

# **Unifying Distributionally Robust Optimization**

---

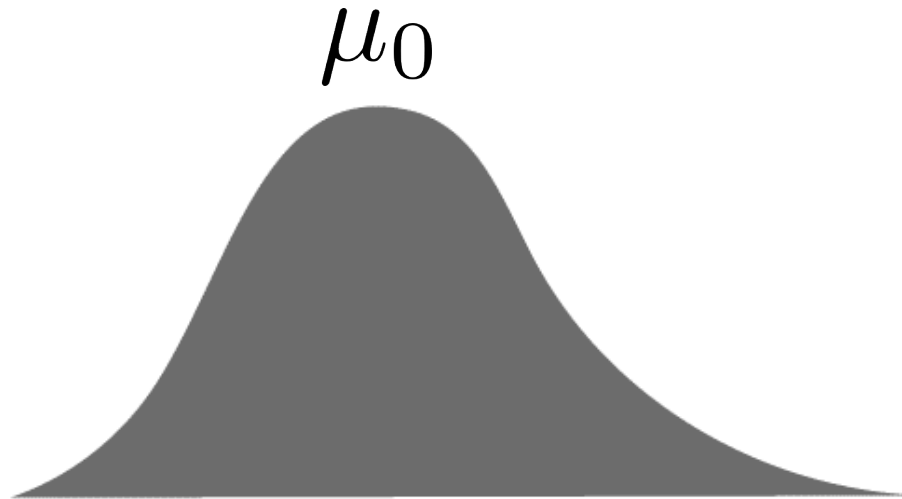
**Optimal transport Approach**

**Jiajin Li**

# Optimization under Uncertainty

---

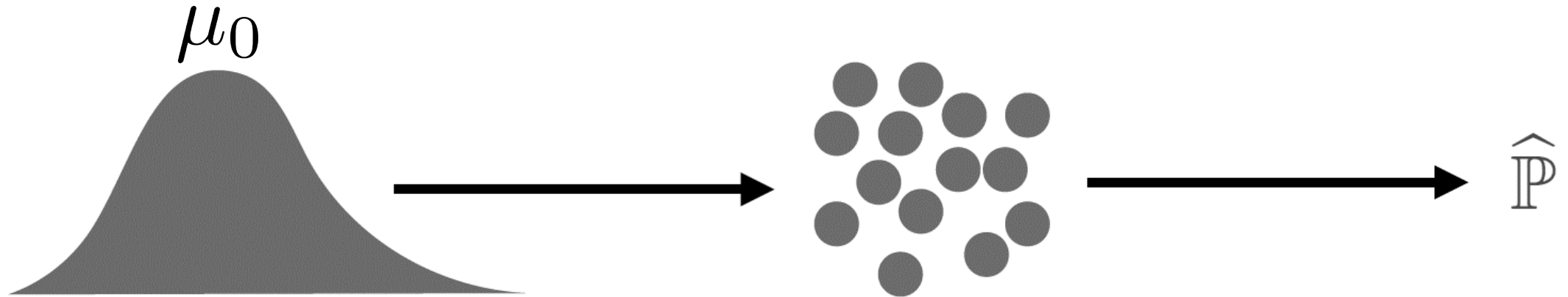
$$\inf_{\ell \in \mathcal{L}} \mathbb{E}_{\mu_0} [\ell(Z)]$$



# Optimization under Uncertainty

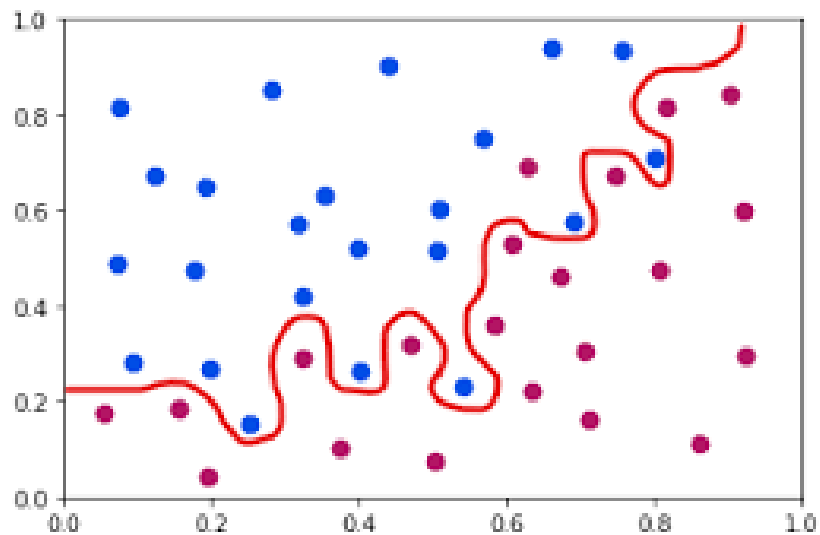
---

$$\inf_{l \in \mathcal{L}} \mathbb{E}_{\hat{\mu}}[\ell(Z)]$$



# SAA Often Fails

---



➤ **Overfitting**

➤ **Adversarial Attack**



$x$

“panda”  
57.7% confidence

+ .007 ×



$\text{sign}(\nabla_x J(\theta, x, y))$

“nematode”  
8.2% confidence

=



$x + \epsilon \text{sign}(\nabla_x J(\theta, x, y))$   
“gibbon”  
99.3 % confidence

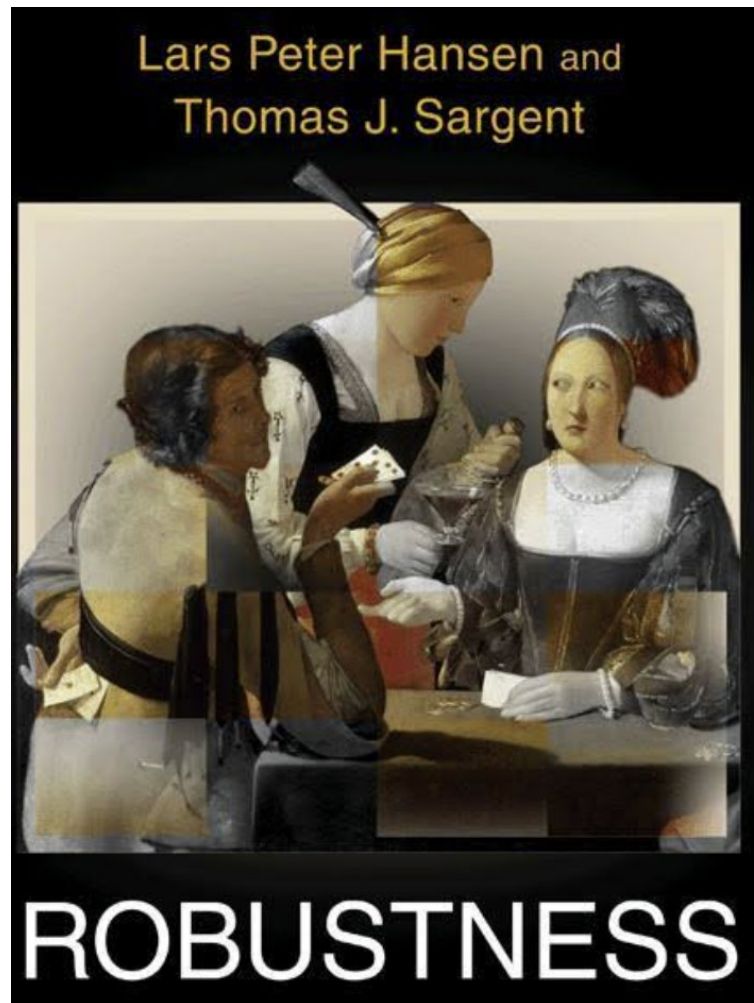
# Model Misspecification

$$\inf_{\ell \in \mathcal{L}} \sup_{\mu \in \mathcal{B}} \mathbb{E}_{\mu} [\ell(Z)]$$



What makes sense for choosing the set  $\mathcal{B}$ ?

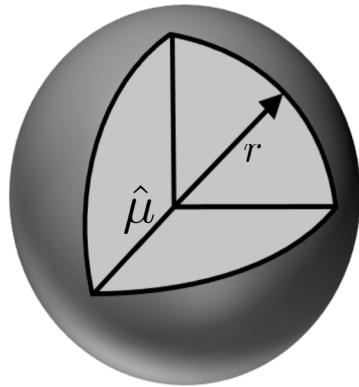
Optimal in Gilboa-Schmeidler (2008) decision theoretic (multiple prior) sense.



# Ambiguity Set

---

- A set of distribution around the baseline distribution  $\hat{\mu}$  :



$$\mathcal{B} := \{\mu \in \mathcal{P}(\mathcal{Z}) : \mathbb{D}(\mu, \hat{\mu}) \leq r\}$$

## Desirable Properties:

- Non-parametric
- Tractable
- Explainable

How to choose the  
“discrepancy”  $D(\cdot, \cdot)$ ?



# How to Choose Probability Metric?

---

**Two** natural ways to model changes in distributions.

A) Our model gets the **likelihood** of outcomes wrong.

B) Our model gets the **outcomes** wrong.

Traditionally A) and B) are seen as separate mechanisms.

**Approach A) Divergence**: Dupuis, James & Peterson '00; Hansen & Sargent '01, '08; Nilim & El Ghaoui '02, '03; Iyengar '05; A. Ben-Tal, L. El Ghaoui, & A. Nemirovski '09; Bertsimas & Sim '04; Bertsimas, Brown, Caramanis '13; Lim & Shanthikumar '04; Lam '13, '17; Csiszár & Breuer '13; Jiang & Guan '12; Hu & Hong '13; Wang, Glynn & Ye '14; Bayrakskan & Love '15; Duchi, Glynn & Namkoong '16; Bertsimas, Gupta & Kallus '13, (LD-) Van Bary et al. '17

**Approach B) Wasserstein**: Scarf '58; Hampel '73; Huber '81; Pflug & Wozabal '07; Mehrotra & Zhang '14; Esfahani & Kuhn '15; Blanchet & Murthy '16; Gao & Kleywegt '16; Duchi & Namkoong '17, (Sinkhorn-) Wang et al. '21.

# How to Choose Probability Metric?

---

**Two** natural ways to model changes in distributions.

Approach A) **Divergence**: Dupuis, James & Peterson '00; Hansen & Sargent '01, '08; Nilim & El Ghaoui '02, '03; Iyengar '05; A. Ben-Tal, L. El Ghaoui, & A. Nemirovski '09; Bertsimas & Sim '04; Bertsimas, Brown, Caramanis '13; Lim

Can we address model misspecification in terms of both **likelihoods** and **actual outcomes**?

**Why is this important?**

A) C  
outc

B) C  
wro

Traditionally A) and B) are seen as separate mechanisms.

er '13;  
Ye '14;  
g '16;  
al. '17

Huber  
Esfahani  
regt '16;  
21.



# Why is this important?

---

A unified theory of DRO will position the area well to address a key question...

How to *practically* choose the distributional uncertainty set?

Even experts see these DRO models as fundamentally different – in some sense, they are not...

# A Unified View of DRO via **Optimal Transport** Approach with **Conditional Moment Constraints**



# Unifying Formulation

---

- OT-DRO with conditional moment constraints

$$\begin{aligned} \sup_{\mu} \{ \mathbb{E}_{\mu}[\ell(Z)] : \mathbb{D}(\mu, \hat{\mu}) \leq r \} & \longrightarrow \sup \mathbb{E}_{\pi}[W \cdot \ell(Z)] \\ & \text{s.t. } \pi \in \mathcal{P}((\mathcal{V} \times \mathcal{W}) \times (\hat{\mathcal{V}} \times \hat{\mathcal{W}})) \\ & \pi_{(\hat{V}, \hat{W})} = \hat{\nu} \\ & \mathbb{E}_{\pi}[W | \mathcal{G}] = 1 \quad \pi\text{-a.s.} \\ & \mathbb{E}_{\pi}[c((V, W), (\hat{V}, \hat{W}))] \leq r. \end{aligned}$$

$Z \rightarrow (V, W)$



Lifting Technique!

# Unifying Formulation

- OT-DRO with conditional moment constraints

$$\sup_{\mu} \{ \mathbb{E}_{\mu}[\ell(Z)] : \mathbb{D}(\mu, \hat{\mu}) \leq r \} \quad \xrightarrow{Z \rightarrow (V, W)} \quad \begin{aligned} & \sup \mathbb{E}_{\pi} [W \cdot \ell(Z)] \\ & \text{s.t. } \pi \in \mathcal{P}((\mathcal{V} \times \mathcal{W}) \times (\mathcal{V} \times \mathcal{W})) \\ & \pi_{(\hat{V}, \hat{W})} = \hat{\nu} \\ & \mathbb{E}_{\pi} [W | \mathcal{G}] = 1 \quad \pi\text{-a.s.} \\ & \mathbb{E}_{\pi} [c((V, W), (\hat{V}, \hat{W}))] \leq r. \end{aligned}$$



Lifting Technique!

If  $\mathcal{G}$  is a trivial sigma field, then

$$\mathbb{E}_{\pi} [W | \mathcal{G}] = \mathbb{E}_{\pi} [W] = 1.$$

# Unifying Formulation

- OT-DRO with conditional moment constraints

$$\begin{aligned} \sup_{\mu} \{ \mathbb{E}_{\mu}[\ell(Z)] : \mathbb{D}(\mu, \hat{\mu}) \leq r \} & \longrightarrow \sup \mathbb{E}_{\pi} [W \cdot \ell(Z)] \\ & \text{s.t. } \pi \in \mathcal{P}((\mathcal{V} \times \mathcal{W}) \times (\hat{\mathcal{V}} \times \hat{\mathcal{W}})) \\ & \pi_{(\hat{V}, \hat{W})} = \hat{\nu} \\ & \mathbb{E}_{\pi} [W | \mathcal{G}] = 1 \quad \pi\text{-a.s.} \\ & \mathbb{E}_{\pi} [c((V, W), (\hat{V}, \hat{W}))] \leq r. \end{aligned}$$

$Z \rightarrow (V, W)$



Lifting Technique!

If  $\mathcal{G}$  is the smallest sigma field generated by  $\hat{\mathcal{V}}$ , then

$$\mathbb{E}_{\pi} [W | \mathcal{G}] = \mathbb{E}_{\pi} [W | \hat{\mathcal{V}}] = 1, \pi\text{-a.s.}$$

# Unifying Formulation

---

- OT-DRO with conditional moment constraints

$$\begin{aligned} \sup_{\mu} \{ \mathbb{E}_{\mu}[\ell(Z)] : \mathbb{D}(\mu, \hat{\mu}) \leq r \} & \xrightarrow{Z \rightarrow (V, W)} \sup \mathbb{E}_{\pi}[W \cdot \ell(Z)] \\ & \text{s.t. } \pi \in \mathcal{P}((\mathcal{V} \times \mathcal{W}) \times (\mathcal{V} \times \mathcal{W})) \\ & \pi_{(\hat{V}, \hat{W})} = \hat{\nu} \\ & \mathbb{E}_{\pi}[W | \mathcal{G}] = 1 \quad \pi\text{-a.s.} \\ & \mathbb{E}_{\pi}[c((V, W), (\hat{V}, \hat{W}))] \leq r. \end{aligned}$$



$$(\ell, \mathcal{Z}, \mathbb{D}, \hat{\mu}, r) \rightarrow (\mathcal{V}, \mathcal{W}, \mathcal{G}, \hat{\nu}, c, r)$$

# Unifying Formulation

---

➤ OT-DRO with “Martingale” Constraint

$$\begin{aligned} & \sup_{\pi \in \Pi((\Omega \times \mathbb{R}_+) \times (\Omega \times \mathbb{R}_+))} \mathbb{E}_\pi [\ell(\theta, \xi) \cdot \zeta] \\ & \text{s.t. } \mathbb{E}_\pi [c_M((\xi, \zeta), (\xi', \zeta'))] \leq \delta, \\ & \mathbb{E}_\pi [\zeta \mid \zeta'] = \zeta', \\ & \Pi_{(\xi', \zeta')} \pi = \hat{\nu}. \end{aligned}$$

This is the “**baseline model**” which is constrained to be  $\hat{\nu}$ .



# We Recover Most DRO Formulations

---



Pick **cost functions & reference measure**  
Recover → Most DRO Formulations!

# Wasserstein DRO

$$\sup_{\pi \in \mathcal{P}(\Omega \times \Omega)} \left\{ \mathbb{E}_{\pi} [\ell(\theta, \xi)] : \mathbb{E}_{\pi} [c(\xi, \xi')] \leq \delta, \Pi_{\xi'} \pi = \hat{\mathbb{P}}_n \right\}$$

$$\sup_{\pi \in \Pi((\Omega \times \mathbb{R}_+) \times (\Omega \times \mathbb{R}_+))} \mathbb{E}_{\pi} [\ell(\theta, \xi) \cdot \zeta]$$

$$\text{s.t. } \mathbb{E}_{\pi} [c(\xi, \xi') + \infty \cdot \mathbb{I}_{\zeta \neq \zeta'}] \leq \delta,$$

$$\mathbb{E}_{\pi} [\zeta \mid \zeta'] = \zeta', \rightarrow \text{Automatically satisfied}$$

$$\Pi_{(\xi', \zeta')} \pi = \hat{\mathbb{P}}_n \times \delta_1.$$



$\zeta = 1, \pi$ -a.s!

# $\phi$ -divergence [Csiszar, 1963, 1967]

---

For two probability measures  $\mathbb{Q}$  and  $\hat{\mathbb{P}}_n \in \mathcal{P}(\Omega)$ , we let  $\rho$  be a dominating measure of  $\mathbb{Q}$  and  $\hat{\mathbb{P}}_n$  (i.e.,  $\mathbb{Q} \ll \rho$  and  $\hat{\mathbb{P}}_n \ll \rho$ ). Then, the  $\phi$  divergence between  $\mathbb{Q}$  and  $\hat{\mathbb{P}}_n$  is defined, independently of  $\rho$ , by

$$D_\phi(\mathbb{Q}, \hat{\mathbb{P}}_n) = \int_{\Omega} \frac{d\hat{\mathbb{P}}_n}{d\rho}(\xi) \phi \left( \frac{d\mathbb{Q}}{d\rho}(\xi) / \frac{d\hat{\mathbb{P}}_n}{d\rho}(\xi) \right) d\rho(\xi),$$

where  $0 \cdot \phi \left( \frac{0}{0} \right) := 0$ , and  $0 \cdot \phi \left( \frac{a}{0} \right) := \lim_{t \rightarrow 0} \frac{\phi(t)}{t} := \alpha \phi'_\infty, \forall \alpha > 0$ .

**The speed of growth of  $\phi$  at infity.**

# $\phi$ -divergence [Csiszar, 1963, 1967]

For two probability measures  $\mathbb{Q}$  and  $\hat{\mathbb{P}}_n \in \mathcal{P}(\Omega)$ , we let  $\rho$  be a dominating measure of  $\mathbb{Q}$  and  $\hat{\mathbb{P}}_n$  (i.e.,  $\mathbb{Q} \ll \rho$  and  $\hat{\mathbb{P}}_n \ll \rho$ ). Then, the  $\phi$  divergence between

## Fact 1 (Decomposition)

$$D_\phi(\mathbb{Q}, \hat{\mathbb{P}}_n) = \int_{\Omega} \phi \left( \frac{d\mathbb{Q}}{d\rho}(\xi) / \frac{d\hat{\mathbb{P}}_n}{d\rho}(\xi) \right) d\hat{\mathbb{P}}_n(\xi), + \phi'_\infty \mathbb{Q} \left( \frac{d\hat{\mathbb{P}}_n}{d\rho}(\xi) = 0 \right).$$

where  $0 \cdot \phi \left( \frac{\infty}{0} \right) := 0$ , and  $0 \cdot \phi \left( \frac{\infty}{0} \right) := \lim_{t \rightarrow 0} \frac{\phi(\infty)}{t} := \alpha \phi'_\infty, \forall \alpha > 0$ .

The speed of growth of  $\phi$  at infity.

# $\phi$ -divergence [Csiszar, 1963,1967]

For two probability measures  $Q$  and  $\hat{P}_n \in \mathcal{P}(\Omega)$ , we let  $\rho$  be a dominating measure of  $Q$  and  $\hat{P}_n$  (i.e.,  $Q \ll \rho$  and  $\hat{P}_n \ll \rho$ ). Then, the  $\phi$  divergence between

**Example 1 (Kullback–Leibler divergence)**

$$\phi'_\infty = \lim_{t \rightarrow \infty} \frac{t \log t}{t} = +\infty$$

$$D_{\text{KL}}(Q, \hat{P}_n) = \begin{cases} \int_{\Omega} \phi \left( \frac{dQ}{d\hat{P}_n} \right) d\hat{P}_n, & Q \ll \hat{P}_n \\ \infty, & \text{Otherwise} \end{cases}$$

The speed of growth of  $\phi$  at infity.

# $\phi$ -divergence [Csiszar, 1963, 1967]

---

For two probability measures  $\mathbb{Q}$  and  $\hat{\mathbb{P}}_n \in \mathcal{P}(\Omega)$ , we let  $\rho$  be a dominating measure of  $\mathbb{Q}$  and  $\hat{\mathbb{P}}_n$  (i.e.,  $\mathbb{Q} \ll \rho$  and  $\hat{\mathbb{P}}_n \ll \rho$ ). Then, the  $\phi$  divergence between

## Fact 2 (Asymmetry)

$$D_\phi(\mathbb{Q}, \hat{\mathbb{P}}_n) = D_\psi(\hat{\mathbb{P}}_n, \mathbb{Q}),$$

where  $\psi(t) = t\phi\left(\frac{1}{t}\right)$  represents the **Ciszar dual** of  $\phi(t)$ .



The speed of growth of  $\phi$  at infity.

# $\phi$ -divergence DRO [ $\phi'_\infty = \infty$ ]

Likelihood ratio  $\zeta$

$$\sup_{\mathbb{Q} \in \mathcal{P}(\Omega)} \left\{ \mathbb{E}_{\mathbb{Q}}[\ell(\theta, \xi)] : \mathbb{E}_{\hat{\mathbb{P}}_n} \left[ \phi \left( \frac{d\mathbb{Q}}{d\hat{\mathbb{P}}_n} \right) \right] \leq \delta \right\}$$

$$\sup_{\pi \in \Pi((\Omega \times \mathbb{R}_+) \times (\Omega \times \mathbb{R}_+))} \mathbb{E}_{\pi} [\ell(\theta, \xi) \cdot \zeta]$$

$$\text{s.t. } \mathbb{E}_{\pi} \left[ \infty \cdot \mathbb{I}_{\xi \neq \xi'} + (\phi(\zeta) - \phi(\zeta'))^+ \right] \leq \delta,$$

$$\mathbb{E}_{\pi} [\zeta] = 1,$$

$$\Pi_{(\xi', \zeta')} \pi = \hat{\mathbb{P}}_n \times \delta_1.$$



$\phi(1) = 0!$

# $\phi$ -divergence DRO [ $\phi'_\infty < \infty$ ]

$$\sup_{\mathbb{Q} \in \mathcal{P}(\Omega)} \left\{ \mathbb{E}_{\mathbb{Q}}[\ell(\theta, \xi)] : \mathbb{E}_{\rho} \left[ \frac{d\hat{\mathbb{P}}_n}{d\rho} \phi \left( \frac{d\mathbb{Q}}{d\rho} / \frac{d\hat{\mathbb{P}}_n}{d\rho} \right) \right] \leq \delta \right\}$$

**Fact**: Suppose that  $\Omega$  is compact, the worst-case distribution will be supported on **N+1 point**, i.e.,  $\xi'_{n+1} \in \arg \max_{\xi \in \Omega} \ell(\theta, \xi)$

$$\hat{\nu}(d\xi', d\zeta') = \frac{1-\epsilon}{n} \sum_i \delta_{(\xi'_i, 1-\epsilon)} + \epsilon \delta_{(\xi'_{n+1}, 0)}$$





# $\phi$ -divergence DRO [ $\phi'_\infty < \infty$ ]

$$\sup_{\pi \in \Pi((\Omega \times \mathbb{R}_+) \times (\Omega \times \mathbb{R}_+))} \mathbb{E}_\pi [\ell(\theta, \xi) \cdot \zeta]$$

$$\text{s.t. } \mathbb{E}_\pi \left[ \infty \cdot \mathbb{I}_{\xi \neq \xi'} + \frac{d\hat{\mu}}{d\rho}(\xi) \phi \left( \zeta / \frac{d\hat{\mu}}{d\rho}(\xi) \right) \right] \leq \delta,$$

$$\mathbb{E}_\pi [\zeta \mid \zeta'] = \zeta',$$

$$\Pi_{(\xi', \zeta')} \pi = \hat{\nu}.$$

$$\rho = (1 - \epsilon) \sum_{i=1}^n \delta_{\xi'_i} + \epsilon \delta_{\xi'_{n+1}}, \text{ for } \epsilon \in (0, 1)$$

$$\hat{\nu}(d\xi', d\zeta') = \frac{1 - \epsilon}{n} \sum_i \delta_{(\xi'_i, 1 - \epsilon)} + \epsilon \delta_{(\xi'_{n+1}, 0)}$$



# $\phi$ -divergence DRO

---

Divergence	$\phi(t), t \geq 0$
Kullback-Leibler	$t \cdot \log(t)$
Burg Entropy	$-\log(t) + t - 1$
$J$ -divergence	$(t - 1) \log(t)$
$\chi^2$ -distance	$\frac{1}{t}(t - 1)^2$
Modified $\chi^2$ -distance	$(t - 1)^2$
Hellinger distance	$(\sqrt{t} - 1)^2$
$\chi$ -divergence of order $n > 1$	$ t - 1 ^n$
Variation distance	$ t - 1 $
Cressie-Read	$\frac{1 - \theta + \theta t - t^\theta}{\theta(1 - \theta)}, \theta \neq 0, 1$

# Sinkhorn DRO

- **Entropic regularization** (Sinkhorn Distance) is popular in Optimal Transport applications in AI [Cuturi. (2013), Peyré & Cuturi. 2017]
- This motivated Wang et al. (2021) to consider the formulation:

$$\begin{aligned} & \sup_{\pi \in \mathcal{P}(\Omega \times \Omega)} \mathbb{E}_{\pi} [\ell(\theta, \xi)] \\ & \text{s.t.}, \mathbb{E}_{\pi} \left[ c(\xi, \xi') + \epsilon \log \left( \frac{d\pi(\xi, \xi')}{d\eta(\xi) d\hat{\mathbb{P}}_n(\xi')} \right) \right] \leq \delta \\ & \Pi_{\xi'} \pi = \hat{\mathbb{P}}_n \end{aligned}$$

Normal distribution

How to recover this formulation?



# Sinkhorn DRO

---

$$\begin{aligned} & \sup_{\pi \in \mathcal{P}(\Omega \times \Omega)} \mathbb{E}_{\pi} [\ell(\theta, \xi)] \\ \text{s.t.}, & \mathbb{E}_{\pi} \left[ c(\xi, \xi') + \epsilon \log \left( \frac{d\pi(\xi, \xi')}{d\eta(\xi) d\hat{\mathbb{P}}_n(\xi')} \right) \right] \leq \delta \\ & \Pi_{\xi'} \pi = \hat{\mathbb{P}}_n \end{aligned}$$

Normal distribution



Again lift the outcome space to  $\Omega \times \Omega \times \mathbb{R}^+$ !

# Sinkhorn DRO $\rightarrow$ KL-DRO

---

Following Wang et al. (2021), we define the kernel distribution as

$$d\nu_{\xi', \epsilon}(\xi) := \frac{\exp\left(-\frac{c(\xi, \xi')}{\epsilon}\right)}{\int_{\Omega} \exp\left(-\frac{c(x, \xi')}{\epsilon}\right) d\eta(x)} d\eta(\xi),$$

and the new reference measure as:

$$d\gamma_0(\xi, \xi') = d\nu_{\xi', \epsilon}(\xi) \times d\hat{\mathbb{P}}_n(\xi').$$



$$\sup_{\pi \in \mathcal{P}(\Omega \times \Omega)} \left\{ \mathbb{E}_{\pi}[\ell(\theta, \xi)] : \epsilon \mathbb{E}_{\gamma_0} \left[ \log \left( \frac{d\pi}{d\gamma_0} \right) \right] \leq \bar{\delta}, \pi_{\xi'} = \hat{\mathbb{P}}_n \right\}$$

**What about new – more powerful formulations?**

# New DRO Model

---

- The adversary has the ability to modify both the **actual outcomes** and the **associated probability**.

$$\begin{aligned} & \sup_{\pi} \mathbb{E}_{\pi} [\ell(\theta, \xi) \cdot \zeta] \\ \text{s.t. } & \mathbb{E}_{\pi} \left[ \gamma_1 \zeta \cdot c(\xi, \xi') + \gamma_2 (\phi(\zeta) - \phi(\zeta'))^+ \right] \leq \delta, \\ & \mathbb{E}_{\pi} [\zeta] = 1, \\ & \Pi_{(\xi', \zeta')} \pi = \hat{\mathbb{P}}_n \times \delta_1. \end{aligned} \tag{P}$$



$\gamma_1 = \infty$ , KL-DRO

# New DRO Model

---

- The adversary has the ability to modify both the **actual outcomes** and the **associated probability**.

$$\begin{aligned} & \sup_{\pi} \mathbb{E}_{\pi} [\ell(\theta, \xi) \cdot \zeta] \\ \text{s.t. } & \mathbb{E}_{\pi} \left[ \gamma_1 \zeta \cdot c(\xi, \xi') + \gamma_2 (\phi(\zeta) - \phi(\zeta'))^+ \right] \leq \delta, \\ & \mathbb{E}_{\pi} [\zeta] = 1, \\ & \Pi_{(\xi', \zeta')} \pi = \hat{\mathbb{P}}_n \times \delta_1. \end{aligned} \tag{P}$$



$\gamma_2 = \infty$ , Wasserstein DRO



# Reformulation Result

---

**Theorem.** (P) is equivalent to

$$\min_{\lambda \geq 0} \lambda \delta + \lambda \gamma_2 \log \left( \mathbb{E}_{\hat{\mathbb{P}}_n} \left[ \exp \left( \frac{\ell_{\lambda \gamma_1}^c(\xi')}{\lambda \gamma_2} \right) \right] \right)$$

where  $\phi(\zeta) = \zeta \log(\zeta) - \zeta + 1$  and the c-transform of  $\ell(\cdot)$  with penalty  $\lambda \gamma_1$  is defined as

$$\ell_{\lambda \gamma_1}^c(\xi') = \sup_{\xi \in \Omega} \{ \ell(\theta, \xi) - \lambda \gamma_1 c(\xi, \xi') \}$$

# Reformulation Result

---

**Theorem.** (P) is equivalent to

$$\min_{\lambda \geq 0} \lambda \delta + \lambda \gamma_2 \log \left( \mathbb{E}_{\hat{\mathbb{P}}_n} \left[ \exp \left( \frac{\ell_{\lambda \gamma_1}^c(\xi')}{\lambda \gamma_2} \right) \right] \right)$$

where  $\phi(\zeta) = \zeta \log(\zeta) - \zeta + 1$  and the c-transform of  $\ell(\cdot)$  with penalty  $\lambda \gamma_1$  is defined as

$$\ell_{\lambda \gamma_1}^c(\xi') = \sup_{\xi \in \Omega} \{ \ell(\theta, \xi) - \lambda \gamma_1 c(\xi, \xi') \}$$

$\gamma_1$  and  $\gamma_2$  play a critical role in controlling the LIKELIHOOD ERROR hedge vs OUTCOME ERROR hedge ← AND STILL TRACTABLE!

# Optimal Transport Plan

---

Structure of the worst case  $\pi^*$ : It must be concentrated on:

$$\left\{ (\xi, \zeta) \in \Omega \times \mathbb{R}^+ : \left[ \begin{array}{l} \xi \in \arg \max_{\xi \in \Omega} [\ell(\theta, \xi) - \lambda^* \gamma_1 c(\xi, \xi'_i)] \\ \zeta = \exp \left( \frac{\ell(\theta, \xi) - \alpha^*}{\lambda^* \gamma_2} - \frac{\gamma_1 c(\xi, \xi'_i)}{\gamma_2} \right) \end{array} \right], \forall i \in [N] \right\},$$

**Perturbation on actual outcomes!**

**Perturbation on data weights!**

Where  $\alpha^*$  is the dual variable of  $\mathbb{E}_\pi[\zeta] = 1$ .

# Tractability

---

**Theorem 2.** Suppose that the loss function  $\ell(\theta, \cdot)$  is a pointwise maximum of concave functions and  $c(\xi, \xi') = \|\xi - \xi'\|_p$ , (P) can be reformulated as a finite convex program.

Can approximate any convex function as the maxima of affine functions...



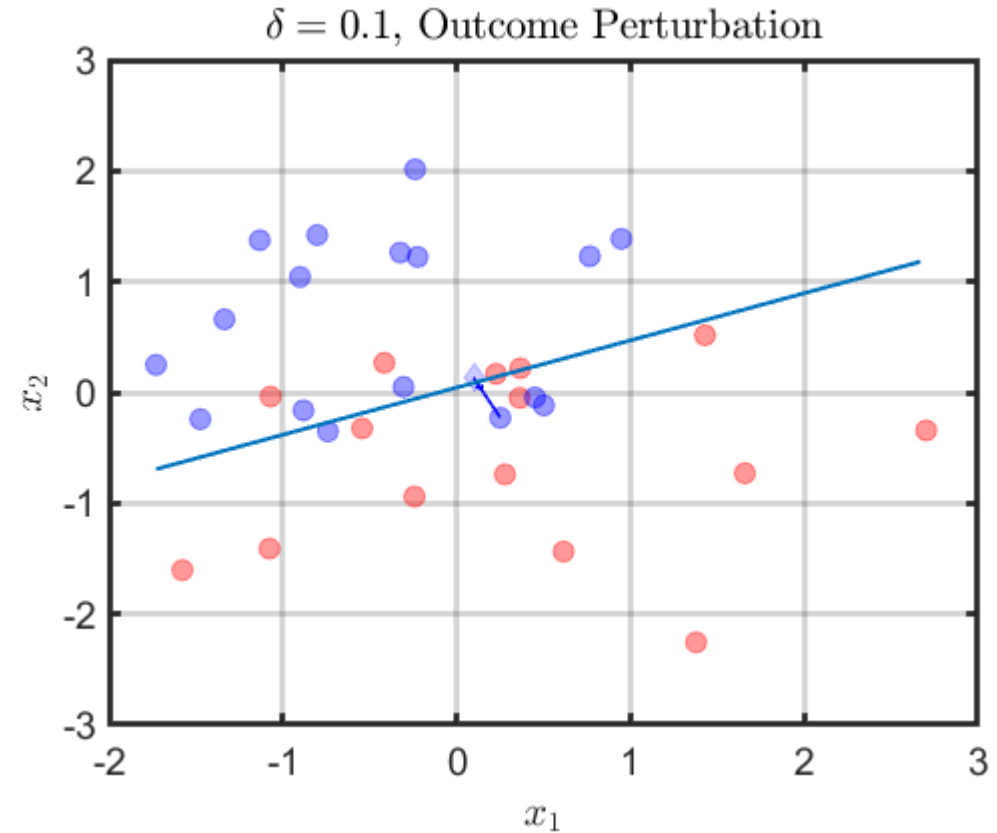
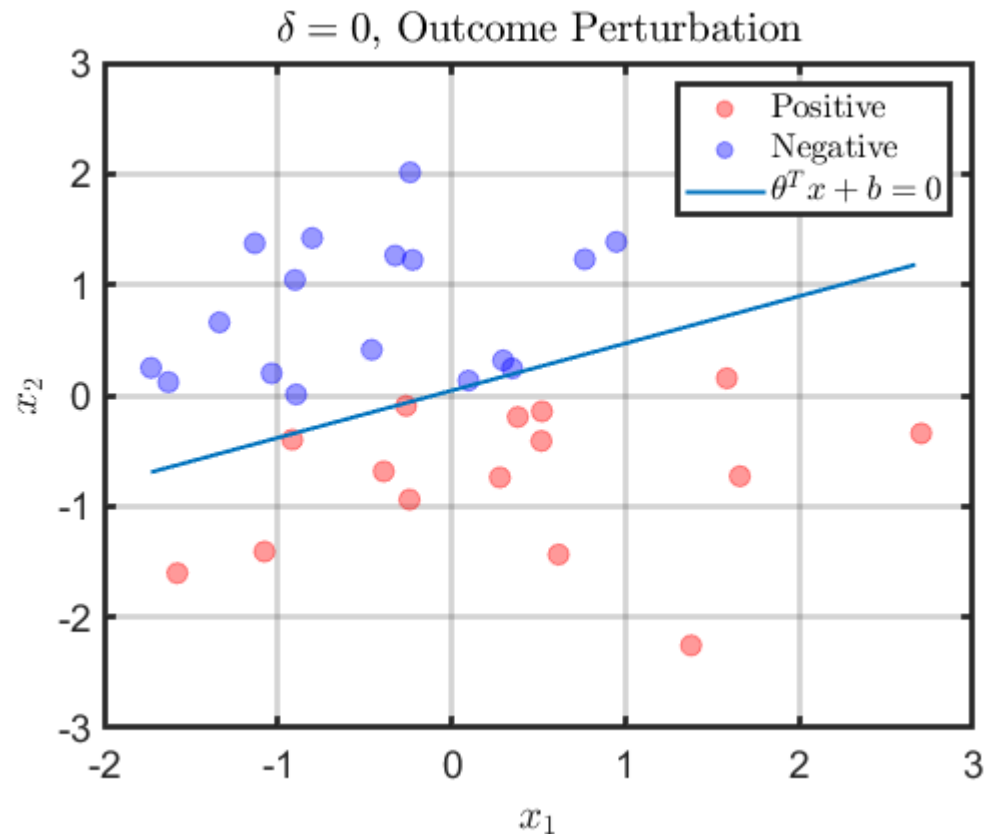
**Mosek & Gurobi help!**

# Support Vector Machine

---

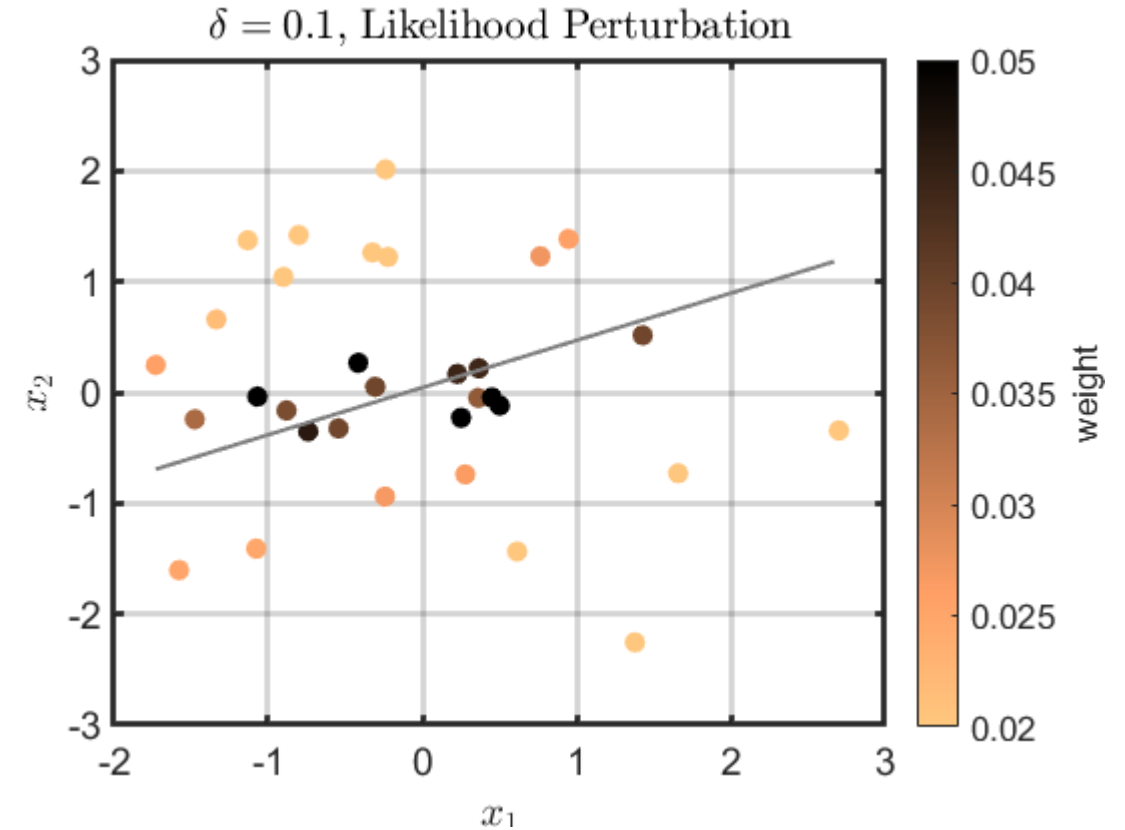
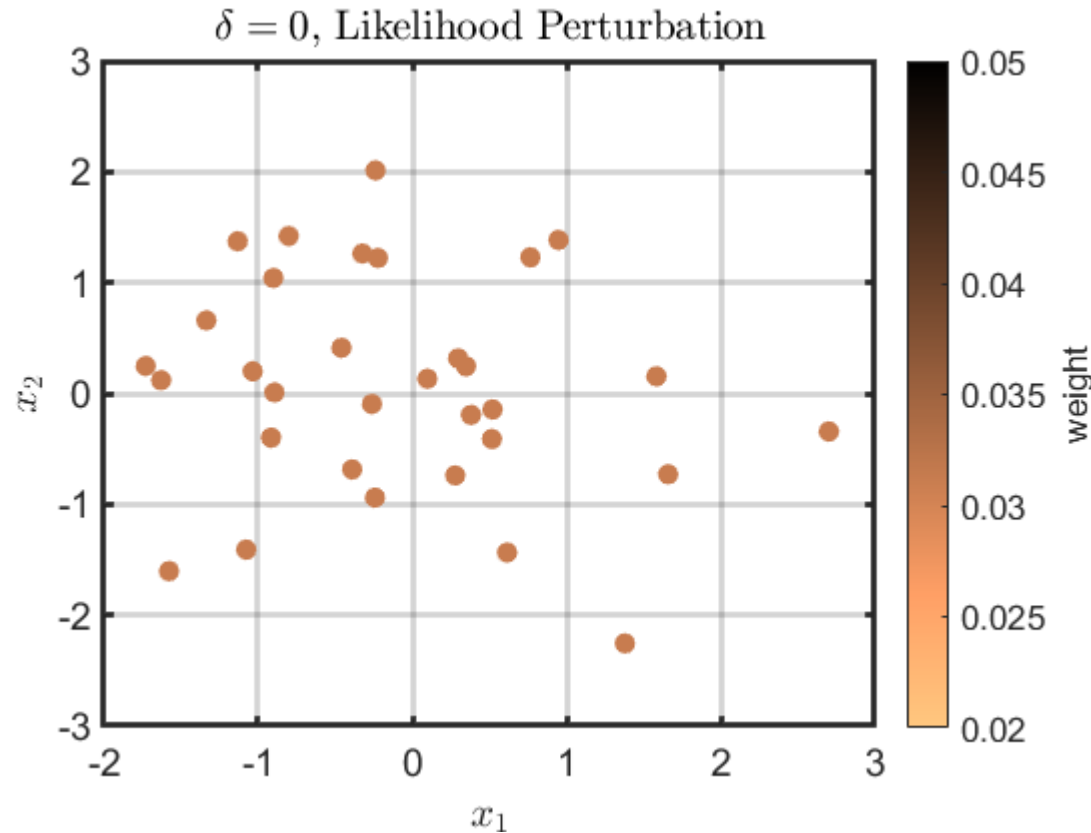
- Binary Classification Problem  $\xi = (x, y), y \in \{+1, -1\}$
- Cost function  $c(\xi, \xi') = \|x - x'\|_2^2 + \infty|y - y'|$
- Loss function  $\ell(\theta, \xi) = \max(1 - y \cdot (\theta^T x + b), 0)$
- Hyperparameter:  $\gamma_1 = \gamma_2 = 1$
- Dimension = 2

# Worst-Case Distribution Visualization



**Outcome Perturbation!**

# Worst-Case Distribution Visualization



**Likelihood Perturbation!**

# Reference:

---

1. Jose Blanchet\*, Daniel Kuhn\*, **Jiajin Li\***, Bahar Tahksen\* (Alphabetical order). Unifying Distributionally Robust Optimization via Optimal Transport Theory. **Working Paper**.
2. **Jiajin Li**, Sirui Lin, Jose Blanchet, Viet Anh Nguyen. Tikhonov Regularization is Optimal Transport Robust under Martingale Constraints, Neural Information Processing Systems (**NeurIPS**), 2022.

Thank you! Q&A?



# Strong Duality Theorem

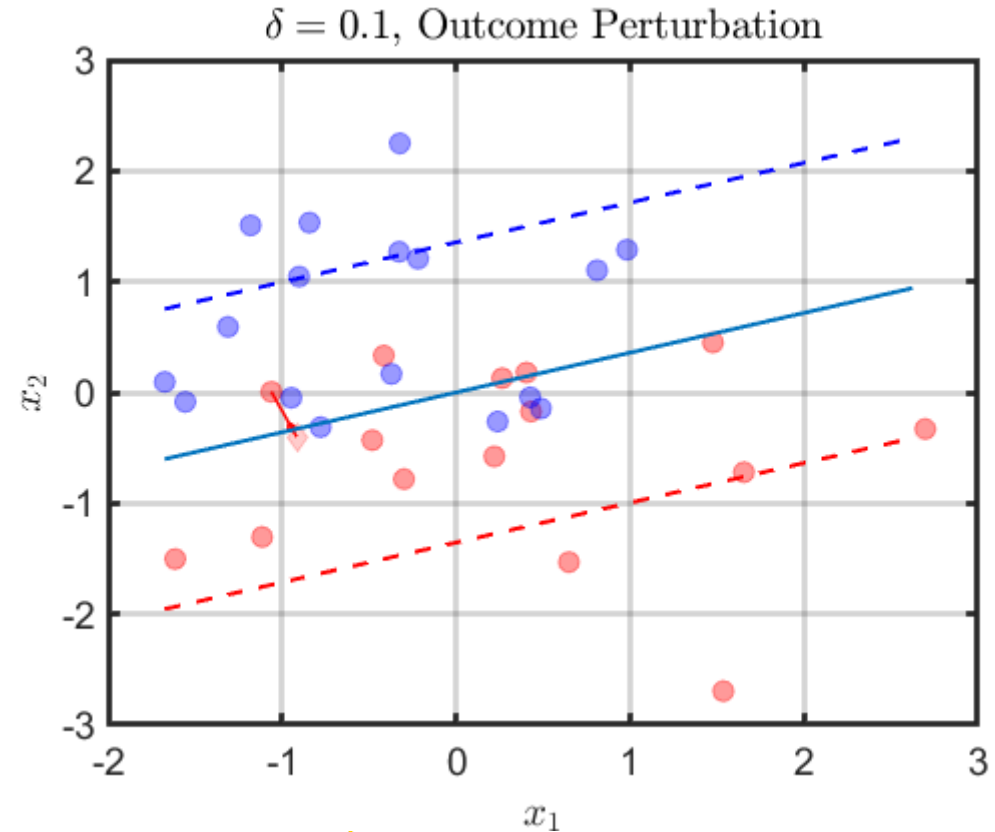
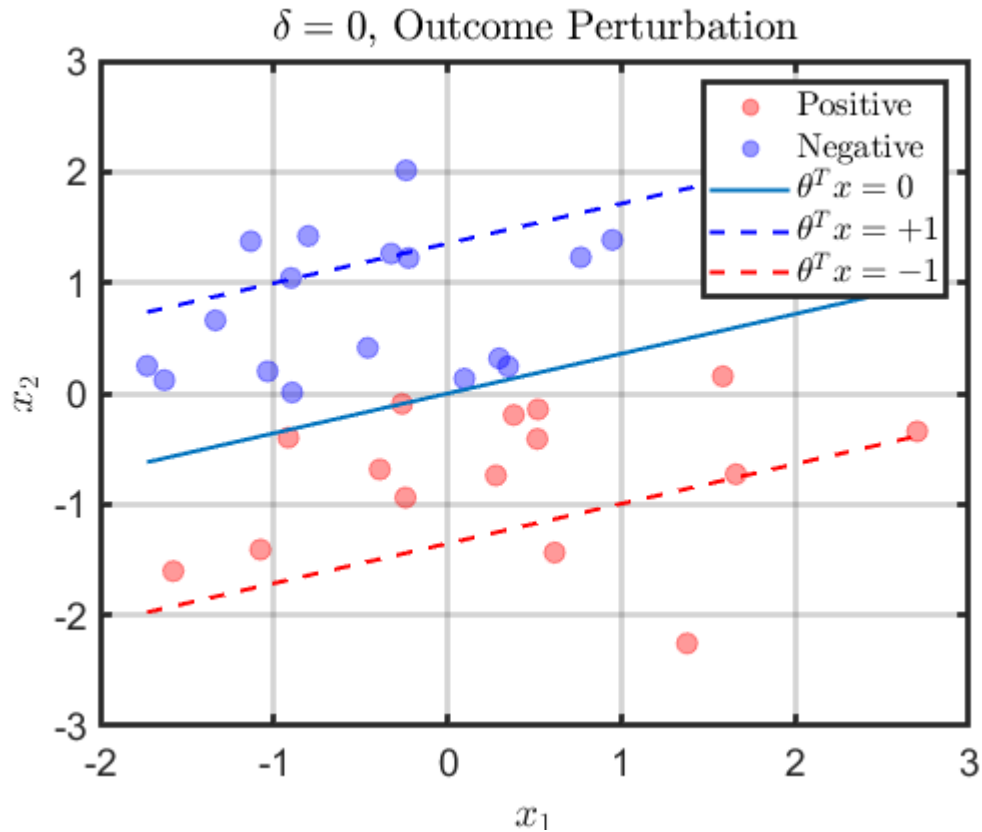
---

If the reference measure  $\hat{\nu}$  is discrete, we have

$$\begin{aligned} & \sup_{\pi \in \Pi((\Omega \times \mathbb{R}_+) \times (\Omega \times \mathbb{R}_+))} \mathbb{E}_{\pi} [\ell(\theta, \xi) \cdot \zeta] \\ & \text{s.t. } \mathbb{E}_{\pi} [c_M((\xi, \zeta), (\xi', \zeta'))] \leq \delta, \\ & \mathbb{E}_{\pi} [\zeta \mid \zeta'] = \zeta', \\ & \Pi_{(\xi', \zeta')} \pi = \hat{\nu} \\ & = \inf_{\lambda \geq 0, \alpha \in \mathbb{R}^N} \lambda \delta + \mathbb{E}_{\hat{\nu}} \left[ \sup_{(\xi, \zeta)} \ell(\theta, \xi) \cdot \zeta - \lambda c_M((\xi, \zeta), (\xi', \zeta')) + \alpha(\zeta')(\zeta - \zeta') \right] \end{aligned}$$

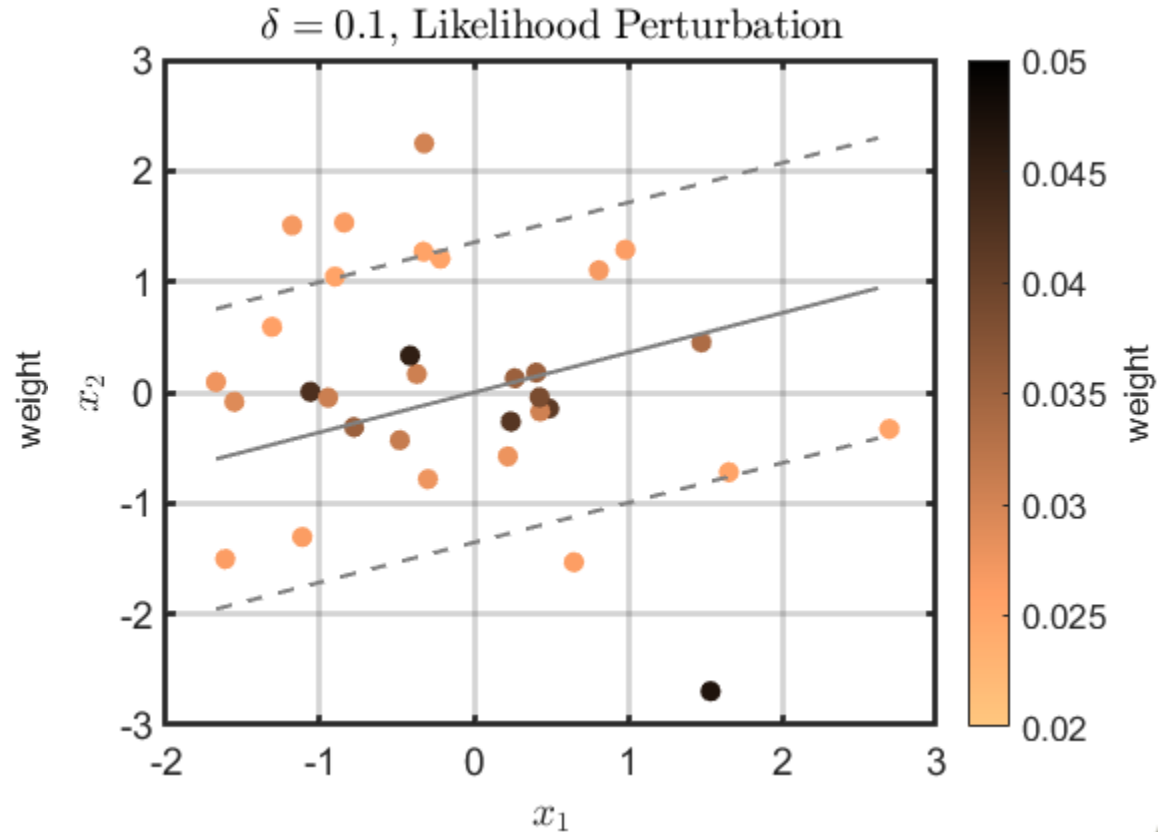
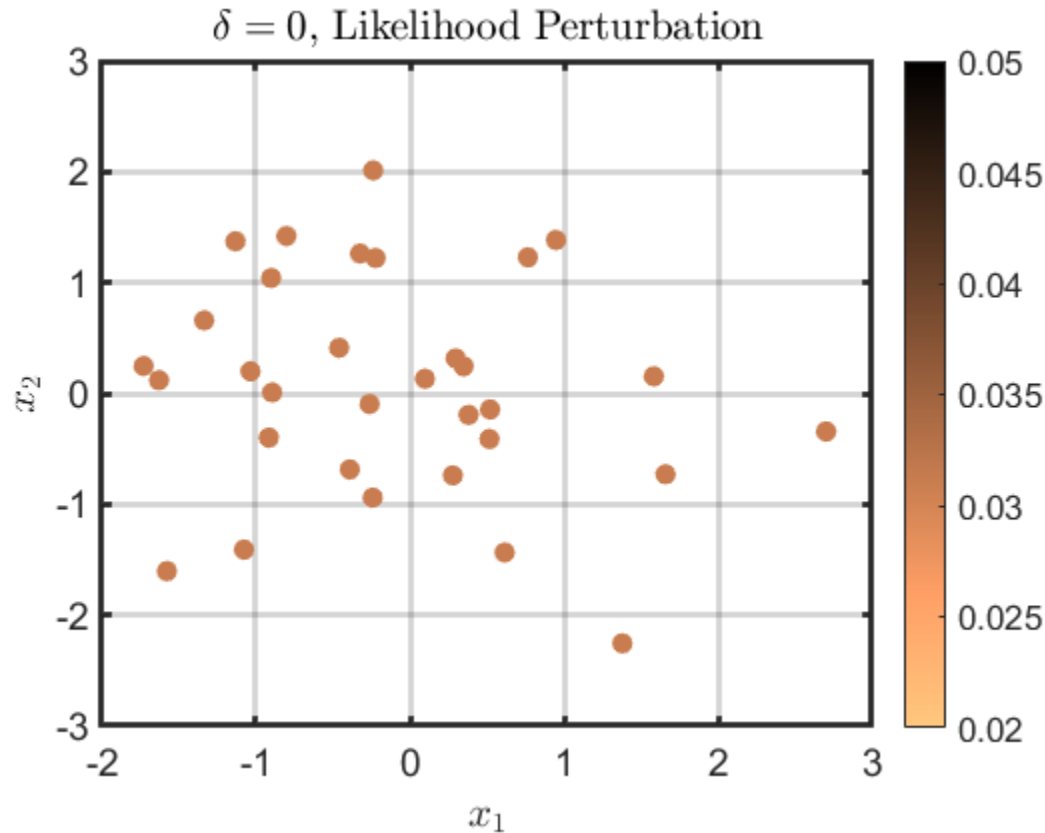
Jiajin Li, Sirui Lin, Jose Blanchet, Viet Anh Nguyen. Tikhonov Regularization is Optimal Transport Robust under Martingale Constraints, Neural Information Processing Systems (**NeurIPS**), 2022.

# Worst-Case Distribution Visualization



**Outcome Perturbation!**

# Worst-Case Distribution Visualization



**Likelihood Perturbation!**