

Causality, Invariance, and Distribution Generalization¹

Sorawit Saengkyongam

ETH Zürich, Seminar for Statistics

¹Some of the slides are provided by Jonas Peters.

Table of Contents

1 Invariance

2 From Invariance to Causality

3 From Causality to Distribution Generalization

Table of Contents

1 Invariance

2 From Invariance to Causality

3 From Causality to Distribution Generalization

Invariance

Setting: Multiple environments $\mathcal{E} = \{e_1, \dots, e_k\}$

$$(X^{e_1}, Y^{e_1}) \sim \mathbb{P}^{e_1}$$

$$(X^{e_2}, Y^{e_2}) \sim \mathbb{P}^{e_2}$$

$$(X^{e_3}, Y^{e_3}) \sim \mathbb{P}^{e_3}$$

Invariance

Setting: Multiple environments $\mathcal{E} = \{e_1, \dots, e_k\}$

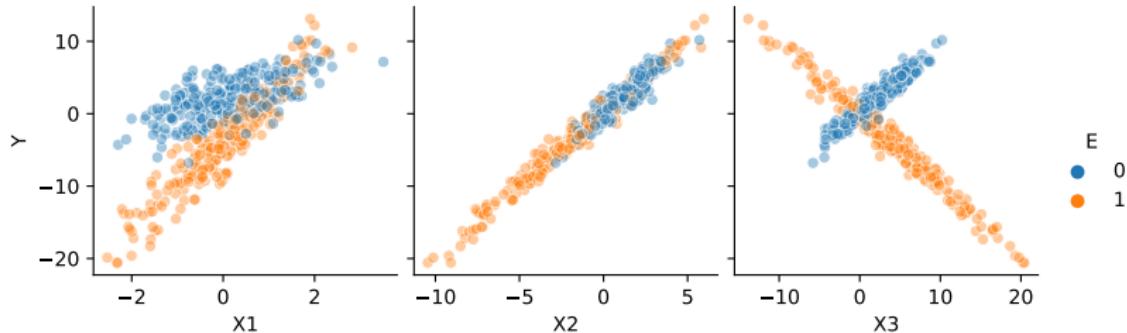
$$(X^{e_1}, Y^{e_1}) \sim \mathbb{P}^{e_1}$$

$$(X^{e_2}, Y^{e_2}) \sim \mathbb{P}^{e_2}$$

$$(X^{e_3}, Y^{e_3}) \sim \mathbb{P}^{e_3}$$

Invariance: $\exists S \subseteq \{1, \dots, d\} : \mathbb{P}(Y^e | X_S^e)$ are identical for all $e \in \mathcal{E}$

Invariance: Example



Here: $\mathbb{P}(Y^0 | X_2^0) = \mathbb{P}(Y^1 | X_2^1)$. We call X^2 is invariant.

Table of Contents

1 Invariance

2 From Invariance to Causality

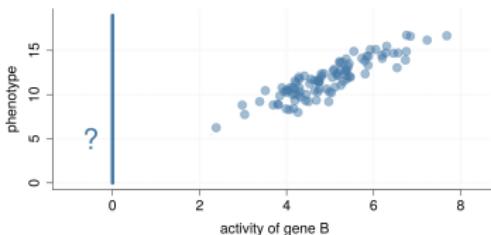
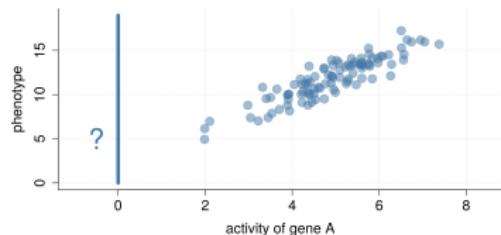
3 From Causality to Distribution Generalization

From Invariance to Causality

Invariance \implies Causality

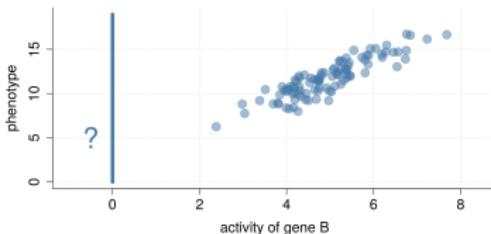
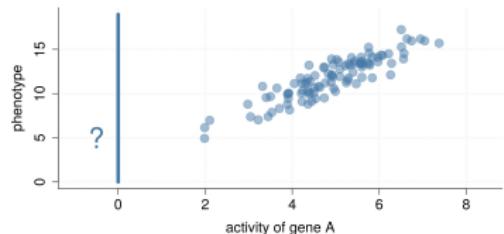
Why does causality matter?

What is the best prediction for the phenotype if we delete gene A (or B)?

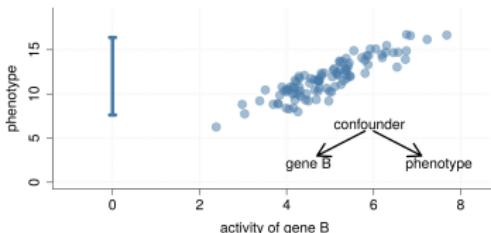
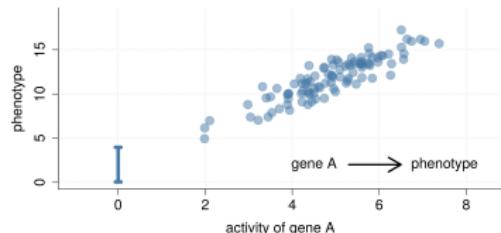


Why does causality matter?

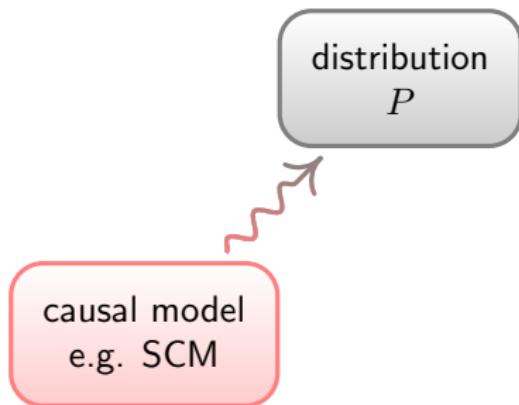
What is the best prediction for the phenotype if we delete gene A (or B)?



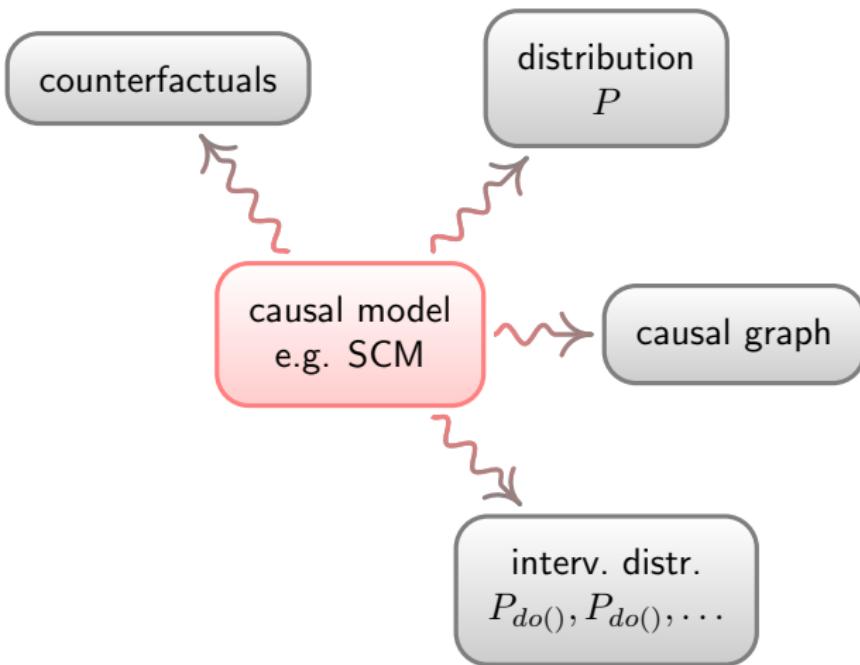
Causality matters: the optimal prediction should depend on the underlying causal structure



What is a causal model?



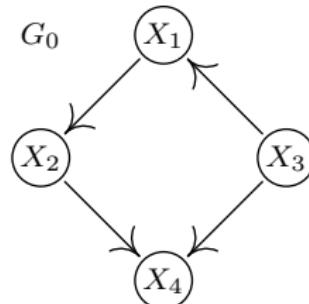
What is a causal model?



SCMs model **observational distributions** over X_1, \dots, X_d . Call it: P .

$$\begin{aligned}X_1 &:= X_3 + N_1 \\X_2 &:= 2X_1 + N_2 \\X_3 &:= N_3 \\X_4 &:= -X_2 - X_3 + N_4\end{aligned}$$

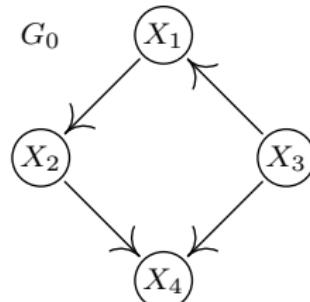
- N_i jointly independent $\mathcal{N}(0, 1)$
- G_0 has no cycles



SCMs model **observational distributions** over X_1, \dots, X_d . Call it: P .

$$\begin{aligned}X_1 &:= X_3 + N_1 \\X_2 &:= 2X_1 + N_2 \\X_3 &:= N_3 \\X_4 &:= -X_2 - X_3 + N_4\end{aligned}$$

- N_i jointly independent $\mathcal{N}(0, 1)$
- G_0 has no cycles

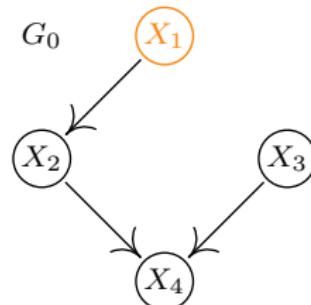


$$\begin{pmatrix} X_1 \\ X_2 \\ X_3 \\ X_4 \end{pmatrix} = \mathcal{N} \left(\begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 2 & 4 & 1 & -5 \\ 4 & 9 & 2 & -11 \\ 1 & 2 & 1 & -3 \\ -5 & -11 & -3 & 15 \end{pmatrix} \right)$$

SCMs model **interventions**, too. Call it: $P_{do(X_1:=0)}$.

$$\begin{aligned}X_1 &:= 0 \\X_2 &:= f_2(X_1, N_2) \\X_3 &:= f_3(N_3) \\X_4 &:= f_4(X_2, X_3, N_4)\end{aligned}$$

- N_i jointly independent
- G_0 has no cycles



Modularity Property

If you intervene only on X_j , you intervene only on X_j .

From Invariance to Causality

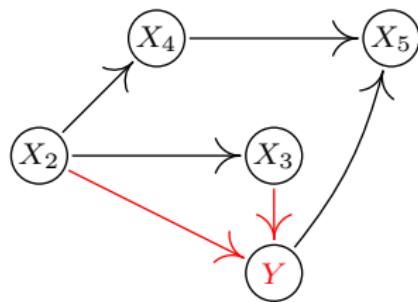
Invariance \implies Causality

The Problem of Causal Discovery:

Y	X_2	X_3	X_4	X_5
3.4	-0.3	5.8	-2.1	2.2
1.7	-0.2	7.0	-1.2	0.4
-2.4	-0.1	4.3	-0.7	3.5
2.3	-0.3	5.5	-1.1	-4.4
3.5	-0.2	3.9	-0.9	-3.9
:	:	:	:	:

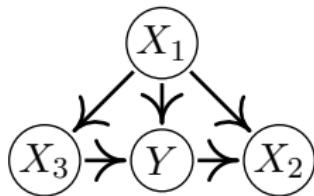
?

causal model

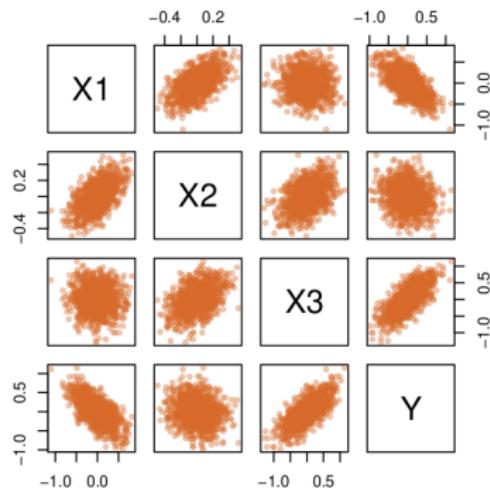
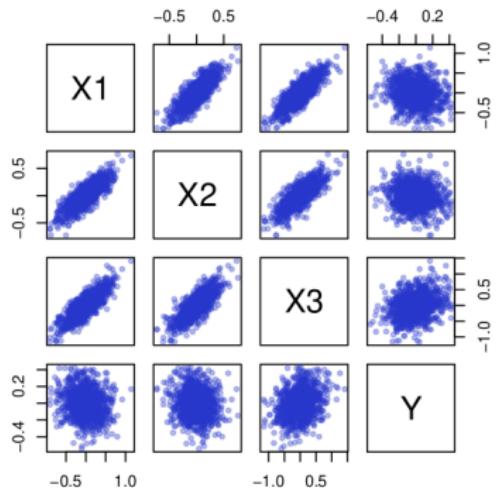


Here: Find the direct causes of Y !

unknown:



known:



linear model

```
> linmod <- lm( Y ~ X)
> summary(linmod)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	8.305e-05	2.067e-03	0.04	0.968
X1	-5.490e-01	9.725e-03	-56.46	<2e-16 ***
X2	-4.078e-01	1.810e-02	-22.52	<2e-16 ***
X3	6.821e-01	6.896e-03	98.91	<2e-16 ***

ICP (R-package InvariantCausalPrediction)

```
> ExpInd
```

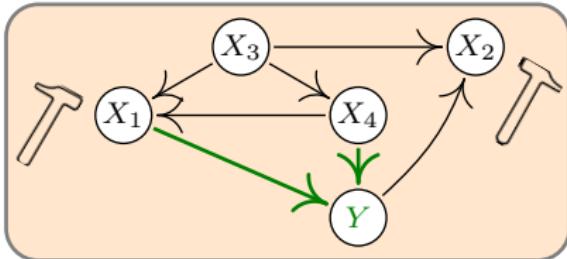
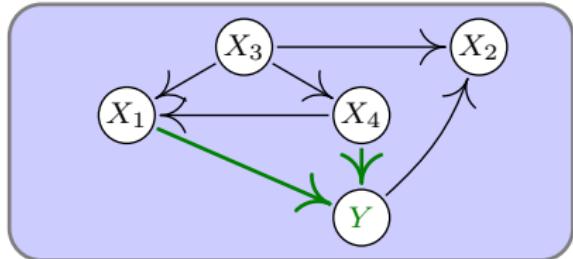
```
[1]1111111111111111111111111111111111111111111111111111111111111111...2222222222222222...
```

```
> icp <- ICP(X,Y,ExpInd)
```

	LOWER BOUND	UPPER BOUND	MAXIMIN EFFECT	P-VALUE	
X1	-0.71	-0.52	-0.52	<1e-09	***
X2	-0.46	0.00	0.00	0.55	
X3	0.58	0.70	0.58	<1e-09	***

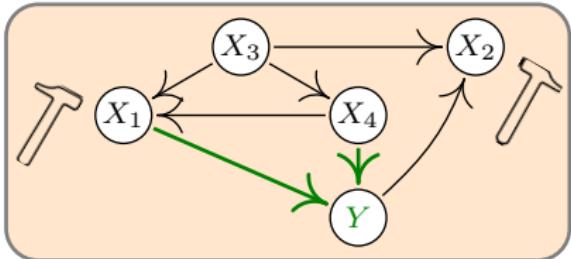
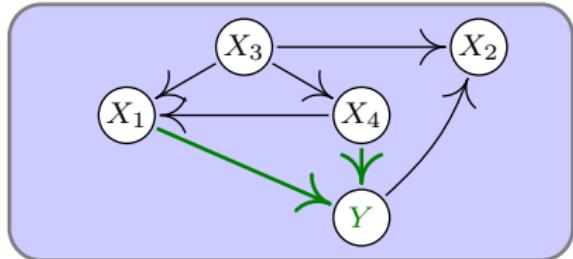
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1					

Fundamental assumption: $X_1, X_4 \rightarrow Y$ is invariant under interventions.



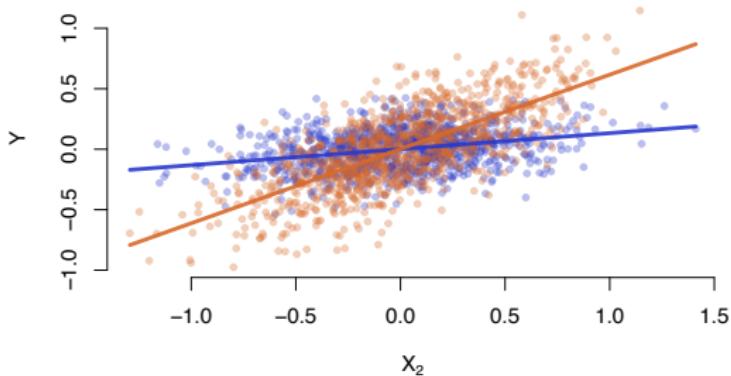
cf. modularity, autonomy, Haavelmo 1944, Aldrich 1989, Pearl 2009, ...

Fundamental assumption: $X_1, X_4 \rightarrow Y$ is invariant under interventions.

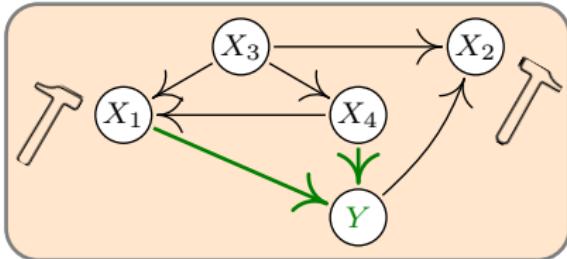
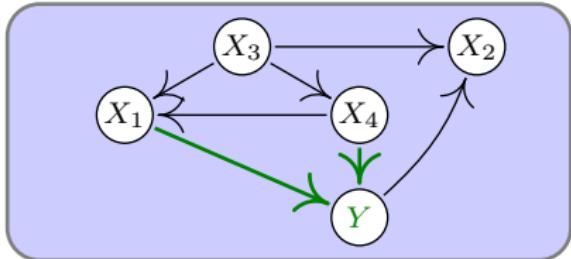


cf. modularity, autonomy, Haavelmo 1944, Aldrich 1989, Pearl 2009, ...

Not all sets of predictors yield an invariant model. Here: $\{2\}$.



Fundamental assumption: $X_1, X_4 \rightarrow Y$ is invariant under interventions.



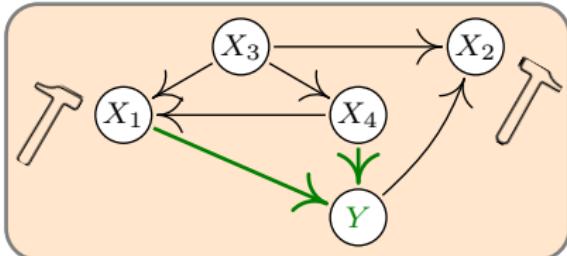
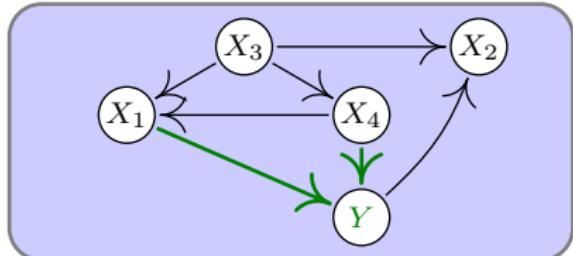
cf. modularity, autonomy, Haavelmo 1944, Aldrich 1989, Pearl 2009, ...

Key idea: Use and data and search for invariant models.

set	\emptyset	$\{1\}$	$\{2\}$	$\{3\}$	\dots	$\{1, 4\}$	$\{2, 4\}$	\dots	$\{1, 3, 4\}$
invariance	\times	\times	\times	\times	\dots	\checkmark	\times	\dots	\checkmark

$$\hat{S} := \bigcap_{S \text{ invariant}} S = \{1, 4\}$$

Fundamental assumption: $X_1, X_4 \rightarrow Y$ is invariant under interventions.



cf. modularity, autonomy, Haavelmo 1944, Aldrich 1989, Pearl 2009, ...

Key idea: Use and data and search for invariant models.

set	\emptyset	$\{1\}$	$\{2\}$	$\{3\}$	\dots	$\{1, 4\}$	$\{2, 4\}$	\dots	$\{1, 3, 4\}$
invariance	\times	\times	\times	\times	\dots	\checkmark	\times	\dots	\checkmark

$$\hat{S} := \bigcap_{S \text{ invariant}} S = \{1, 4\}$$

JP, Bühlmann, Meinshausen, JRSS-B 2016 (with discussion): $P(\hat{S} \subseteq S^*) \geq 1 - \alpha$.

Invariant Causal Prediction

Setting: Multiple environments $\mathcal{E} = \{e_1, \dots, e_k\}$

Given: For each $e \in \mathcal{E}$, an i.i.d. sample of (X^e, Y^e) .

Invariant Causal Prediction

Setting: Multiple environments $\mathcal{E} = \{e_1, \dots, e_k\}$

Given: For each $e \in \mathcal{E}$, an i.i.d. sample of (X^e, Y^e) .

Invariance $H_{0,S}$:

- The conditional distributions $\mathbb{P}(Y^e | X_S^e)$ are identical for all $e \in \mathcal{E}$.
- (Linear setting) $Y^e = X_S^e \gamma + \epsilon^e$, where $\epsilon^e \perp\!\!\!\perp X^e$ and $\epsilon^e \sim F$
- (Non-linear setting) $Y \perp\!\!\!\perp E | X_S$

Invariant Causal Prediction

Setting: Multiple environments $\mathcal{E} = \{e_1, \dots, e_k\}$

Given: For each $e \in \mathcal{E}$, an i.i.d. sample of (X^e, Y^e) .

Invariance $H_{0,S}$:

- The conditional distributions $\mathbb{P}(Y^e | X_S^e)$ are identical for all $e \in \mathcal{E}$.
- (Linear setting) $Y^e = X_S^e \gamma + \epsilon^e$, where $\epsilon^e \perp\!\!\!\perp X^e$ and $\epsilon^e \sim F$
- (Non-linear setting) $Y \perp\!\!\!\perp E | X_S$

Environments: Different interventions (not on Y). Then, H_{0,PA_Y} holds.

Invariant Causal Prediction

Setting: Multiple environments $\mathcal{E} = \{e_1, \dots, e_k\}$

Given: For each $e \in \mathcal{E}$, an i.i.d. sample of (X^e, Y^e) .

Invariance $H_{0,S}$:

- The conditional distributions $\mathbb{P}(Y^e | X_S^e)$ are identical for all $e \in \mathcal{E}$.
- (Linear setting) $Y^e = X_S^e \gamma + \epsilon^e$, where $\epsilon^e \perp\!\!\!\perp X^e$ and $\epsilon^e \sim F$
- (Non-linear setting) $Y \perp\!\!\!\perp E | X_S$

Environments: Different interventions (not on Y). Then, H_{0,PA_Y} holds.

Theorem (Peters et al. 2016)

Assume H_{0,PA_Y} holds. For any test level α we obtain

$$P(\hat{S} \subseteq \text{PA}_Y) \geq 1 - \alpha$$

The guarantee holds without identifiability assumptions (if PA_Y is not identifiable, then ICP returns an empty set).

Table of Contents

1 Invariance

2 From Invariance to Causality

3 From Causality to Distribution Generalization

From Causality to Distribution Generalization

Causality \implies Invariance \implies Generalization

response variable	Y
covariates	$X := X_1, \dots, X_d$
training model	M
training data	i.i.d. from $P_M^{(X,Y)}$
wanted	

response variable	Y
covariates	$X := X_1, \dots, X_d$
training model	M
training data	i.i.d. from $P_M^{(X,Y)}$
wanted	prediction model $Y \approx \hat{f}(X)$ that performs well on test data (which may be different from training data)

domain generalization, out-of-distribution prediction, covariate shift, ...

response variable	Y
covariates	$X := X_1, \dots, X_d$
training model	M
training data	i.i.d. from $P_M^{(X,Y)}$
wanted	prediction model $Y \approx \hat{f}(X)$ that performs well on test data (which may be different from training data)

domain generalization, out-of-distribution prediction, covariate shift, ...

$$\text{minimax solution: } \arg\min_{f_\diamond \in \mathcal{F}} \sup_{\tilde{M} \in \mathcal{N}(M)} E_{\tilde{M}} [(Y - f_\diamond(X))^2].$$

response variable	Y
covariates	$X := X_1, \dots, X_d$
training model	M
training data	i.i.d. from $P_M^{(X,Y)}$
wanted	prediction model $Y \approx \hat{f}(X)$ that performs well on test data (which may be different from training data)

domain generalization, out-of-distribution prediction, covariate shift, ...

$$\text{minimax solution: } \underset{f_\diamond \in \mathcal{F}}{\operatorname{argmin}} \sup_{\tilde{M} \in \mathcal{N}(M)} E_{\tilde{M}} [(Y - f_\diamond(X))^2].$$

A causal model

response variable	Y
covariates	$X := X_1, \dots, X_d$
training model	M
training data	i.i.d. from $P_M^{(X,Y)}$
wanted	prediction model $Y \approx \hat{f}(X)$ that performs well on test data (which may be different from training data)

domain generalization, out-of-distribution prediction, covariate shift, ...

$$\text{minimax solution: } \arg\min_{f_\diamond \in \mathcal{F}} \sup_{\tilde{M} \in \mathcal{N}(M)} E_{\tilde{M}} [(Y - f_\diamond(X))^2].$$

A **causal model** satisfies

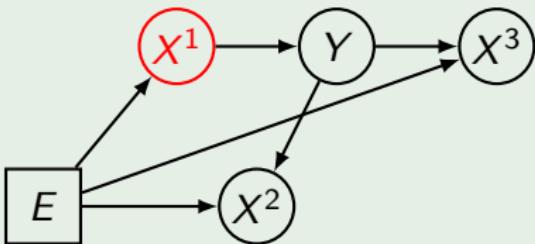
$$f_{causal} = \arg\min_{f_\diamond \in \mathcal{F}} \sup_{i \in \mathcal{I}} E_{M(i)} [(Y - f_\diamond(X))^2],$$

where $M(i)$: an intervention model and \mathcal{I} : all interventions on X .

Distribution Generalization

Example

$$\mathcal{S}(e) : \begin{cases} X^1 := \mu_e + \epsilon_{X_1} \\ Y := \beta_1 X^1 + \epsilon_Y \\ X^2 := \gamma_e Y + \epsilon_{X_2} \\ X^3 := \eta_e Y + \epsilon_{X_3} \end{cases}$$

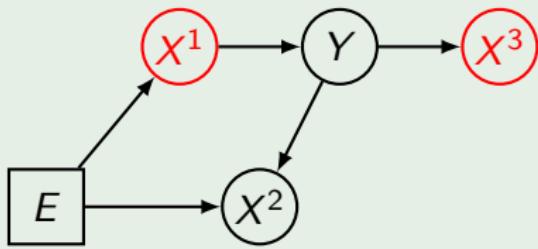


In some environments, X^2 and X^3 may be positively correlated with Y , while in others they can be negatively correlated. This poses a threat to generalization.

Distribution Generalization

Example (What if...)

$$\mathcal{S}(\epsilon) : \begin{cases} X^1 := \mu_e + \epsilon_{X_1} \\ Y := \beta_1 X^1 + \epsilon_Y \\ X^2 := \gamma_e Y + \epsilon_{X_2} \\ X^3 := \eta Y + \epsilon_{X_3} \end{cases}$$



Distribution Generalization²

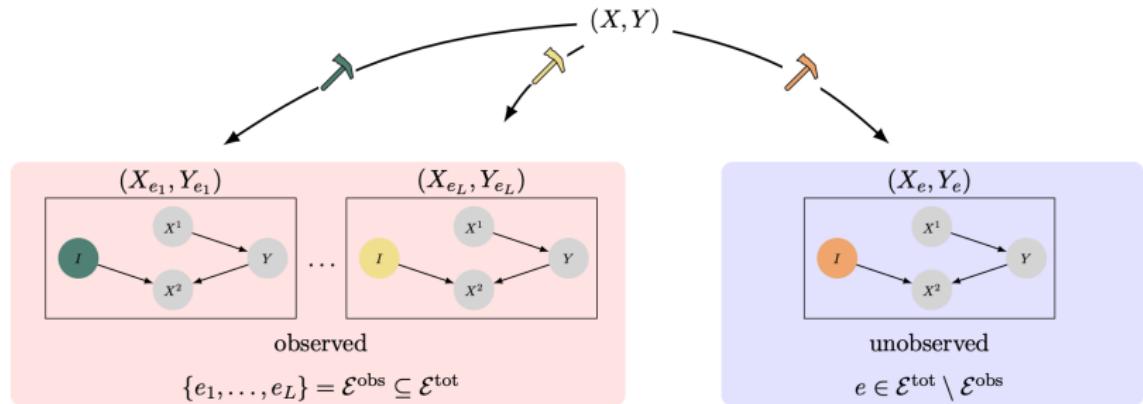


Figure 3.: Illustration of multi-environment data generation setting. Only some environments are observed, but one would like to be able to make predictions on any further potentially unobserved environment.

²Figure 3 is taken from Pfister et al. 2021

Distribution Generalization

Given: Training data: (Y^{tr}, X^{tr}, E^{tr}) , Test data: (X^{test}, E^{test})

Distribution Generalization

Given: Training data: (Y^{tr}, X^{tr}, E^{tr}) , Test data: (X^{test}, E^{test})

Invariance: $H_{0,S} : Y^{tr} \perp\!\!\!\perp E^{tr} \mid X_S^{tr}$

Invariant and most predictive: Let \mathcal{S}^{inv} be a set of all invariant set. For each $S \in \mathcal{S}^{\text{inv}}$,

- ① Find an estimate $\hat{\mathbb{E}}[Y^{tr} \mid X_S^{tr}]$
- ② Compute the MSE $\mathbb{E}[(Y^{tr} - \hat{\mathbb{E}}[Y^{tr} \mid X_S^{tr}])^2]$

Return the subset S^* (and its corresponding estimate $\hat{\mathbb{E}}[Y^{tr} \mid X_{S^*}^{tr}]$) that yields the lowest MSE.

Distribution Generalization

Given: Training data: (Y^{tr}, X^{tr}, E^{tr}) , Test data: (X^{test}, E^{test})

Invariance: $H_{0,S} : Y^{tr} \perp\!\!\!\perp E^{tr} \mid X_S^{tr}$

Invariant and most predictive: Let \mathcal{S}^{inv} be a set of all invariant set. For each $S \in \mathcal{S}^{\text{inv}}$,

- ① Find an estimate $\hat{\mathbb{E}}[Y^{tr} \mid X_S^{tr}]$
- ② Compute the MSE $\mathbb{E}[(Y^{tr} - \hat{\mathbb{E}}[Y^{tr} \mid X_S^{tr}])^2]$

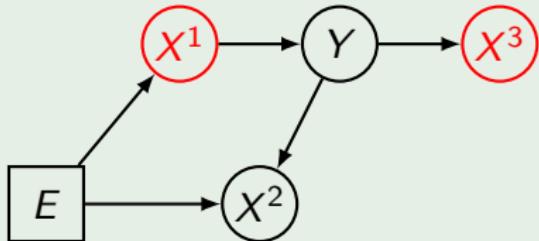
Return the subset S^* (and its corresponding estimate $\hat{\mathbb{E}}[Y^{tr} \mid X_{S^*}^{tr}]$) that yields the lowest MSE.

Environments: If $Y^{test} \perp\!\!\!\perp E^{test} \mid X_{S^*}^{test}$ holds then $\hat{\mathbb{E}}[Y^{tr} \mid X_{S^*}^{tr}]$ is expected to perform well on the test data.

Distribution Generalization

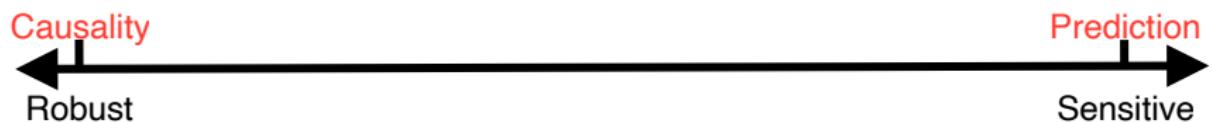
Example

$$\mathcal{S}(\epsilon) : \begin{cases} X^1 := \mu_e + \epsilon_{X_1} \\ Y := \beta_1 X^1 + \epsilon_Y \\ X^2 := \gamma_e Y + \epsilon_{X_2} \\ X^3 := \eta Y + \epsilon_{X_3} \end{cases}$$

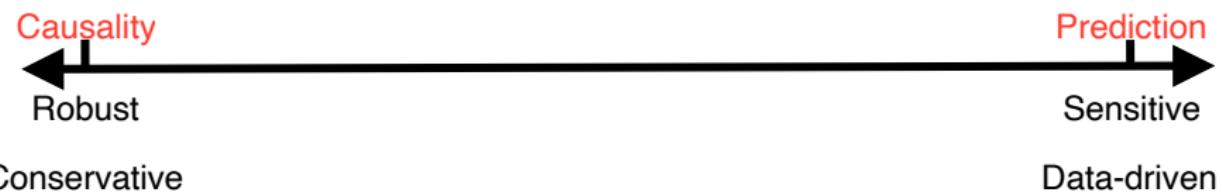


We have that $\{X^1\}$ and $\{X^1, X^3\}$ are the invariant sets. Although the causal parent is $\{X^1\}$, we may want to use both X^1, X^3 for predicting Y since the conditional expectation $\mathbb{E}[Y | X^1, X^3]$ is stable across environments and more predictive.

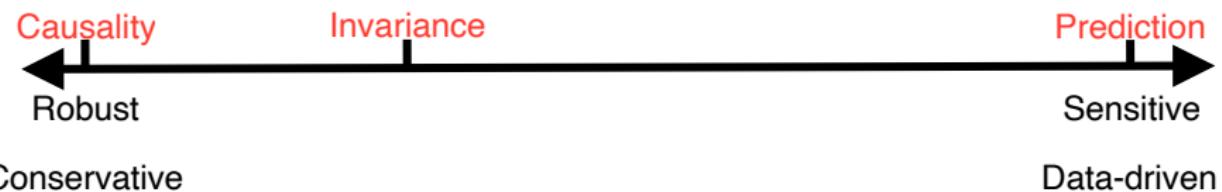
Conclusion and Open Questions



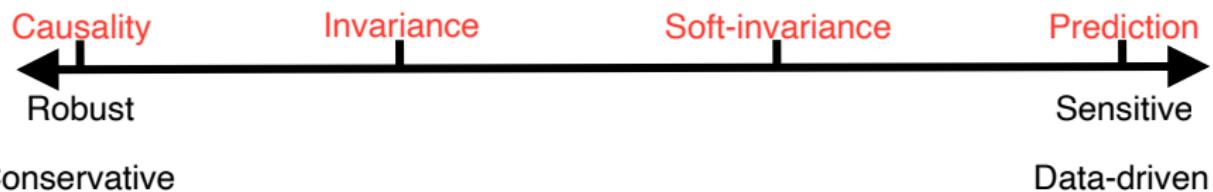
Conclusion and Open Questions



Conclusion and Open Questions



Conclusion and Open Questions



$$b_{AR}^\gamma := \underset{b}{\operatorname{argmin}} \underbrace{\mathbf{E}_M(Y - Xb)^2}_{\text{prediction}} + \gamma \underbrace{\|\mathbf{E}_M A^t(Y - Xb)\|_2^2}_{\text{invariance}}$$

References I

-  Bühlmann, Peter (2020). "Invariance, causality and robustness". In.
-  Heinze-Deml, Christina, Jonas Peters, and Nicolai Meinshausen (2018). "Invariant Causal Prediction for Nonlinear Models". In: *Journal of Causal Inference* 6.2.
-  Pearl, J. (2009). *Causality: Models, Reasoning, and Inference*. 2nd. New York, USA: Cambridge University Press.
-  Peters, Jonas, Peter Bühlmann, and Nicolai Meinshausen (2016). "Causal inference by using invariant prediction: identification and confidence intervals". In: *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, pp. 947–1012.
-  Pfister, N., E. G. William, J. Peters, R. Aebersold, and P. Bühlmann (2021). "Stabilizing Variable Selection and Regression". In: *Annals of Applied Statistics* 15.3, pp. 1220–1246.

References II

-  Rothenhäusler, Dominik, Nicolai Meinshausen, Peter Bühlmann, and Jonas Peters (2021). "Anchor regression: Heterogeneous data meet causality". In: *Journal of the Royal Statistical Society Series B: Statistical Methodology* 83.2, pp. 215–246.
-  Saengkyongam, Sorawit, Nikolaj Thams, Jonas Peters, and Niklas Pfister (2023). "Invariant policy learning: A causal perspective". In: *IEEE Transactions on Pattern Analysis and Machine Intelligence*.