# Causality with Robust Machine Learning

清华大学 刘家硕

THU Trustworthy-Al Group

2020年11月11日

- 1 Invariant Causal Prediction
- ② Distributionally Robust Learning
- 3 Relationships
- 4 Conclusions

# Paper: Causal inference using invariant prediction: identification and confidence intervals<sup>1</sup>

- Main idea: use the invariance of the causal relationships under different environments for causal inference.
- Problem setting: data  $(X^e, Y^e)$  from different environments  $e \in \mathcal{E}$  with unknown experimental conditions or type of interventions.
- Invariance Assumption:

Assumption 1 (Invariant prediction) There exists a vector of coefficients  $\gamma^* = (\gamma_1^*, \dots, \gamma_p^*)^t$  with support  $S^* := \{k : \gamma_k^* \neq 0\} \subseteq \{1, \dots, p\}$  that satisfies

for all 
$$e \in \mathcal{E}$$
:  $X^e$  has an arbitrary distribution and 
$$Y^e = \mu + X^e \gamma^* + \varepsilon^e, \quad \varepsilon^e \sim F_{\varepsilon} \text{ and } \varepsilon^e \perp \!\!\! \perp X^e_{S^*}, \tag{3}$$

where  $\mu \in \mathbb{R}$  is an intercept term,  $\varepsilon^e$  is random noise with mean zero, finite variance and the same distribution  $F_{\varepsilon}$  across all  $e \in \mathcal{E}$ .

- Equivalent to  $P(Y^e|X_{S*}^e)$  identical for all  $e \in \mathcal{E}$ .
- S\* corresponds to what kind of covariates?

<sup>&</sup>lt;sup>1</sup>Peters, J. , Bühlmann, Peter, & Meinshausen, N. . (2015). Causal inference using invariant prediction: identification and confidence intervals. Stats, 78(5), 947-1012.

#### Invariant Causal Prediction

• Parents of Y satisfies the invariance assumption:

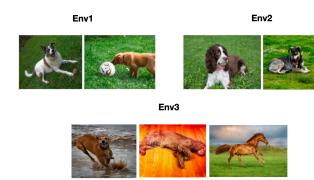
Proposition 1 Consider a linear structural equation model, as formally defined in Section [4.1], for the variables  $(X_1 = Y, X_2, ..., X_p, X_{p+1})$ , with coefficients  $(\beta_{jk})_{j,k=1,...,p+1}$ , whose structure is given by a directed acyclic graph. The independence assumption on the noise variables in Section [4.1] can here be replaced by the strictly weaker assumption that  $\varepsilon_1^e \perp \{\varepsilon_j^e; j \in AN(1)\}$  for all environments  $e \in \mathcal{E}$ , where  $extbf{AN}(1)$  are the ancestors of  $extbf{Y}$ . Then Assumption [1] holds for the parents of  $extbf{Y}$ , namely  $extbf{S}^* = extbf{PA}(1)$ , and  $extbf{Y}^* = extbf{S}_1$ , as defined in Section [4.1] under the following assumption:

for each  $e \in \mathcal{E}$ : the experimental setting e arises by one or several interventions on variables from  $\{X_2, \ldots, X_{p+1}\}$  but interventions on Y are not allowed; here, we allow for do-interventions [Pearl, [2009]] (see also Section [4.2.1]] and note that the assigned values can be random, too), or soft-interventions [Eberhardt] and Scheines, [2007] (see also Sections [4.2.2] and [4.2.3]).

Practically, S\* is not unique.

## An Example for Demonstration

Invariant Causal Prediction 0000000



- When environment set  $\mathcal{E}$  contains Env1 and Env2: grass is invariant.
- When environment set  $\mathcal{E}$  contains Env1 and Env3: grass is variant.

# Plausible causal predictors and coefficients

$$H_{0,\gamma,S}(\mathcal{E}): \quad \gamma_k = 0 \text{ if } k \not \in S \quad \text{and} \quad \left\{ \begin{array}{l} \exists F_\varepsilon \text{ such that for all } e \in \mathcal{E} \\ Y^e = X^e \gamma + \varepsilon^e, \text{ where } \varepsilon^e \perp \!\!\! \perp X^e_S \text{ and } \varepsilon^e \sim F_\varepsilon. \end{array} \right.$$

- Plausible causal predictors:
  - We call the variables S ⊆ {1,...,p} plausible causal predictors under E if the following null hypothesis holds true:

$$H_{0,S}(\mathcal{E}): \exists \gamma \in \mathbb{R}^p \text{ such that } H_{0,\gamma,S}(\mathcal{E}) \text{ is true.}$$
 (5)

(ii) The identifiable causal predictors under interventions ε are defined as the following subset of plausible causal predictors

$$S(\mathcal{E}) \ := \bigcap_{S: H_{0,S}(\mathcal{E}) \ is \ true} S \ = \bigcap_{\gamma \in \Gamma(\mathcal{E})} \{k: \gamma_k \neq 0\}. \tag{6}$$

- *S*(*E*) ⊆ *S*\*
- ullet larger  ${\mathcal E}$ , the more  $S({\mathcal E})$
- Plausible causal coefficients

**Definition 2 (Plausible causal coefficients)** We define the set  $\Gamma_S(\mathcal{E})$  of plausible causal coefficients for the set  $S \subseteq \{1, ..., p\}$  and the global set  $\Gamma(\mathcal{E})$  of plausible causal coefficients under  $\mathcal{E}$  as

$$\Gamma_S(\mathcal{E}) := \{ \gamma \in \mathbb{R}^p : H_{0,\gamma,S}(\mathcal{E}) \text{ is true} \},$$
(7)

$$\Gamma(\mathcal{E}) := \bigcup_{S \subseteq \{1, \dots, p\}} \Gamma_S(\mathcal{E}).$$
 (8)

6 / 21

## Estimation of identifiable causal predictors

#### Generic method for invariant prediction

- 1) For each set  $S \subseteq \{1, ..., p\}$ , test whether  $H_{0,S}(\mathcal{E})$  holds at level  $\alpha$  (we will discuss later concrete examples).
- 2) Set  $\hat{S}(\mathcal{E})$  as

Invariant Causal Prediction

$$\hat{S}(\mathcal{E}) := \bigcap_{S:H_{0,S}(\mathcal{E}) \text{ not rejected}} S.$$
 (12)

For the confidence sets, define

$$\hat{\Gamma}(\mathcal{E}) := \bigcup_{S \subseteq \{1, \dots, p\}} \hat{\Gamma}_S(\mathcal{E}), \tag{13}$$

where

$$\hat{\Gamma}_S(\mathcal{E}) := \begin{cases} \emptyset & H_{0,S}(\mathcal{E}) \text{ can be rejected at level } \alpha \\ \hat{C}(S) & \text{otherwise.} \end{cases}$$
 (14)

Here,  $\hat{C}(S)$  is a  $(1 - \alpha)$ -confidence set for the regression vector  $\beta^{\text{pred}}(S)$  that is obtained by pooling the data.

## Identifiability results

**Theorem 2** Consider a (linear) Gaussian SEM as in (19) and (20) with interventions. Then, with  $S(\mathcal{E})$  as in (6), all causal predictors are identifiable, that is

$$S(\mathcal{E}) = \mathbf{PA}(Y) = \mathbf{PA}(1) \tag{22}$$

if one of the following three assumptions is satisfied:

- i) The interventions are do-interventions (Section [4.2.1]) with a<sup>e</sup><sub>j</sub> ≠ E(X<sup>1</sup><sub>j</sub>) and there is at least one single intervention on each variable other than Y, that is for each j ∈ {2,..., p + 1} there is an experiment e with A<sup>e</sup> = {j}.
- ii) The interventions are noise interventions (Section [4.2.2]) with 1 ≠ E(A<sub>j</sub><sup>e</sup>)<sup>2</sup> < ∞, and again, there is at least one single intervention on each variable other than Y. If the interventions act additively rather than multiplicatively, we require EC<sub>j</sub><sup>e</sup> ≠ 0 or 0 < Var C<sub>j</sub><sup>e</sup> < ∞.</p>
- iii) The interventions are simultaneous noise interventions (Section  $\boxed{4.2.3}$ ). This result still holds if we allow changing linear coefficients  $\beta_{j,k}^{e=2} \neq \beta_{j,k}^{e=1}$  in  $\boxed{21}$  with (possibly random) coefficients  $\beta_{j,k}^{e=2}$ .

The statements remain correct if we replace the null hypothesis (10) with its weaker version (16).

8 / 21

- 1 Invariant Causal Prediction
- 2 Distributionally Robust Learning
- 3 Relationships
- 4 Conclusions

## Robust Learning

Robust learning takes the form:

$$\theta = \arg\min_{\theta \in \Theta} \sup_{P \in \mathcal{P}} \mathbb{E}_{X,Y \sim P}[\ell(\theta; X, Y)] \tag{1}$$

where  $\mathcal{P}$  denotes the uncertainty set.

Adversarial Robustness:

$$\min_{\theta} \mathbb{E}_{X,Y \sim P_{tr}} \left[ \sup_{\|\Delta\| \le \rho} \ell(\theta; X + \Delta, Y) \right] \tag{2}$$

Distributional Robustness:

$$\min_{\theta} \sup_{P:D(P,P_{tr}) \le \rho} \mathbb{E}_{X,Y \sim P}[\ell(\theta; X, Y)] \tag{3}$$

- Remarks:
  - Adversarial robustness is a special case of distributional robustness.
  - The choice of the uncertainty set is important.

Recently, there are mainly two kinds of DRO based on the chosen distance metric:

f-divergence DRO<sup>2</sup>:

$$D_f(P||P_{tr}) = \int f(\frac{dP}{dP_{tr}})dP_{tr}$$
 (4)

Wasserstein DRO<sup>3</sup>:

$$W_c(P; P_{tr}) = \inf_{M \in \Pi(P; P_{tr})} \mathbb{E}_{Z_1, Z_2 \sim M}[c(Z_1, Z_2)]$$
 (5)

 $<sup>^2</sup>$ Duchi, J. C. . (2018). Learning models with uniform performance via distributionally robust optimization.

<sup>&</sup>lt;sup>3</sup>Sinha, A. , Namkoong, H. , Volpi, R. , & Duchi, J. . (2017). Certifying some distributional robustness with principled adversarial training.

### f-divergence DRO

## Paper: Learning Models with Uniform Performance via Distributionally Robust Optimization<sup>4</sup>

Objective function:

$$\widehat{\theta}_n \in \operatorname*{argmin}_{\theta \in \Theta} \left\{ \mathcal{R}_f(\theta; \widehat{P}_n) := \sup_{Q \ll \widehat{P}_n} \left\{ \mathbb{E}_Q[\ell(\theta; X)] : D_f(Q \| \widehat{P}_n) \leq \rho \right\} \right\}.$$

Optimization:

**Proposition 1.** Let P be an arbitrary probability measure on (X, A). Then, for any  $\rho > 0$ , we have for all  $\theta \in \Theta$ 

$$\sup_{Q \ll P} \left\{ \mathbb{E}_{Q}[\ell(\theta; x)] : D_{f}(Q \| P) \le \rho \right\} = \inf_{\lambda \ge 0, \eta \in \mathbb{R}} \left\{ \mathbb{E}_{P} \left[ \lambda f^{*} \left( \frac{\ell(\theta; X) - \eta}{\lambda} \right) \right] + \lambda \rho + \eta \right\}. \tag{5}$$

Moreover, if the supremum on the left hand side is finite, there are finite  $\lambda(\theta) \geq 0$  and  $\eta(\theta) \in \mathbb{R}$ attaining the infimum on the right hand side.

Simplified dual formulation for the Cressie-Read family:

**Lemma 1.** Let P be an arbitrary probability measure on  $(\mathcal{X}, \mathcal{A})$ . Then, for  $k \in (1, \infty)$  and  $k_* =$  $k/(k-1) \in (1,\infty)$ , and any  $\rho > 0$ , we have for all  $\theta \in \Theta$ 

$$\mathcal{R}_{k}(\theta; P) = \inf_{\eta \in \mathbb{R}} \left\{ c_{k}(\rho) \mathbb{E}_{P} \left[ \left( \ell(\theta; X) - \eta \right)_{+}^{k_{*}} \right]^{\frac{1}{k_{*}}} + \eta \right\}. \tag{8}$$

where  $c_k(\rho) := (k(k-1)\rho + 1)^{\frac{1}{k}}$ .

12 / 21

<sup>&</sup>lt;sup>4</sup>Duchi, J. C. . (2018). Learning models with uniform performance via distributionally robust optimization.

# Paper: Certifying Some Distributional Robustness with Principled Adversarial Training<sup>5</sup>

- Why Wasserstein distance?
  - More flexible: do not require the same support.
  - More difficult to optimize.
- Objective function:

$$\min_{\theta} \sup_{P:W_c(P,P_{tr}) \le \rho} \mathbb{E}_{X,Y \sim P}[\ell(\theta; X, Y)]$$
 (6)

Optimization:

Proposition 1. Let  $\ell: \Theta \times Z \to \mathbb{R}$  and  $c: Z \times Z \to \mathbb{R}_+$  be continuous. Let  $\phi_{\gamma}(\theta; z_0) =$  $\sup_{z\in\mathcal{Z}} \{\ell(\theta;z) - \gamma c(z,z_0)\}\$  be the robust surrogate (2b). For any distribution Q and any  $\rho > 0$ ,

$$\sup_{P:W_c(P,Q) \le \rho} \mathbb{E}_P[\ell(\theta;Z)] = \inf_{\gamma \ge 0} \left\{ \gamma \rho + \mathbb{E}_Q[\phi_{\gamma}(\theta;Z)] \right\}, \tag{5}$$

and for any  $\gamma \geq 0$ , we have

$$\sup_{P} \{ \mathbb{E}_{P}[\ell(\theta; Z)] - \gamma W_{c}(P, Q) \} = \mathbb{E}_{Q}[\phi_{\gamma}(\theta; Z)]. \tag{6}$$

<sup>&</sup>lt;sup>5</sup>Sinha, A. , Namkoong, H. , Volpi, R. , & Duchi, J. . (2017). Certifying some distributional robustness with principled adversarial training.

#### Other works in DRO

- Group DRO: consider the group-level DRO<sup>6</sup>
- Marginal DRO: consider the DRO on the marginal distributions of covariates<sup>7</sup>
- f-divergence DRO with varaince regularizer<sup>8</sup>
- WDRO for logistic regression <sup>9</sup>
- WDRO for linear regression <sup>10</sup>

<sup>&</sup>lt;sup>6</sup>Shiori Sagawa et al. DISTRIBUTIONALLY ROBUST NEURAL NETWORKS FOR GROUP SHIFTS: ON THE IMPORTANCE OF REGULARIZATION FOR WORST-CASE **GENERALIZATION** 

<sup>&</sup>lt;sup>7</sup>John C. Duchi et al. Distributionally Robust Losses Against Mixture Covariate Shifts

<sup>&</sup>lt;sup>8</sup>Hongseok Namkoong et al. Variance-based Regularization with Convex Objectives

<sup>&</sup>lt;sup>9</sup>Soroosh Shafieezadeh-Abadeh et al. Distributionally Robust Logistic Regression <sup>10</sup>Ruidi Chen et al. A Robust Learning Approach for Regression Models Based on

Distributionally Robust Optimization

- 3 Relationships
- 4 Conclusions

0000

Causality can be viewed as robust learning with certain uncertainty set<sup>11</sup>:

$$\theta^{causal} = \arg\min_{\theta} \sup_{Q \in \mathcal{Q}^{do}} \mathbb{E}_{X,Y \sim Q}[\ell(\theta; X, Y)]$$
 (7)

Understanding:

- ullet ightarrow :  $heta^{causal}$  obviously minimizes the worst case loss
- ← : consider the converse-negative proposition:

Not 
$$\theta^{causal} \to \text{Not } \arg\min_{\theta} \sup_{Q \in \mathcal{Q}^{do}}$$
 (8)

◆: another perspective:

$$\sup_{Q \in \mathcal{Q}^{do}} \mathbb{E}_{X,Y \sim Q}[\ell(\theta; X, Y)] = \begin{cases} \infty & \theta \neq \theta^{causal} \\ \operatorname{Var}(\epsilon_{y}) & \theta = \theta^{causal} \end{cases}$$
(9)

<sup>&</sup>lt;sup>11</sup>Meinshausen, N. (2018). Causality from a distributional robustness point of view. 6-10.

## A Toy Illustration of $\leftarrow$

Consider a simple data generation mechanism:

$$S \to Y \dashrightarrow V \tag{10}$$

or in equation:

$$S \leftarrow \mathcal{N}(0, \epsilon_S) \tag{11}$$

$$Y \leftarrow S + \epsilon_Y \tag{12}$$

$$V \leftarrow \alpha_e Y + \epsilon_V \tag{13}$$

where  $\alpha_e$  is changing across environments/distributions, which means V should not be used to predict Y. Then  $\mathcal{Q}^{do}$  contains:

$$P_{S=a_1}^{do}, \dots, P_{S=a_m}^{do}, P_{V=b_1}^{do}, \dots, P_{V=b_k}^{do}$$
 (14)

Actually, from  $P_{V=b}^{do}$ , we can know that V is not the direct cause of Y.

## Difference between Causality, Robustness and Distributional Robustness

#### The goals are different:

- Causality: to estimate causal coefficients, to identify causal covaraites
- Traditional Robustness: the predictive accuracy for a reference distribution is optimal:

$$\arg\min_{\theta} \sup_{Q \in \mathcal{Q}} \mathbb{E}_{X,Y \sim P^{obs}}[\ell(\theta; X, Y)] \tag{15}$$

Distributional Robustness: robust over a set of distributions

$$\arg\min_{\theta} \sup_{Q \in \mathcal{Q}} \mathbb{E}_{X,Y \sim Q}[\ell(\theta; X, Y)] \tag{16}$$

Difference between Causality and DRO:

- Causality: seeks for the true causal parameters
- DRO: only seeks for robustness over distributions but does not care the true parameters

- Invariant Causal Prediction
- ② Distributionally Robust Learning
- Relationships
- 4 Conclusions

#### Conclusion

- Distributionally robust learning is closely related to causality.
- How to approach causality by the way of DRO?
- How to build the uncertainty set with the assistance of causality?



Conclusions



清华大学 刘家硕 THU Trustworthy-AI Group
Causality with Robust Machine Learning 21 / 21