

# Distributionally robust risk evaluation with shape constraints

Yu Gui

Department of Statistics and Data Science, the Wharton School



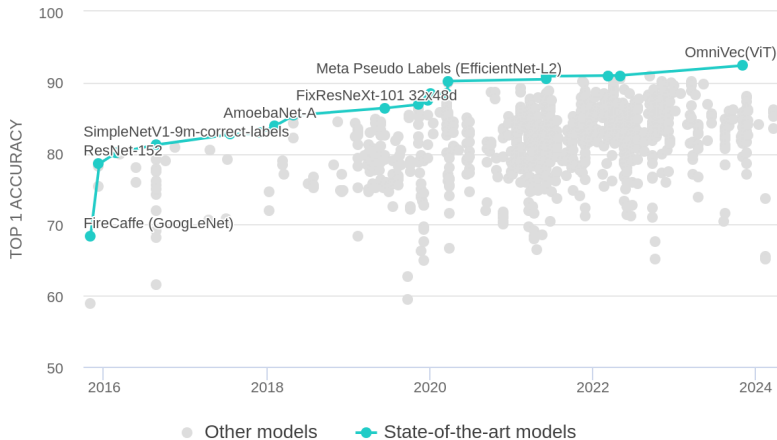


Rina Foygel Barber @UChicago



Cong Ma @UChicago

# ImageNet Dataset



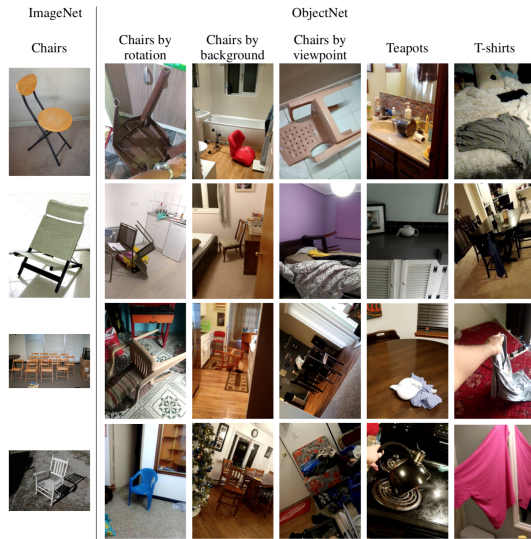
Leaderboard: image classification on ImageNet\*

\*Deng et al. (2009)

### Question

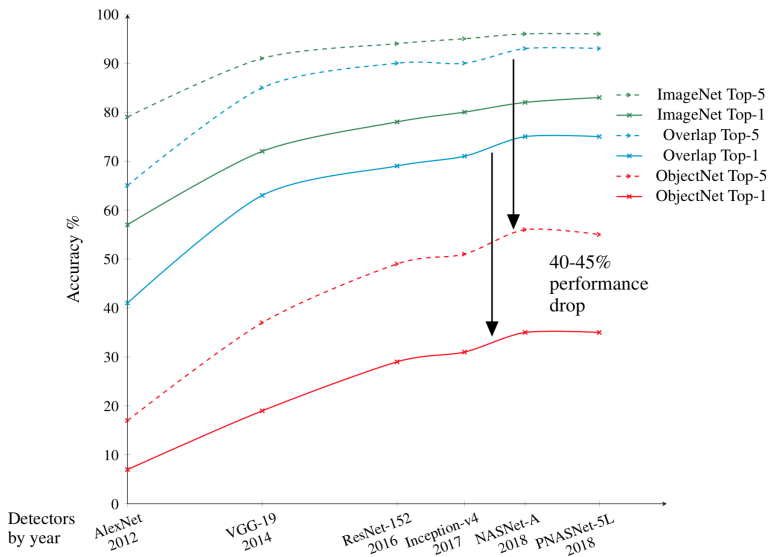
What if test distribution  $\neq$  training distribution?

# AN EXAMPLE: OBJECTNET<sup>†</sup>



<sup>†</sup>Barbu et al. (2019)

# PERFORMANCE ON OBJECTNET



### Question

How to quantify the out-of-sample performance?

# STATISTICAL INFERENCE WITH DISTRIBUTION SHIFT

$$\mathbb{E}_P[R_\alpha(X)] \leq \alpha \quad \xRightarrow{P^{\text{test}} \neq P} \quad \mathbb{E}_{P^{\text{test}}}[R_\alpha(X)] = ?$$

# STATISTICAL INFERENCE WITH DISTRIBUTION SHIFT

$$\mathbb{E}_P[R_\alpha(X)] \leq \alpha \quad \xRightarrow{P^{\text{test}} \neq P} \quad \mathbb{E}_{P^{\text{test}}}[R_\alpha(X)] = ?$$

**Example:** hypothesis test for  $P \in \mathcal{H}_0$  with data from  $P^{\text{test}}$  (Thams et al., 2023)

- Risk function

$$R_\alpha(X) = \phi_\alpha(X)$$

- Valid type-I error control with data from  $P$

$$\mathbb{P}_P(\phi_\alpha(X) = 1) \leq \alpha \quad \Longleftrightarrow \quad \mathbb{E}_P[R_\alpha(X)] \leq \alpha$$

# STATISTICAL INFERENCE WITH DISTRIBUTION SHIFT

$$\mathbb{E}_P[R_\alpha(X)] \leq \alpha \quad \xRightarrow{P^{\text{test}} \neq P} \quad \mathbb{E}_{P^{\text{test}}}[R_\alpha(X)] = ?$$

**A concrete example:** predictive inference under covariate shift<sup>‡</sup>

---

<sup>‡</sup>(Vovk et al., 2005; Tibshirani et al., 2019).

# STATISTICAL INFERENCE WITH DISTRIBUTION SHIFT

$$\mathbb{E}_P[R_\alpha(X)] \leq \alpha \quad \xRightarrow{P^{\text{test}} \neq P} \quad \mathbb{E}_{P^{\text{test}}}[R_\alpha(X)] = ?$$

**A concrete example:** predictive inference under covariate shift<sup>‡</sup>

- Prediction set  $\hat{C}_{1-\alpha}$  constructed with a dataset  $\mathcal{D}$  drawn from  $P$
- Risk function

$$R_\alpha(X) = \mathbb{P}\left(Y \notin \hat{C}_{1-\alpha}(X) \mid X\right)$$

- *Conformal prediction* CP: validity when  $\{(X, Y)\} \cup \mathcal{D}$  is exchangeable (implies  $X \sim P$ )

$$\text{for any } \alpha \in (0, 1) \quad \mathbb{P}(Y \notin \hat{C}_{1-\alpha}(X)) \leq \alpha \quad \Longleftrightarrow \quad \mathbb{E}_P[R_\alpha(X)] \leq \alpha$$

<sup>‡</sup>(Vovk et al., 2005; Tibshirani et al., 2019).

**“Estimable” distribution shift**

**“Estimable” distribution shift**

$$\hat{\mathbf{w}} \approx \frac{dP^{\text{test}}}{dP}$$

## “Estimable” distribution shift

- Covariate shift<sup>†</sup>: choose  $\beta$

$$\mathbb{E}_{P^{\text{test}}}[R_\beta(X)] = \mathbb{E}_P \left[ \frac{dP^{\text{test}}}{dP}(X) R_\beta(X) \right] \approx \mathbb{E}_P[\hat{\mathbf{w}}(X) R_\beta(X)] \leq \alpha$$

---

<sup>†</sup>(Sugiyama, 2011)

## “Estimable” distribution shift

- An example: missing at random (MAR)<sup>†</sup>

$\mathbf{M} \in \mathbb{R}^{d_1 \times d_2}$       $M_{ij}$  is observed independently with probability  $p_{ij} \in (0, 1)$

- $\mathcal{S} = \{(i, j) : M_{i,j} \text{ is observed}\}$  and  $(i_*, j_*) \mid \mathcal{S} \sim \text{Unif}(\mathcal{S}^c)$

$$\mathbb{P}\left((i_*, j_*) = (i_k, j_k) \mid \mathcal{S} \cup \{(i_*, j_*)\} = \{(i_l, j_l)\}_{l \leq n+1}\right) = \frac{(1 - p_{i_k j_k})/p_{i_k j_k}}{\sum_{l \leq n+1} (1 - p_{i_l j_l})/p_{i_l j_l}}$$

$$\text{importance sampling with “density ratio”} = \frac{1 - p_{i,j}}{p_{i,j}}$$

missingness  $\approx$  distribution shift between sampled and unsampled populations

<sup>†</sup>Gui, Yu, Rina Barber, and Cong Ma. “Conformalized matrix completion.” Advances in Neural Information Processing Systems 36 (2023): 4820-4844.

## “Estimable” distribution shift

- Covariate shift: choose  $\beta$

$$\mathbb{E}_{P^{\text{test}}}[R_\beta(X)] = \mathbb{E}_P \left[ \frac{dP^{\text{test}}}{dP}(X) R_\beta(X) \right] \approx \mathbb{E}_P[\hat{\mathbf{w}}(X) R_\beta(X)] \leq \alpha$$

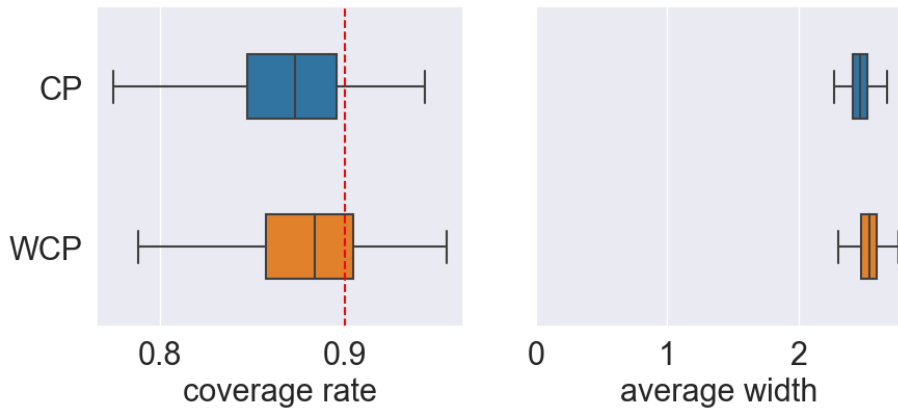
## “Estimable” distribution shift

- Covariate shift: choose  $\beta$

$$\mathbb{E}_{P^{\text{test}}}[R_\beta(X)] = \mathbb{E}_P \left[ \frac{dP^{\text{test}}}{dP}(X) R_\beta(X) \right] \approx \mathbb{E}_P[\hat{\mathbf{w}}(X) R_\beta(X)] \leq \alpha$$

An inevitable error term  $\|\mathbf{w} - \hat{\mathbf{w}}\|_1$ !

**An example with a wine quality dataset<sup>§</sup>:** white wine (4898) vs red wine (1599)



# DISTRIBUTIONALLY ROBUST LEARNING (DRL)<sup>†</sup>

**Worst-case control:** choose  $\beta$

$$\mathbb{E}_{P^{\text{test}}}[R_{\beta}(X)] \leq \sup_{Q' \in \mathcal{Q}} \mathbb{E}_{Q'}[R_{\beta}(X)] \leq \alpha \quad \text{if } P^{\text{test}} \in \mathcal{Q} \quad (\text{DRL})$$

---

<sup>†</sup>El Ghaoui and Lebrete (1997); Ben-Tal and Nemirovski (1998); Lam (2016); Duchi and Namkoong (2019); Blanchet et al. (2019)

# DISTRIBUTIONALLY ROBUST LEARNING (DRL)<sup>†</sup>

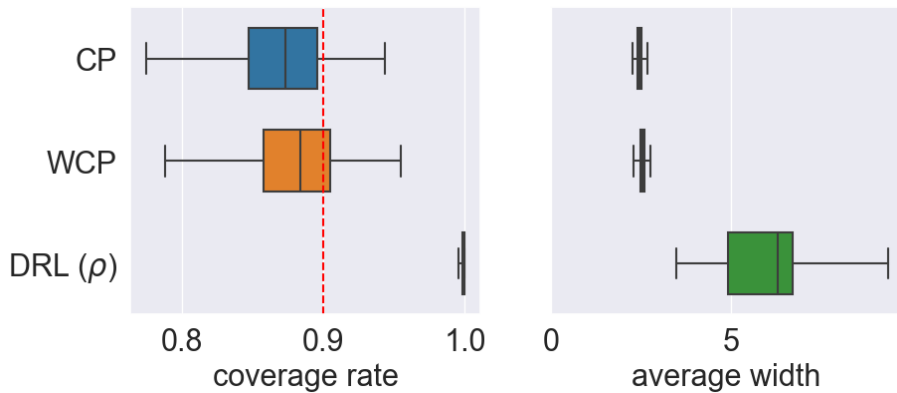
**Worst-case control:** choose  $\beta$

$$\mathbb{E}_{P^{\text{test}}}[R_{\beta}(X)] \leq \sup_{Q' \in \mathcal{Q}} \mathbb{E}_{Q'}[R_{\beta}(X)] \leq \alpha \quad \text{if } P^{\text{test}} \in \mathcal{Q} \quad (\text{DRL})$$

Too conservative/pessimistic!

---

<sup>†</sup>El Ghaoui and Le Bret (1997); Ben-Tal and Nemirovski (1998); Lam (2016); Duchi and Namkoong (2019); Blanchet et al. (2019)



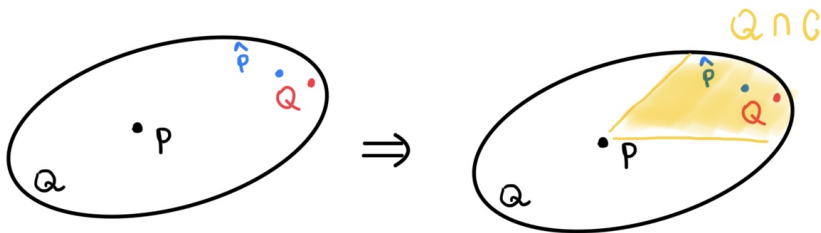
$$\rho \approx D_{\text{KL}}(P^{\text{test}} \parallel P)$$

### A middle ground?

*Misspecification of reweighting methods* VS *Overly pessimism of (DRL)*

## A middle ground?

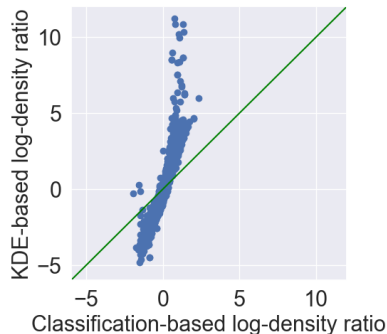
*Misspecification of reweighting methods* VS *Overly pessimism of (DRL)*



## A middle ground?

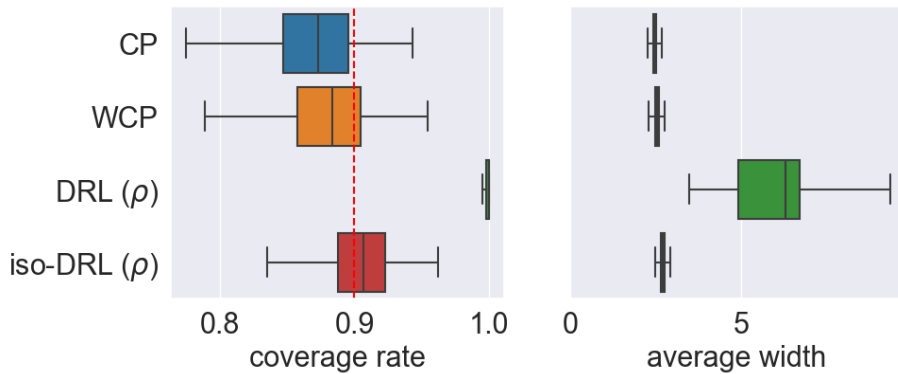
*Misspecification of reweighting methods* VS *Overly pessimism of (DRL)*

**Fitted density ratio  $\hat{w}$  vs  $\frac{dP^{\text{test}}}{dP}$  proxy:** an illustrative example with a wine quality dataset



- Biased but exhibits an approximately isotonic trend
- Under(Over)-represented regions in  $P^{\text{test}}$  are revealed by the under(over)-represented regions in  $\hat{P}$
- Use the side information to construct an additional cone constraint

$$\mathcal{Q}_{\hat{w}}^{\text{iso}} = \{Q' : dQ'/dP \text{ is isotonic in } \hat{w}\}$$



$$\rho \approx D_{\text{KL}}(P^{\text{test}} \parallel P)$$

- Under any fixed partial order  $\preccurlyeq$  on  $\mathcal{X} \subseteq \mathbb{R}^d$

$$\mathcal{Q}_{\preccurlyeq}^{\text{iso}} = \{Q' : dQ'/dP \text{ is isotonic under } \preccurlyeq\}$$

# ISO-DRL UNDER GENERAL PARTIAL ORDERS

- Under any fixed partial order  $\preccurlyeq$  on  $\mathcal{X} \subseteq \mathbb{R}^d$

$$\mathcal{Q}_{\preccurlyeq}^{\text{iso}} = \{Q' : dQ'/dP \text{ is isotonic under } \preccurlyeq\}$$

- iso-DRL chooses  $\beta$  such that

$$\sup_{Q \in \mathcal{Q} \cap \mathcal{Q}_{\preccurlyeq}^{\text{iso}}} \mathbb{E}_Q [R_\beta(X)] \leq \alpha \quad (\text{iso-DRL})$$

## Question

How to solve the cone-constrained optimization problem (iso-DRL)?

## Question

How to solve the cone-constrained optimization problem (iso-DRL)?

- At the population level: a cone-constrained optimization problem in function space?
- With a finite sample: efficient computation? consistent estimate?

## Question

How to solve the cone-constrained optimization problem (iso-DRL)?

- At the population level: a cone-constrained optimization problem in function space?
- With a finite sample: efficient computation? consistent estimate?

**Improvements over DRL?**

# AN EQUIVALENT FORMULATION

Gui et al, 2024 (Theorem 3.1)

Under regularity **conditions** on  $\mathcal{Q}$ , it holds that

$$\sup_{Q \in \mathcal{Q} \cap \mathcal{Q}_{\preceq}^{\text{iso}}} \mathbb{E}_Q [R_\beta(X)] = \sup_{Q \in \mathcal{Q}} \mathbb{E}_Q [R_\beta^{\text{iso}}(X)] \quad (\text{Equiv})$$

$$R_\beta^{\text{iso}}(X) = \operatorname{argmin}_{a \in \mathcal{C}_{\preceq}^{\text{iso}}} \int (a - R_\beta)^2 dP$$

$\mathcal{C}_{\preceq}^{\text{iso}}$  = cone of isotonic functions under  $\preceq$

# AN EQUIVALENT FORMULATION

Gui et al, 2024 (Theorem 3.1)

Under regularity **conditions** on  $\mathcal{Q}$ , it holds that

$$\sup_{Q \in \mathcal{Q} \cap \mathcal{Q}_{\approx}^{\text{iso}}} \mathbb{E}_Q [R_{\beta}(X)] = \sup_{Q \in \mathcal{Q}} \mathbb{E}_Q [\textcolor{red}{R}_{\beta}^{\text{iso}}(X)] \quad (\text{Equiv})$$

## ► Examples of $\mathcal{Q}$

- $\Gamma$ -marginal selection model in sensitivity analysis ([Rosenbaum, 1987](#); [Tan, 2006](#))

$$\mathcal{Q} = \left\{ Q : \Gamma^{-1} \leq \frac{dQ}{dP}(X) \leq \Gamma \text{ almost surely} \right\} \quad (\Gamma\text{-MS})$$

- $f$ -divergence constrained distribution shift ([Ben-Tal and Nemirovski, 1998](#); [El Ghaoui and Le Bret, 1997](#); [Duchi and Namkoong, 2019](#))

$$\mathcal{Q} = \{ Q : D_f(Q \parallel P) \leq \rho \} \quad (f\text{-Div})$$

# AN EQUIVALENT FORMULATION

Gui et al, 2024 (Theorem 3.1)

Under regularity **conditions** on  $\mathcal{Q}$ , it holds that

$$\sup_{Q \in \mathcal{Q} \cap \mathcal{Q}_{\mathcal{Y}}^{\text{iso}}} \mathbb{E}_Q [R_{\beta}(X)] = \sup_{Q \in \mathcal{Q}} \mathbb{E}_Q [\textcolor{red}{R}_{\beta}^{\text{iso}}(X)] \quad (\text{Equiv})$$

► **Examples of  $\mathcal{Q}$**

- $\Gamma$ -marginal selection model in sensitivity analysis
- $f$ -divergence constrained distribution shift

► **Two sources of computational costs are separated:**

# AN EQUIVALENT FORMULATION

Gui et al, 2024 (Theorem 3.1)

Under regularity **conditions** on  $\mathcal{Q}$ , it holds that

$$\sup_{Q \in \mathcal{Q} \cap \mathcal{Q}_{\mathbb{X}}^{\text{iso}}} \mathbb{E}_Q [R_{\beta}(X)] = \sup_{Q \in \mathcal{Q}} \mathbb{E}_Q [\textcolor{red}{R}_{\beta}^{\text{iso}}(X)] \quad (\text{Equiv})$$

► **Examples of  $\mathcal{Q}$**

- $\Gamma$ -marginal selection model in sensitivity analysis
- $f$ -divergence constrained distribution shift

► **Two sources of computational costs are separated:**

- $\mathcal{Q} \longrightarrow$  computational cost in solving (**DRL**) with  $\textcolor{red}{R}_{\beta}^{\text{iso}}$

# AN EQUIVALENT FORMULATION

Gui et al, 2024 (Theorem 3.1)

Under regularity **conditions** on  $\mathcal{Q}$ , it holds that

$$\sup_{Q \in \mathcal{Q} \cap \mathcal{Q}_{\preceq}^{\text{iso}}} \mathbb{E}_Q [R_{\beta}(X)] = \sup_{Q \in \mathcal{Q}} \mathbb{E}_Q [R_{\beta}^{\text{iso}}(X)] \quad (\text{Equiv})$$

## ► Examples of $\mathcal{Q}$

- $\Gamma$ -marginal selection model in sensitivity analysis
- $f$ -divergence constrained distribution shift

## ► Two sources of computational costs are separated:

- $\mathcal{Q} \longrightarrow$  computational cost in solving (**DRL**) with  $R_{\beta}^{\text{iso}}$
- $\mathcal{Q}_{\preceq}^{\text{iso}} \longrightarrow$  isotonic projection of  $R$

# AN EQUIVALENT FORMULATION

Gui et al, 2024 (Theorem 3.1)

Under regularity **conditions** on  $\mathcal{Q}$ , it holds that

$$\sup_{Q \in \mathcal{Q} \cap \mathcal{Q}_{\preceq}^{\text{iso}}} \mathbb{E}_Q [R_{\beta}(X)] = \sup_{Q \in \mathcal{Q}} \mathbb{E}_Q [R_{\beta}^{\text{iso}}(X)] \quad (\text{Equiv})$$

## ► Examples of $\mathcal{Q}$

- $\Gamma$ -marginal selection model in sensitivity analysis
- $f$ -divergence constrained distribution shift

## ► Two sources of computational costs are separated:

- $\mathcal{Q} \longrightarrow$  computational cost in solving (**DRL**) with  $R_{\beta}^{\text{iso}}$
- $\mathcal{Q}_{\preceq}^{\text{iso}} \longrightarrow$  isotonic projection of  $R$

- (**Equiv**) holds at **both population and sample levels**: *reference measure can be  $P$  or  $\hat{P}_n$*

# AN EQUIVALENT FORMULATION

Gui et al, 2024 (Theorem 3.1)

Under regularity **conditions** on  $\mathcal{Q}$ , it holds that

$$\sup_{Q \in \mathcal{Q} \cap \mathcal{Q}_{\infty}^{\text{iso}}} \mathbb{E}_Q [R_{\beta}(X)] = \sup_{Q \in \mathcal{Q}} \mathbb{E}_Q [\textcolor{red}{R}_{\beta}^{\text{iso}}(X)] \quad (\text{Equiv})$$

Shape constraints protect against “nonsmooth” or adversarial distribution shifts

# FINITE-SAMPLE ESTIMATE

$$\Delta^{\text{iso}}(R; \mathcal{Q}) = \sup_{Q \in \mathcal{Q} \cap \mathcal{Q}_{\leq}^{\text{iso}}} \mathbb{E}_Q [R_{\beta}(X)] = \sup_{w \# P \in \mathcal{Q} \cap \mathcal{Q}_{\leq}^{\text{iso}}} \mathbb{E}_P [w(X) \cdot R_{\beta}(X)]$$

# FINITE-SAMPLE ESTIMATE

$$\Delta^{\text{iso}}(R; \mathcal{Q}) = \sup_{Q \in \mathcal{Q} \cap \mathcal{Q}_{\preceq}^{\text{iso}}} \mathbb{E}_Q [R_{\beta}(X)] = \sup_{w \# P \in \mathcal{Q} \cap \mathcal{Q}_{\preceq}^{\text{iso}}} \mathbb{E}_P [w(X) \cdot R_{\beta}(X)]$$



# FINITE-SAMPLE ESTIMATE

$$\Delta^{\text{iso}}(R; \mathcal{Q}) = \sup_{Q \in \mathcal{Q} \cap \mathcal{Q}_{\preceq}^{\text{iso}}} \mathbb{E}_Q [R_\beta(X)] = \sup_{w_{\#} P \in \mathcal{Q} \cap \mathcal{Q}_{\preceq}^{\text{iso}}} \mathbb{E}_P [w(X) \cdot R_\beta(X)]$$



$$\begin{aligned} \hat{\Delta}^{\text{iso}}(\mathcal{Q}) &= \sup_{w_{\#} P \in \mathcal{Q} \cap \mathcal{Q}_{\preceq}^{\text{iso}}} \mathbb{E}_{\hat{P}_n} [w(X) \cdot r_\beta(X)] \\ &\stackrel{(\text{Equiv})}{=} \sup_{w_{\#} P \in \mathcal{Q}} \mathbb{E}_{\hat{P}_n} [w(X) \cdot \hat{r}_\beta^{\text{iso}}(X)] \end{aligned}$$

# FINITE-SAMPLE ESTIMATE

$$\Delta^{\text{iso}}(R; \mathcal{Q}) = \sup_{Q \in \mathcal{Q} \cap \mathcal{Q}_{\leq}^{\text{iso}}} \mathbb{E}_Q [R_{\beta}(X)] = \sup_{w_{\#} P \in \mathcal{Q} \cap \mathcal{Q}_{\leq}^{\text{iso}}} \mathbb{E}_P [w(X) \cdot R_{\beta}(X)]$$



$$\begin{aligned} \hat{\Delta}^{\text{iso}}(\mathcal{Q}) &= \sup_{w_{\#} P \in \mathcal{Q} \cap \mathcal{Q}_{\leq}^{\text{iso}}, \|w\|_{\infty} \leq \Omega} \mathbb{E}_{\hat{P}_n} [w(X) \cdot r_{\beta}(X)] \\ &\stackrel{(\text{Equiv})}{=} \sup_{w_{\#} P \in \mathcal{Q}, \|w\|_{\infty} \leq \Omega} \mathbb{E}_{\hat{P}_n} [w(X) \cdot \hat{r}_{\beta}^{\text{iso}}(X)] \end{aligned}$$

# FINITE-SAMPLE ESTIMATE

$$\begin{aligned}\hat{\Delta}^{\text{iso}}(\mathcal{Q}) &= \sup_{w \# P \in \mathcal{Q} \cap \mathcal{Q}_{\preceq}^{\text{iso}}, \|w\|_{\infty} \leq \Omega} \mathbb{E}_{\hat{P}_n} [w(X) \cdot r_{\beta}(X)] \\ &\stackrel{(\text{Equiv})}{=} \sup_{w \# P \in \mathcal{Q}, \|w\|_{\infty} \leq \Omega} \mathbb{E}_{\hat{P}_n} [w(X) \cdot \hat{r}_{\beta}^{\text{iso}}(X)]\end{aligned}$$

- $r_{\beta}(X)$  is a noisy observation of  $R_{\beta}(X)$
- $\hat{r}_{\beta}^{\text{iso}}(X)$  is the isotonic projection of  $r_{\beta}(X)$  w.r.t.  $\hat{P}_n$

# FINITE-SAMPLE ESTIMATE

Gui et al, 2024 (Theorem 4.4, informal)

For both ( $\Gamma$ -MS) and ( $f$ -Div) with adequately large  $\Omega$ ,

$$\left| \Delta^{\text{iso}}(R; \mathcal{Q}) - \hat{\Delta}^{\text{iso}}(\mathcal{Q}) \right| \lesssim \mathcal{R}_n(\mathcal{C}_{\leq, \Omega}^{\text{iso}}) + \sqrt{\frac{\log n}{n}}$$

# FINITE-SAMPLE ESTIMATE

Gui et al, 2024 (Theorem 4.4, informal)

For both ( $\Gamma$ -MS) and ( $f$ -Div) with adequately large  $\Omega$ ,

$$\left| \Delta^{\text{iso}}(R; \mathcal{Q}) - \hat{\Delta}^{\text{iso}}(\mathcal{Q}) \right| \lesssim \mathcal{R}_n(\mathcal{C}_{\preccurlyeq, \Omega}^{\text{iso}}) + \sqrt{\frac{\log n}{n}}$$

Bounding the Rademacher complexity

►  $d = 1$  (Chatterjee and Lafferty, 2019)

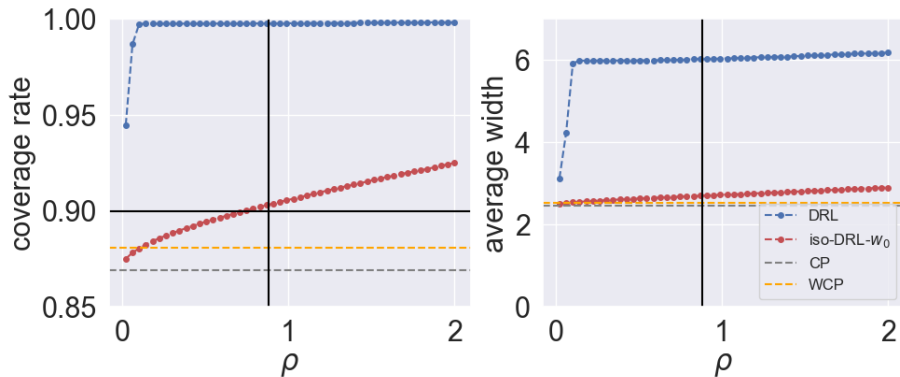
$$\mathcal{R}_n(\mathcal{C}_{\preccurlyeq, \Omega}^{\text{iso}}) \lesssim n^{-1/2}$$

►  $d \geq 2$  with componentwise order, i.e.  $\mathbf{x} \preccurlyeq \mathbf{z}$  iff  $x_i \leq z_i$  for all  $i \in [d]$  (Han et al., 2019)

$$\mathcal{R}_n(\mathcal{C}_{\preccurlyeq, \Omega}^{\text{iso}}) \lesssim n^{-1/d}$$

# EMPIRICAL PERFORMANCE

Wine quality data set with varying  $\rho$



► Conditional distribution

$$Y \mid X \sim \mathcal{N}(X^\top \beta + \sin(X_1) + 0.2X_3^2, 1)$$

► Marginal distributions

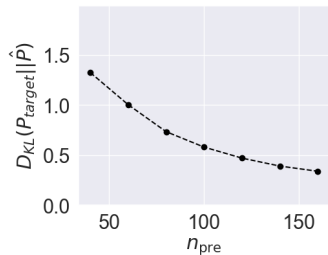
$$\begin{cases} \text{training distribution} & P : X \sim \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d), \\ \text{test distribution} & P^{\text{test}} : X \sim \mathcal{N}(\mu, \mathbf{I}_d + \zeta \cdot \mathbf{\Omega}), \end{cases}$$

►  $d = 5$ ,  $\mathbf{\Omega} = \mathbf{11}^\top$ , and  $\mu = (2/\sqrt{d}) \cdot (1, \dots, 1)^\top$

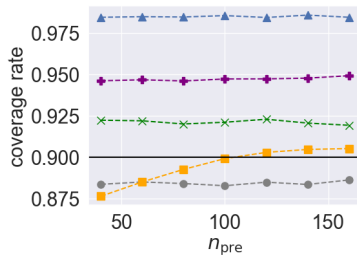
►  $\zeta = 0$ : well-specified  $\hat{w}$  via logistic regression;  $\zeta > 0$ : misspecified  $\hat{w}$

# VARYING SPLITTING RATIO $\eta$ : WELL-SPECIFIED DENSITY RATIO

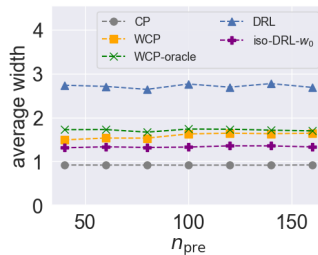
Estimated density ratio  $\hat{w}$  via logistic regression using  $\eta \times 100\%$  data



(a)  $D_{\text{KL}}(P^{\text{test}} || \hat{P})$  versus  $n_{\text{pre}}$



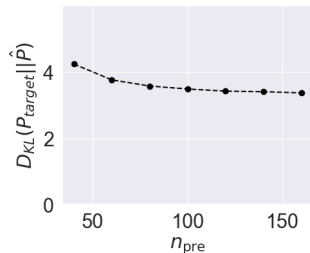
(b) Comparison with varying  $n_{\text{pre}}$



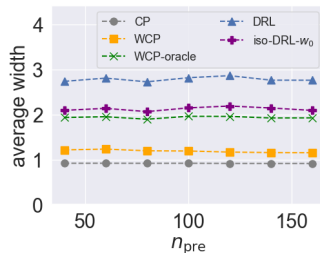
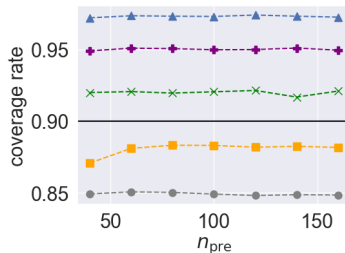
Results with well-specified density ratio ( $\zeta = 0$ )<sup>¶</sup>

<sup>¶</sup> $\rho = \rho^* = D_{\text{KL}}(P^{\text{test}} || P)$

# VARYING SPLITTING RATIO $\eta$ : MISSPECIFIED DENSITY RATIO



(a)  $D_{\text{KL}}(P^{\text{test}} || \hat{P})$  versus  $n_{\text{pre}}$



(b) Comparison with varying  $n_{\text{pre}}$

Results with misspecified density ratio ( $\zeta = 1$ )

- Distribution shift can harm the validity of statistical inference
- By incorporating shape constraints, (iso-DRL) offers one way to balance the misspecification of reweighting methods and the pessimism of DRL

Thank you!

# CONDITION ON $\mathcal{Q}$

- Change of “variable”

$$Q \in \mathcal{Q} \quad \text{if and only if} \quad w_{\#}P \in \mathcal{B}$$

- Convex ordering ( $\preceq^{cvx}$ ): for two distributions  $Q$  and  $P$ ,

$$Q \preceq^{cvx} P \quad \text{if and only if} \quad \mathbb{E}_Q[\psi(X)] \leq \mathbb{E}_P[\psi(X)] \quad \text{for any convex function } \psi$$

## Condition (Closedness under convex ordering)

The set  $\mathcal{B}$  is closed under convex ordering such that

$$\text{if } Q' \in \mathcal{B}, \text{ then } Q \in \mathcal{B} \text{ for any } Q \preceq^{cvx} Q' \quad (\text{conditions})$$

## A DETOUR: CONFORMAL PREDICTION

- Any distribution  $P_{X,Y}$  (completely unknown)
- $\{(X_i, Y_i)\}_{i \leq n+1} \sim P_{X,Y}$  are exchangeable with unobserved  $Y_{n+1}$

### Finite-sample validity

Construct marginal confidence intervals any  $\alpha \in (0, 1)$

$$\mathbb{P}(Y_{n+1} \in C_{1-\alpha}(X_{n+1})) \geq 1 - \alpha$$

# SPLIT CONFORMAL PREDICTION

- Split dataset into a training set and a calibration set  $\mathcal{D}_{\text{calib}} = \{(X_i, Y_i)\}_{i \leq n}$
- Prefit  $\hat{\mu} : \mathcal{X} \rightarrow \mathcal{Y}$  on the training set  $\implies$  nonconformity score  $R(x, y)$

# SPLIT CONFORMAL PREDICTION

- Split dataset into a training set and a calibration set  $\mathcal{D}_{\text{calib}} = \{(X_i, Y_i)\}_{i \leq n}$
- Prefit  $\hat{\mu} : \mathcal{X} \rightarrow \mathcal{Y}$  on the training set  $\implies$  nonconformity score  $R(x, y)$
- Exchangeability of  $\mathcal{D}_{\text{calib}} \cup \{(X_{n+1}, Y_{n+1})\}$

$$\left( R(X_{n+1}, Y_{n+1}) \mid \{R(x_i, y_i)\}_{i \leq n+1} \right) \sim \frac{1}{n+1} \sum_{i=1}^{n+1} \delta_{R(x_i, y_i)}$$

Calculate the quantile

$$q_{1-\alpha} = \text{Quantile}_{1-\alpha} \left( \frac{1}{n+1} \sum_{i=1}^{n+1} \delta_{R(x_i, y_i)} + \frac{1}{n+1} \delta_{\infty} \right)$$

Construct the prediction interval

$$C_{1-\alpha}(X_{n+1}) = \left\{ y : R(X_{n+1}, y) \leq q_{1-\alpha} \right\}$$

# TWO INGREDIENTS OF CONFORMAL PREDICTION

- Exchangeable data  $\{(X_i, Y_i)\}_{i \leq n+1}$
- Symmetric algorithm  $\mathcal{A}$  (not required in split conformal prediction)

# TWO INGREDIENTS OF CONFORMAL PREDICTION

- Exchangeable data  $\{(X_i, Y_i)\}_{i \leq n+1}$
- Symmetric algorithm  $\mathcal{A}$  (not required in split conformal prediction)

**Question:** What if  $\{(X_i, Y_i)\}_{i \leq n+1}$  are not exchangeable? How can we fix this?

# CP UNDER WEIGHTED EXCHANGEABILITY

## ► Weighted exchangeability

### Definition (Tibshirani et al., 2019)

Random variables  $\{V_i\}_{i \leq n+1}$  are said to be weighted exchangeable with weight functions  $\{w_i\}_{i \leq n+1}$  if the joint density can be factorized by

$$f(v_1, \dots, v_{n+1}) = \left\{ \prod_{i \leq n+1} w_i(v_i) \right\} \cdot g(v_1, \dots, v_{n+1})$$

where  $g$  is any function that does not depend on the ordering of its inputs.

# CP UNDER WEIGHTED EXCHANGEABILITY

- If  $\{Z_i = (X_i, Y_i)\}_{i \leq n+1}$  are weighted exchangeable with weight functions  $w_i$

$$\left\{ R(Z_{n+1}) \middle| \{R(z_i)\}_{i \leq n+1} \right\} \sim \sum_{i \leq n+1} p_i(Z_1, \dots, Z_{n+1}) \delta_{R(Z_i)}$$

where  $p_i$ 's are standardized weights

$$p_i^w(z_1, \dots, z_{n+1}) = \frac{\sum_{\sigma: \sigma(n+1)=i} \prod_{j \leq n+1} w_j(z_{\sigma(j)})}{\sum_{\sigma} \prod_{j \leq n+1} w_j(z_{\sigma(j)})}, \quad i = 1, \dots, n+1$$

# CP UNDER WEIGHTED EXCHANGEABILITY

- If  $\{Z_i = (X_i, Y_i)\}_{i \leq n+1}$  are weighted exchangeable with weight functions  $w_i$

$$\left\{ R(Z_{n+1}) \middle| \{R(z_i)\}_{i \leq n+1} \right\} \sim \sum_{i \leq n+1} p_i(Z_1, \dots, Z_{n+1}) \delta_{R(Z_i)}$$

where  $p_i$ 's are standardized weights

$$p_i^w(z_1, \dots, z_{n+1}) = \frac{\sum_{\sigma: \sigma(n+1)=i} \prod_{j \leq n+1} w_j(z_{\sigma(j)})}{\sum_{\sigma} \prod_{j \leq n+1} w_j(z_{\sigma(j)})}, \quad i = 1, \dots, n+1$$

- Construct the prediction interval

$$\hat{C}_{1-\alpha}(X_{n+1}) = \{y \in \mathcal{Y} : R(X_{n+1}, y) \leq q_{1-\alpha}^w\}$$

with the threshold

$$q_{1-\alpha}^w = \text{Quantile}_{1-\alpha} \left( \sum_{i \leq n} p_i^w \delta_{R(Z_i)} + p_{n+1}^w \delta_{\infty} \right)$$

- Barbu, A., D. Mayo, J. Alverio, W. Luo, C. Wang, D. Gutfreund, J. Tenenbaum, and B. Katz (2019). Objectnet: A large-scale bias-controlled dataset for pushing the limits of object recognition models. *Advances in neural information processing systems* 32.
- Ben-Tal, A. and A. Nemirovski (1998). Robust convex optimization. *Mathematics of operations research* 23(4), 769–805.
- Blanchet, J., Y. Kang, and K. Murthy (2019). Robust wasserstein profile inference and applications to machine learning. *Journal of Applied Probability* 56(3), 830–857.
- Chatterjee, S. and J. Lafferty (2019). Adaptive risk bounds in unimodal regression. *Bernoulli*.
- Cortez, P., A. Cerdeira, F. Almeida, T. Matos, and J. Reis (2009). Modeling wine preferences by data mining from physicochemical properties. *Decision support systems* 47(4), 547–553.
- Deng, J., W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei (2009). Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee.
- Duchi, J. and H. Namkoong (2019). Variance-based regularization with convex objectives. *Journal of Machine Learning Research* 20(68), 1–55.
- El Ghaoui, L. and H. Lebret (1997). Robust solutions to least-squares problems with uncertain data. *SIAM Journal on matrix analysis and applications* 18(4), 1035–1064.
- Han, Q., T. Wang, S. Chatterjee, and R. J. Samworth (2019). Isotonic regression in general dimensions. *The Annals of Statistics*.
- Lam, H. (2016). Robust sensitivity analysis for stochastic systems. *Mathematics of Operations Research* 41(4), 1248–1275.

- Rosenbaum, P. R. (1987). Sensitivity analysis for certain permutation inferences in matched observational studies. *Biometrika* 74(1), 13–26.
- Sugiyama, M. (2011). Learning under non-stationarity: Covariate shift adaptation by importance weighting. In *Handbook of Computational Statistics: Concepts and Methods*, pp. 927–952. Springer.
- Tan, Z. (2006). A distributional approach for causal inference using propensity scores. *Journal of the American Statistical Association* 101(476), 1619–1637.
- Thams, N., S. Saengkyongam, N. Pfister, and J. Peters (2023). Statistical testing under distributional shifts. *Journal of the Royal Statistical Society Series B: Statistical Methodology* 85(3), 597–663.
- Tibshirani, R. J., R. F. Barber, E. J. Candès, and A. Ramdas (2019). Conformal prediction under covariate shift. In *NeurIPS*.
- Vovk, V., A. Gammerman, and G. Shafer (2005). *Algorithmic learning in a random world*, Volume 29. Springer.