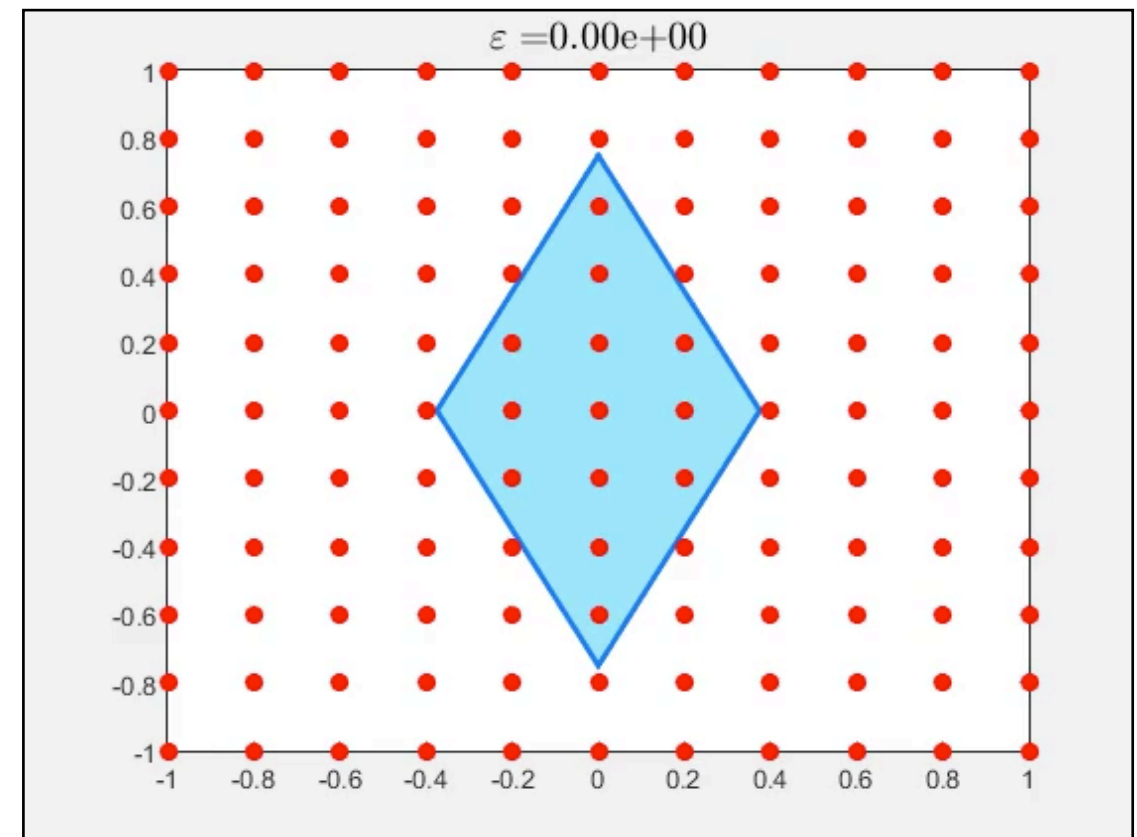
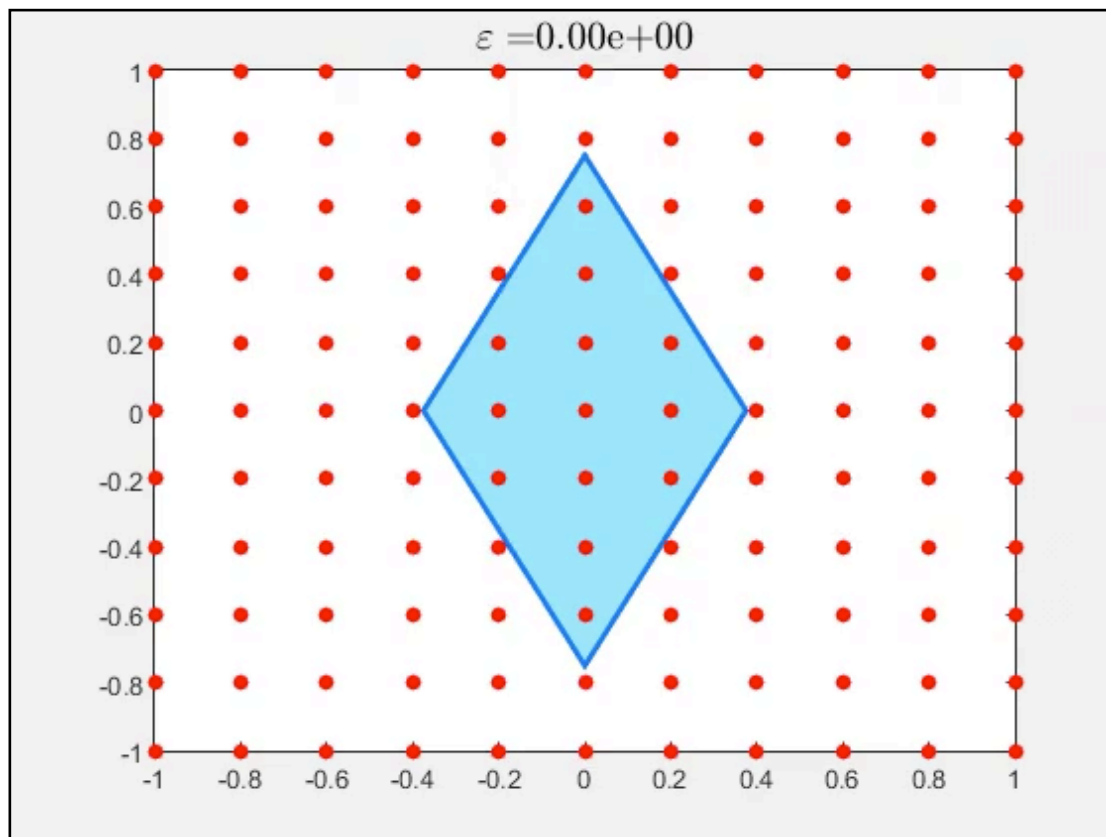
$$\sup_{Q \in \hat{\mathcal{P}}_N} Q[\xi \notin \mathbb{A}] = \sup_{Q \in \hat{\mathcal{P}}_N} \mathbb{E}^Q[1 - \mathbb{I}_{\mathbb{A}}(\xi)] = \sup_{Q \in \hat{\mathcal{P}}_N} \mathbb{E}^Q\left[\max_{j \leq J} \{l_j(\xi)\}\right]$$



Uncertainty Quantification

Probability of being inside a polytope:

Probability of being outside a polytope:



Application 1: Portfolio Selection

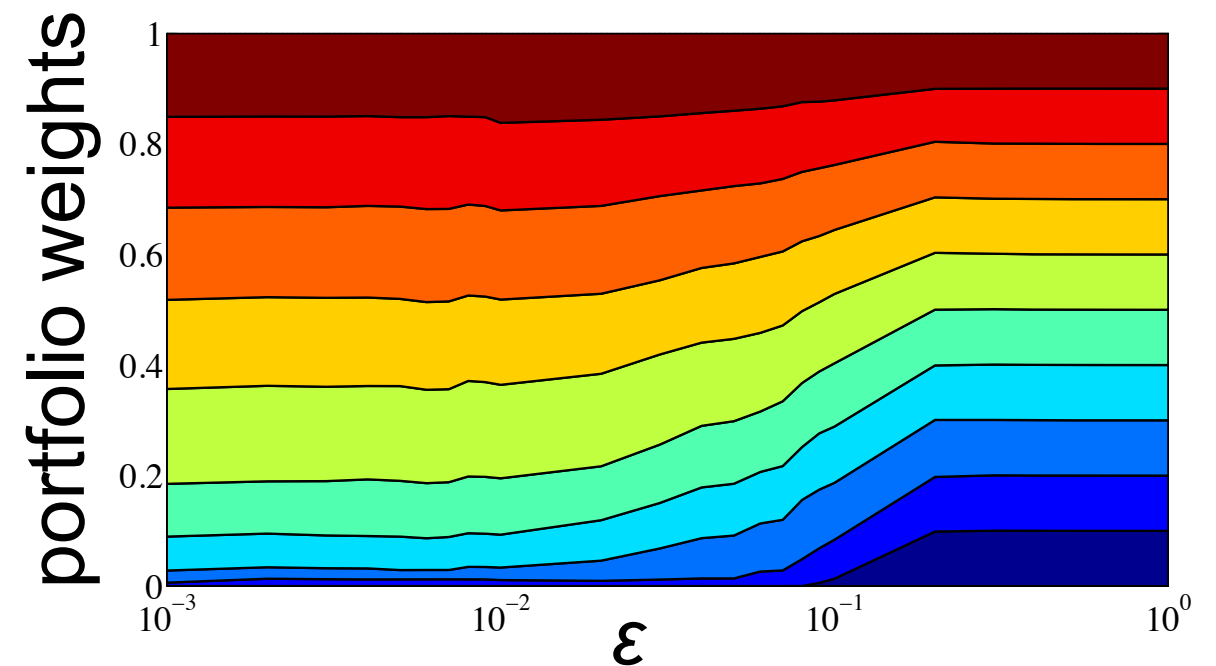
$$\min_{x \in \mathcal{X}} \left\{ \mathbb{E}^{\mathbb{P}} \left[-x^{\top} \xi \right] + \rho \mathbb{P}\text{-CVaR}_{\alpha}(-x^{\top} \xi) \right\}$$

▶ 10 assets

▶ $\rho = 10$

▶ $\alpha = 20\%$

▶ $\xi_i = \psi + \zeta_i$ where $\psi \sim \mathcal{N}(0, 2\%)$
and $\zeta_i \sim \mathcal{N}(i \times 3\%, i \times 2.5\%)$



Fact: The 1/n portfolio is hard to beat out of sample.⁸⁾

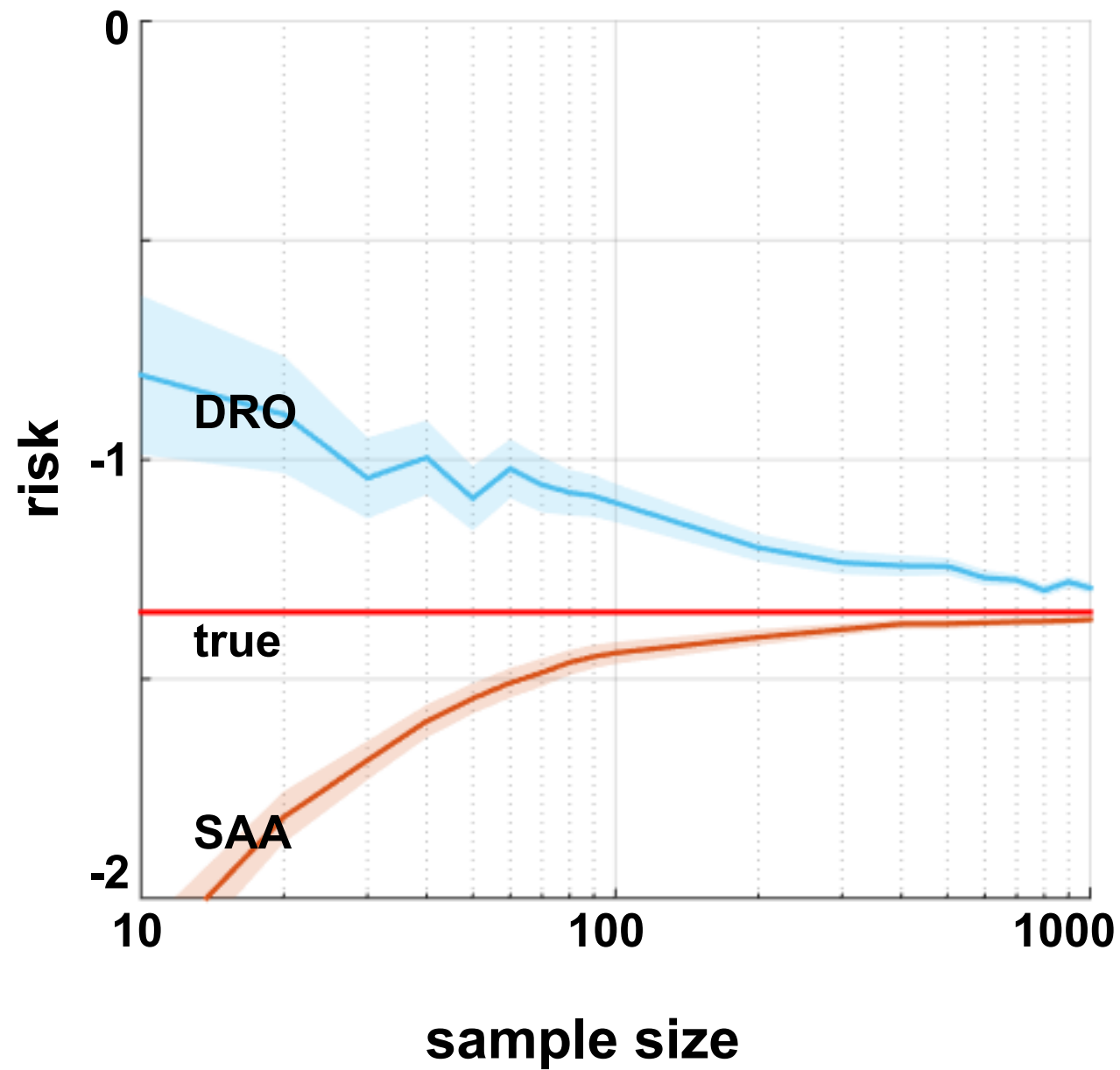
Possible Explanation: It is optimal for $\epsilon \rightarrow \infty$.⁹⁾

⁸⁾ DeMiguel, Garlappi & Uppal, *Rev. Financ. Stud.*, 2009;

⁹⁾ Pflug, Pichler & Wozabal, *J. Bank. Financ.*, 2012.

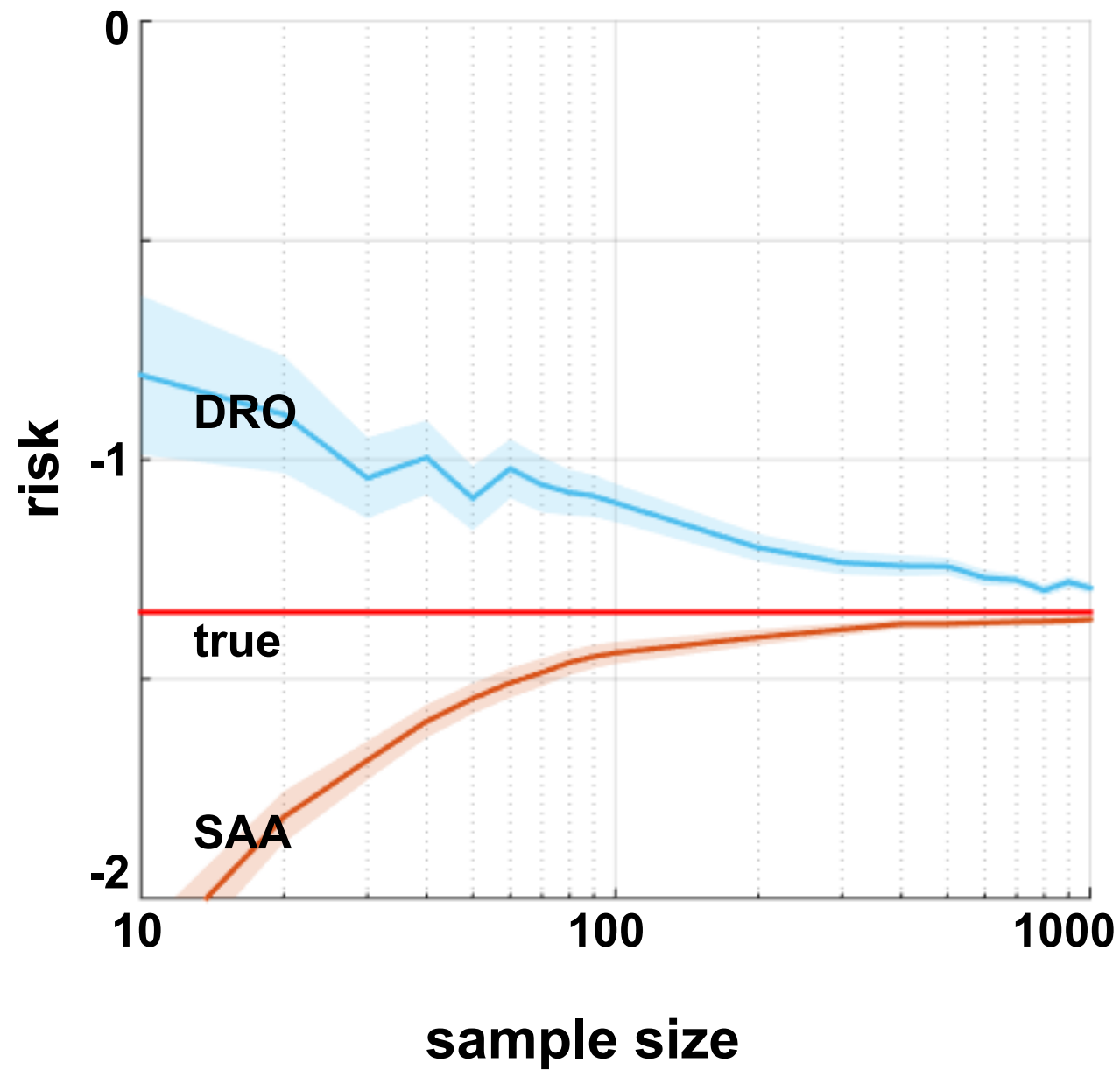
Learning Curves

what we **think** to get ...



Learning Curves

what we **think** to get ...



what we **actually** get ...

