

Wasserstein Distributionally Robust Optimization: Theory and Applications in Machine Learning

Daniel Kuhn, Peyman Mohajerin Esfahani,
Viet Anh Nguyen, Soroosh Shafieezadeh-Abadeh

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

Decision-Making under Uncertainty

Risk:

$$\mathcal{R}(\mathbb{P}, \ell) = \mathbb{E}^{\mathbb{P}}[\ell(\xi)]$$

Optimal risk:

$$\mathcal{R}(\mathbb{P}, \mathcal{L}) = \inf_{\ell \in \mathcal{L}} \mathbb{E}^{\mathbb{P}}[\ell(\xi)]$$

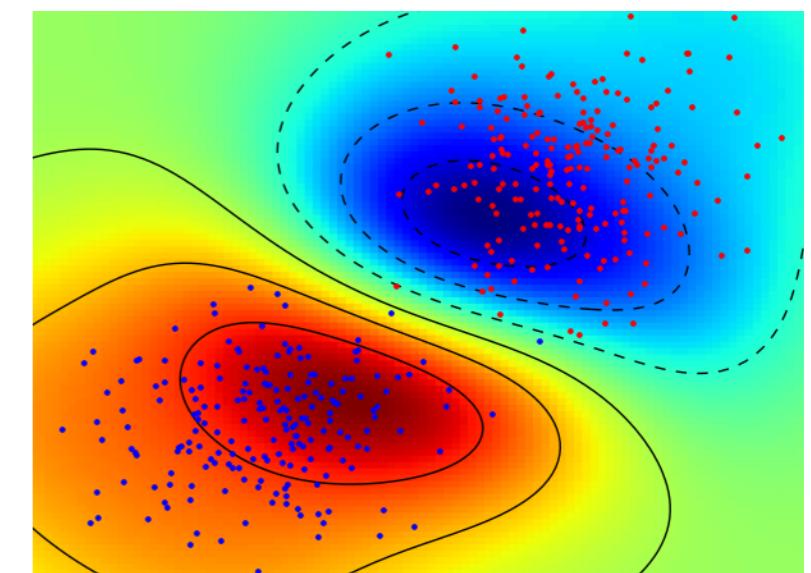
Applications:



Supply Chain Mgmt.



Portfolio Mgmt.



Machine Learning

Data-Driven Decision-Making

Available information:

Structural: $\textcolor{red}{P}$ supported on $\Xi \subseteq \mathbb{R}^m$

Statistical: $\hat{\xi}_1, \hat{\xi}_2, \dots, \hat{\xi}_N \sim \textcolor{red}{P}^N$

Nominal problem:

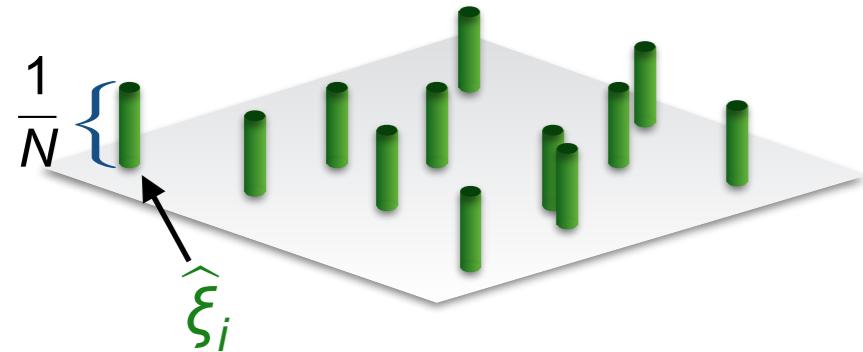


Nominal risk: $\mathcal{R}(\hat{P}_N, \ell)$

Optimal nominal risk: $\mathcal{R}(\hat{P}_N, \mathcal{L})$

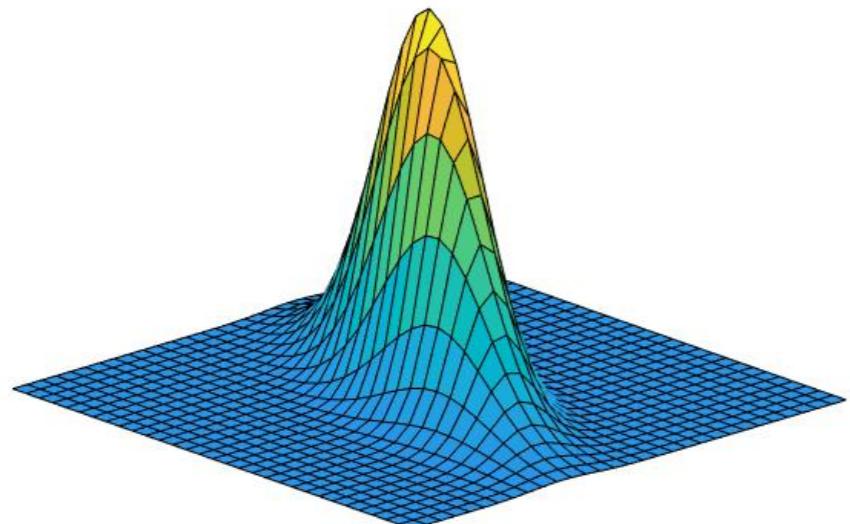
Nominal Distribution

Non-parametric estimators:



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

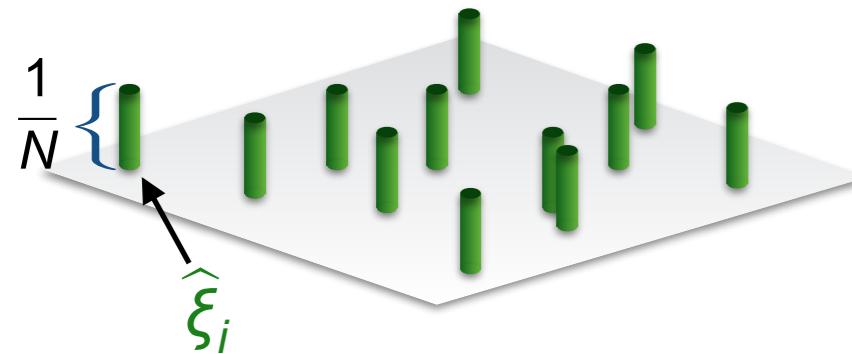
Parametric estimators:



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

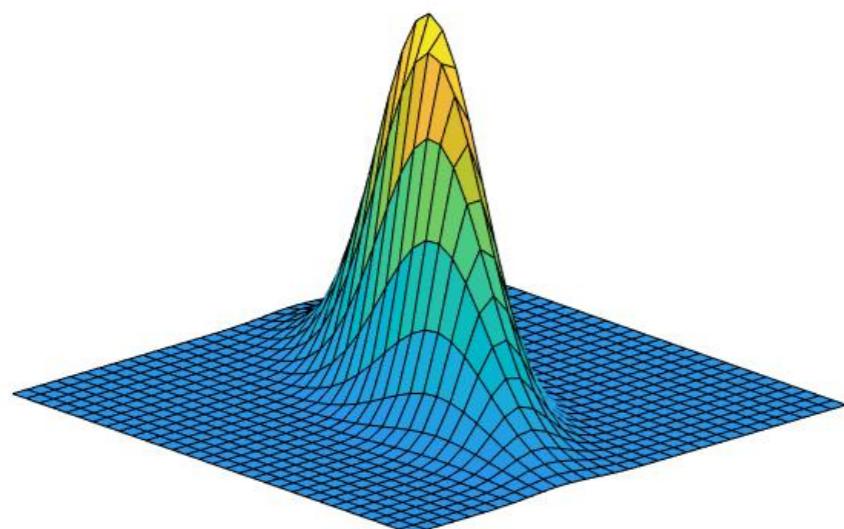
Nominal Distribution

Non-parametric estimators:



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

Parametric estimators:



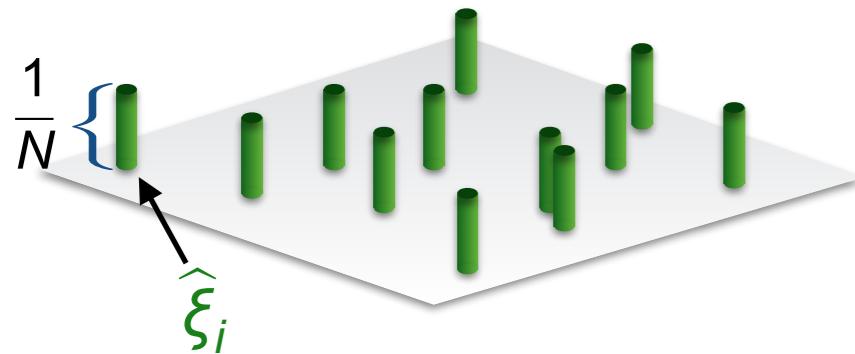
Elliptical distribution: $\hat{\mathbb{P}}_N = \mathcal{E}_g(\hat{\mu}_N, \hat{\Sigma}_N)$

density generator

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

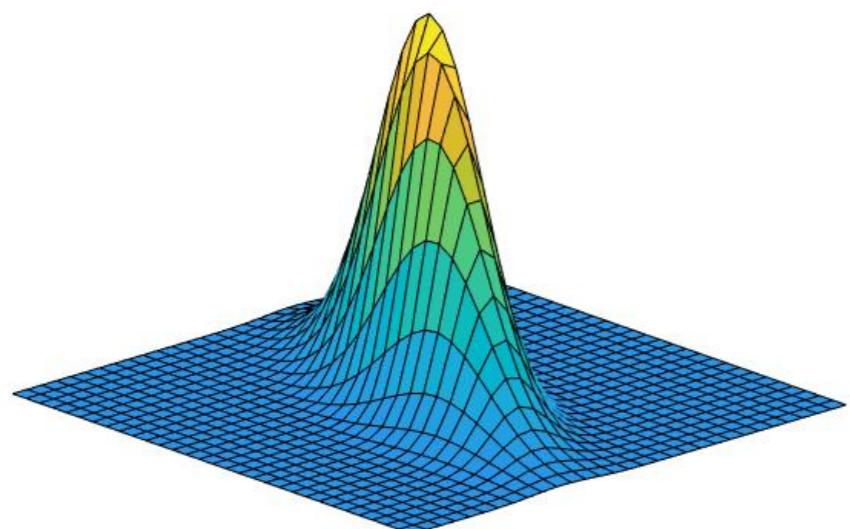
Nominal Distribution

Non-parametric estimators:



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

Parametric estimators:

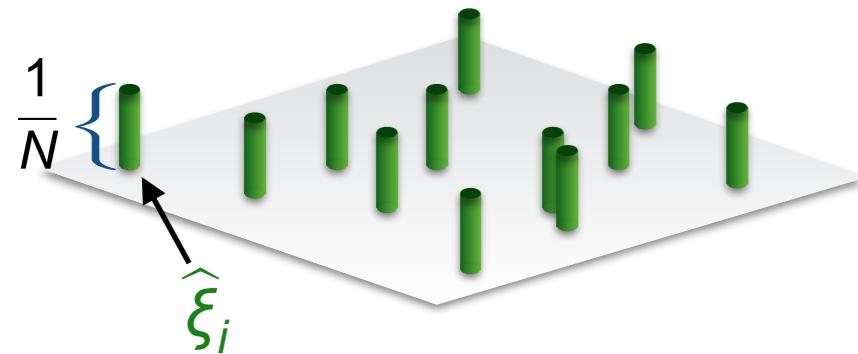


Elliptical distribution: $\widehat{\mathbb{P}}_N = \mathcal{E}_g(\hat{\mu}_N, \hat{\Sigma}_N)$

mean

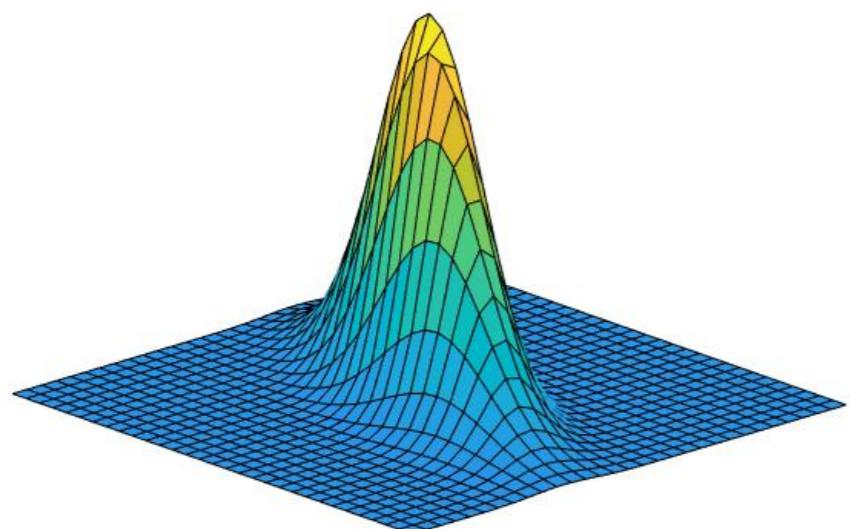
Nominal Distribution

Non-parametric estimators:



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

Parametric estimators:



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

covariance matrix

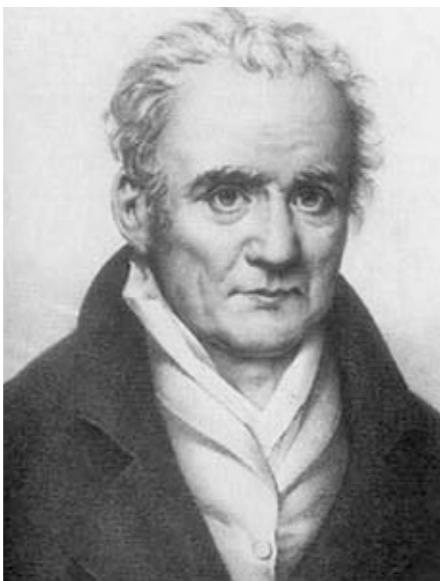
Estimation Errors

The nominal distribution $\hat{\mathbb{P}}_N$

- ▶ depends on $\hat{\xi}_1, \hat{\xi}_2, \dots, \hat{\xi}_N$ and is thus random;
- ▶ differs from the data-generating distribution \mathbb{P} .

Q: How to measure estimation errors?

A: Use the Monge / Kantorovich / Wasserstein distance!



Gaspard Monge
(1746 – 1818)



Leonid Kantorovich
(1912 – 1986)

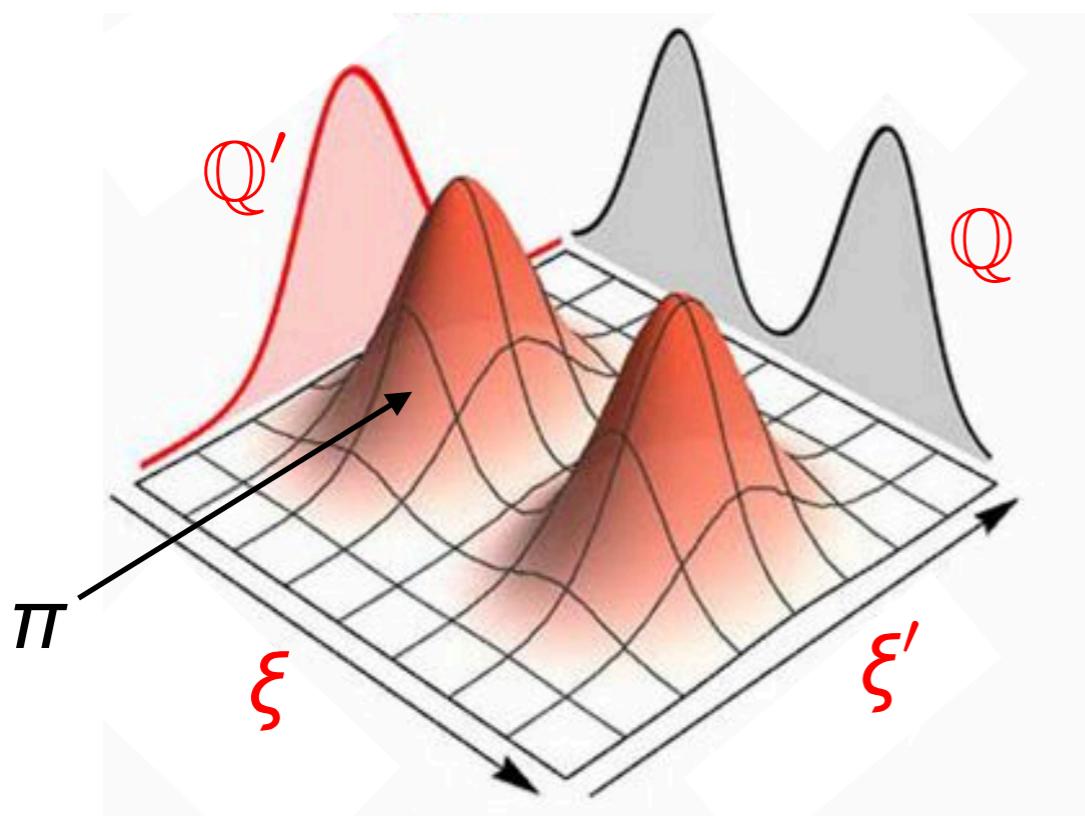


Leonid Vaserštejn
(*1944)

Wasserstein Distance

Definition:

$$W_p(Q, Q') = \left(\inf_{\pi \in \Pi(Q, Q')} \int_{\mathbb{R}^m \times \mathbb{R}^m} \|\xi - \xi'\|^p \pi(d\xi, d\xi') \right)^{\frac{1}{p}}$$

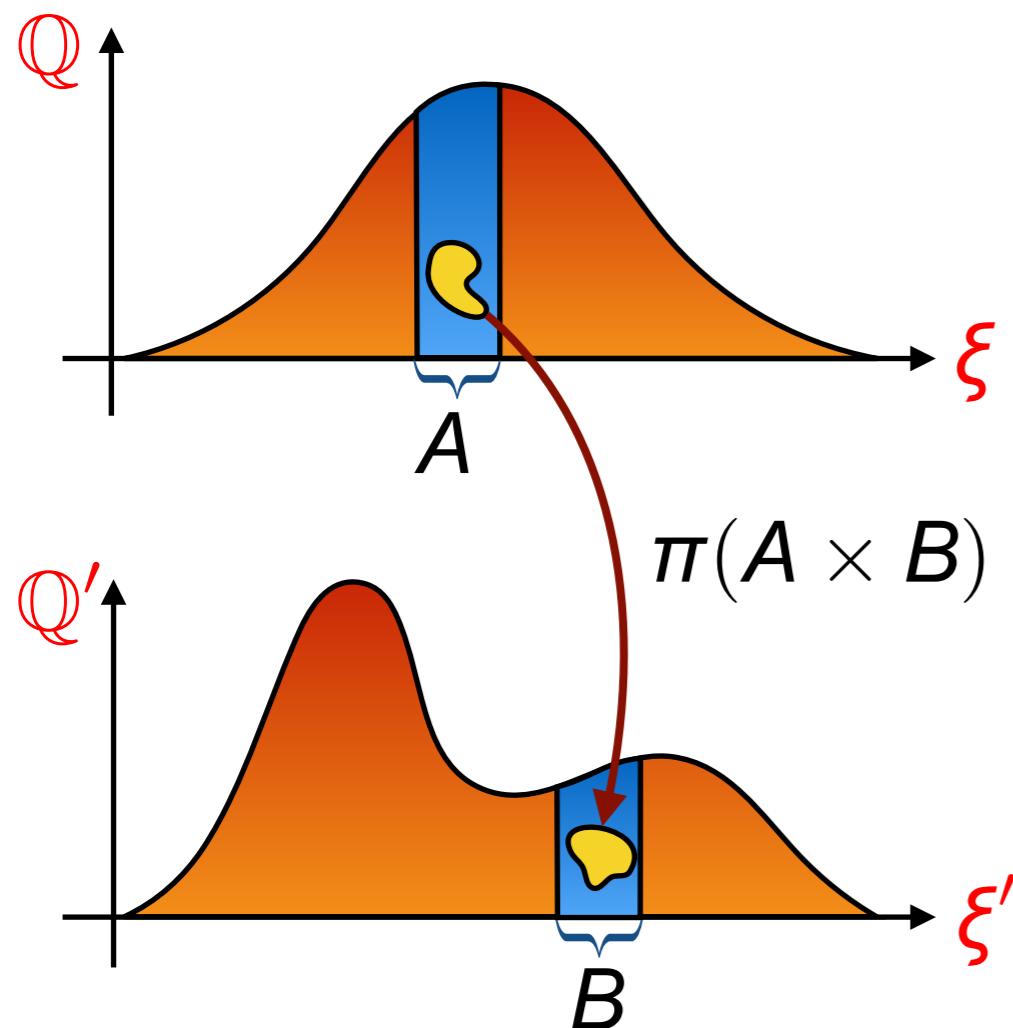


$\Pi(Q, Q')$ = set of couplings with
marginals Q and Q'

Wasserstein Distance

Definition:

$$W_p(Q, Q') = \left(\inf_{\pi \in \Pi(Q, Q')} \int_{\mathbb{R}^m \times \mathbb{R}^m} \|\xi - \xi'\|^p \pi(d\xi, d\xi') \right)^{\frac{1}{p}}$$



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

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

Stability Theory

Assume:

- ▶ $\ell(\xi)$ is Lipschitz continuous with $\text{Lip}(\ell) = L$;
- ▶ the estimation error in $\hat{\mathbb{P}}_N$ is small, $W_1(\hat{\mathbb{P}}_N, \mathbb{P}) \leq \varepsilon$.

Question: How big is the estimation error in $\mathcal{R}(\hat{\mathbb{P}}_N, \ell)$?

Stability Theory

Assume:

- ▶ $\ell(\xi)$ is Lipschitz continuous with $\text{Lip}(\ell) = L$;
- ▶ the estimation error in $\widehat{\mathbb{P}}_N$ is small, $W_1(\widehat{\mathbb{P}}_N, \mathbb{P}) \leq \varepsilon$.

Question: How big is the estimation error in $\mathcal{R}(\widehat{\mathbb{P}}_N, \ell)$?

$$\left| \mathcal{R}(\widehat{\mathbb{P}}_N, \ell) - \mathcal{R}(\mathbb{P}, \ell) \right| = L \cdot \left| \mathbb{E}^{\widehat{\mathbb{P}}_N}[\ell(\xi)/L] - \mathbb{E}^{\mathbb{P}}[\ell(\xi)/L] \right|$$

$$\leq L \cdot W_1(\widehat{\mathbb{P}}_N, \mathbb{P}) \leq L \cdot \varepsilon$$

↑

K-R Theorem

⇒ Estimation error at most amplified by L

Stability Theory

Assume:

- ▶ $\ell(\xi)$ is Lipschitz continuous with $\text{Lip}(\ell) \leq L \quad \forall \ell \in \mathcal{L};$
- ▶ the estimation error in $\hat{\mathbb{P}}_N$ is small, $W_1(\hat{\mathbb{P}}_N, \mathbb{P}) \leq \varepsilon.$

Question: How big is the estimation error in $\mathcal{R}(\hat{\mathbb{P}}_N, \mathcal{L})?$

Stability Theory

Assume:

- ▶ $\ell(\xi)$ is Lipschitz continuous with $\text{Lip}(\ell) \leq L \quad \forall \ell \in \mathcal{L}$;
- ▶ the estimation error in $\widehat{\mathbb{P}}_N$ is small, $W_1(\widehat{\mathbb{P}}_N, \mathbb{P}) \leq \varepsilon$.

Question: How big is the estimation error in $\mathcal{R}(\widehat{\mathbb{P}}_N, \mathcal{L})$?

$$\left| \mathcal{R}(\widehat{\mathbb{P}}_N, \mathcal{L}) - \mathcal{R}(\mathbb{P}, \mathcal{L}) \right| \leq L \cdot \varepsilon$$

⇒ Estimation error at most amplified by L

Distributionally Robust Optimization (DRO)

Nominal risk: $\mathcal{R}(\hat{\mathbb{P}}_N, \ell)$

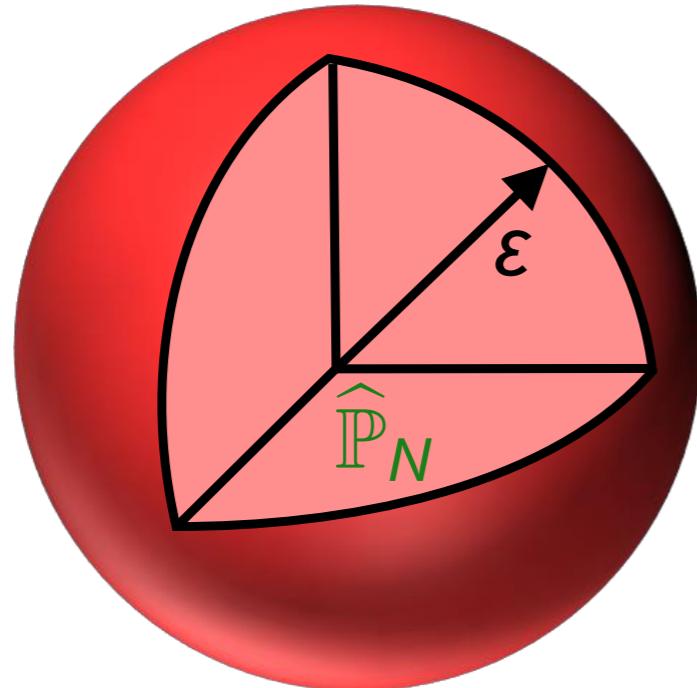
Optimal nominal risk: $\mathcal{R}(\hat{\mathbb{P}}_N, \mathcal{L})$

Q: If we have found a $\hat{\mathbb{P}}_N$ that approximates \mathbb{P} best in Wasserstein distance, how can we further reduce the out-of-sample risk?

A: Robustify the risk w.r.t. all distributions \mathbb{Q} with $W_p(\mathbb{Q}, \hat{\mathbb{P}}_N) \leq \varepsilon$

Wasserstein DRO

Definition: $\mathbb{B}_{\varepsilon,p}(\hat{\mathbb{P}}_N) = \left\{ \mathbb{Q} \in \mathcal{P}(\Xi) : W_p(\mathbb{Q}, \hat{\mathbb{P}}_N) \leq \varepsilon \right\}$



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

Worst-case risk:

$$\mathcal{R}_{\varepsilon,p}(\hat{\mathbb{P}}_N, \ell) = \sup_{\mathbb{Q} \in \mathbb{B}_{\varepsilon,p}(\hat{\mathbb{P}}_N)} \mathbb{E}^{\mathbb{Q}}[\ell(\xi)]$$

Worst-case optimal risk:

$$\mathcal{R}_{\varepsilon,p}(\hat{\mathbb{P}}_N, \mathcal{L}) = \inf_{\ell \in \mathcal{L}} \mathcal{R}_{\varepsilon,p}(\hat{\mathbb{P}}_N, \ell)$$

Three Simple Bounds on the Worst-Case Expectation

Robust Lower Bound (any p)

Theorem: If $\widehat{\mathbb{P}}_N$ is the empirical distribution, then

$$\mathcal{R}_{\varepsilon,p}(\widehat{\mathbb{P}}_N, \ell) \geq \left\{ \begin{array}{ll} \sup & \frac{1}{N} \sum_{i=1}^N \ell(\widehat{\xi}_i + \theta_i) \\ \text{s.t.} & \theta_i \in \mathbb{R}^m \quad \forall i \in [N] \\ & \widehat{\xi}_i + \theta_i \in \Xi \quad \forall i \in [N] \\ & \frac{1}{N} \sum_{i=1}^N \|\theta_i\|^p \leq \varepsilon^p. \end{array} \right.$$

Lipschitz Regularization (any p)

Theorem: The worst-case risk is bounded above by the Lipschitz-regularized nominal risk, i.e.,

$$\mathcal{R}_{\varepsilon,p}(\widehat{\mathbb{P}}_N, \ell) \leq \mathcal{R}(\widehat{\mathbb{P}}_N, \ell) + \varepsilon \cdot \text{Lip}(\ell).$$

Gelbrich Bound ($p = 2$)

Theorem: If $\mathbb{Q} \sim (\mu, \Sigma)$ and $\mathbb{Q}' \sim (\mu', \Sigma')$, then¹⁾

$$W_2(\mathbb{Q}, \mathbb{Q}') \geq \sqrt{\|\mu - \mu'\|_2^2 + \text{Tr} \left(\Sigma + \Sigma' - 2 \left(\Sigma^{\frac{1}{2}} \Sigma' \Sigma^{\frac{1}{2}} \right)^{\frac{1}{2}} \right)}.$$

¹⁾ Gelbrich, *Mathematische Nachrichten*, 1990.

Gelbrich Bound ($p = 2$)

Theorem: If $\mathbb{Q} \sim (\mu, \Sigma)$ and $\mathbb{Q}' \sim (\mu', \Sigma')$, then¹⁾

$$W_2(\mathbb{Q}, \mathbb{Q}') \geq \underbrace{\sqrt{\|\mu - \mu'\|_2^2 + \text{Tr} \left(\Sigma + \Sigma' - 2 \left(\Sigma^{\frac{1}{2}} \Sigma' \Sigma^{\frac{1}{2}} \right)^{\frac{1}{2}} \right)}}_{= d_G((\mu, \Sigma), (\mu', \Sigma'))}.$$

¹⁾ Gelbrich, *Mathematische Nachrichten*, 1990.

Gelbrich Bound ($p = 2$)

Theorem: If $\mathbb{Q} \sim (\mu, \Sigma)$ and $\mathbb{Q}' \sim (\mu', \Sigma')$, then¹⁾

$$W_2(\mathbb{Q}, \mathbb{Q}') \geq \underbrace{\sqrt{\|\mu - \mu'\|_2^2 + \text{Tr} \left(\Sigma + \Sigma' - 2 \left(\Sigma^{\frac{1}{2}} \Sigma' \Sigma^{\frac{1}{2}} \right)^{\frac{1}{2}} \right)}}_{= d_G((\mu, \Sigma), (\mu', \Sigma'))}.$$

The bound is exact if \mathbb{Q}, \mathbb{Q}' are elliptical with the same generator.

¹⁾ Gelbrich, *Mathematische Nachrichten*, 1990.

Gelbrich Bound ($p = 2$)

Uncertainty set for mean vectors and covariance matrices:

$$\mathcal{U}_\varepsilon(\hat{\mu}, \hat{\Sigma}) = \left\{ (\mu, \Sigma) \in \mathbb{R}^m \times \mathbb{S}_+^m : d_G \left((\hat{\mu}, \hat{\Sigma}), (\mu, \Sigma) \right) \leq \varepsilon^2 \right\}$$

Gelbrich Bound ($p = 2$)

Uncertainty set for mean vectors and covariance matrices:

$$\mathcal{U}_\varepsilon(\hat{\mu}, \hat{\Sigma}) = \left\{ (\mu, \Sigma) \in \mathbb{R}^m \times \mathbb{S}_+^m : d_G \left((\hat{\mu}, \hat{\Sigma}), (\mu, \Sigma) \right) \leq \varepsilon^2 \right\}$$

Mean-covariance “projection”:

$$T: \mathbb{Q} \mapsto \left(\mathbb{E}^\mathbb{Q}[\xi], \mathbb{E}^\mathbb{Q}[\xi \xi^\top] - \mathbb{E}^\mathbb{Q}[\xi] \mathbb{E}^\mathbb{Q}[\xi]^\top \right)$$

Gelbrich Bound ($p = 2$)

Uncertainty set for mean vectors and covariance matrices:

$$\mathcal{U}_\varepsilon(\hat{\mu}, \hat{\Sigma}) = \left\{ (\mu, \Sigma) \in \mathbb{R}^m \times \mathbb{S}_+^m : d_G \left((\hat{\mu}, \hat{\Sigma}), (\mu, \Sigma) \right) \leq \varepsilon^2 \right\}$$

Mean-covariance “projection”:

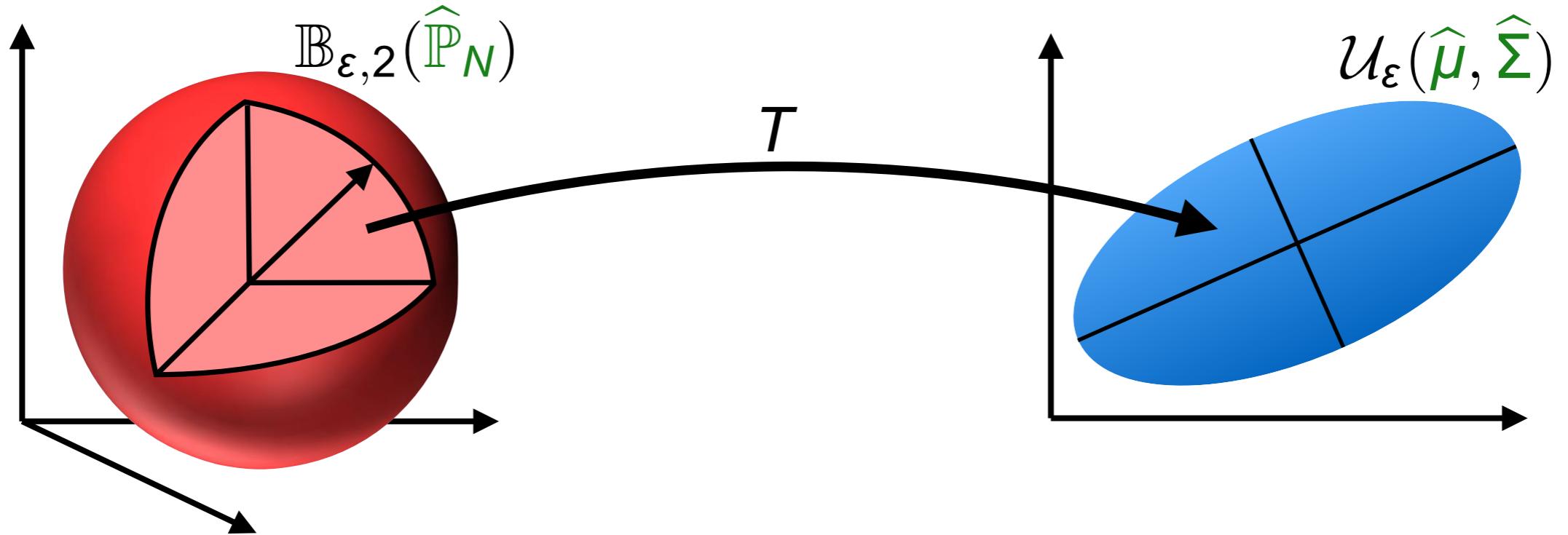
$$T: \mathbb{Q} \mapsto \left(\underbrace{\mathbb{E}^\mathbb{Q}[\xi]}_{\mu}, \underbrace{\mathbb{E}^\mathbb{Q}[\xi \xi^\top] - \mathbb{E}^\mathbb{Q}[\xi] \mathbb{E}^\mathbb{Q}[\xi]^\top}_{\Sigma} \right)$$

Gelbrich Bound ($p = 2$)

Uncertainty set for mean vectors and covariance matrices:

$$\mathcal{U}_\varepsilon(\hat{\mu}, \hat{\Sigma}) = \left\{ (\mu, \Sigma) \in \mathbb{R}^m \times \mathbb{S}_+^m : d_G \left((\hat{\mu}, \hat{\Sigma}), (\mu, \Sigma) \right) \leq \varepsilon^2 \right\}$$

Mean-covariance “projection”:

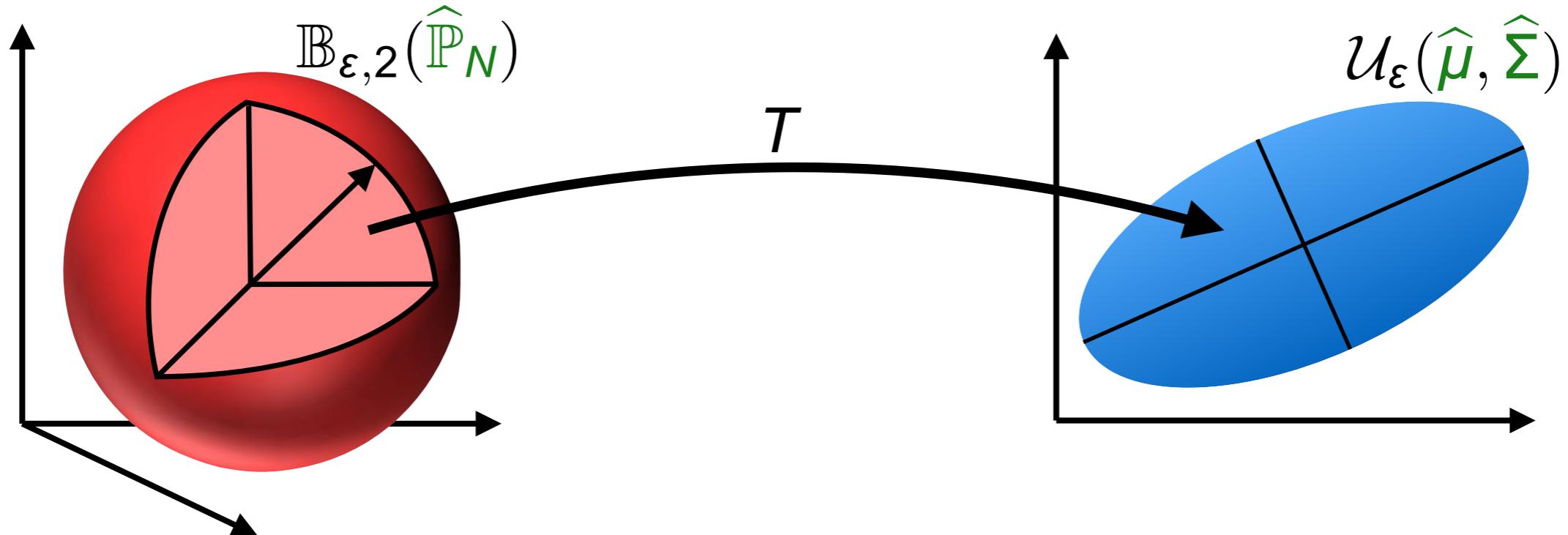


Gelbrich Bound ($p = 2$)

Uncertainty set for mean vectors and covariance matrices:

$$\mathcal{U}_\varepsilon(\hat{\mu}, \hat{\Sigma}) = \left\{ (\mu, \Sigma) \in \mathbb{R}^m \times \mathbb{S}_+^m : d_G \left((\hat{\mu}, \hat{\Sigma}), (\mu, \Sigma) \right) \leq \varepsilon^2 \right\}$$

Mean-covariance “projection”:



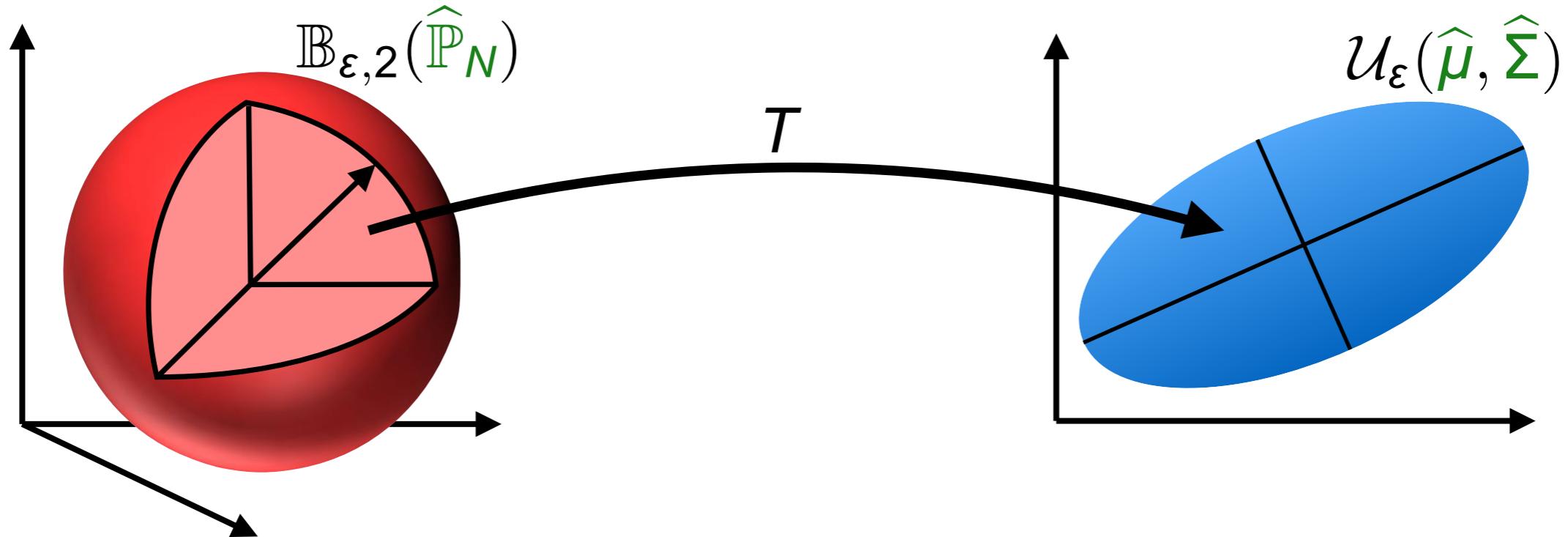
Note: $T(\mathbb{B}_{\varepsilon,2}(\hat{P}_N)) \subseteq \mathcal{U}_\varepsilon(\hat{\mu}, \hat{\Sigma})$

Gelbrich Bound ($p = 2$)

Uncertainty set for mean vectors and covariance matrices:

$$\mathcal{U}_\varepsilon(\hat{\mu}, \hat{\Sigma}) = \left\{ (\mu, \Sigma) \in \mathbb{R}^m \times \mathbb{S}_+^m : d_G \left((\hat{\mu}, \hat{\Sigma}), (\mu, \Sigma) \right) \leq \varepsilon^2 \right\}$$

Mean-covariance “projection”:



Note: $T(\mathbb{B}_{\varepsilon,2}(\hat{\mathbb{P}}_N)) = \mathcal{U}_\varepsilon(\hat{\mu}, \hat{\Sigma})$ if $\hat{\mathbb{P}}_N$ is elliptical!

Gelbrich Bound ($p = 2$)

The Gelbrich hull:

$$\mathbb{G}_\varepsilon(\hat{\mu}, \hat{\Sigma}) = T^{-1}(\mathcal{U}_\varepsilon(\hat{\mu}, \hat{\Sigma}))$$

Gelbrich Bound ($p = 2$)

The Gelbrich hull:

$$\mathbb{G}_\varepsilon(\hat{\mu}, \hat{\Sigma}) = T^{-1}(\mathcal{U}_\varepsilon(\hat{\mu}, \hat{\Sigma}))$$

Note: $\mathbb{B}_{\varepsilon,2}(\hat{\mathbb{P}}_N) \subseteq \mathbb{G}_\varepsilon(\hat{\mu}, \hat{\Sigma})$

Gelbrich Bound ($p = 2$)

The Gelbrich hull:

$$\mathbb{G}_\varepsilon(\hat{\mu}, \hat{\Sigma}) = T^{-1}(\mathcal{U}_\varepsilon(\hat{\mu}, \hat{\Sigma}))$$

Note: $\mathbb{B}_{\varepsilon,2}(\hat{\mathbb{P}}_N) \subseteq \mathbb{G}_\varepsilon(\hat{\mu}, \hat{\Sigma})$

Theorem: If $\hat{\mathbb{P}}_N$ has mean $\hat{\mu}$ and covariance matrix $\hat{\Sigma}$, then

$$\mathcal{R}_{\varepsilon,2}(\hat{\mathbb{P}}_N, \ell) \leq \sup_{\mathbb{Q} \in \mathbb{G}_\varepsilon(\hat{\mu}, \hat{\Sigma})} \mathbb{E}^{\mathbb{Q}}[\ell(\xi)].$$

Gelbrich Bound ($p = 2$)

The Gelbrich hull:

$$\mathbb{G}_\varepsilon(\hat{\mu}, \hat{\Sigma}) = T^{-1}(\mathcal{U}_\varepsilon(\hat{\mu}, \hat{\Sigma}))$$

Note: $\mathbb{B}_{\varepsilon,2}(\hat{\mathbb{P}}_N) \subseteq \mathbb{G}_\varepsilon(\hat{\mu}, \hat{\Sigma})$

Theorem: If $\hat{\mathbb{P}}_N$ has mean $\hat{\mu}$ and covariance matrix $\hat{\Sigma}$, then

$$\begin{aligned}\mathcal{R}_{\varepsilon,2}(\hat{\mathbb{P}}_N, \ell) &\leq \sup_{\mathbb{Q} \in \mathbb{G}_\varepsilon(\hat{\mu}, \hat{\Sigma})} \mathbb{E}^{\mathbb{Q}}[\ell(\xi)]. \\ &= \overline{\mathcal{R}}_\varepsilon(\hat{\mu}, \hat{\Sigma}, \ell)\end{aligned}$$

Efficient Computation of the Worst-Case Risk

Strong Duality

Theorem: If ℓ is u.s.c. and integrable under $\widehat{\mathbb{P}}_N$, then¹⁾

$$\mathcal{R}_{\varepsilon,p}(\widehat{\mathbb{P}}_N, \ell) = \inf_{\gamma \geq 0} \mathbb{E}^{\widehat{\mathbb{P}}_N} [\ell_\gamma(\xi)] + \gamma \varepsilon^p,$$

where $\ell_\gamma(\xi) = \sup_{z \in \Xi} \ell(z) - \gamma \|z - \xi\|^p$.

¹⁾ Mohajerin Esfahani & Kuhn, *Math. Program.*, 2018; Guan & Zhao, *Oper. Res. Lett.*, 2018; Blanchet & Murthy, *Math. Oper. Res.*, 2019; Gao & Kleywegt, *arxiv*, 2017.

Strong Duality

Theorem: If ℓ is u.s.c. and integrable under $\widehat{\mathbb{P}}_N$, then¹⁾

$$\mathcal{R}_{\varepsilon,p}(\widehat{\mathbb{P}}_N, \ell) = \inf_{\gamma \geq 0} \mathbb{E}^{\widehat{\mathbb{P}}_N} [\ell_\gamma(\xi)] + \gamma \varepsilon^p,$$

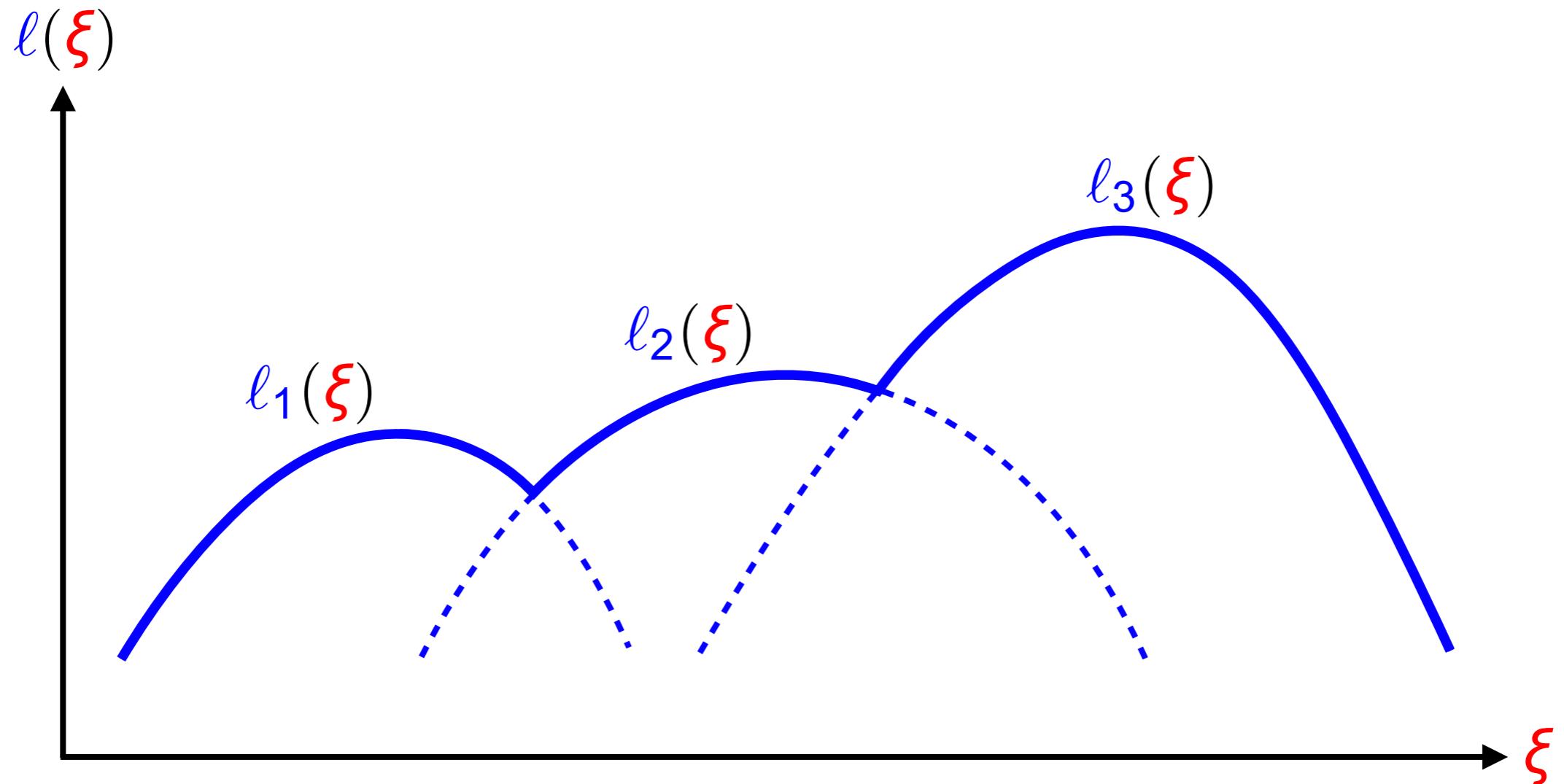
where $\ell_\gamma(\xi) = \sup_{z \in \Xi} \ell(z) - \gamma \|z - \xi\|^p$.



Moreau envelope

¹⁾ Mohajerin Esfahani & Kuhn, *Math. Program.*, 2018; Guan & Zhao, *Oper. Res. Lett.*, 2018; Blanchet & Murthy, *Math. Oper. Res.*, 2019; Gao & Kleywegt, *arxiv*, 2017.

Piecewise Concave Loss



$\ell(\xi) = \max_j \ell_j(\xi)$, where each ℓ_j is proper, concave and u.s.c.

Piecewise Concave Loss

Theorem: If Ξ is convex, $\widehat{\mathbb{P}}_N$ is the empirical distribution and $\ell(\xi)$ is piecewise concave, then

$$\begin{aligned} \mathcal{R}_{\varepsilon,p}(\widehat{\mathbb{P}}_N, \ell) &= \\ \inf \quad &\gamma \varepsilon^p + \frac{1}{N} \sum_{i=1}^N s_i \\ \text{s.t. } &\gamma \in \mathbb{R}_+, \quad s_i \in \mathbb{R}, \quad u_{ij} \in \mathbb{R}^m, \quad v_{ij} \in \mathbb{R}^m \quad \forall i, j \\ &[-\ell_j]^*(u_{ij} - v_{ij}) + \sigma_\Xi(v_{ij}) - u_{ij}^\top \widehat{\xi}_i + \varphi(q) \left\| \frac{u_{ij}}{\gamma} \right\|_*^q \leq s_i \quad \forall i, j, \end{aligned}$$

where $\frac{1}{p} + \frac{1}{q} = 1$ and $\varphi(q) = (q-1)^{q-1}/q^q$.

Piecewise Concave Loss

Finite convex optimization problem!
Size polynomial in # samples and # pieces!

$$\mathcal{R}_{\varepsilon,p}(\widehat{\mathbb{P}}_N, \ell) =$$

$$\inf \quad \gamma \varepsilon^p + \frac{1}{N} \sum_{i=1}^N s_i$$

$$\text{s.t. } \gamma \in \mathbb{R}_+, \quad s_i \in \mathbb{R}, \quad u_{ij} \in \mathbb{R}^m, \quad v_{ij} \in \mathbb{R}^m \quad \forall i, j$$

$$[-\ell_j]^*(u_{ij} - v_{ij}) + \sigma_\Xi(v_{ij}) - u_{ij}^\top \widehat{\xi}_i + \varphi(q) \underbrace{\gamma \left\| \frac{u_{ij}}{\gamma} \right\|_*^q}_{\text{convex conjugate of } -\ell_j} \leq s_i \quad \forall i, j$$

support function of Ξ
perspective of $\|\cdot\|_^q$*

Piecewise Concave Loss

Theorem: If Ξ is convex, $\widehat{\mathbb{P}}_N$ is the empirical distribution and $\ell(\xi)$ is piecewise concave, then

$$\mathcal{R}_{\varepsilon,p}(\widehat{\mathbb{P}}_N, \ell) = \max \quad \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^J a_{ij} \ell_j \left(\widehat{\xi}_i + \frac{\theta_{ij}}{a_{ij}} \right)$$

$$\text{s.t. } a_{ij} \in \mathbb{R}_+, \quad \theta_{ij} \in \mathbb{R}^m \quad \forall i \in [N], \quad \forall j \in [J]$$

$$\widehat{\xi}_i + \frac{\theta_{ij}}{a_{ij}} \in \Xi \quad \forall i \in [N], \quad \forall j \in [J]$$

$$\sum_{j=1}^J a_{ij} = 1 \quad \forall i \in [N]$$

$$\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^J a_{ij} \left\| \frac{\theta_{ij}}{a_{ij}} \right\|^p \leq \varepsilon^p.$$

Main Takeaways

- ▶ if ℓ is piecewise concave and $\widehat{\mathbb{P}}_N$ empirical, the worst-case risk can be computed by convex optimization for any $p \in [1, \infty]$;
 - ⇒ generalizes earlier results for $p = 1$ ¹⁾
- ▶ if ℓ is concave and $\widehat{\mathbb{P}}_N$ empirical, the robust lower bound is exact;
- ▶ the dual convex program
 - ▶ provides a finite reduction of the worst-case risk problem;
 - ▶ is always solvable;
 - ▶ can be used to find a worst-case distribution (may not exist for $p = 1$, always exists for $p > 1$).

¹⁾ Mohajerin Esfahani & Kuhn, *Math. Program.*, 2018.

Convex Loss

Theorem: If $\Xi = \mathbb{R}^m$ and $\ell(\xi)$ is convex, then

$$\mathcal{R}_{\varepsilon,1}(\widehat{\mathbb{P}}_N, \ell) = \mathcal{R}(\widehat{\mathbb{P}}_N, \ell) + \varepsilon \cdot \text{Lip}(\ell).$$

Worst-Case Risk for $p = 1$

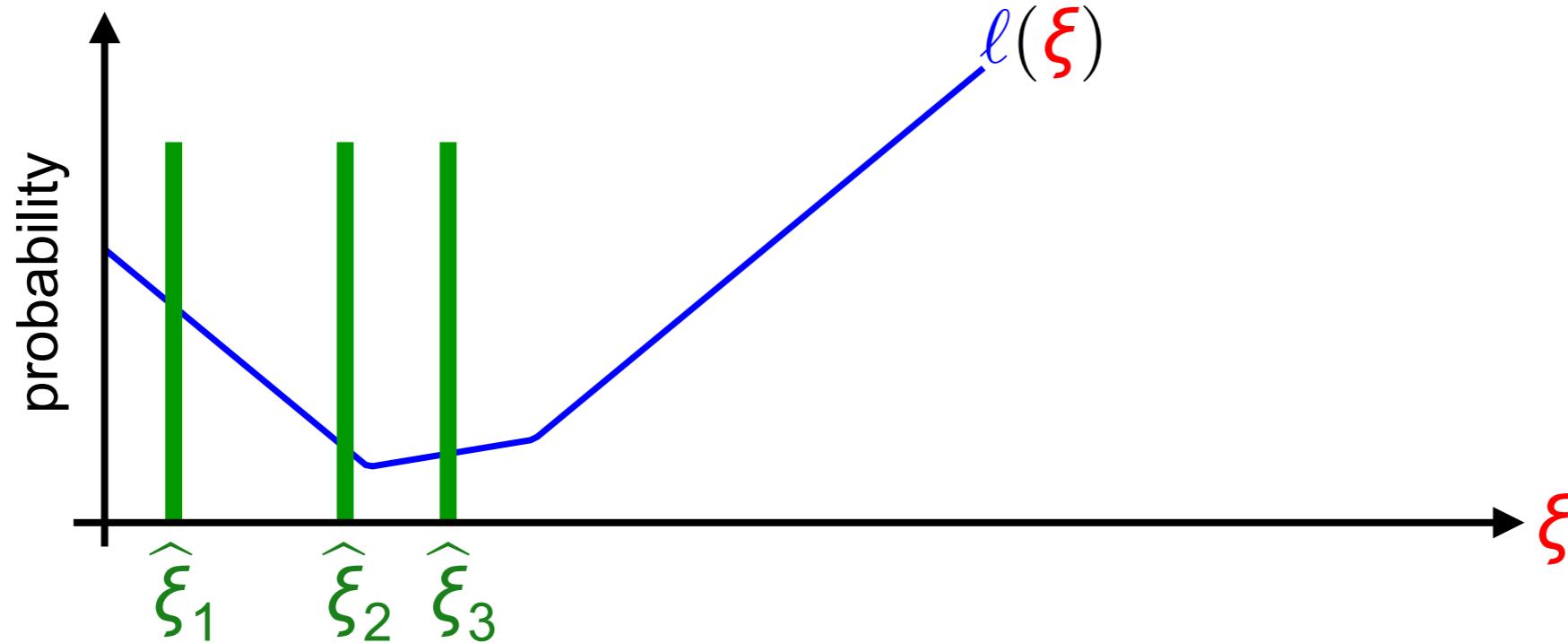


nominal distribution = empirical distribution

support set $\Xi = \mathbb{R}$

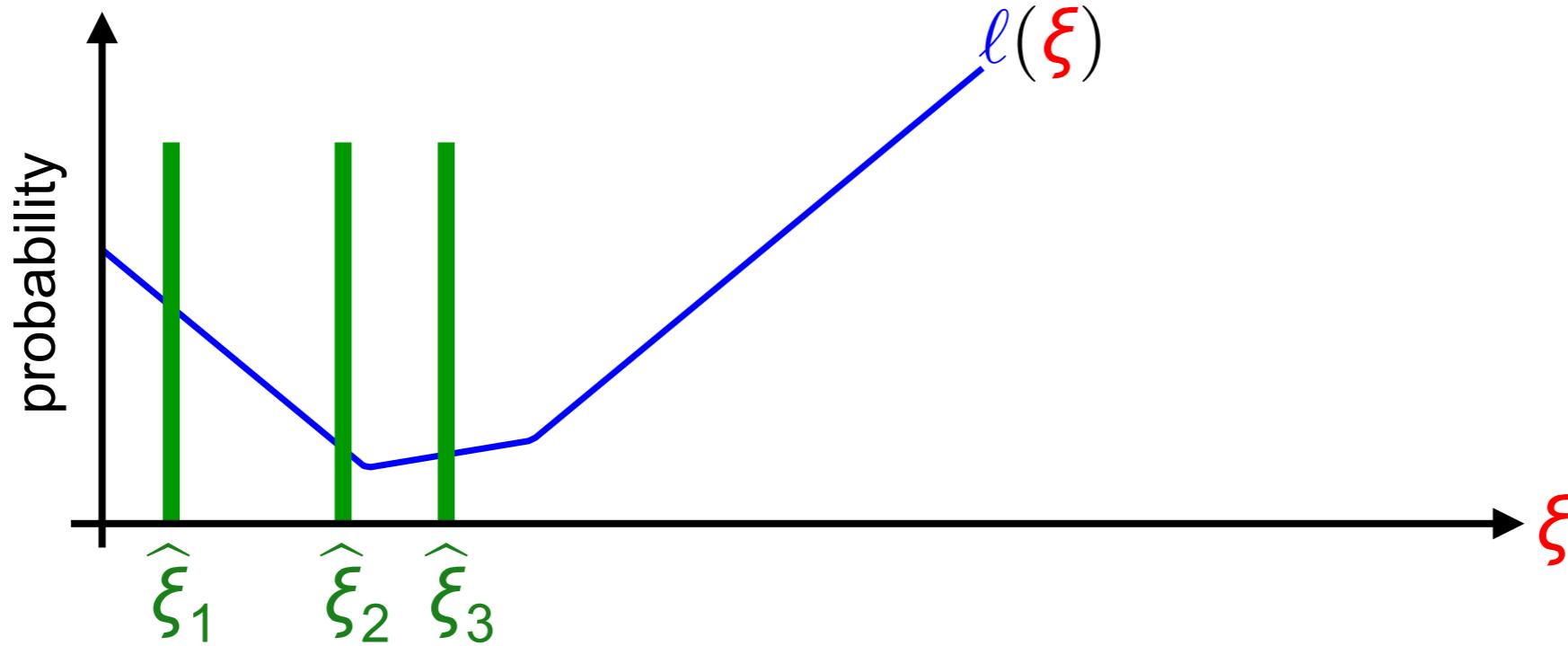
Wasserstein radius $\varepsilon = 0.7$

Worst-Case Risk for $p = 1$



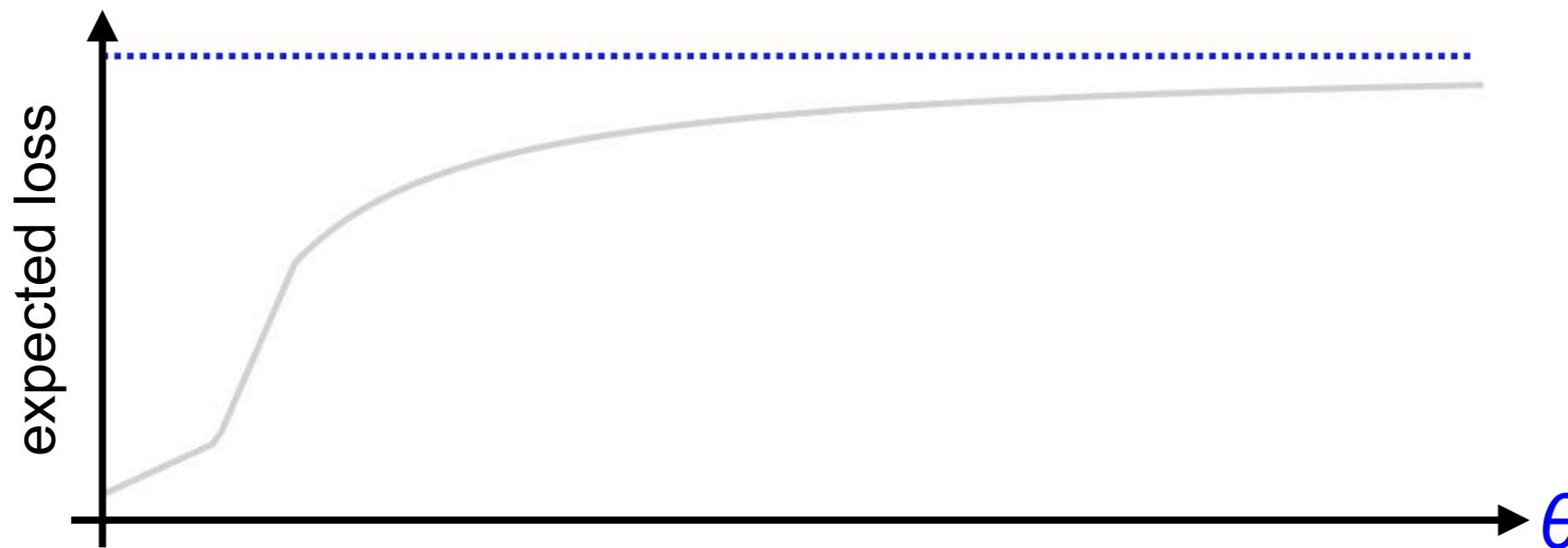
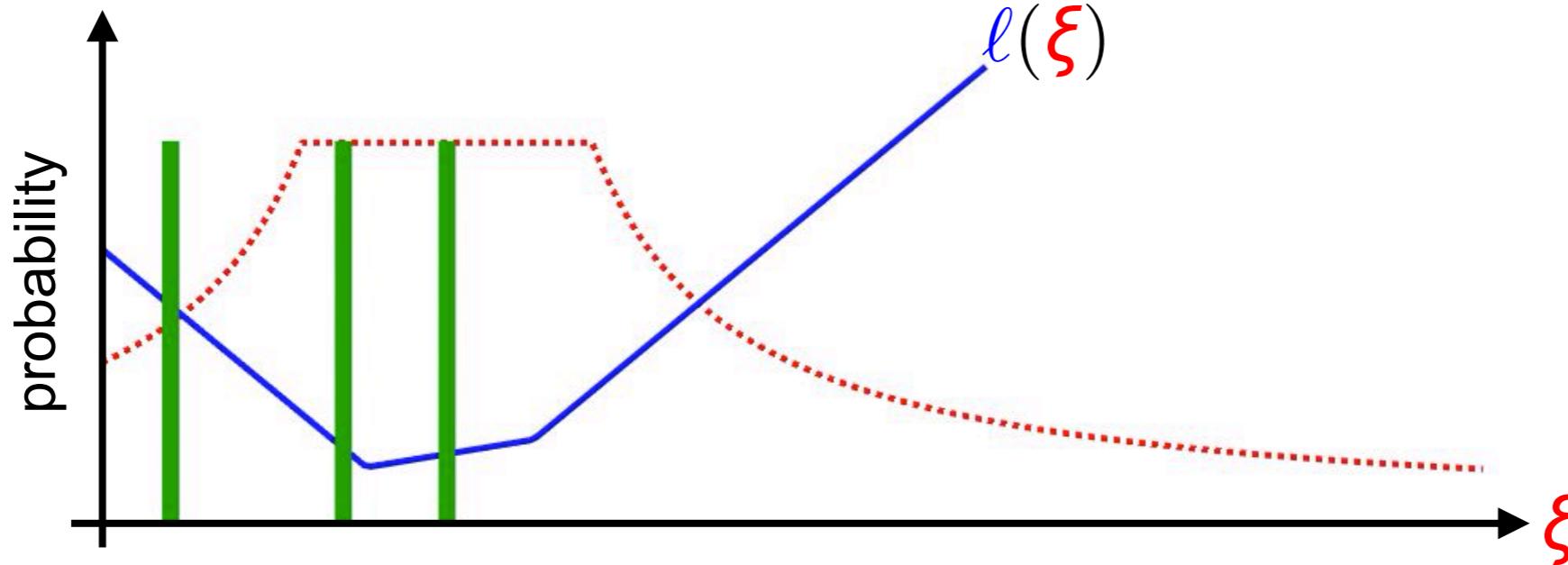
Convex and Lipschitz continuous loss

Worst-Case Risk for $p = 1$

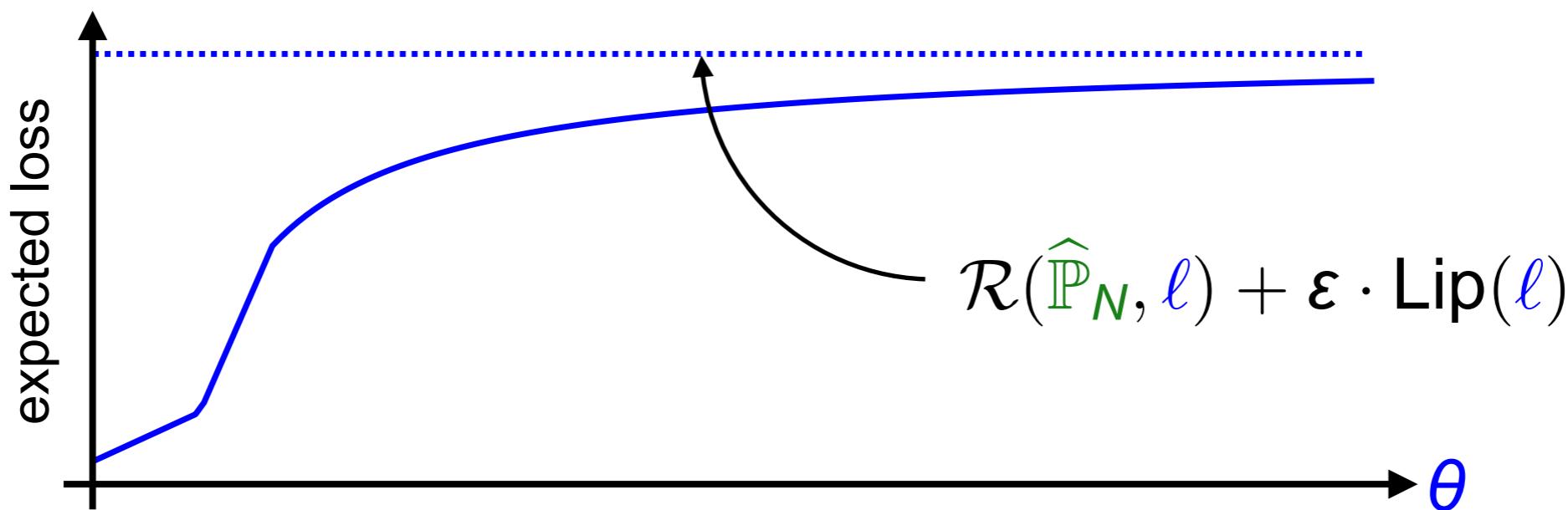
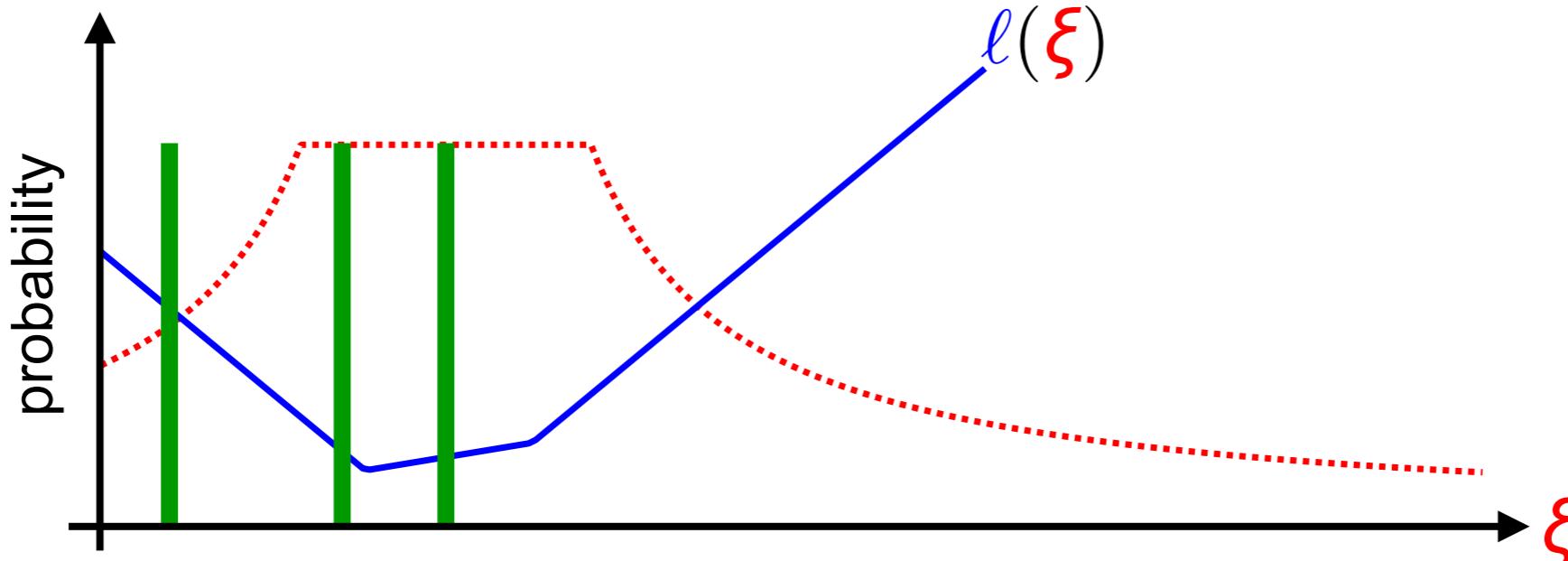


Increase $\hat{\xi}_3$ by θ (and reduce its probability if necessary)

Worst-Case Risk for $p = 1$



Worst-Case Risk for $p = 1$



Main Takeaways

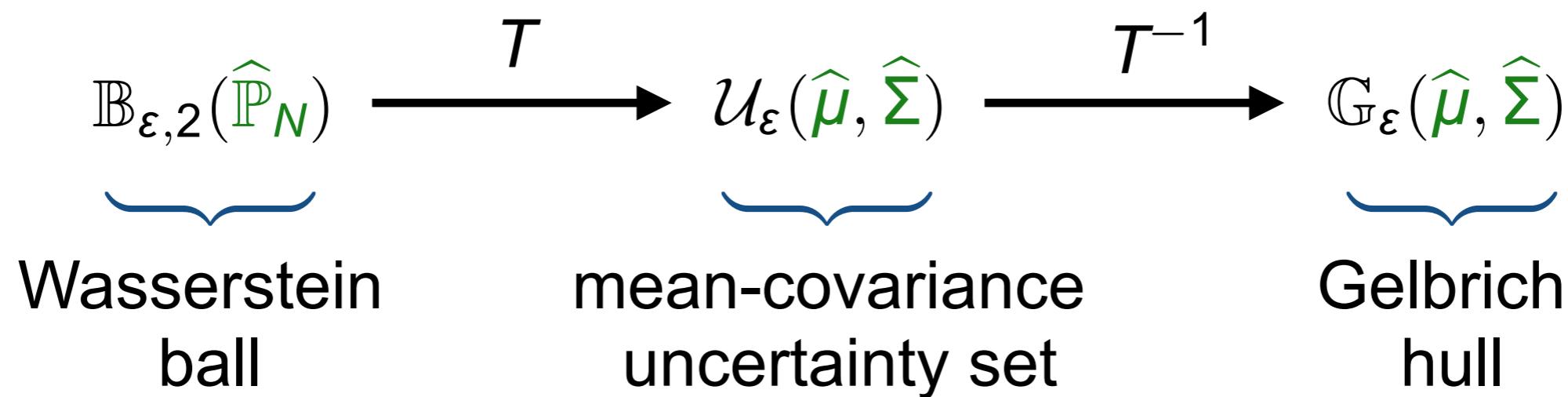
- ▶ if ℓ is convex, $\Xi = \mathbb{R}^m$ and $p = 1$, the worst case risk coincides with the Lipschitz-regularized nominal risk;
- ▶ no distribution attains the worst-case risk;
- ▶ the worst-case risk is attained asymptotically: push a linearly decreasing mass to infinity in the direction of steepest ascent;
- ▶ computing the worst-case risk is generally hard but tractable in special cases.

Computing the Gelbrich Bound

Recall: $\mathcal{R}_{\varepsilon,2}(\widehat{\mathbb{P}}_N, \ell) \leq \overline{\mathcal{R}}_{\varepsilon}(\widehat{\mu}, \widehat{\Sigma}, \ell)$

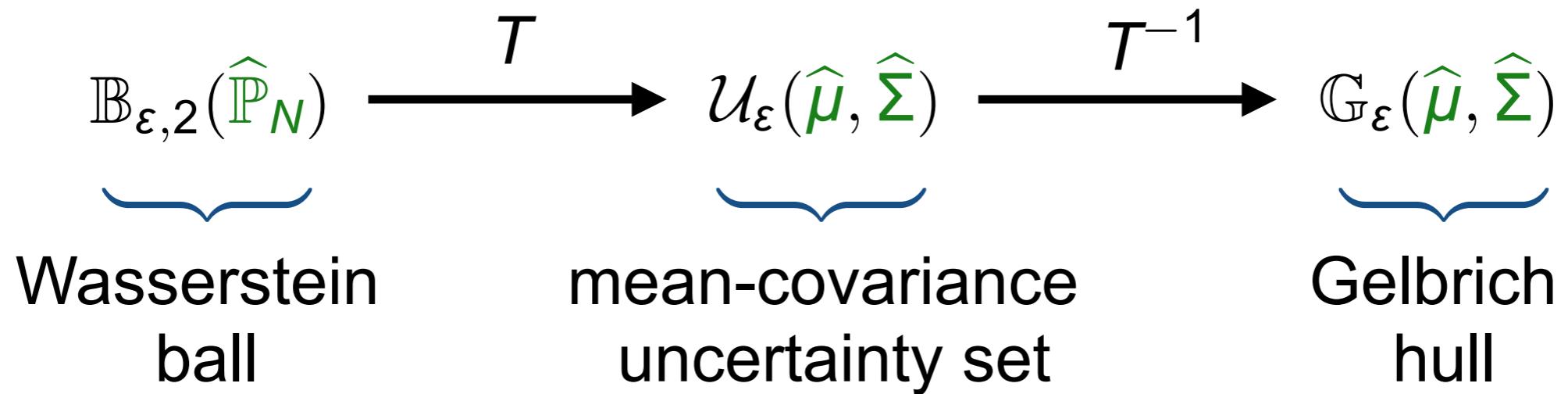
Computing the Gelbrich Bound

Recall: $\mathcal{R}_{\varepsilon,2}(\widehat{\mathbb{P}}_N, \ell) \leq \overline{\mathcal{R}}_{\varepsilon}(\widehat{\mu}, \widehat{\Sigma}, \ell)$



Computing the Gelbrich Bound

Recall: $\mathcal{R}_{\varepsilon,2}(\widehat{\mathbb{P}}_N, \ell) \leq \overline{\mathcal{R}}_{\varepsilon}(\widehat{\mu}, \widehat{\Sigma}, \ell)$



$$\mathbb{G}_{\varepsilon}(\widehat{\mu}, \widehat{\Sigma}) = \bigcup_{(\mu, \Sigma) \in \mathcal{U}_{\varepsilon}(\widehat{\mu}, \widehat{\Sigma})} \underbrace{\mathcal{P}(\mu, \Sigma)}_{\text{Chebyshev ambiguity set}}$$

Computing the Gelbrich Bound

Recall: $\mathcal{R}_{\varepsilon,2}(\widehat{\mathbb{P}}_N, \ell) \leq \overline{\mathcal{R}}_{\varepsilon}(\widehat{\mu}, \widehat{\Sigma}, \ell)$

$$\overline{\mathcal{R}}_{\varepsilon}(\widehat{\mu}, \widehat{\Sigma}, \ell) = \sup_{\mathbb{Q} \in \mathbb{G}_{\varepsilon}(\widehat{\mu}, \widehat{\Sigma})} \mathbb{E}^{\mathbb{Q}}[\ell(\xi)]$$

Computing the Gelbrich Bound

Recall: $\mathcal{R}_{\varepsilon,2}(\widehat{\mathbb{P}}_N, \ell) \leq \overline{\mathcal{R}}_{\varepsilon}(\widehat{\mu}, \widehat{\Sigma}, \ell)$

$$\overline{\mathcal{R}}_{\varepsilon}(\widehat{\mu}, \widehat{\Sigma}, \ell) = \sup_{(\mu, \Sigma) \in \mathcal{U}_{\varepsilon}(\widehat{\mu}, \widehat{\Sigma})} \sup_{\mathbb{Q} \in \mathcal{P}(\mu, \Sigma)} \mathbb{E}^{\mathbb{Q}}[\ell(\xi)]$$

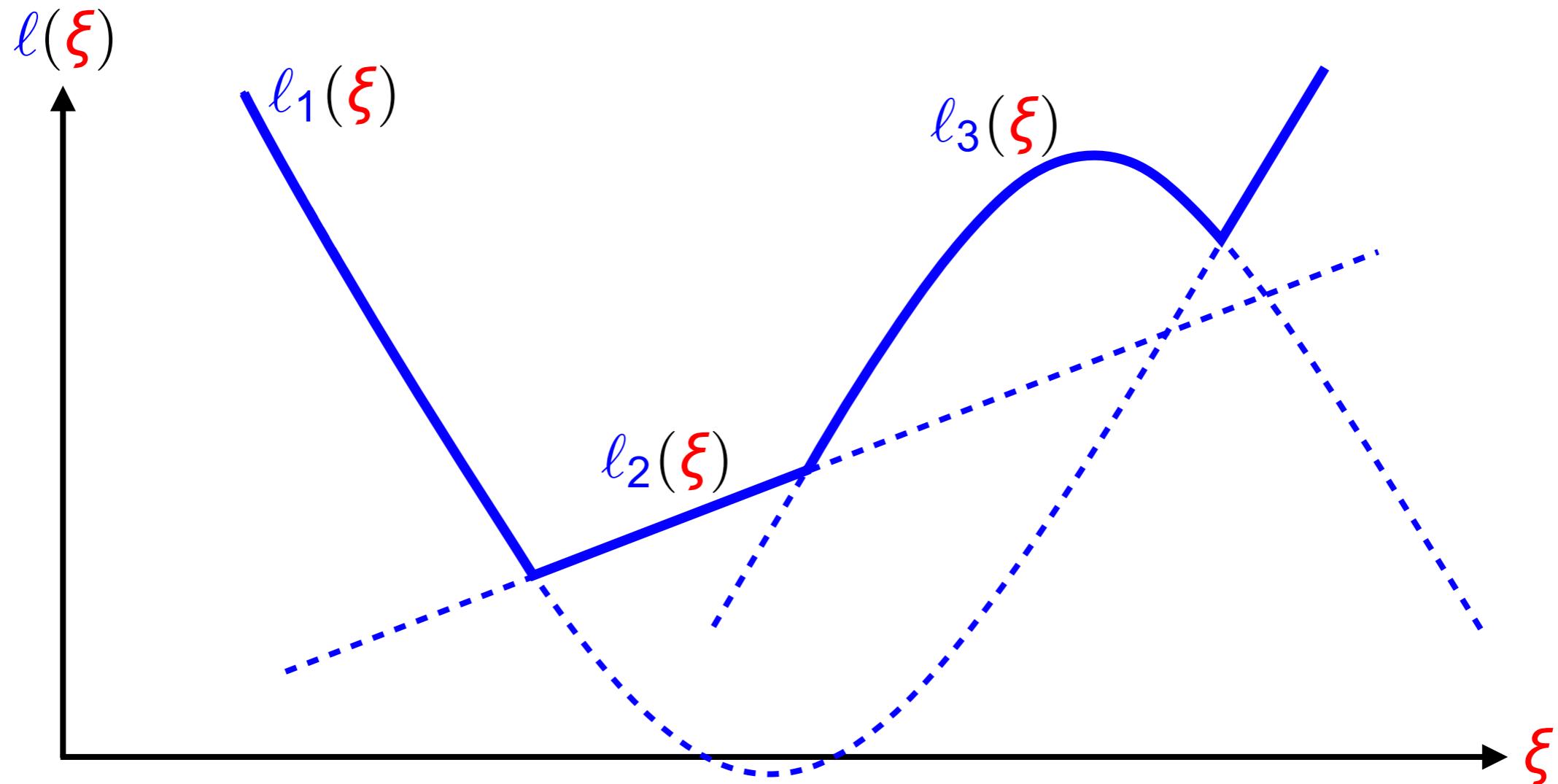
Computing the Gelbrich Bound

Recall: $\mathcal{R}_{\varepsilon,2}(\widehat{\mathbb{P}}_N, \ell) \leq \overline{\mathcal{R}}_{\varepsilon}(\widehat{\mu}, \widehat{\Sigma}, \ell)$

$$\overline{\mathcal{R}}_{\varepsilon}(\widehat{\mu}, \widehat{\Sigma}, \ell) = \sup_{(\mu, \Sigma) \in \mathcal{U}_{\varepsilon}(\widehat{\mu}, \widehat{\Sigma})} \sup_{\mathbb{Q} \in \mathcal{P}(\mu, \Sigma)} \mathbb{E}^{\mathbb{Q}}[\ell(\xi)]$$

2nd layer Chebyshev risk
of robustness

Piecewise Quadratic Loss



$$\ell(\xi) = \max_j \ell_j(\xi), \text{ where } \ell_j(\xi) = \xi^\top Q_j \xi + 2q_j^\top \xi + q_j^0$$

Piecewise Quadratic Loss

Theorem: If $\Xi = \mathbb{R}^m$, and $\ell(\xi)$ is piecewise quadratic, then the Gelbrich risk satisfies

$$\begin{aligned} \overline{\mathcal{R}}_\varepsilon(\hat{\mu}, \hat{\Sigma}, \ell) &= \\ \inf \quad & y_0 + \gamma \left(\varepsilon^2 - \|\hat{\mu}\|_2^2 - \text{Tr}[\hat{\Sigma}] \right) + z + \text{Tr}[Z] \\ \text{s.t. } & y \in \mathbb{R}_+, \quad y_0 \in \mathbb{R}, \quad y \in \mathbb{R}^m, \quad Y \in \mathbb{S}^m, \quad x \in \mathbb{R}_+, \quad Z \in \mathbb{S}_+^m \\ & \begin{bmatrix} yI - Y & y + \gamma\hat{\mu} \\ y^\top + \gamma\hat{\mu}^\top & z \end{bmatrix} \succeq 0, \quad \begin{bmatrix} yI - Y & y\hat{\Sigma}^{\frac{1}{2}} \\ \gamma\hat{\Sigma}^{\frac{1}{2}} & Z \end{bmatrix} \succeq 0 \\ & \begin{bmatrix} Y - Q_j & y - q_j \\ y^\top - q_j^\top & y_0 - q_j^0 \end{bmatrix} \succeq 0 \quad \forall j \in [J]. \end{aligned}$$

Piecewise Quadratic Loss

Theorem: If $\Xi = \mathbb{R}^m$, and $\ell(\xi)$ is piecewise quadratic, then the Gelbrich risk satisfies

$$\overline{\mathcal{R}}_\varepsilon(\hat{\mu}, \hat{\Sigma}, \ell) =$$

$$\max \sum_{j=1}^J \text{Tr}[Q_j \Theta_j] + 2q_j^\top \theta_j + q_j^0 a_j$$

$$\text{s.t. } \mu \in \mathbb{R}^m, \Sigma \in \mathbb{S}_+^m, a_j \in \mathbb{R}_+, \theta_j \in \mathbb{R}^m, \Theta_j \in \mathbb{S}_+^m \quad \forall j \in [J]$$

$$\begin{bmatrix} \Theta_j & \theta_j \\ \theta_j^\top & a_j \end{bmatrix} \succeq 0 \quad \forall j \in [J]$$

$$\sum_{j=1}^J a_j = 1, \quad \sum_{j=1}^J \theta_j = \mu, \quad \sum_{j=1}^J \Theta_j = \Sigma + \mu\mu^\top, \quad (\mu, \Sigma) \in \mathcal{U}_\varepsilon(\hat{\mu}, \hat{\Sigma}).$$

Properties of the Gelbrich Bound

- ▶ bounds the Wasserstein risk and inherits its statistical guarantees;
- ▶ can be computed via convex optimization;
- ▶ induces the uncertainty set $\mathcal{U}_\varepsilon(\hat{\mu}, \hat{\Sigma})$ for μ and Σ
 \implies outer robust program *convex* for *any* loss function;
- ▶ exact if ℓ is quadratic and $\hat{\mathbb{P}}_N$ elliptical;
- ▶ worst-case distribution available in closed form if ℓ is quadratic and $\hat{\mathbb{P}}_N$ elliptical; hard to compute if ℓ is piecewise quadratic.

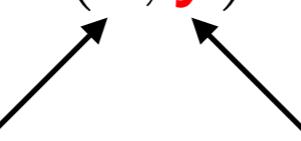
Applications in Machine Learning

Classification

Random vector: $\xi = (x, y) \sim \mathbb{P}$

input $\in \mathbb{R}^n$

output $\in \{-1, +1\}$



Classification

Random vector: $\xi = (x, y) \sim \mathbb{P}$

Goal: predict output from input



Classification

Random vector: $\xi = (x, y) \sim \mathbb{P}$

Goal: predict output from input



Ideal prediction:
$$\begin{cases} +1 & \text{if } \mathbb{P}[y = +1 | x] \geq 50\% \\ -1 & \text{if } \mathbb{P}[y = -1 | x] > 50\% \end{cases}$$

Classification

Random vector: $\xi = (x, y) \sim \mathbb{P}$

Reality:

- ▶ \mathbb{P} unknown
- ▶ we are given training samples $\hat{\xi}_1, \hat{\xi}_2, \dots, \hat{\xi}_N \sim \mathbb{P}^N$

Classification

Random vector: $\xi = (\mathbf{x}, \mathbf{y}) \sim \mathbb{P}$

Reality:

- ▶ \mathbb{P} unknown
- ▶ we are given training samples $\hat{\xi}_1, \hat{\xi}_2, \dots, \hat{\xi}_N \sim \mathbb{P}^N$

Practical approach: predict output as $\text{sign}(\mathbf{w}^\top \mathbf{x})$

\implies prediction correct if $\mathbf{y} \cdot \mathbf{w}^\top \mathbf{x} > 0$

Classification

Random vector: $\xi = (\mathbf{x}, \mathbf{y}) \sim \mathbb{P}$

Reality:

- ▶ \mathbb{P} unknown
- ▶ we are given training samples $\hat{\xi}_1, \hat{\xi}_2, \dots, \hat{\xi}_N \sim \mathbb{P}^N$

Practical approach: predict output as $\text{sign}(\mathbf{w}^\top \mathbf{x})$

\implies prediction error $L(\mathbf{y} \cdot \mathbf{w}^\top \mathbf{x})$

Classification

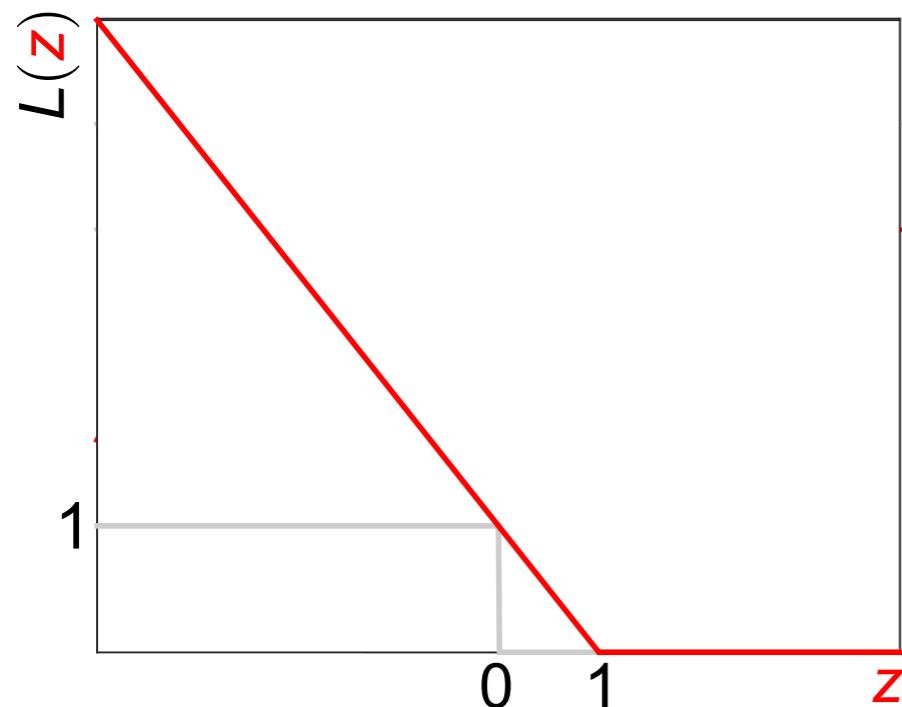
Random vector: $\xi = (\mathbf{x}, y) \sim \mathbb{P}$

Reality:

- ▶ \mathbb{P} unknown
- ▶ we are given training samples $\hat{\xi}_1, \hat{\xi}_2, \dots, \hat{\xi}_N \sim \mathbb{P}^N$

Practical approach: predict output as $\text{sign}(\mathbf{w}^\top \mathbf{x})$

⇒ prediction error $L(y \cdot \mathbf{w}^\top \mathbf{x})$



Support vector machine:

Hinge loss:

$$L(z) = \max\{0, 1 - z\}$$

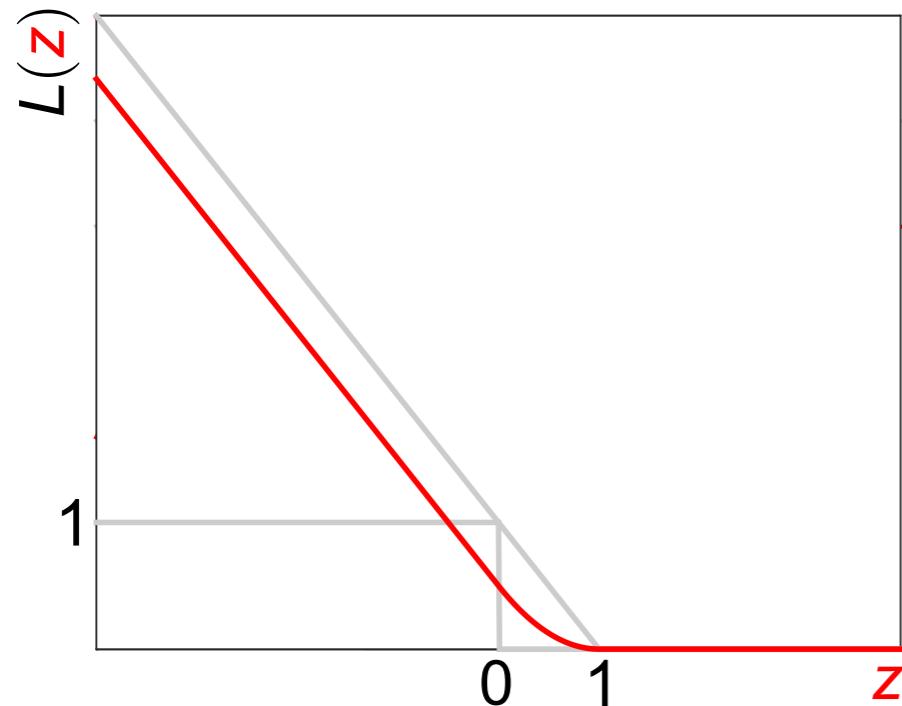
Classification

Random vector: $\xi = (x, y) \sim \mathbb{P}$

Reality:

- ▶ \mathbb{P} unknown
- ▶ we are given training samples $\hat{\xi}_1, \hat{\xi}_2, \dots, \hat{\xi}_N \sim \mathbb{P}^N$

Practical approach: predict output as $\text{sign}(w^\top x)$
 \implies prediction error $L(y \cdot w^\top x)$



Smooth support vector machine:

Smooth hinge loss:

$$L(z) = \begin{cases} \frac{1}{2} - z & \text{if } z \leq 0 \\ \frac{1}{2}(1 - z)^2 & \text{if } 0 < z < 1 \\ 0 & \text{else} \end{cases}$$

Classification

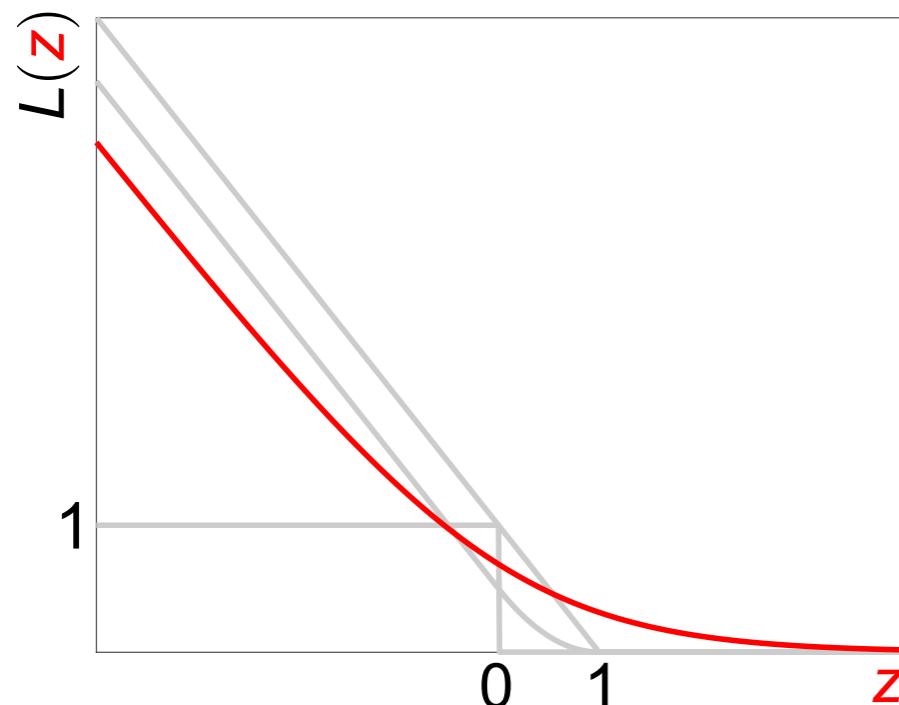
Random vector: $\xi = (x, y) \sim \mathbb{P}$

Reality:

- ▶ \mathbb{P} unknown
- ▶ we are given training samples $\hat{\xi}_1, \hat{\xi}_2, \dots, \hat{\xi}_N \sim \mathbb{P}^N$

Practical approach: predict output as $\text{sign}(w^\top x)$

⇒ prediction error $L(y \cdot w^\top x)$



Logistic regression:

Logloss:

$$L(z) = \log(1 + e^{-z})$$

Classification

Empirical risk minimization:

$$\inf_{\mathbf{w} \in \mathbb{R}^n} \mathbb{E}^{\widehat{\mathbb{P}}_N} [L(\mathbf{y} \cdot \mathbf{w}^\top \mathbf{x})]$$

Classification

Empirical risk minimization:

$$\inf_{w \in \mathbb{R}^n} \underbrace{\mathbb{E}^{\widehat{\mathbb{P}}_N}[L(y \cdot w^\top x)]}_{\text{nominal risk}}$$

- ▶ replace true risk with nominal risk
- ▶ $\widehat{\mathbb{P}}_N$ = empirical distribution
- ▶ susceptible to overfitting

Classification

Distributionally robust classification:

$$\inf_{\mathbf{w} \in \mathbb{R}^n} \sup_{\mathbb{Q} \in \mathbb{B}_{\varepsilon, 1}(\hat{\mathbb{P}}_N)} \mathbb{E}^{\mathbb{Q}}[L(\mathbf{y} \cdot \mathbf{w}^\top \mathbf{x})]$$

Classification

Distributionally robust classification:

$$\inf_{\mathbf{w} \in \mathbb{R}^n} \sup_{Q \in \mathbb{B}_{\varepsilon, 1}(\hat{\mathbb{P}}_N)} \mathbb{E}^Q [L(y \cdot \mathbf{w}^\top \mathbf{x})]$$

worst-case risk

- replace true risk with worst-case risk
- norm on input-output space: $\|\xi\| = \|\mathbf{x}\| + \frac{\kappa}{2}|y|$
- κ = cost of flipping a label

Classification

Distributionally robust classification:

$$\inf_{\mathbf{w} \in \mathbb{R}^n} \sup_{\mathbb{Q} \in \mathbb{B}_{\varepsilon, 1}(\widehat{\mathbb{P}}_N)} \mathbb{E}^{\mathbb{Q}}[L(\mathbf{y} \cdot \mathbf{w}^\top \mathbf{x})] \quad (*)$$

Theorem: If $\kappa = \infty$ and $p = 1$ then¹⁾

$$(*) \equiv \inf_{\mathbf{w} \in \mathbb{R}^n} \mathbb{E}^{\widehat{\mathbb{P}}_N}[L(\mathbf{y} \cdot \mathbf{w}^\top \mathbf{x})] + \varepsilon \cdot \text{Lip}(L) \cdot \|\mathbf{w}\|_*.$$

¹⁾ Shafeezadeh-Abadeh, Kuhn & Mohajerin Esfahani, *J. Mach. Learn. Res.*, 2019; Blanchet, Kang & Murthy, *arXiv*, 2016; Gao, Chen & Kleywegt, *arXiv*, 2017.

Distributionally robust classification:

$$\inf_{\mathbf{w} \in \mathbb{R}^n} \sup_{Q \in \mathbb{B}_{\varepsilon, 1}(\hat{\mathbb{P}}_N)} \mathbb{E}^Q [L(\mathbf{y} \cdot \mathbf{w}^\top \mathbf{x})] \quad (*)$$

Theorem: If $\kappa = \infty$ and $L(\mathbf{z})$ is a standard loss function, then

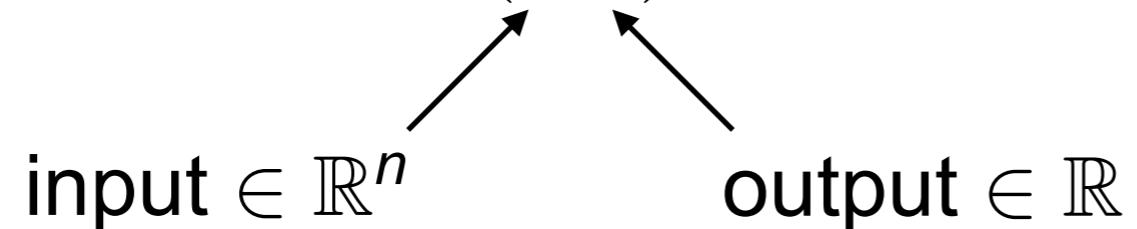
$$(*) \equiv \underbrace{\inf_{\mathbf{w} \in \mathbb{R}^n} \mathbb{E}^{\hat{\mathbb{P}}_N} [L(\mathbf{y} \cdot \mathbf{w}^\top \mathbf{x})]}_{\text{empirical risk}} + \underbrace{\varepsilon \cdot \text{Lip}(L) \cdot \|\mathbf{w}\|_*}_{\text{regularization term}}.$$

empirical risk regularization
term

¹⁾ Shafeezadeh-Abadeh, Kuhn & Mohajerin Esfahani, *J. Mach. Learn. Res.*, 2019; Blanchet, Kang & Murthy, *arXiv*, 2016; Gao, Chen & Kleywegt, *arXiv*, 2017.

Regression

Random vector: $\xi = (x, y) \sim \mathbb{P}$



Regression

Random vector: $\xi = (x, y) \sim \mathbb{P}$

Goal: predict output from input



Regression

Random vector: $\xi = (x, y) \sim \mathbb{P}$

Goal: predict output from input



Ideal prediction: $\mathbb{E}^{\mathbb{P}}[y | x]$

Regression

Random vector: $\xi = (x, y) \sim \mathbb{P}$

Reality:

- ▶ \mathbb{P} unknown
- ▶ we are given training samples $\hat{\xi}_1, \hat{\xi}_2, \dots, \hat{\xi}_N \sim \mathbb{P}^N$

Regression

Random vector: $\xi = (\mathbf{x}, \mathbf{y}) \sim \mathbb{P}$

Reality:

- ▶ \mathbb{P} unknown
- ▶ we are given training samples $\hat{\xi}_1, \hat{\xi}_2, \dots, \hat{\xi}_N \sim \mathbb{P}^N$

Practical approach: predict output as $\mathbf{w}^\top \mathbf{x}$

\implies prediction error $L(\mathbf{w}^\top \mathbf{x} - \mathbf{y})$

Regression

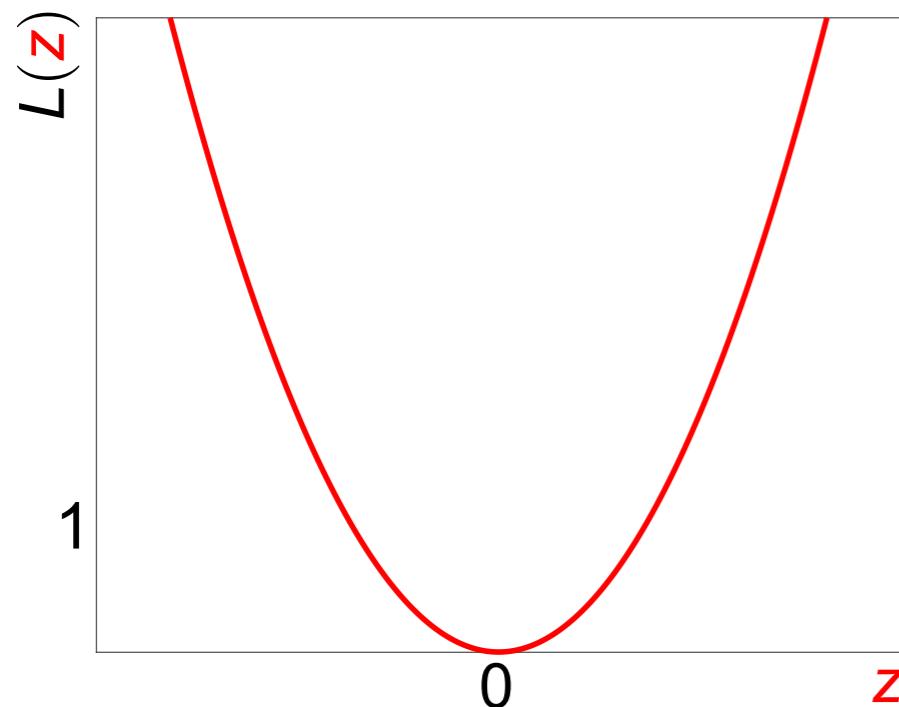
Random vector: $\xi = (\mathfrak{x}, \mathfrak{y}) \sim \mathbb{P}$

Reality:

- ▶ \mathbb{P} unknown
- ▶ we are given training samples $\hat{\xi}_1, \hat{\xi}_2, \dots, \hat{\xi}_N \sim \mathbb{P}^N$

Practical approach: predict output as $\mathfrak{w}^\top \mathfrak{x}$

⇒ prediction error $L(\mathfrak{w}^\top \mathfrak{x} - \mathfrak{y})$



Ordinary least squares:

Square loss:

$$L(\mathfrak{z}) = \frac{1}{2}\mathfrak{z}^2$$

Regression

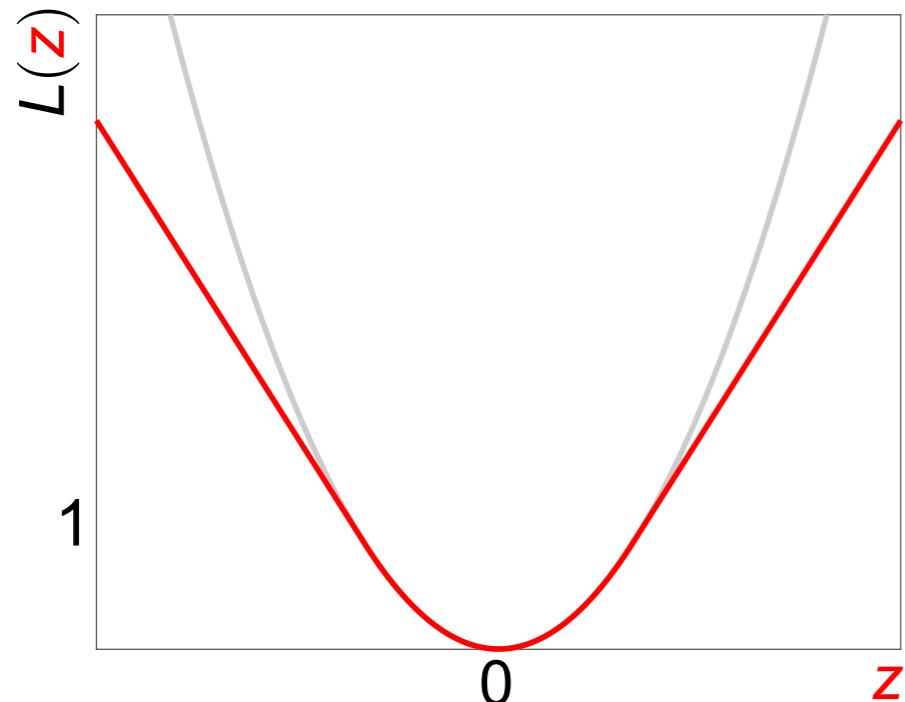
Random vector: $\xi = (x, y) \sim \mathbb{P}$

Reality:

- ▶ \mathbb{P} unknown
- ▶ we are given training samples $\hat{\xi}_1, \hat{\xi}_2, \dots, \hat{\xi}_N \sim \mathbb{P}^N$

Practical approach: predict output as $w^\top x$

⇒ prediction error $L(w^\top x - y)$



Robust regression:

Huber loss:

$$L(z) = \begin{cases} \frac{1}{2}z^2 & \text{if } |z| \leq \delta \\ \delta(|z| - \frac{1}{2}\delta) & \text{else} \end{cases}$$

Regression

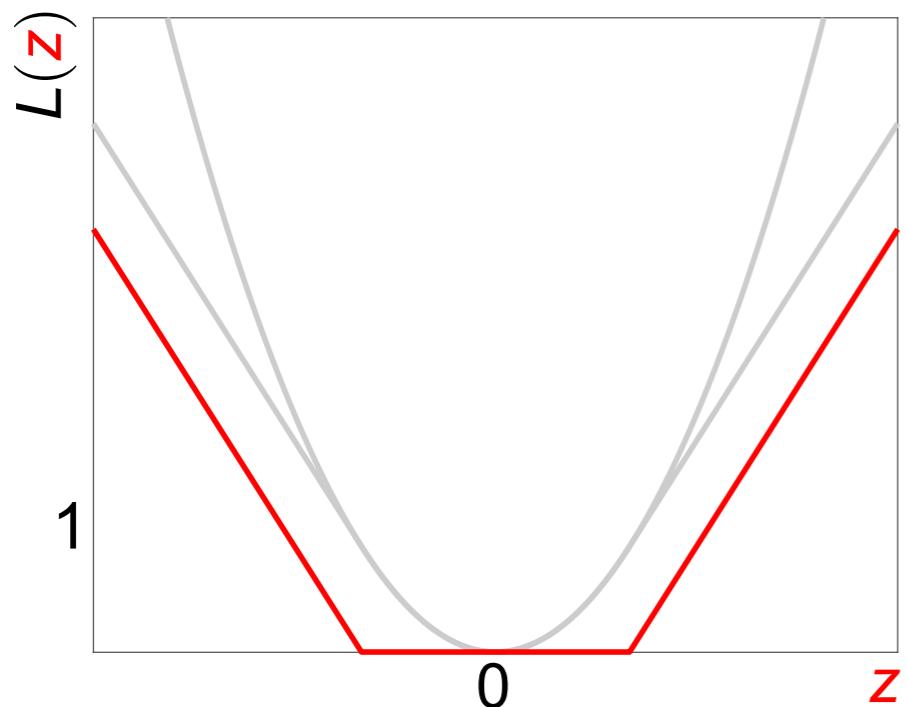
Random vector: $\xi = (x, y) \sim \mathbb{P}$

Reality:

- ▶ \mathbb{P} unknown
- ▶ we are given training samples $\hat{\xi}_1, \hat{\xi}_2, \dots, \hat{\xi}_N \sim \mathbb{P}^N$

Practical approach: predict output as $w^\top x$

⇒ prediction error $L(w^\top x - y)$



Support vector regression:

δ -insensitive loss:

$$L(z) = \max\{0, |z| - \varepsilon\}$$

Regression

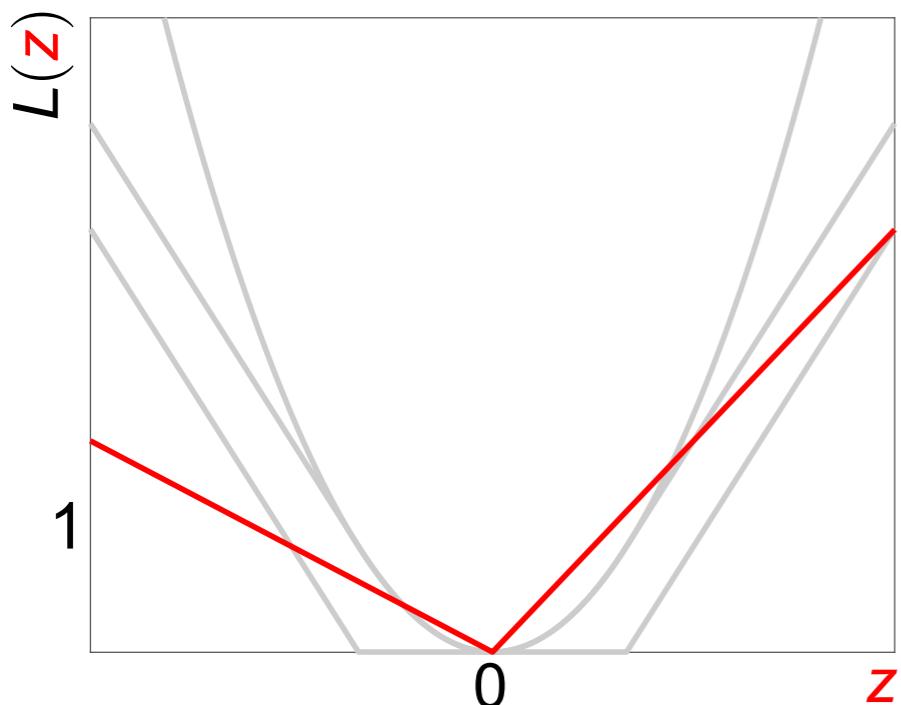
Random vector: $\xi = (x, y) \sim \mathbb{P}$

Reality:

- ▶ \mathbb{P} unknown
- ▶ we are given training samples $\hat{\xi}_1, \hat{\xi}_2, \dots, \hat{\xi}_N \sim \mathbb{P}^N$

Practical approach: predict output as $w^\top x$

⇒ prediction error $L(w^\top x - y)$



Quantile regression:

Pinball loss:

$$L(z) = \max\{-\tau z, (1-\tau)z\}$$

Regression

Empirical risk minimization:

$$\inf_{\mathbf{w} \in \mathbb{R}^n} \mathbb{E}^{\widehat{\mathbb{P}}_N} [L(\mathbf{w}^\top \mathbf{x} - y)]$$

Regression

Distributionally robust regression:

$$\inf_{\mathbf{w} \in \mathbb{R}^n} \sup_{\mathbb{Q} \in \mathbb{B}_{\varepsilon, p}(\widehat{\mathbb{P}}_N)} \mathbb{E}^{\mathbb{Q}}[L(\mathbf{w}^\top \mathbf{x} - y)]$$

Regression

Distributionally robust regression:

$$\inf_{\mathbf{w} \in \mathbb{R}^n} \sup_{\mathbb{Q} \in \mathbb{B}_{\varepsilon, p}(\widehat{\mathbb{P}}_N)} \mathbb{E}^{\mathbb{Q}}[L(\mathbf{w}^\top \mathbf{x} - y)]$$

- ▶ norm on input-output space: $\|\xi\| = \|\mathbf{x}\| + \frac{\kappa}{2}|y|$

Regression

Distributionally robust regression:

$$\inf_{\mathbf{w} \in \mathbb{R}^n} \sup_{\mathbb{Q} \in \mathbb{B}_{\varepsilon, p}(\widehat{\mathbb{P}}_N)} \mathbb{E}^{\mathbb{Q}}[L(\mathbf{w}^\top \mathbf{x} - y)] \quad (*)$$

Theorem: If $\kappa = \infty$ and $p = 1$, then¹⁾

$$(*) \equiv \inf_{\mathbf{w} \in \mathbb{R}^n} \mathbb{E}^{\widehat{\mathbb{P}}_N}[L(\mathbf{w}^\top \mathbf{x} - y)] + \varepsilon \cdot \text{Lip}(L) \cdot \|\mathbf{w}\|_*.$$

¹⁾ Shafieezadeh-Abadeh, Kuhn & Mohajerin Esfahani, *J. Mach. Learn. Res.*, 2019; Blanchet, Kang & Murthy, *arXiv*, 2016; Gao, Chen & Kleywegt, *arXiv*, 2017.

Regression

Distributionally robust regression:

$$\inf_{\mathbf{w} \in \mathbb{R}^n} \sup_{\mathbb{Q} \in \mathbb{B}_{\varepsilon, p}(\widehat{\mathbb{P}}_N)} \mathbb{E}^{\mathbb{Q}}[L(\mathbf{w}^\top \mathbf{x} - y)] \quad (*)$$

Theorem: If $\kappa = \infty$ and $p = 1$, then¹⁾

$$(*) \equiv \underbrace{\inf_{\mathbf{w} \in \mathbb{R}^n} \mathbb{E}^{\widehat{\mathbb{P}}_N}[L(\mathbf{w}^\top \mathbf{x} - y)]}_{\text{empirical risk}} + \underbrace{\varepsilon \cdot \text{Lip}(L) \cdot \|\mathbf{w}\|_*}_{\text{regularization term}}.$$

¹⁾ Shafeezadeh-Abadeh, Kuhn & Mohajerin Esfahani, *J. Mach. Learn. Res.*, 2019; Blanchet, Kang & Murthy, *arXiv*, 2016; Gao, Chen & Kleywegt, *arXiv*, 2017.

Regression

Distributionally robust regression:

$$\inf_{\mathbf{w} \in \mathbb{R}^n} \sup_{\mathbb{Q} \in \mathbb{B}_{\varepsilon, p}(\widehat{\mathbb{P}}_N)} \mathbb{E}^{\mathbb{Q}}[L(\mathbf{w}^\top \mathbf{x} - y)] \quad (*)$$

Theorem: If $\kappa = \infty$, $p = 2$ and $L(z)$ is the square error, then¹⁾

$$(*) \equiv \inf_{\mathbf{w} \in \mathbb{R}^n} \left[\left(\mathbb{E}^{\widehat{\mathbb{P}}_N}[L(\mathbf{w}^\top \mathbf{x} - y)] \right)^{\frac{1}{2}} + \varepsilon \cdot \|\mathbf{w}\|_* \right]^2.$$

¹⁾ Blanchet, Kang & Murthy, arXiv, 2016.

Regression

Distributionally robust regression:

$$\inf_{\mathbf{w} \in \mathbb{R}^n} \sup_{\mathbb{Q} \in \mathbb{B}_{\varepsilon, p}(\widehat{\mathbb{P}}_N)} \mathbb{E}^{\mathbb{Q}}[L(\mathbf{w}^\top \mathbf{x} - y)] \quad (*)$$

Theorem: If $\kappa = \infty$, $p = 2$ and $L(z)$ is the square error, then¹⁾

$$(*) \equiv \inf_{\mathbf{w} \in \mathbb{R}^n} \left[\underbrace{\left(\mathbb{E}^{\widehat{\mathbb{P}}_N}[L(\mathbf{w}^\top \mathbf{x} - y)] \right)^{\frac{1}{2}}}_{\text{square root of empirical risk}} + \varepsilon \cdot \|\mathbf{w}\|_* \underbrace{\vphantom{\left(\mathbb{E}^{\widehat{\mathbb{P}}_N}[L(\mathbf{w}^\top \mathbf{x} - y)] \right)^{\frac{1}{2}}}^2}_{\text{regularization term}} \right]^2.$$

⇒ generalized LASSO estimation problem

¹⁾ Blanchet, Kang & Murthy, arXiv, 2016.

Maximum Likelihood Estimation

Given: Training samples $\hat{\xi}_1, \hat{\xi}_2, \dots, \hat{\xi}_N \in \mathbb{R}^m$ i.i.d. from $\mathcal{N}(\mu, \Sigma)$

Maximum Likelihood Estimation

Given: Training samples $\hat{\xi}_1, \hat{\xi}_2, \dots, \hat{\xi}_N \in \mathbb{R}^m$ i.i.d. from $\mathcal{N}(\mu, \Sigma)$

Input: Covariance matrix Σ

Output: Precision matrix $X = \Sigma^{-1}$

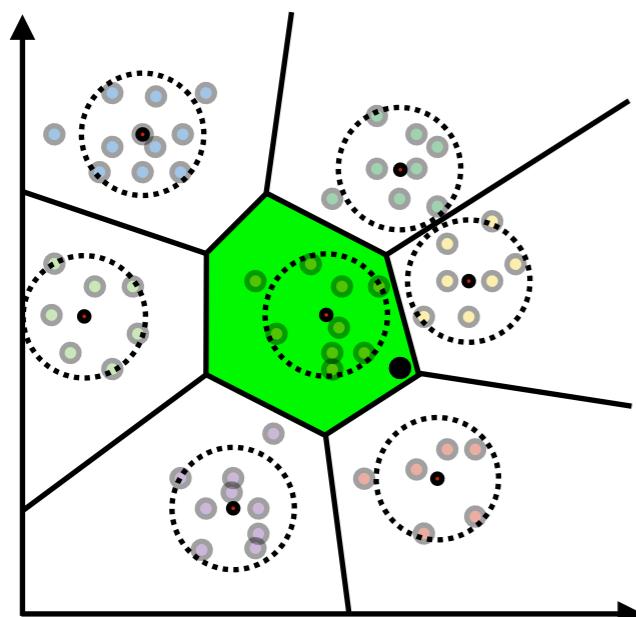
Maximum Likelihood Estimation

Given: Training samples $\hat{\xi}_1, \hat{\xi}_2, \dots, \hat{\xi}_N \in \mathbb{R}^m$ i.i.d. from $\mathcal{N}(\mu, \Sigma)$

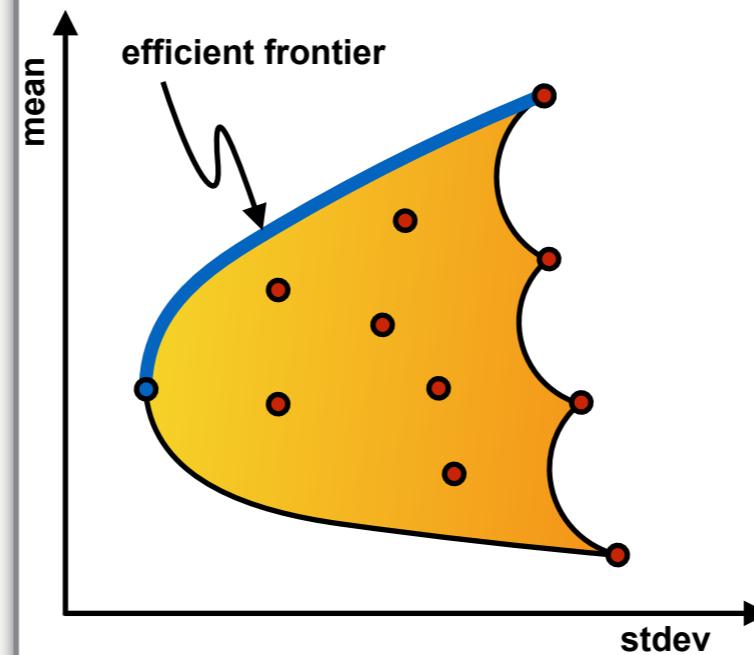
Input: Covariance matrix Σ

Output: Precision matrix $X = \Sigma^{-1}$

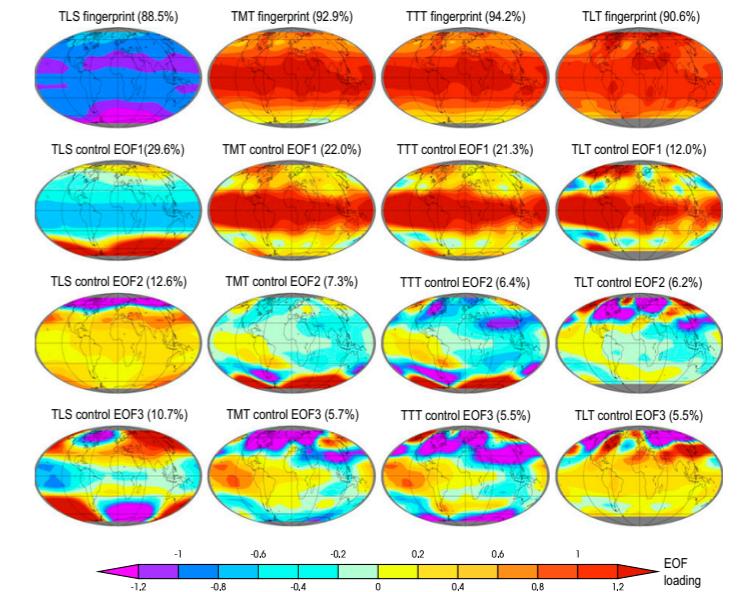
Linear discriminant analysis¹⁾



Markowitz portfolio analysis²⁾



Optimal fingerprint method³⁾



¹⁾ Fisher, *Ann. Eugen.*, 1936.

²⁾ Markowitz, *J. Financ.*, 1952.

³⁾ Ribes et al., *Clim. Dynam.*, 2009.

Maximum Likelihood Estimation

Given: Training samples $\hat{\xi}_1, \hat{\xi}_2, \dots, \hat{\xi}_N \in \mathbb{R}^m$ i.i.d. from $\mathcal{N}(\mu, \Sigma)$

Empirical estimators:

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N \hat{\xi}_i \quad \text{and} \quad \hat{\Sigma} = \frac{1}{N} \sum_{i=1}^N (\hat{\xi}_i - \hat{\mu})(\hat{\xi}_i - \hat{\mu})^\top$$

Maximum Likelihood Estimation

Given: Training samples $\hat{\xi}_1, \hat{\xi}_2, \dots, \hat{\xi}_N \in \mathbb{R}^m$ i.i.d. from $\mathcal{N}(\mu, \Sigma)$

Empirical estimators:

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N \hat{\xi}_i \quad \text{and} \quad \hat{\Sigma} = \frac{1}{N} \sum_{i=1}^N (\hat{\xi}_i - \hat{\mu})(\hat{\xi}_i - \hat{\mu})^\top$$

► $m > N \implies \hat{\Sigma}$ not invertible

Maximum Likelihood Estimation

Given: Training samples $\hat{\xi}_1, \hat{\xi}_2, \dots, \hat{\xi}_N \in \mathbb{R}^m$ i.i.d. from $\mathcal{N}(\mu, \Sigma)$

Maximum likelihood estimators:

$$\inf_{\mu \in \mathbb{R}^m, \mathbf{X} \in \mathbb{S}_+^m} \left\{ -\log \det \mathbf{X} + \frac{1}{N} \sum_{i=1}^N (\hat{\xi}_i - \mu) \mathbf{X} (\hat{\xi}_i - \mu)^\top \right\}$$

Maximum Likelihood Estimation

Given: Training samples $\hat{\xi}_1, \hat{\xi}_2, \dots, \hat{\xi}_N \in \mathbb{R}^m$ i.i.d. from $\mathcal{N}(\mu, \Sigma)$

Maximum likelihood estimators:

$$\inf_{\mu \in \mathbb{R}^m, \mathbf{X} \in \mathbb{S}_+^m} \left\{ -\log \det \mathbf{X} + \frac{1}{N} \sum_{i=1}^N (\hat{\xi}_i - \mu) \mathbf{X} (\hat{\xi}_i - \mu)^\top \right\}$$

- ▶ $m > N \implies$ problem unbounded!
- ▶ $m \leq N \implies \mu^* = \hat{\mu}, \mathbf{X}^* = \hat{\Sigma}^{-1}$

Maximum Likelihood Estimation

Given: Training samples $\hat{\xi}_1, \hat{\xi}_2, \dots, \hat{\xi}_N \in \mathbb{R}^m$ i.i.d. from $\mathcal{N}(\mu, \Sigma)$

Robust maximum likelihood estimators:¹⁾

$$\inf_{\mu \in \mathbb{R}^m, X \in \mathbb{S}_+^m} \left\{ -\log \det X + \sup_{Q \in \mathbb{B}_{\varepsilon, 2}(\hat{\mathbb{P}}_N)} \mathbb{E}^Q [(\xi - \mu) X (\xi - \mu)^T] \right\}$$

¹⁾ Nguyen, Kuhn & Mohajerin Esfahani, arxiv, 2018.

Maximum Likelihood Estimation

Given: Training samples $\hat{\xi}_1, \hat{\xi}_2, \dots, \hat{\xi}_N \in \mathbb{R}^m$ i.i.d. from $\mathcal{N}(\mu, \Sigma)$

Robust maximum likelihood estimators:¹⁾

$$\inf_{\mu \in \mathbb{R}^m, X \in \mathbb{S}_+^m} \left\{ -\log \det X + \sup_{Q \in \mathbb{B}_{\varepsilon, 2}(\hat{\mathbb{P}}_N)} \mathbb{E}^Q [(\xi - \mu) X (\xi - \mu)^T] \right\}$$

- ▶ normal nominal distribution $\hat{\mathbb{P}}_N = \mathcal{N}(\hat{\mu}, \hat{\Sigma})$
- ▶ recover ML problem for $\varepsilon = 0$
- ▶ regularization by robustification

¹⁾ Nguyen, Kuhn & Mohajerin Esfahani, arxiv, 2018.

Wasserstein Shrinkage Estimator¹⁾

Theorem: If $\varepsilon > 0$ and $\widehat{\Sigma} = \sum_{i=1}^m \lambda_i \cdot v_i v_i^\top$, then

$$\mu^* = \widehat{\mu} \quad \text{and} \quad X^* = \sum_{i=1}^m x_i^* \cdot v_i v_i^\top,$$

where

$$x_i^* = y^* \left[1 - \frac{1}{2} \left(\sqrt{\lambda_i^2 (y^*)^2 + 4\lambda_i y^*} - \lambda_i y^* \right) \right]$$

and y^* is the unique positive solution of

$$(\varepsilon^2 - \frac{1}{2} \sum_{i=1}^m \lambda_i) y - m + \frac{1}{2} \sum_{i=1}^m \sqrt{\lambda_i^2 y^2 + 4\lambda_i y} = 0.$$

¹⁾ Nguyen, Kuhn & Mohajerin Esfahani, arxiv, 2018.

Wasserstein Shrinkage Estimator¹⁾

Theorem: If $\varepsilon > 0$ and $\widehat{\Sigma} = \sum_{i=1}^m \lambda_i \cdot v_i v_i^\top$, then

$$\mathbf{X}^* = \sum_{i=1}^m \mathbf{x}_i^* \cdot v_i v_i^\top$$

where

$$\mathbf{x}_i^* = \gamma^* \left[1 - \frac{1}{2} \left(\sqrt{\lambda_i^2 (\gamma^*)^2 + 4\lambda_i \gamma^*} - \lambda_i \gamma^* \right) \right]$$

and γ^* is the unique positive solution of

$$(\varepsilon^2 - \frac{1}{2} \sum_{i=1}^m \lambda_i) \gamma - m + \frac{1}{2} \sum_{i=1}^m \sqrt{\lambda_i^2 \gamma^2 + 4\lambda_i \gamma} = 0.$$

- ▶ \mathbf{X}^* and $\widehat{\Sigma}$ commute $\implies \mathbf{X}^*$ is **rotation-equivariant**
- ▶ $x_i^* > 0$ for all $i \implies \mathbf{X}^*$ is **invertible**
- ▶ \mathbf{X}^* is **easy to compute** (spectral decomposition + bisection)

¹⁾ Nguyen, Kuhn & Mohajerin Esfahani, arxiv, 2018.

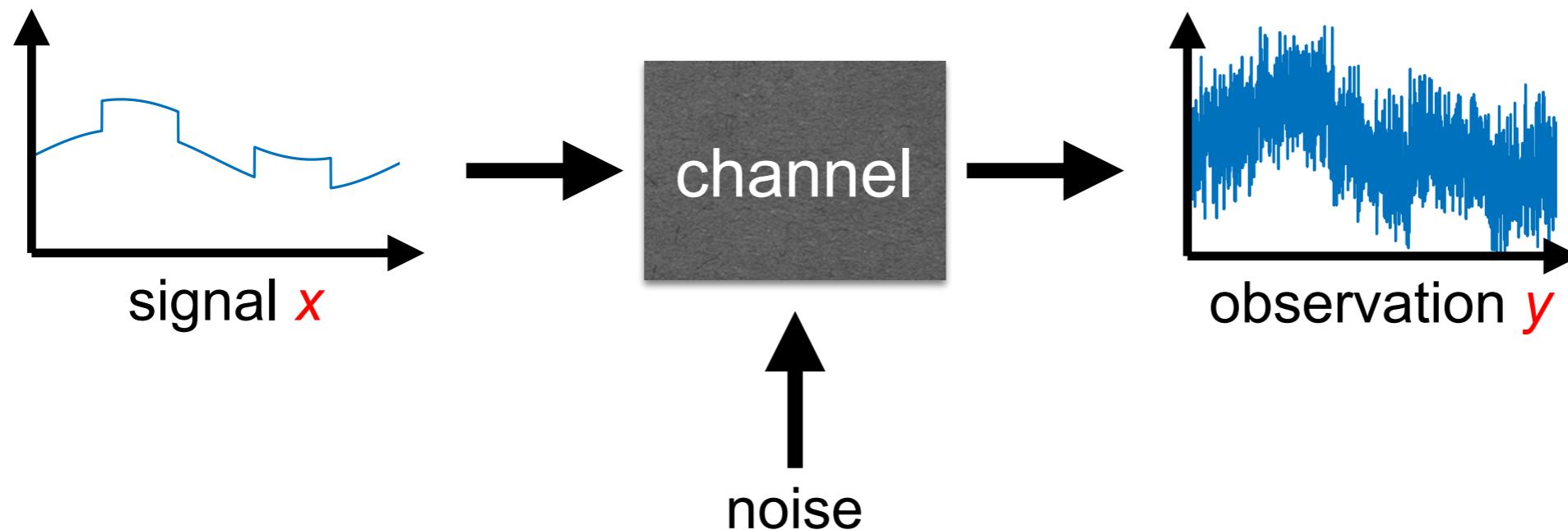
Wasserstein Shrinkage Estimator

Properties of the Wasserstein shrinkage estimator:

- ▶ exists for $m > N$ if $\varepsilon > 0$
- ▶ rotation-equivariant
- ▶ preserves the order of the eigenvalues of $\hat{\Sigma}$
- ▶ condition number improves as $\varepsilon \rightarrow \infty$
- ▶ same complexity as spectral decomposition of $\hat{\Sigma}$
- ▶ offers out-of-sample guarantee on Stein's loss
- ▶ follows from robustifying the ML estimation problem
- ▶ similar performance as graphical lasso at computational cost of linear shrinkage estimator

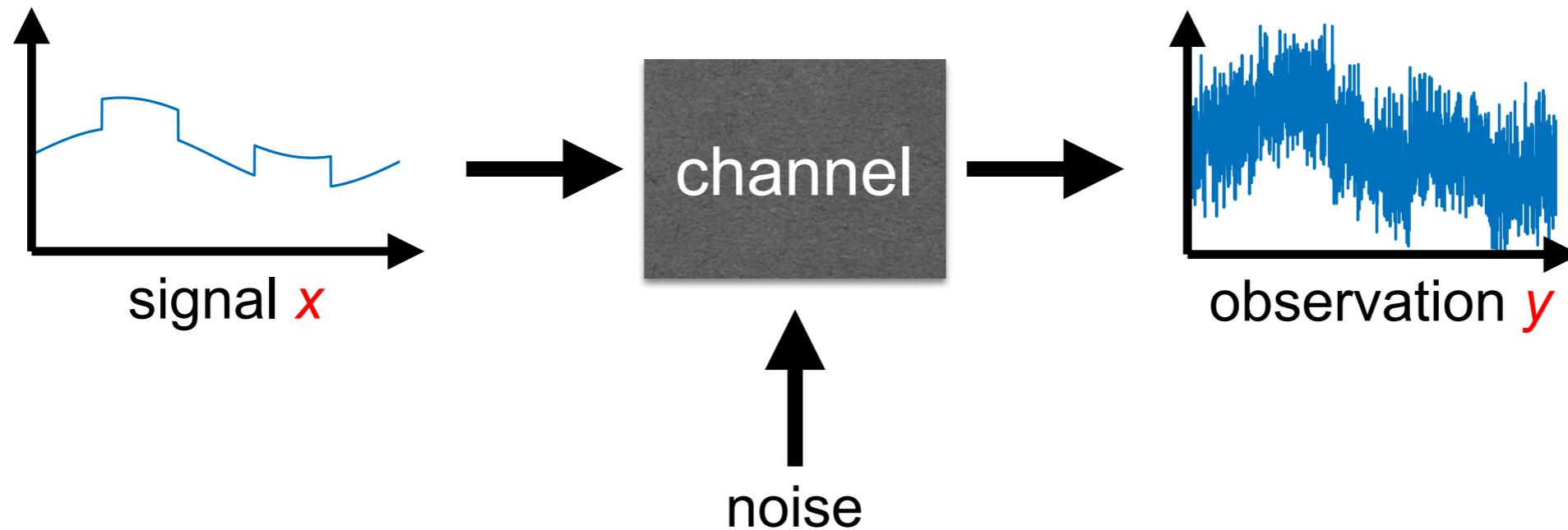
Minimum Mean Square Error Estimation

Problem: Infer signal $x \in \mathbb{R}^{m_x}$ from noisy observation $y \in \mathbb{R}^{m_y}$



Minimum Mean Square Error Estimation

Problem: Infer signal $\mathbf{x} \in \mathbb{R}^{m_x}$ from noisy observation $\mathbf{y} \in \mathbb{R}^{m_y}$



The distribution \mathbb{P} of $\xi = (\mathbf{x}, \mathbf{y}) \in \mathbb{R}^m$, $m = m_x + m_y$, is unknown.

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

Minimum Mean Square Error Estimation

Problem: Infer signal $\mathbf{x} \in \mathbb{R}^{m_x}$ from noisy observation $\mathbf{y} \in \mathbb{R}^{m_y}$

Robust MMSE estimation problem:¹⁾

$$\inf_{\psi(\cdot)} \sup_{\mathbb{Q} \in \mathbb{B}_{\varepsilon, 2}(\hat{\mathbb{P}}_N)} \mathbb{E}^{\mathbb{Q}} [\|\mathbf{x} - \psi(\mathbf{y})\|_2^2]$$

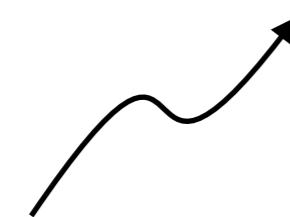
¹⁾ Shafieezadeh-Abadeh, Nguyen, Kuhn & Mohajerin Esfahani, NeuRIPS, 2018.

Minimum Mean Square Error Estimation

Problem: Infer signal $\mathbf{x} \in \mathbb{R}^{m_x}$ from noisy observation $\mathbf{y} \in \mathbb{R}^{m_y}$

Robust MMSE estimation problem:¹⁾

$$\inf_{\psi(\cdot)} \sup_{Q \in \mathbb{B}_{\varepsilon, 2}(\hat{\mathbb{P}}_N)} \mathbb{E}^Q [\|\mathbf{x} - \psi(\mathbf{y})\|_2^2]$$


estimator
(any function of \mathbf{y})

¹⁾ Shafieezadeh-Abadeh, Nguyen, Kuhn & Mohajerin Esfahani, NeuRIPS, 2018.

Minimum Mean Square Error Estimation

Problem: Infer signal $\mathbf{x} \in \mathbb{R}^{m_x}$ from noisy observation $\mathbf{y} \in \mathbb{R}^{m_y}$

Robust MMSE estimation problem:¹⁾

$$\inf_{\psi(\cdot)} \sup_{Q \in \mathbb{B}_{\varepsilon, 2}(\hat{\mathbb{P}}_N)} \mathbb{E}^Q [\|\mathbf{x} - \psi(\mathbf{y})\|_2^2]$$

MSE

¹⁾ Shafieezadeh-Abadeh, Nguyen, Kuhn & Mohajerin Esfahani, NeuRIPS, 2018.

Minimum Mean Square Error Estimation

Theorem:¹⁾ If $\widehat{\Sigma} \succ 0$, then the estimation problem is equivalent to

$$\max \text{Tr} [\Sigma_{xx} - \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{yx}]$$

$$\text{s.t. } \Sigma_{xx} \in \mathbb{S}_+^{m_x}, \quad \Sigma_{yy} \in \mathbb{S}_+^{m_y}, \quad \Sigma_{xy} = \Sigma_{yx}^\top \in \mathbb{R}^{m_x \times m_y}$$

$$\Sigma = \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix}, \quad \text{Tr} \left[\Sigma + \widehat{\Sigma} - 2 \left(\widehat{\Sigma}^{\frac{1}{2}} \Sigma \widehat{\Sigma}^{\frac{1}{2}} \right)^{\frac{1}{2}} \right] \leq \varepsilon^2.$$

If Σ_{xx}^* , Σ_{yy}^* and Σ_{xy}^* are optimal, then

$$\psi^*(y) = \Sigma_{xy}^* (\Sigma_{yy}^*)^{-1} (y - \widehat{\mu}_y) + \widehat{\mu}_x$$

is a robust MMSE estimator.

¹⁾ Shafieezadeh-Abadeh, Nguyen, Kuhn & Mohajerin Esfahani, NeuRIPS, 2018.

Minimum Mean Square Error Estimation

Theorem:¹⁾ If $\widehat{\Sigma} \succ 0$, then the estimation problem is equivalent to

$$\max \text{Tr} [\Sigma_{xx} - \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{yx}]$$

$$\text{s.t. } \Sigma_{xx} \in \mathbb{S}_+^{m_x}, \quad \Sigma_{yy} \in \mathbb{S}_+^{m_y}, \quad \Sigma_{xy} = \Sigma_{yx}^\top \in \mathbb{R}^{m_x \times m_y}$$

$$\Sigma = \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix}, \quad \text{Tr} \left[\Sigma + \widehat{\Sigma} - 2 \left(\widehat{\Sigma}^{\frac{1}{2}} \Sigma \widehat{\Sigma}^{\frac{1}{2}} \right)^{\frac{1}{2}} \right] \leq \varepsilon^2.$$

If Σ_{xx}^* , Σ_{yy}^* and Σ_{xy}^* are optimal, then

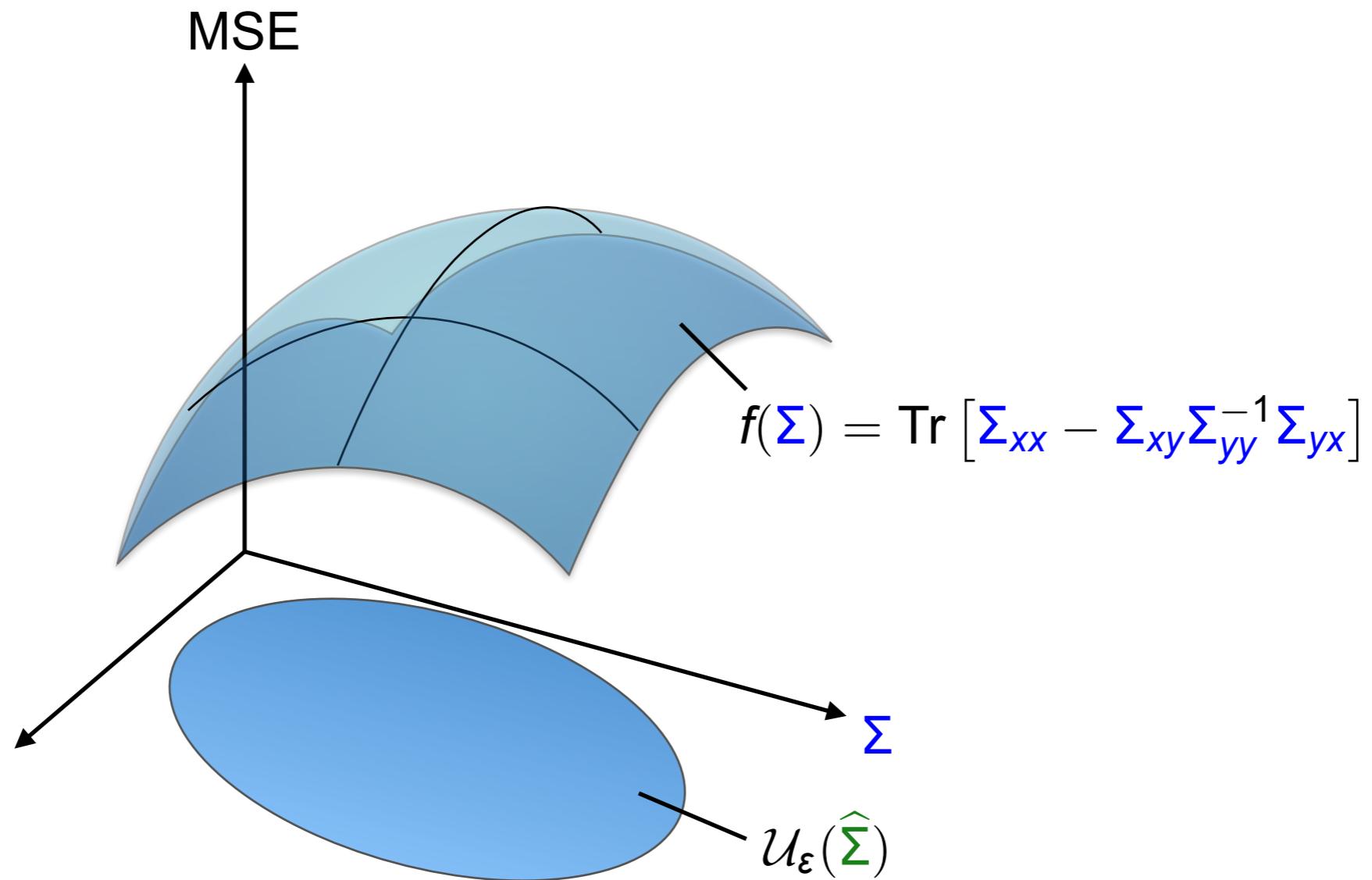
$$\psi^*(\mathbf{y}) = \Sigma_{xy}^* (\Sigma_{yy}^*)^{-1} (\mathbf{y} - \widehat{\mu}_y) + \widehat{\mu}_x$$

is a robust MMSE estimator.

- ▶ ψ^* is an **affine decision rule**
- ▶ ψ^* and $\mathcal{Q}^* = \mathcal{E}_g(\widehat{\mu}, \Sigma^*)$ form a **Nash equilibrium**

¹⁾ Shafieezadeh-Abadeh, Nguyen, Kuhn & Mohajerin Esfahani, NeuRIPS, 2018.

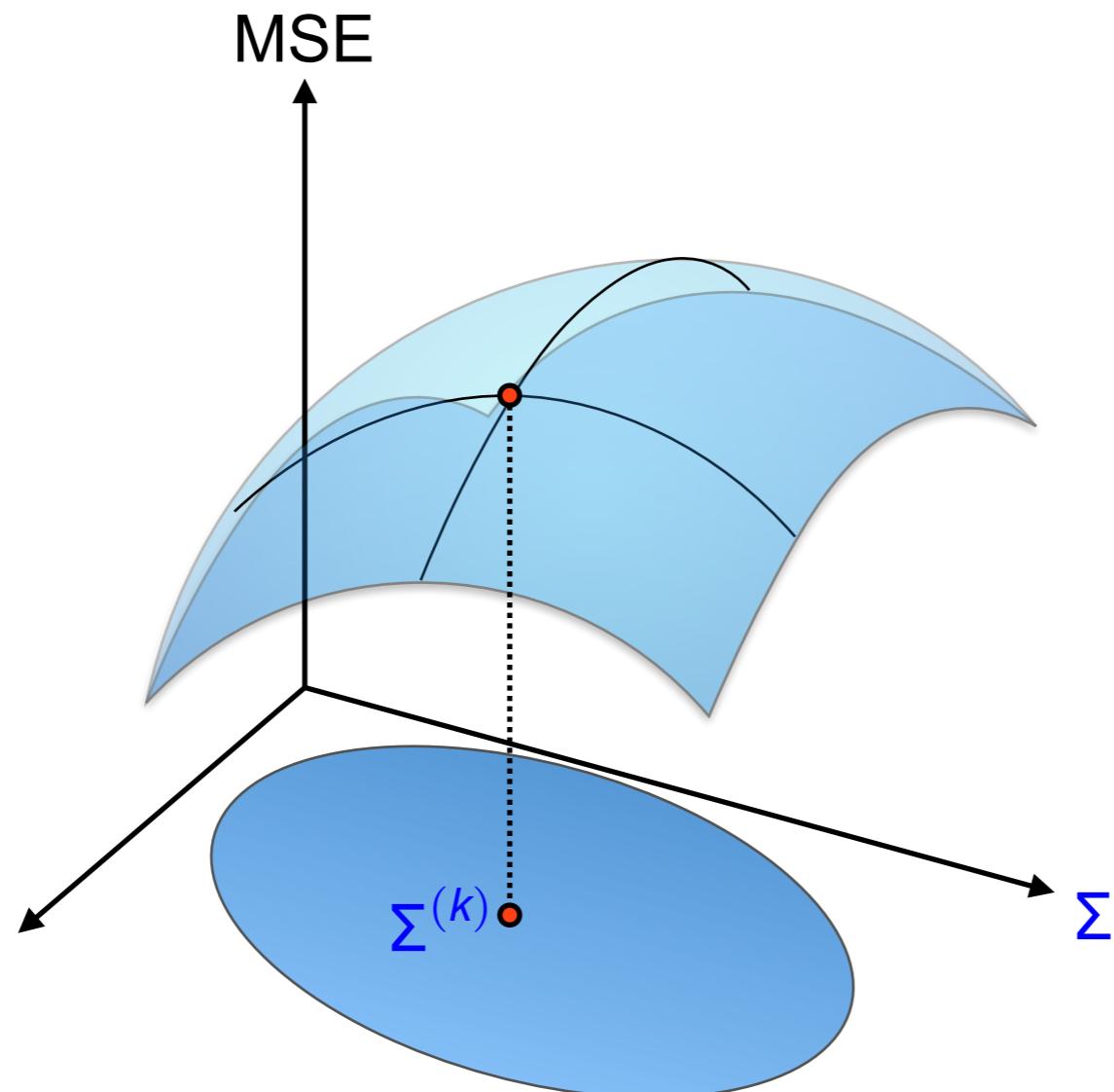
Frank-Wolfe Algorithm



Task: Solve the nonlinear SDP

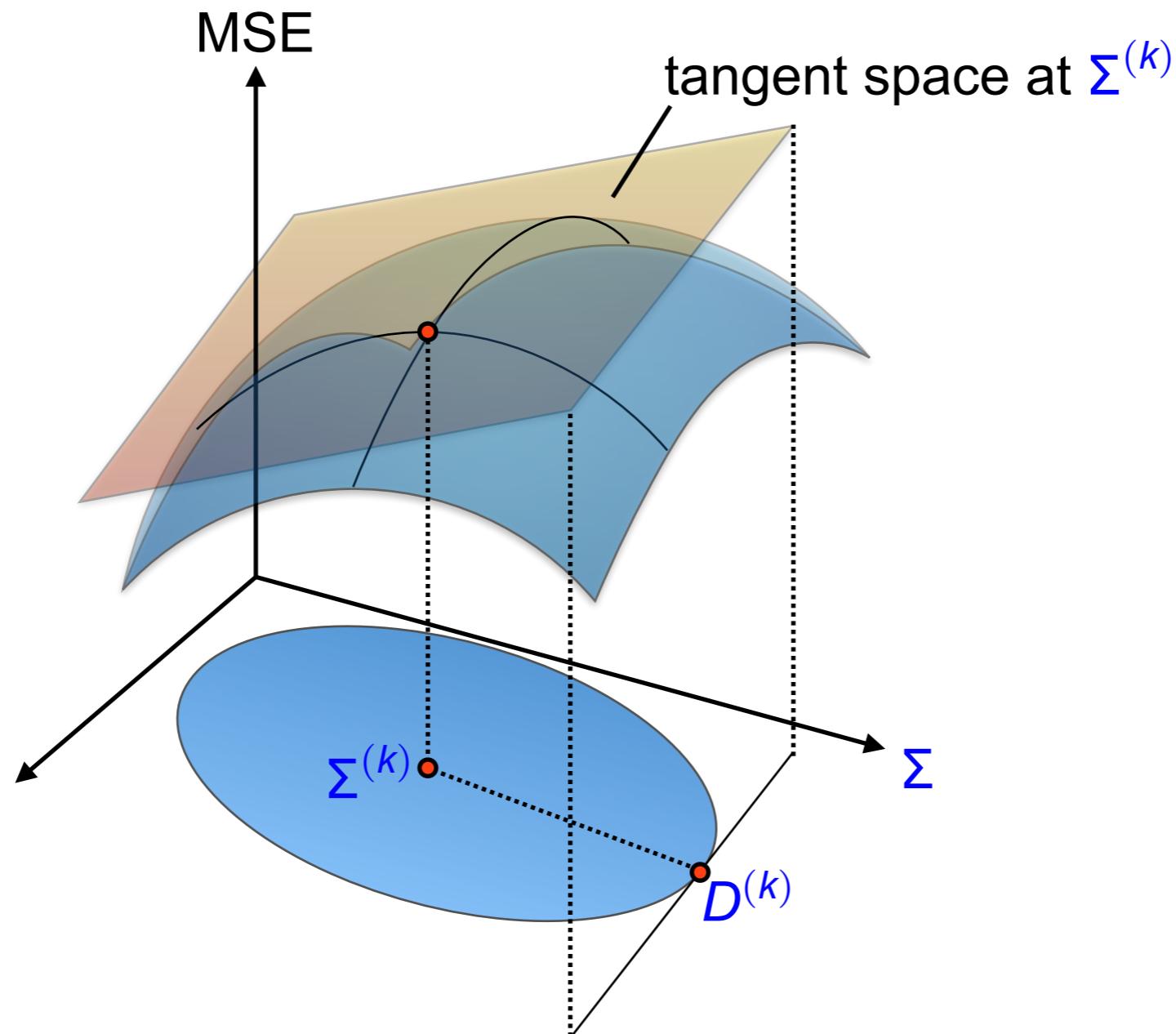
$$\max \left\{ f(\Sigma) : \Sigma \in \mathcal{U}_\varepsilon(\hat{\Sigma}), \Sigma \succeq 0 \right\}$$

Frank-Wolfe Algorithm



Step 0: Pick initial iterate $\Sigma^{(k)}$

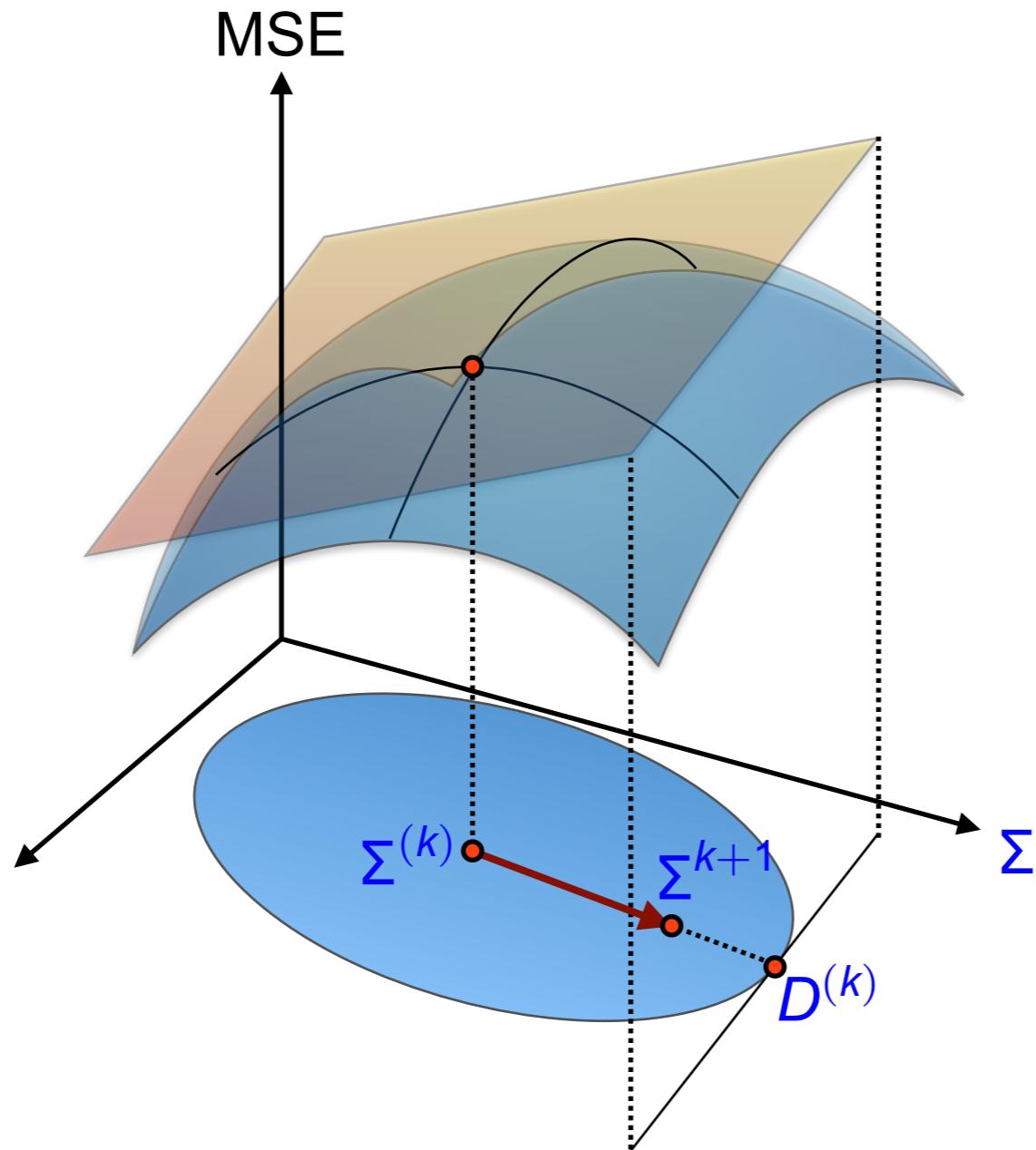
Frank-Wolfe Algorithm



Step 1: Solve direction-finding subproblem

$$D^{(k)} \in \arg \max \left\{ \text{Tr} \left[D \nabla f(\Sigma^{(k)}) \right] : D \in \mathcal{U}_\varepsilon(\widehat{\Sigma}), D \succeq 0 \right\}$$

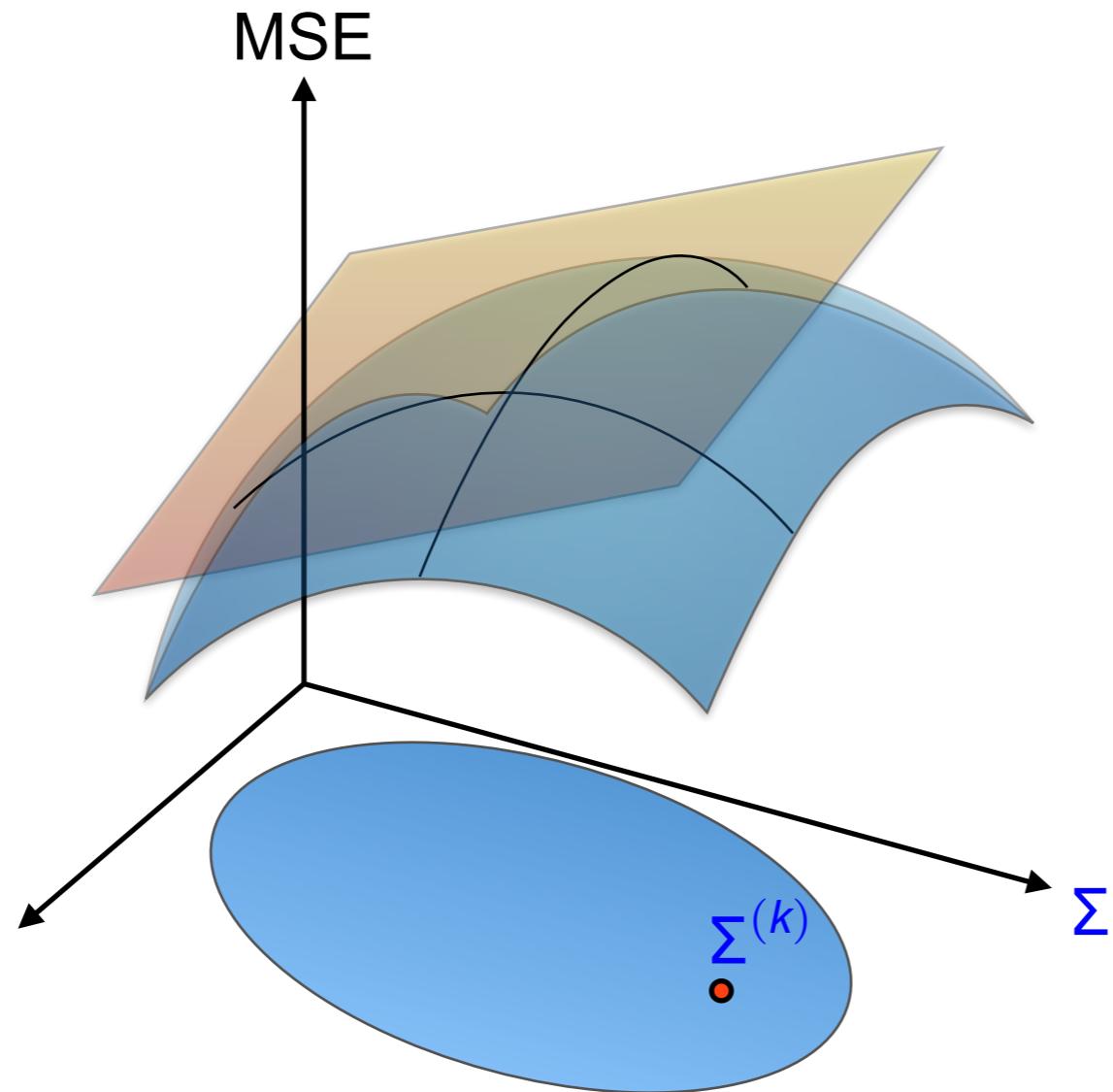
Frank-Wolfe Algorithm



Step 2: Construct new iterate

$$\Sigma^{(k+1)} = \alpha_k \cdot D^{(k)} + (1 - \alpha_k) \cdot \Sigma^{(k)}$$

Frank-Wolfe Algorithm



Step 3: Set $k \leftarrow k + 1$ and go to Step 1.

Frank-Wolfe Algorithm

Key benefits:¹⁾

1. Nonlinear SDP is “nice”:

- ▶ f is concave and smooth, ∇f is bounded away from 0
- ▶ $\mathcal{U}_\varepsilon(\widehat{\Sigma})$ is strongly convex

¹⁾ Shafieezadeh-Abadeh, Nguyen, Kuhn & Mohajerin Esfahani, NeuRIPS, 2018.

Frank-Wolfe Algorithm

Key benefits:¹⁾

1. Nonlinear SDP is “nice”:

- ▶ f is concave and smooth, ∇f is bounded away from 0
- ▶ $\mathcal{U}_\varepsilon(\widehat{\Sigma})$ is strongly convex

⇒ **convergence of FW is linear!**

¹⁾ Shafieezadeh-Abadeh, Nguyen, Kuhn & Mohajerin Esfahani, NeuRIPS, 2018.

Frank-Wolfe Algorithm

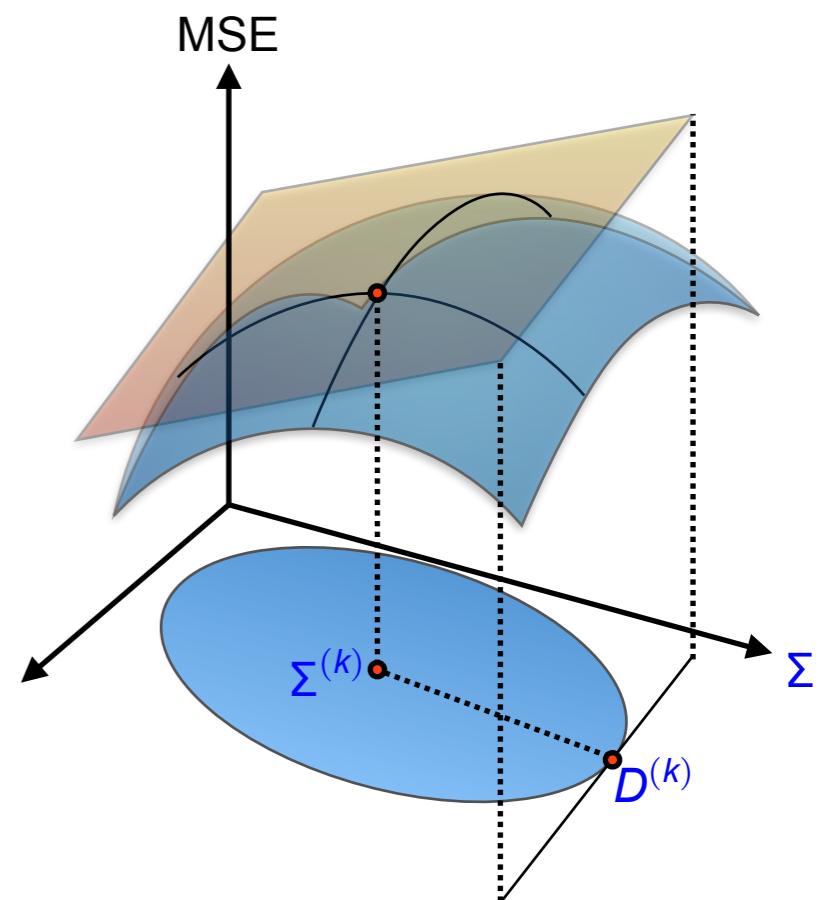
Key benefits:¹⁾

2. Direction-finding subproblem can be solved in closed form:

$$D^{(k)} = (\gamma^*)^2 \left(\gamma^* I - \nabla f(\Sigma^{(k)}) \right)^{-1} \widehat{\Sigma} \left(\gamma^* I - \nabla f(\Sigma^{(k)}) \right)^{-1},$$

where γ^* solves

$$\text{Tr} \left[\widehat{\Sigma} \left(I - \gamma \left[\gamma I - \nabla f(\Sigma^{(k)}) \right]^{-1} \right)^2 \right] = \varepsilon^2.$$



¹⁾ Shafieezadeh-Abadeh, Nguyen, Kuhn & Mohajerin Esfahani, NeuRIPS, 2018.

Frank-Wolfe Algorithm

Key benefits:¹⁾

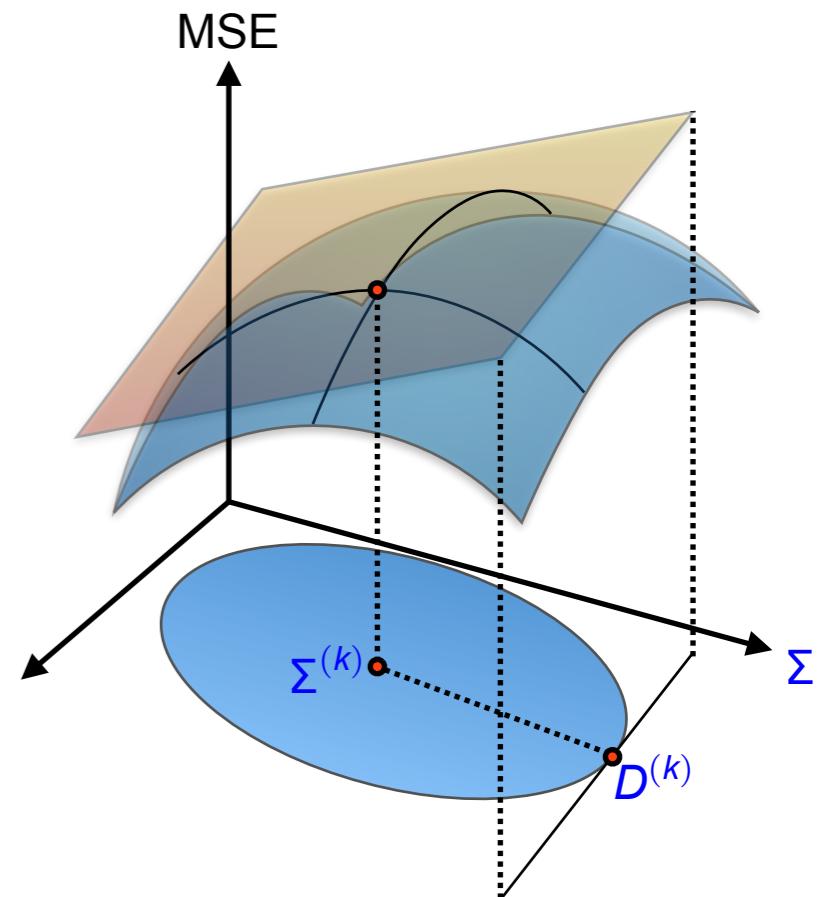
2. Direction-finding subproblem can be solved in closed form:

$$D^{(k)} = (\gamma^*)^2 \left(\gamma^* I - \nabla f(\Sigma^{(k)}) \right)^{-1} \widehat{\Sigma} \left(\gamma^* I - \nabla f(\Sigma^{(k)}) \right)^{-1},$$

where γ^* solves

$$\text{Tr} \left[\widehat{\Sigma} \left(I - \gamma \left[\gamma I - \nabla f(\Sigma^{(k)}) \right]^{-1} \right)^2 \right] = \varepsilon^2.$$

⇒ **use bisection; extremely fast!**



¹⁾ Shafieezadeh-Abadeh, Nguyen, Kuhn & Mohajerin Esfahani, NeuRIPS, 2018.

Summary

Wasserstein DRO machine learning models:

- ▶ Explain commonly used regularization terms
- ▶ Explain the benefits of shrinkage estimators
- ▶ Motivate new regularization schemes
- ▶ Equivalent to efficiently solvable convex programs
- ▶ Sometimes solvable in quasi-closed form (bisection!)
- ▶ Sometimes solvable via fast FW algorithms

This Talk is Based on...

- [1] D. Kuhn, P. Mohajerin Esfahani, V. Nguyen and S. Shafieezadeh Abadeh. **Wasserstein Distributionally Robust Optimization: Theory and Applications in Machine Learning.** INFORMS TutORials in Operations Research. 2019.
- [2] P. Mohajerin Esfahani and D. Kuhn. **Data-Driven Distributionally Robust Optimization using the Wasserstein Metric: Performance Guarantees and Tractable Reformulations.** *Mathematical Programming* 171(1–2), 115–166, 2018.
- [3] V. Nguyen, D. Kuhn & P. Mohajerin Esfahani. **Distributionally Robust Inverse Covariance Estimation: The Wasserstein Shrinkage Estimator.** Working Paper, 2018.
- [4] S. Shafieezadeh Abadeh, P. Mohajerin Esfahani and D. Kuhn. **Distributionally Robust Logistic Regression.** *Neural Information Processing Systems*, 2015.
- [5] S. Shafieezadeh Abadeh, P. Mohajerin Esfahani and D. Kuhn. **Regularization via Mass Transportation.** *Journal of Machine Learning Research* 20(103), 1–68, 2019.
- [6] S. Shafieezadeh Abadeh, P. Mohajerin Esfahani, V. Nguyen and D. Kuhn. **Wasserstein Distributionally Robust Kalman Filtering.** *Neural Information Processing Systems*, 2018.
- [7] J. (T.) Zhen, D. Kuhn & W. Wiesemann. **Nonlinear Distributionally Robust Optimization.** Working Paper, 2019.

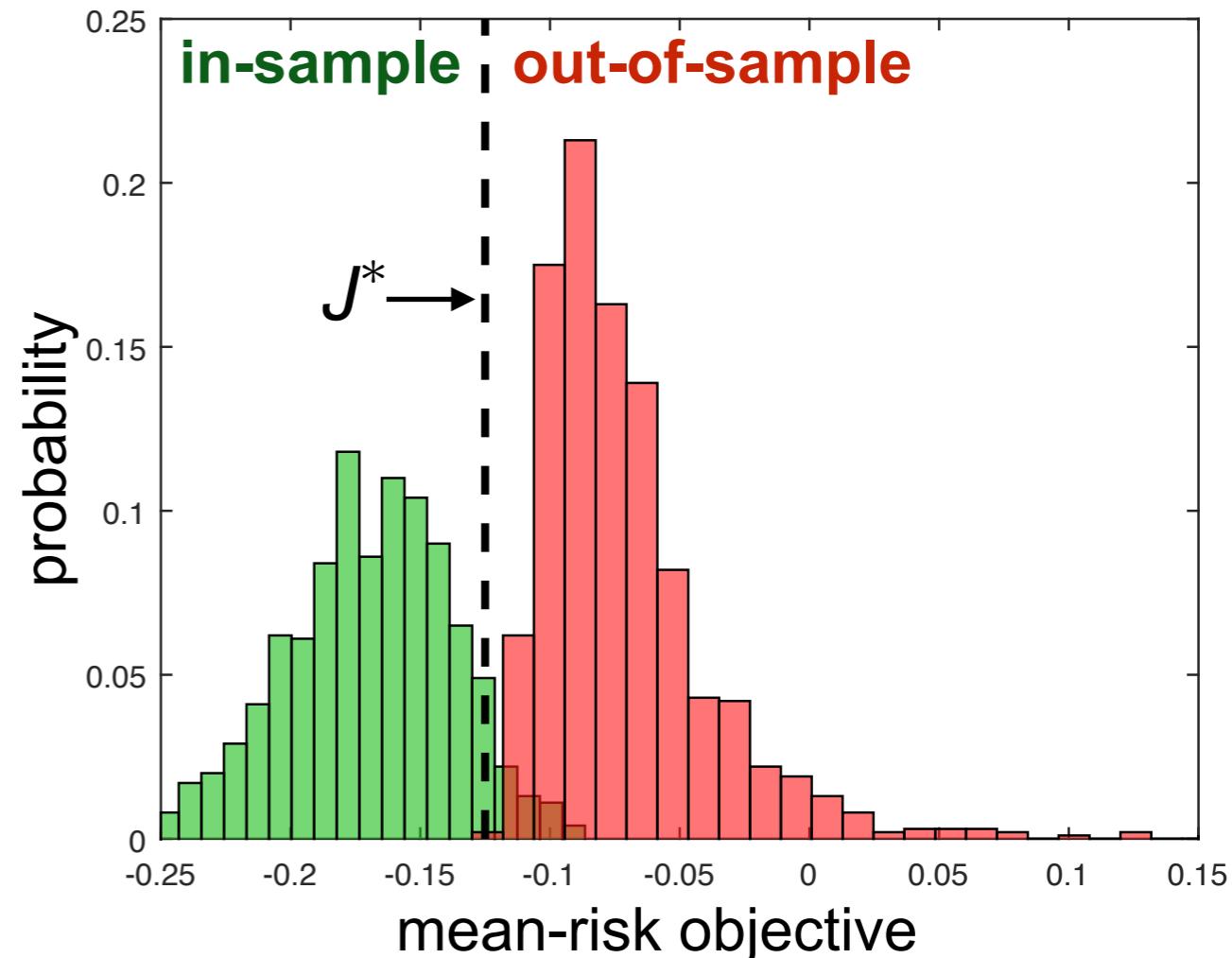
SAA with Scarce Data

Mean-risk portfolio problem

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

- ▶ 10 assets
- ▶ $\rho = 10$
- ▶ $\alpha = 20\%$
- ▶ $\boldsymbol{\xi}_i = \boldsymbol{\psi} + \boldsymbol{\zeta}_i$ where $\boldsymbol{\psi} \sim \mathcal{N}(0, 2\%)$
and $\boldsymbol{\zeta}_i \sim \mathcal{N}(i \times 3\%, i \times 2.5\%)$

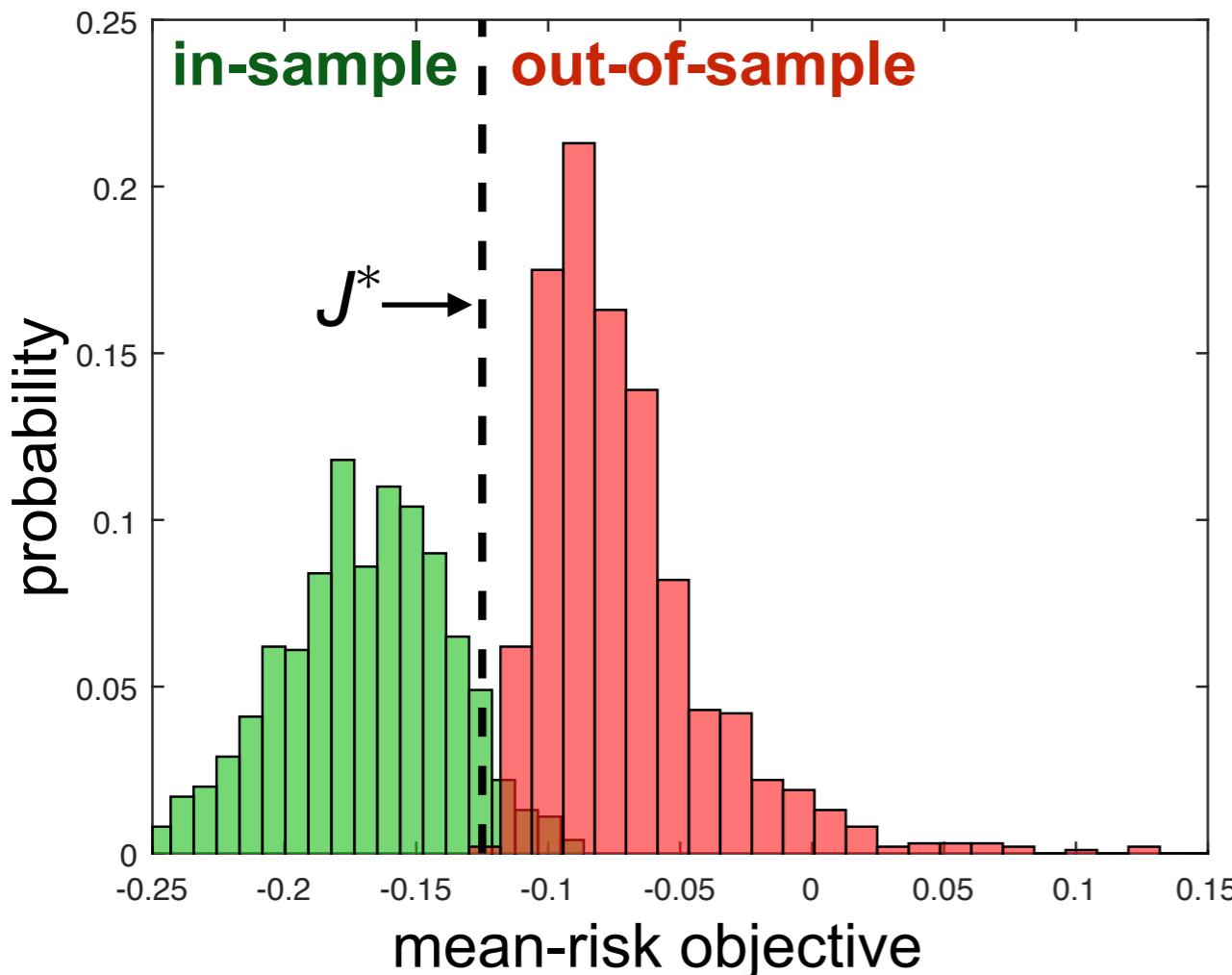
Performance of SAA solution



- ▶ 30 training samples
- ▶ in-sample: optimistic bias
- ▶ out-of-sample: pessimistic bias

DRO with Scarce Data

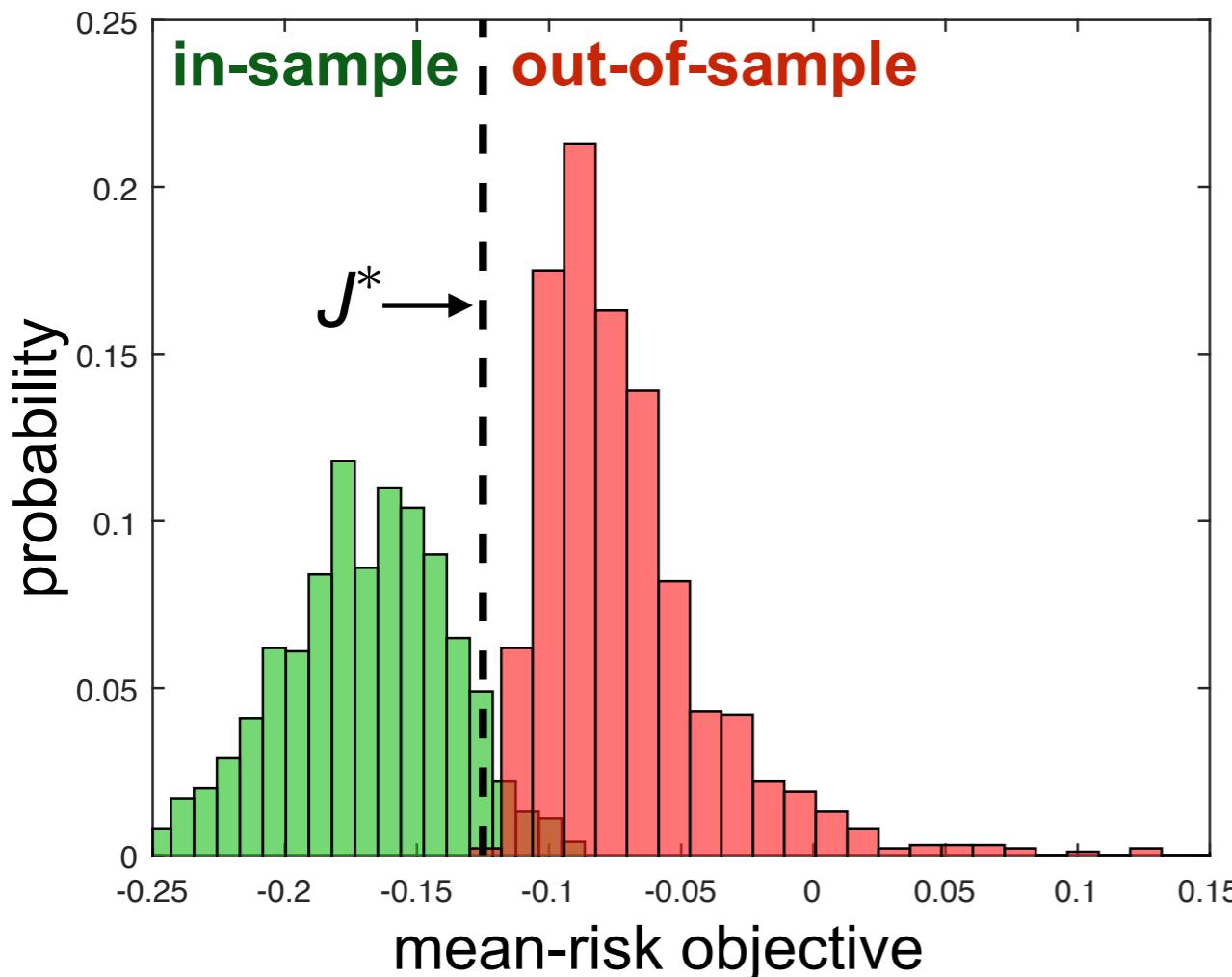
Performance of SAA solution



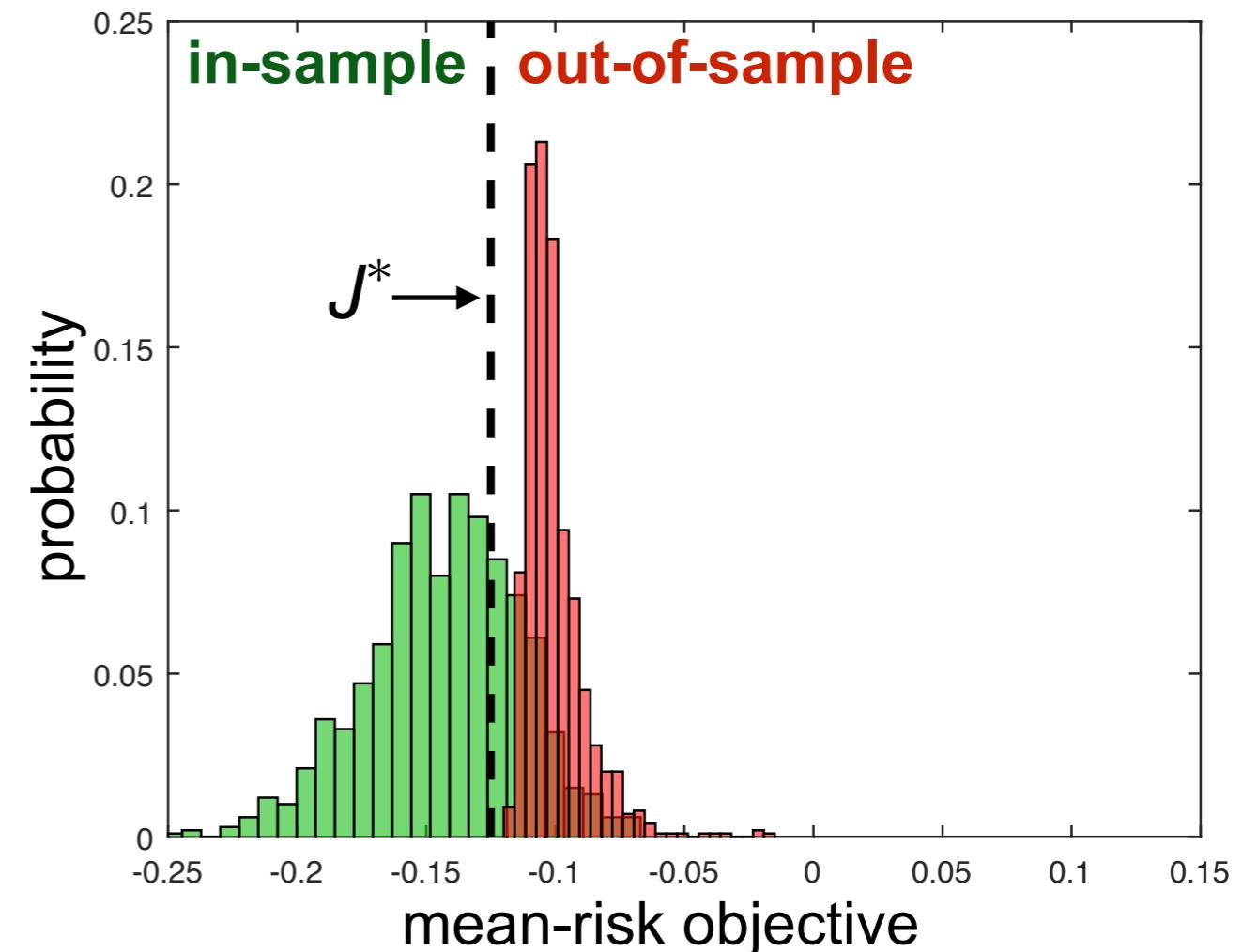
- ▶ in-sample: optimistic bias
- ▶ out-of-sample: pessimistic bias

DRO with Scarce Data

Performance of SAA solution



Performance of DRO solution

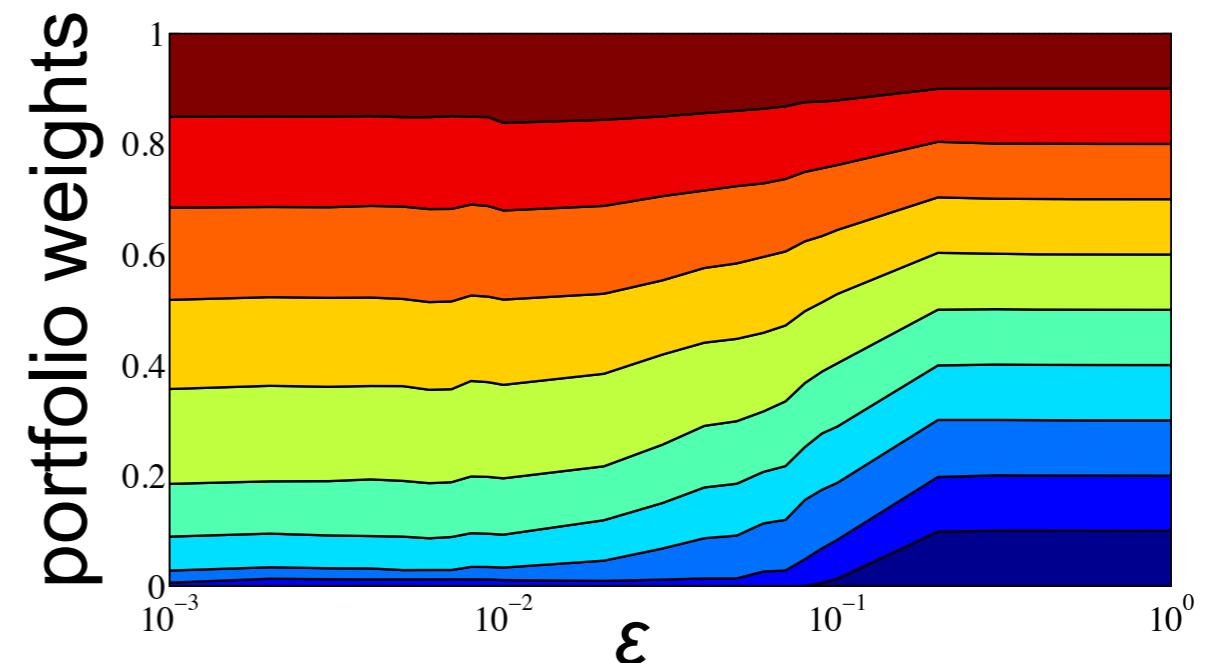


- ▶ in-sample: optimistic bias
- ▶ out-of-sample: pessimistic bias
- ▶ DRO reduces bias & post-decision disappointment

Application 1: Portfolio Selection

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

- ▶ 10 assets
- ▶ $\rho = 10$
- ▶ $\alpha = 20\%$
- ▶ $\boldsymbol{\xi}_i = \boldsymbol{\psi} + \boldsymbol{\zeta}_i$ where $\boldsymbol{\psi} \sim \mathcal{N}(0, 2\%)$
and $\boldsymbol{\zeta}_i \sim \mathcal{N}(i \times 3\%, i \times 2.5\%)$



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

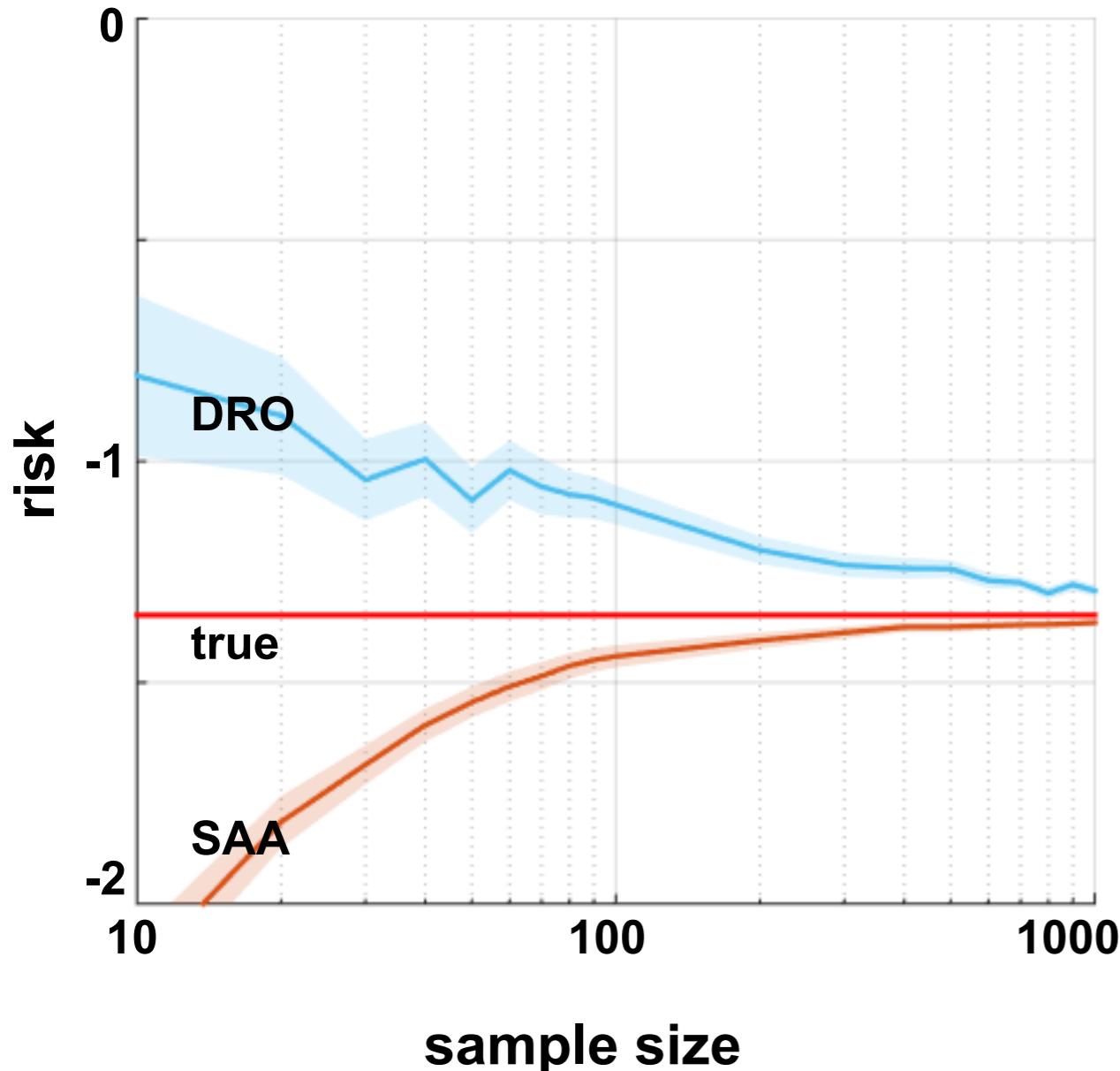
Possible Explanation: It is optimal for $\varepsilon \rightarrow \infty$.²⁾

¹⁾ DeMiguel, Garlappi & Uppal, *Rev. Financ. Stud.*, 2009;

²⁾ Pflug, Pichler & Wozabal, *J. Bank. Financ.*, 2012.

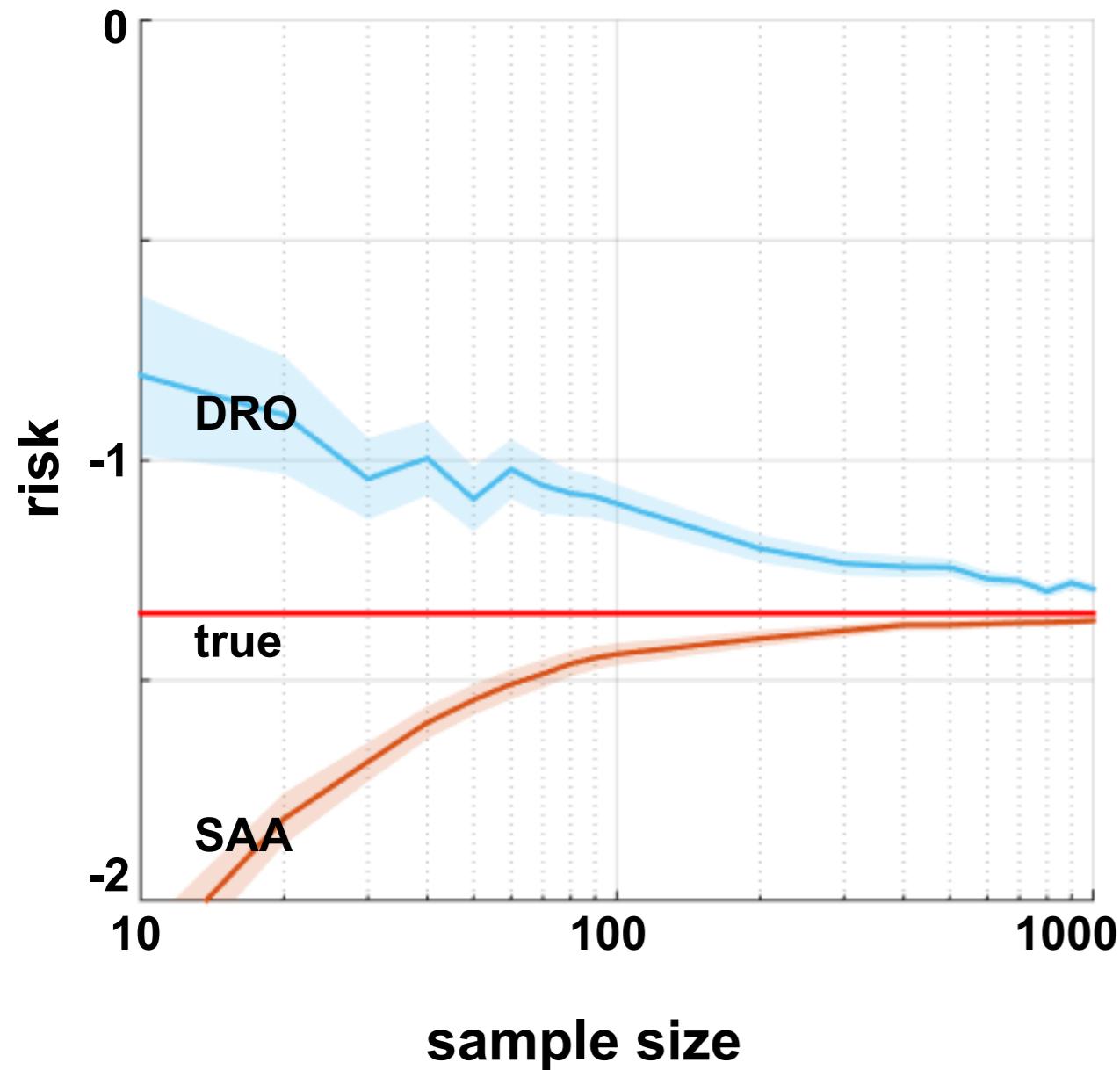
Learning Curves

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

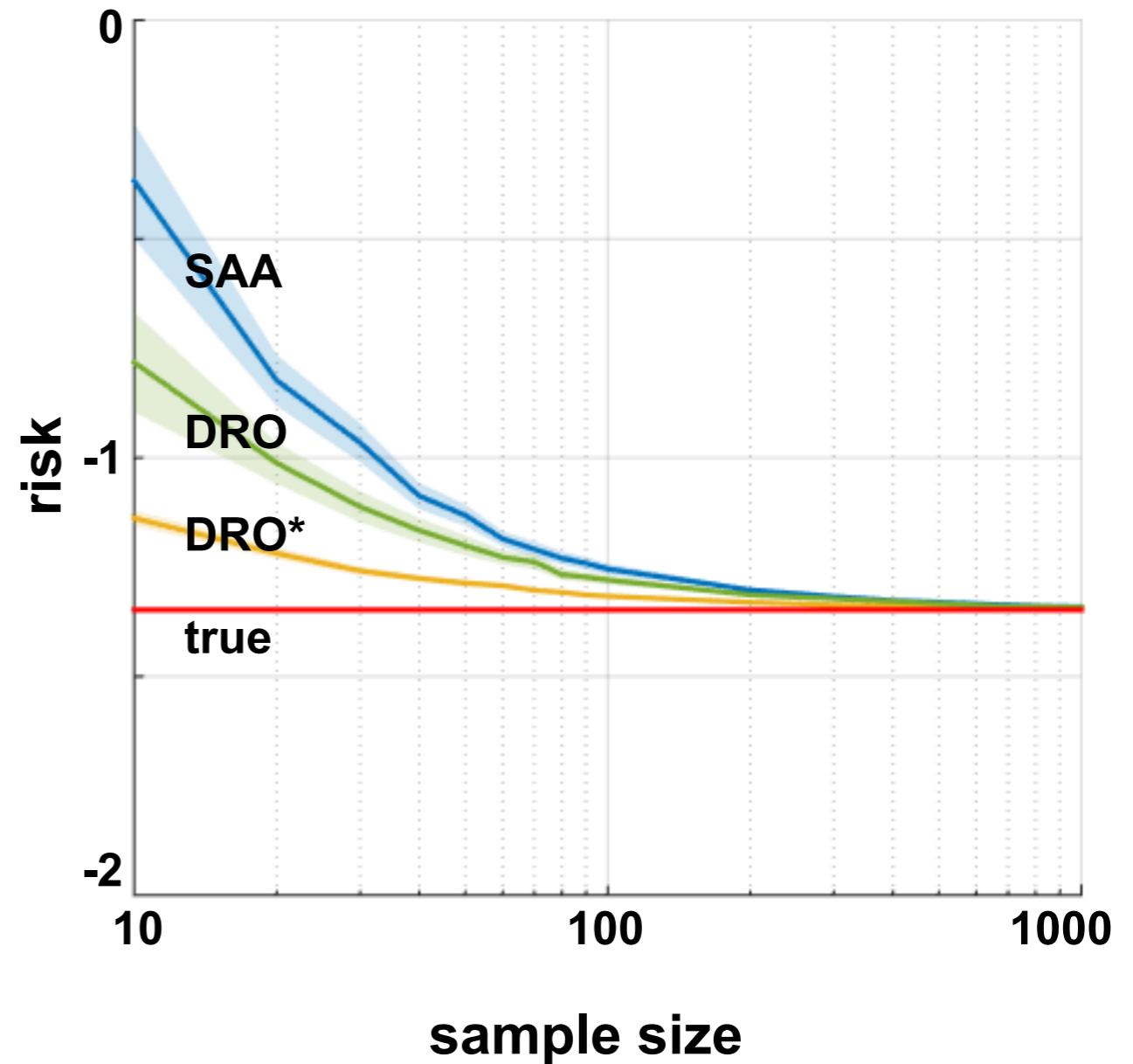


Learning Curves

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



what we **actually** get ...



Tractability

Worst-case expectation problem:

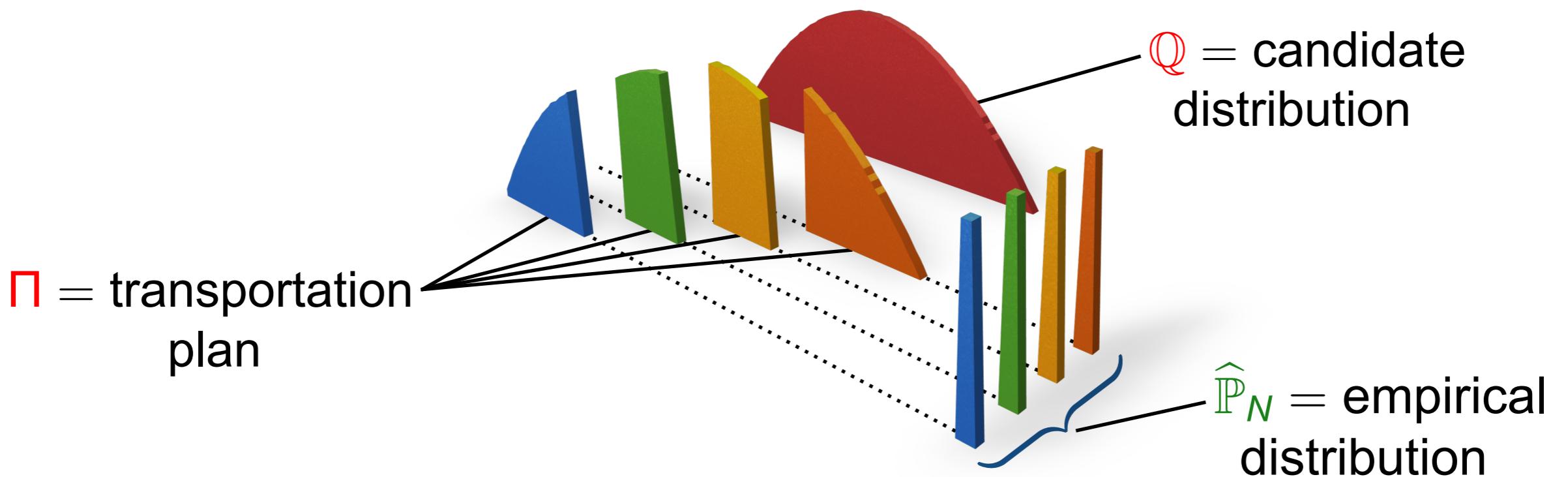
$$\sup_{Q \in \widehat{\mathcal{P}}_N} \mathbb{E}^Q [\ell(\xi)]$$

Use the Kantorovich-Rubinstein theorem:

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

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

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

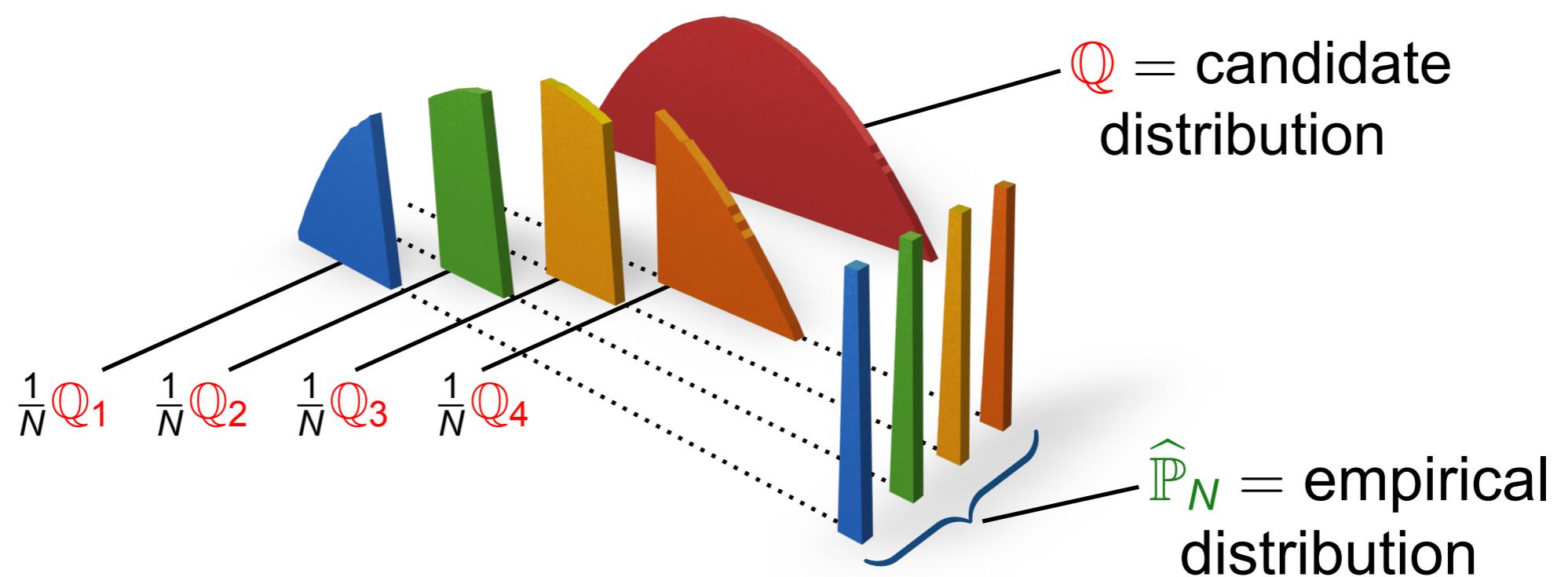


Tractability

Decompose Π into Q_1, \dots, Q_N :

$$\sup_{Q_i} \frac{1}{N} \sum_{i=1}^N \int_{\Xi} \ell(\xi) Q_i(d\xi)$$

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



Tractability

Dual of the moment problem is a robust program:

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

Tractability

Introduce the indicator function of Ξ :

$$\inf_{\lambda, s_i} \quad \lambda \varepsilon + \frac{1}{N} \sum_{i=1}^N s_i$$

$$\text{s.t. } \ell(\xi) - \lambda \|\xi - \hat{\xi}_i\| - \delta_\Xi(\xi) \leq s_i \quad \forall \xi \in \mathbb{R}^m, \quad \forall i \leq N$$

$$\lambda \geq 0$$

Tractability

Reformulate robust program as bilevel program:

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

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

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

\$[-\ell]^*(z_i - v_i)\$
 \$\sigma_\Xi(v_i)\$
 \$\|z_i\|_* \leq \lambda\$

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**