

# Neural Causal AI

## Adversarial Invariance Learning from Multiple Environments

**Jianqing Fan**

Princeton University

with **Yihong Gu**, **Cong Fang**, and **Peter Buehlmann**



# Outlines

- 1 Introduction
- 2 Endogeneity in High Dimension
- 3 Multi-Environment Linear Regression
- 4 Neural Causal Learning
- 5 Causality under SCM
- 6 Implementation and Numerical Studies



Neural Causal AI

# Outlines

- 1 Introduction
- 2 Endogeneity in High Dimension
- 3 Multi-Environment Linear Regression
- 4 Neural Causal Learning
- 5 Causality under SCM
- 6 Implementation and Numerical Studies



Neural Causal AI



Yihong Gu



Cong Fang



Peter Buehlmann

# Outlines

- 1 Introduction
- 2 Endogeneity in High Dimension
- 3 Multi-Environment Linear Regression
- 4 Neural Causal Learning
- 5 Causality under SCM
- 6 Implementation and Numerical Studies



Neural Causal AI



Yihong Gu



Cong Fang



Peter Buehlmann



# Introduction

# What Is Causality?

Wikipedia: one event, process, or state contributes to **production of another**:

- ★ relations hold in the past must hold in the future
- ★ relations hold in one environment must hold in another.

## Invariance

Phil. of Sci.: Phenomenon that no evidence against is regarded a truth.

Causality  $\approx$  Invariance under MEs

# What Is Causality?

Wikipedia: one event, process, or state contributes to **production of another**:

- ★ relations hold in the past must hold in the future
- ★ relations hold in one environment must hold in another.

## Invariance

Phil. of Sci.: Phenomenon that no evidence against is regarded a truth.

Causality  $\approx$  Invariance under MEs

# What Is Causality?

Wikipedia: one event, process, or state contributes to **production of another**:

- ★ relations hold in the past must hold in the future
- ★ relations hold in one environment must hold in another.

## Invariance

Phil. of Sci.: Phenomenon that no evidence against is regarded a truth.

Causality  $\approx$  Invariance under MEs

# What Is Causality?

Wikipedia: one event, process, or state contributes to **production of another**:

- ★ relations hold in the past must hold in the future
- ★ relations hold in one environment must hold in another.

## Invariance

Phil. of Sci.: Phenomenon that no evidence against is regarded a truth.

**Causality  $\approx$  Invariance under MEs**

★ Prediction

★ Attribution

★ Inferences

★ Causality

## Typical Processes:

- ★ Collect response variable  $Y$  and its associated variables  $X \in \mathbb{R}^p$ .
- ★ Use statistical machine algorithms to select important variables.

What can be wrong?

★ Prediction

★ Attribution

★ Inferences

★ Causality

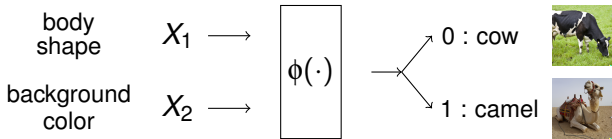
## Typical Processes:

- ★ Collect response variable  $Y$  and its associated variables  $X \in \mathbb{R}^p$ .
- ★ Use statistical machine algorithms to select important variables.

## What can be wrong?

# An illustrative example

## ★ Classification uses two features:



★ Standard SML : Data  $\mathcal{D}$ : ●70% cows on grass ( $X_2$  green),

●80% camels on sand ( $X_2$  yellow)

Get  $\mathcal{D}_{train} + \mathcal{D}_{test}$   $\implies$   $\hat{\phi}(\cdot)$  works well on  $\mathcal{D}_{test}$ ,  
but relies on  $X_2$  (spurious)

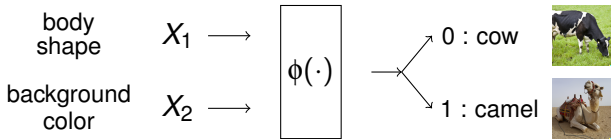
★ Prediction: Not robust in other environments (marketing).

★ Attribution: Wrong mechanism or treatment targets!



# An illustrative example

## ★ Classification uses two features:



★ Standard SML : Data  $\mathcal{D}$ : ●70% cows on grass ( $X_2$  green),

●80% camels on sand ( $X_2$  yellow)

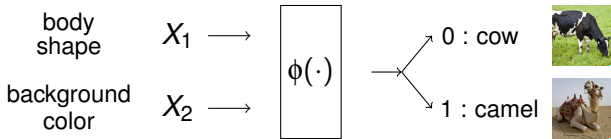
Get  $\mathcal{D}_{train} + \mathcal{D}_{test}$   $\implies$   $\hat{\phi}(\cdot)$  works well on  $\mathcal{D}_{test}$ ,  
but relies on  $X_2$  (spurious)

★ Prediction: Not robust in other environments (marketing).

★ Attribution: Wrong mechanism or treatment targets!

# An illustrative example

## ★ Classification uses two features:



- ★ Standard SML : Data  $\mathcal{D}$ : • 70% cows on grass ( $X_2$  green),  
• 80% camels on sand ( $X_2$  yellow)

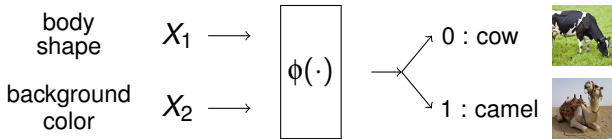
Get  $\mathcal{D}_{train} + \mathcal{D}_{test}$   $\implies$   $\hat{\phi}(\cdot)$  works well on  $\mathcal{D}_{test}$ ,  
but relies on  $X_2$  (spurious)

## What is wrong?

- ★ Prediction: Not robust in other environments (marketing).
- ★ Attribution: Wrong mechanism or treatment targets!

# An illustrative example

## ★ Classification uses two features:



★ Standard SML : Data  $\mathcal{D}$ : ● 70% cows on grass ( $X_2$  green),

● 80% camels on sand ( $X_2$  yellow)

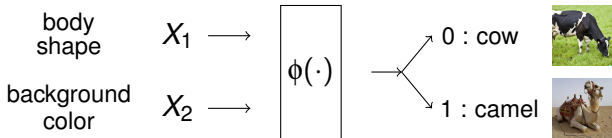
Get  $\mathcal{D}_{train} + \mathcal{D}_{test}$   $\implies$   $\hat{\phi}(\cdot)$  works well on  $\mathcal{D}_{test}$ ,  
but relies on  $X_2$  (spurious)

★ Prediction: Not robust in other environments (marketing).

★ Attribution: Wrong mechanism or treatment targets!

# An illustrative example

## ★ Classification uses two features:



★ Standard SML : Data  $\mathcal{D}$ : ● 70% cows on grass ( $X_2$  green),

● 80% camels on sand ( $X_2$  yellow)

Get  $\mathcal{D}_{train} + \mathcal{D}_{test}$   $\implies$   $\hat{\phi}(\cdot)$  works well on  $\mathcal{D}_{test}$ ,  
but relies on  $X_2$  (spurious)

★ Prediction: Not robust in other environments (marketing).

★ Attribution: Wrong mechanism or treatment targets!

# Can Machine Learn Causality?

Eliminate endogeneity?

Train a Causal AI + What Data?

**Can Machine Learn Causality?**

**Eliminate endogeneity?**

**Train a Causal AI + What Data?**

**Can Machine Learn Causality?**

**Eliminate endogeneity?**

**Train a Causal AI + What Data?**

## Use data heterogeneity

$\mathcal{D}$ : 70% cows on grass  
80% camels on sand

$\tilde{\mathcal{D}}$ : 50% cows on grass  
60% camels on sand

assoc. of  $X_2$  and  $Y$  varies in  $\mathcal{D}$  and  $\tilde{\mathcal{D}} \implies X_2$  is spurious variable

$X_2$  endogenous spurious  $\implies$  inconsistency

Today's Talk: Variable Selection (Causality Learning) from Invariance

$X_1$  ✓,  $X_2$  ✗,  $(X_1, X_2)$  ✗,  $X_3 = \text{temperature}$  ✓



## Use data heterogeneity

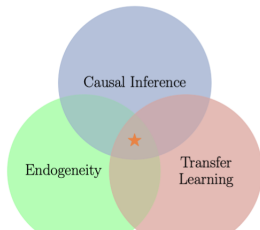
$\mathcal{D}$ : 70% cows on grass  
80% camels on sand

$\tilde{\mathcal{D}}$ : 50% cows on grass  
60% camels on sand

assoc. of  $X_2$  and  $Y$  varies in  $\mathcal{D}$  and  $\tilde{\mathcal{D}} \implies X_2$  is spurious variable

$X_2$  **endogenous spurious**  $\implies$  **inconsistency**

Today's Talk: Variable Selection (Causality Learning) from **Invariance**



## Use data heterogeneity

$\mathcal{D}$ : 70% cows on grass  
80% camels on sand

$\tilde{\mathcal{D}}$ : 50% cows on grass  
60% camels on sand

assoc. of  $X_2$  and  $Y$  varies in  $\mathcal{D}$  and  $\tilde{\mathcal{D}} \implies X_2$  is spurious variable

$X_2$  **endogenous spurious**  $\implies$  **inconsistency**

Today's Talk: Variable Selection (Causality Learning) from **Invariance**

$X_1$  ✓,  $X_2$  ✗,  $(X_1, X_2)$  ✗,  $X_3 = \text{temperature}$  ✓

## Use data heterogeneity

$\mathcal{D}$ : 70% cows on grass  
80% camels on sand

$\tilde{\mathcal{D}}$ : 50% cows on grass  
60% camels on sand

assoc. of  $X_2$  and  $Y$  varies in  $\mathcal{D}$  and  $\tilde{\mathcal{D}} \implies X_2$  is spurious variable

$X_2$  endogenous spurious  $\implies$  inconsistency

Today's Talk: Variable Selection (Causality Learning) from **Invariance**

$X_1$  ✓,  $X_2$  ✗,  $(X_1, X_2)$  ✗,  $X_3 = \text{temperature}$  ✓

## Use data heterogeneity

$\mathcal{D}$ : 70% cows on grass  
80% camels on sand

$\tilde{\mathcal{D}}$ : 50% cows on grass  
60% camels on sand

assoc. of  $X_2$  and  $Y$  varies in  $\mathcal{D}$  and  $\tilde{\mathcal{D}} \implies X_2$  is spurious variable

$X_2$  **endogenous spurious**  $\implies$  **inconsistency**

Today's Talk: Variable Selection (Causality Learning) from **Invariance**

$X_1$  ✓,  $X_2$  ✗,  $(X_1, X_2)$  ✗,  $X_3$  =temperature ✓

## Use data heterogeneity

$\mathcal{D}$ : 70% cows on grass  
80% camels on sand

$\tilde{\mathcal{D}}$ : 50% cows on grass  
60% camels on sand

assoc. of  $X_2$  and  $Y$  varies in  $\mathcal{D}$  and  $\tilde{\mathcal{D}} \implies X_2$  is spurious variable

$X_2$  endogenous spurious  $\implies$  inconsistency

Today's Talk: Variable Selection (Causality Learning) from **Invariance**

$X_1$  ✓,  $X_2$  ✗,  $(X_1, X_2)$  ✗,  $X_3 = \text{temperature}$  ✓

↑ exogeneous spurious

# Spurious variables

{ **Endogeneous:**

{ **Exogeneous:**

# Spurious variables

- Endogeneous:** background colors
- Exogeneous:** time photo taken

# Spurious variables

{	<b>Endogeneous:</b>	background colors	<b>harmful-bias</b>
	<b>Exogeneous:</b>	time photo taken	<b>unbiased-var.</b>

$\ell_1$ , SCAD, SIS



# Spurious variables

{	<b>Endogeneous:</b>	background colors	<b>harmful-bias</b>
	<b>Exogeneous:</b>	time photo taken	<b>unbiased-var.</b>

$\ell_1$ , SCAD, SIS

## Eliminate endogeneity by FAIR-NN

{	<b>Endogeneous:</b>	background colors	<b>harmful-bias</b>
	<b>Exogeneous:</b>	time photo taken	<b>unbiased-var.</b>

$\ell_1$ , SCAD, SIS

## Eliminate endogeneity by FAIR-NN

### What only one environment?

# Endogeneity in High Dimension

- ★ Fan, J. and Liao, Y. (2014). Endogeneity in ultrahigh dimension. *Ann. Statist.*, **42**, 872-917.
- ★ Fan, J., Han, F., and Liu, H. (2014). Challenges of Big Data analysis. *Natl. Sci. Rev.*, **1**, 293-314.

# Assumptions in Variable Selection

Stylized Model:  $Y = \mathbf{X}^T \beta_0 + \varepsilon$ ,  $\mathbb{E}\varepsilon\mathbf{X} = 0$  or stronger,

$\beta_0$  sparse.

★ Tons of equations!

★ can not be validated!

# Assumptions in Variable Selection

Stylized Model:  $Y = \mathbf{X}^T \beta_0 + \varepsilon$ ,  $\mathbb{E}\varepsilon\mathbf{X} = 0$  or stronger,

★ Tons of equations!

$\beta_0$  sparse.

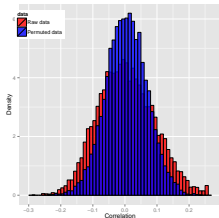
★ can not be validated!

## Prostate cancer study

Data: 148 microarrays from GEO data

Response: Expressions of gene DDR

Covariates: remaining 12,718 genes



# Assumptions in Variable Selection

Stylized Model:  $Y = \mathbf{X}^T \beta_0 + \varepsilon$ ,  $\mathbb{E}\varepsilon\mathbf{X} = 0$  or stronger,

★ Tons of equations!

$\beta_0$  sparse.

★ can not be validated!

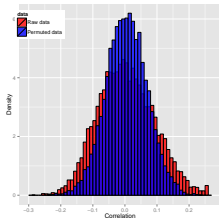
## Prostate cancer study

Data: 148 microarrays from GEO data

Response: Expressions of gene DDR

Covariates: remaining 12,718 genes

Example:  $Y = 2X_1 + X_2 + \varepsilon$ ,



$$\mathbb{E}(\varepsilon|X_1) = 0, \mathbb{E}(\varepsilon|X_2) = 0$$

# Assumptions in Variable Selection

Stylized Model:  $Y = \mathbf{X}^T \beta_0 + \varepsilon$ ,  $\mathbb{E}\varepsilon\mathbf{X} = 0$  or stronger,

★ Tons of equations!

$\beta_0$  sparse.

★ can not be validated!

## Prostate cancer study

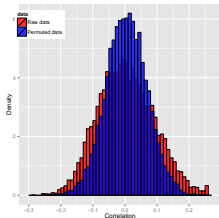
Data: 148 microarrays from GEO data

Response: Expressions of gene DDR

Covariates: remaining 12,718 genes

Example:  $Y = 2X_1 + X_2 + \varepsilon$ ,

Netting: Collecting many variables  $\{X_j\}_{j=1}^p$ .



$\mathbb{E}(\varepsilon|X_1) = 0, \mathbb{E}(\varepsilon|X_2) = 0$



# Assumptions in Variable Selection

**Stylized Model:**  $Y = \mathbf{X}^T \beta_0 + \varepsilon$ ,  $\mathbb{E}\varepsilon\mathbf{X} = 0$  or stronger,

★ Tons of equations!

$\beta_0$  sparse.

★ can not be validated!

## Prostate cancer study

**Data:** 148 microarrays from GEO data

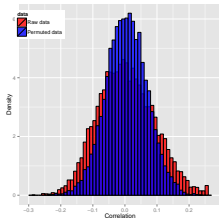
**Response:** Expressions of gene DDR

**Covariates:** remaining 12,718 genes

**Example:**  $Y = 2X_1 + X_2 + \varepsilon$ ,

**Netting:** Collecting many variables  $\{X_j\}_{j=1}^p$ .

■ Many  $X_j$ 's related to  $Y$ ,



$\mathbb{E}(\varepsilon|X_1) = 0, \mathbb{E}(\varepsilon|X_2) = 0$





# Assumptions in Variable Selection

Stylized Model:  $Y = \mathbf{X}^T \beta_0 + \varepsilon$ ,  $\mathbb{E}\varepsilon\mathbf{X} = 0$  or stronger,

★ Tons of equations!

$\beta_0$  sparse.

★ can not be validated!

## Prostate cancer study

Data: 148 microarrays from GEO data

Response: Expressions of gene DDR

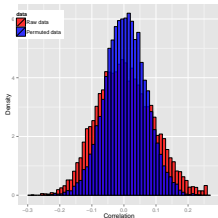
Covariates: remaining 12,718 genes

Example:  $Y = 2X_1 + X_2 + \varepsilon$ ,

Netting: Collecting many variables  $\{X_j\}_{j=1}^p$ .

■ Many  $X_j$ 's related to  $Y$ , hence to  $\varepsilon = Y - 2X_1 - X_2$  for large  $p$ :

$\text{corr}(X_j, \varepsilon) \neq 0$ , for some  $j$ . **Endogeneity**



$\mathbb{E}(\varepsilon|X_1) = 0, \mathbb{E}(\varepsilon|X_2) = 0$



# Solutions

Model:  $Y = \mathbf{X}_{S_0}^T \beta_{S_0} + \varepsilon$  with  $\mathbb{E}(\varepsilon | \mathbf{X}_{S_0}) = 0$  or **weaker**. **more realistic**

Example:  $\mathbb{E}\varepsilon \mathbf{X}_{S_0} = 0$ ,

# Solutions

Model:  $Y = \mathbf{X}_{S_0}^T \beta_{S_0} + \varepsilon$  with  $\mathbb{E}(\varepsilon | \mathbf{X}_{S_0}) = 0$  or **weaker**. **more realistic**

Example:  $\mathbb{E} \varepsilon \mathbf{X}_{S_0} = 0$ ,

# Solutions

Model:  $Y = \mathbf{X}_{S_0}^T \beta_{S_0} + \varepsilon$  with  $\mathbb{E}(\varepsilon | \mathbf{X}_{S_0}) = 0$  or **weaker**. **more realistic**

Example:  $\mathbb{E}\varepsilon \mathbf{X}_{S_0} = 0$ ,  $\mathbb{E}\varepsilon \mathbf{X}_{S_0}^2 = 0$ ,

★ Variables  $\mathbf{X}_{S_0}$  are special or causal, as more equations than unknowns.

# Solutions

Model:  $Y = \mathbf{X}_{S_0}^T \beta_{S_0} + \varepsilon$  with  $\mathbb{E}(\varepsilon | \mathbf{X}_{S_0}) = 0$  or **weaker**. **more realistic**

Example:  $\mathbb{E}\varepsilon \mathbf{X}_{S_0} = 0$ ,  $\mathbb{E}\varepsilon \mathbf{X}_{S_0}^2 = 0$ ,  $\mathbb{E}\varepsilon \mathbf{X}_{S_0}^3 = 0$

★ Variables  $\mathbf{X}_{S_0}$  are special or causal, as more equations than unknowns.

**Invariant**

# Solutions

Model:  $Y = \mathbf{X}_{S_0}^T \beta_{S_0} + \varepsilon$  with  $\mathbb{E}(\varepsilon | \mathbf{X}_{S_0}) = 0$  or **weaker**. **more realistic**

Example:  $\mathbb{E}\varepsilon \mathbf{X}_{S_0} = 0$ ,  $\mathbb{E}\varepsilon \mathbf{X}_{S_0}^2 = 0$ ,  $\mathbb{E}\varepsilon \mathbf{X}_{S_0}^3 = 0$

★ Variables  $\mathbf{X}_{S_0}$  are special or causal, as more equations than unknowns. **Invariant**

Generalization:  $\mathbb{E}(Y - \mathbf{X}_{S_0}^T \beta_{S_0}) f(\mathbf{X}_{S_0}) = 0$  for  $f \in \mathcal{F}$  **GM constraints**

# Solutions

Model:  $Y = \mathbf{X}_{S_0}^T \beta_{S_0} + \varepsilon$  with  $\mathbb{E}(\varepsilon | \mathbf{X}_{S_0}) = 0$  or **weaker**. **more realistic**

Example:  $\mathbb{E}\varepsilon \mathbf{X}_{S_0} = 0$ ,  $\mathbb{E}\varepsilon \mathbf{X}_{S_0}^2 = 0$ ,  $\mathbb{E}\varepsilon \mathbf{X}_{S_0}^3 = 0$

★ Variables  $\mathbf{X}_{S_0}$  are special or causal, as more equations than unknowns. **Invariant**

Generalization:  $\mathbb{E}(Y - \mathbf{X}_{S_0}^T \beta_{S_0}) f(\mathbf{X}_{S_0}) = 0$  for  $f \in \mathcal{F}$  **GM constraints**

## Constrained LS

# Solutions

Model:  $Y = \mathbf{X}_{S_0}^T \beta_{S_0} + \varepsilon$  with  $\mathbb{E}(\varepsilon | \mathbf{X}_{S_0}) = 0$  or **weaker**. **more realistic**

Example:  $\mathbb{E}\varepsilon \mathbf{X}_{S_0} = 0$ ,  $\mathbb{E}\varepsilon \mathbf{X}_{S_0}^2 = 0$ ,

★ Variables  $\mathbf{X}_{S_0}$  are special or causal, as more equations than unknowns. **Invariant**

Generalization:  $\mathbb{E}(Y - \mathbf{X}_{S_0}^T \beta_{S_0}) f(\mathbf{X}_{S_0}) = 0$  for  $f \in \mathcal{F}$  **GM constraints**

## Soft Constrained LS

$$\min_{\beta} \sum_{i=1}^n \varepsilon_i^2 + \lambda (\|\sum_{i=1}^n \varepsilon_i \mathbf{X}_{i,S_0}\|^2 + \|\sum_{i=1}^n \varepsilon_i \mathbf{X}_{i,S_0}^2\|^2), \quad \varepsilon_i = Y_i - \mathbf{X}_{i,S_0}^T \beta_{S_0}.$$



# Solutions

Model:  $Y = \mathbf{X}_{S_0}^T \beta_{S_0} + \varepsilon$  with  $\mathbb{E}(\varepsilon | \mathbf{X}_{S_0}) = 0$  or **weaker**. **more realistic**

Example:  $\mathbb{E}\varepsilon \mathbf{X}_{S_0} = 0$ ,  $\mathbb{E}\varepsilon \mathbf{X}_{S_0}^2 = 0$ ,  $\mathbb{E}\varepsilon \mathbf{X}_{S_0}^3 = 0$

★ Variables  $\mathbf{X}_{S_0}$  are special or causal, as more equations than unknowns. **Invariant**

Generalization:  $\mathbb{E}(Y - \mathbf{X}_{S_0}^T \beta_{S_0}) f(\mathbf{X}_{S_0}) = 0$  for  $f \in \mathcal{F}$  **GM constraints**

## Soft Constrained LS

$$\min_{\beta} \sum_{i=1}^n \varepsilon_i^2 + \lambda \left( \max_{f \in \mathcal{F}} \sum_{i=1}^n \varepsilon_i f(\mathbf{X}_{i,S_0}) \right), \quad \varepsilon_i = Y_i - \mathbf{x}_{i,S_0}^T \beta_{S_0}$$

# Multi-Environment Linear Reg

- ★ Fan, J., Fang, C., Gu, Y., and Zhang, T. (2024+). Environment Invariant Linear Least Squares. *Ann. Statist.*

★ Multi-environment regression: For each  $e$ ,  $(X_i^{(e)}, Y_i^{(e)})_{i=1}^n \sim i.i.d. \mu^{(e)} \in \mathcal{U}_{\beta^*}$ :

$$Y^{(e)} = (\beta_{S^*}^*)^\top X_{S^*}^{(e)} + \varepsilon^{(e)} \quad \text{with} \quad \mathbb{E}[\varepsilon^{(e)} X_{S^*}^{(e)}] = 0.$$

◆  $S^*, \beta^*$  are **invariant**.

◆ More realistic and weaker than  $\mathbb{E}[\varepsilon^{(e)} X^{(e)}] = 0$  for regression.

★ Heterogeneous: Each environment does not provide a consistent estimator.

★ Multi-environment regression: For each  $e$ ,  $(X_i^{(e)}, Y_i^{(e)})_{i=1}^n \sim i.i.d. \mu^{(e)} \in \mathcal{U}_{\beta^*}$ :

$$Y^{(e)} = (\beta_{S^*}^*)^\top X_{S^*}^{(e)} + \varepsilon^{(e)} \quad \text{with} \quad \mathbb{E}[\varepsilon^{(e)} X_{S^*}^{(e)}] = 0.$$

◆  $S^*, \beta^*$  are **invariant**. ← learning object

◆ More realistic and weaker than  $\mathbb{E}[\varepsilon^{(e)} X^{(e)}] = 0$  for regression.

★ Heterogeneous: Each environment does not provide a consistent estimator.

★ Multi-environment regression: For each  $e$ ,  $(X_i^{(e)}, Y_i^{(e)})_{i=1}^n \sim i.i.d. \mu^{(e)} \in \mathcal{U}_{\beta^*}$ :

$$Y^{(e)} = (\beta_{S^*}^*)^\top X_{S^*}^{(e)} + \varepsilon^{(e)} \quad \text{with} \quad \mathbb{E}[\varepsilon^{(e)} X_{S^*}^{(e)}] = 0.$$

◆  $S^*, \beta^*$  are **invariant**. ← learning object

◆ More realistic and weaker than  $\mathbb{E}[\varepsilon^{(e)} X^{(e)}] = 0$  for regression.

★ Heterogeneous: Each environment does not provide a consistent estimator.

# Focused linear invariance regularizer

## Population-level penalty:

★ delete endogenous variables

$$J(\beta) = \sum_{j \in \mathcal{S}(\beta)} \sum_{e=1}^K \left| \mathbb{E} \left[ \underbrace{\left( Y^{(e)} - \beta_{\mathcal{S}(\beta)}^\top X_{\mathcal{S}(\beta)}^{(e)} \right)}_{\varepsilon^{(e)}} X_j^{(e)} \right] \right|^2$$

# Focused linear invariance regularizer

## Population-level penalty:

★ delete endogenous variables

$$J(\beta) = \sum_{j \in S(\beta)} \sum_{e=1}^K \left| \mathbb{E} \left[ \underbrace{\left( Y^{(e)} - \beta_{S(\beta)}^\top X_{S(\beta)}^{(e)} \right)}_{\varepsilon^{(e)}} X_j^{(e)} \right] \right|^2$$

- ★ If  $S$  is selected, minimizing  $J(\beta)$  encourages  $X_j^{(e)}$  and  $\varepsilon^{(e)}$  uncorrelated across for all  $j \in S$  and all environments.

# A multi-environment version of linear least squares

## ★ Population-level EILLS:

environment-invariant linear least-squares

$$Q(\beta; \gamma) = R(\beta) + \gamma J(\beta)$$

★ EILLS estimator  $\hat{\beta}_Q = \operatorname{argmin}_{\beta} \hat{Q}(\beta; \gamma)$ . ( $\mathbb{E} \rightsquigarrow \hat{\mathbb{E}}$ )

★ Regularized EILLS estimator:  $\hat{\beta}_L = \operatorname{argmin}_{\beta} \hat{Q}(\beta; \gamma) + \lambda \|\beta\|_0$ .




# A multi-environment version of linear least squares

## ★ Population-level EILLS:

environment-invariant linear least-squares

$$Q(\beta; \gamma) = R(\beta) + \gamma J(\beta) = \sum_{e=1}^K \mathbb{E}[|Y^{(e)} - \beta^\top X^{(e)}|^2] \\ + \gamma \sum_{j=1}^p 1_{\{\beta_j \neq 0\}} \times \sum_{e=1}^K |\mathbb{E}[(Y^{(e)} - \beta^\top X^{(e)})X_j^{(e)}]|^2$$

 delete endogenous variables

★ EILLS estimator  $\hat{\beta}_Q = \operatorname{argmin}_{\beta} \hat{Q}(\beta; \gamma)$ . ( $\mathbb{E} \rightsquigarrow \hat{\mathbb{E}}$ )


★ Regularized EILLS estimator:  $\hat{\beta}_L = \operatorname{argmin}_{\beta} \hat{Q}(\beta; \gamma) + \lambda \|\beta\|_0$ .

# A multi-environment version of linear least squares

## ★ Population-level EILLS:

environment-invariant linear least-squares

$$Q(\beta; \gamma) = R(\beta) + \gamma J(\beta) = \sum_{e=1}^K \mathbb{E}[|Y^{(e)} - \beta^\top X^{(e)}|^2] \\ + \gamma \sum_{j=1}^p 1_{\{\beta_j \neq 0\}} \times \sum_{e=1}^K |\mathbb{E}[(Y^{(e)} - \beta^\top X^{(e)})X_j^{(e)}]|^2$$

 delete endogenous variables

## ★ EILLS estimator $\hat{\beta}_Q = \operatorname{argmin}_{\beta} \hat{Q}(\beta; \gamma)$ . ( $\mathbb{E} \rightsquigarrow \hat{\mathbb{E}}$ )


## ★ Regularized EILLS estimator: $\hat{\beta}_L = \operatorname{argmin}_{\beta} \hat{Q}(\beta; \gamma) + \lambda \|\beta\|_0$ .

# A multi-environment version of linear least squares

## ★ Population-level EILLS:


environment-invariant linear least-squares

$$Q(\beta; \gamma) = R(\beta) + \gamma J(\beta) = \sum_{e=1}^K \mathbb{E}[|Y^{(e)} - \beta^\top X^{(e)}|^2] \\ + \gamma \sum_{j=1}^p 1_{\{\beta_j \neq 0\}} \times \sum_{e=1}^K |\mathbb{E}[(Y^{(e)} - \beta^\top X^{(e)})X_j^{(e)}]|^2$$

 delete endogenous variables

## ★ EILLS estimator $\hat{\beta}_Q = \operatorname{argmin}_{\beta} \hat{Q}(\beta; \gamma)$ . ( $\mathbb{E} \rightsquigarrow \hat{\mathbb{E}}$ )

## ★ Regularized EILLS estimator: $\hat{\beta}_L = \operatorname{argmin}_{\beta} \hat{Q}(\beta; \gamma) + \lambda \|\beta\|_0$ .

 delete exog. var.

## Two type of Spurious Variables

**Exogenous**  $\Rightarrow$  variance

## Two type of Spurious Variables

**Exogenous**  $\Rightarrow$  variance

**Endogenous**  $\Rightarrow$  biases + incon.

## Two type of Spurious Variables

**Exogenous**  $\Rightarrow$  variance

↑ reduced by  $\ell_0(\beta)$  or  $\|\beta\|_1$  or SCAD

**Endogenous**  $\Rightarrow$  biases + incon.

## Two type of Spurious Variables

**Exogenous**  $\Rightarrow$  variance

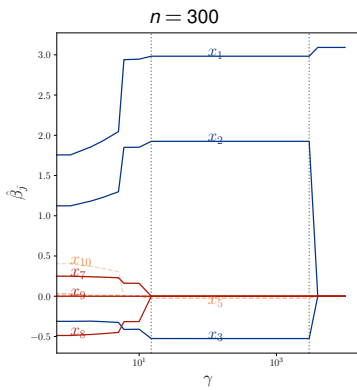
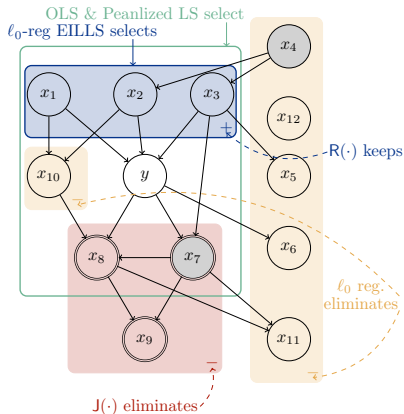
reduced by  $\ell_0(\beta)$  or  $\|\beta\|_1$  or SCAD

**Endogenous**  $\Rightarrow$  biases + incon.

reduced by  $J(\beta)$

# How is $S^*$ selected in SCM?

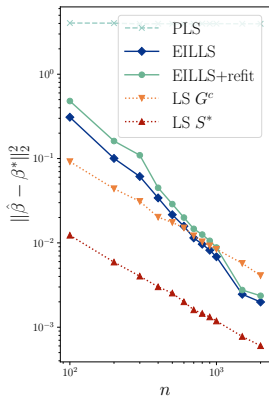
- ★  $p = 12$ ,  $S^* = \{1, 2, 3\}$ ,  $\mathcal{G} = \{7, 8, 9\}$  (double circled).
- ★  $e = 1$  observational env
- ★  $e = 2$  interventional env: intervene on  $x_4, x_7$  (shaded)



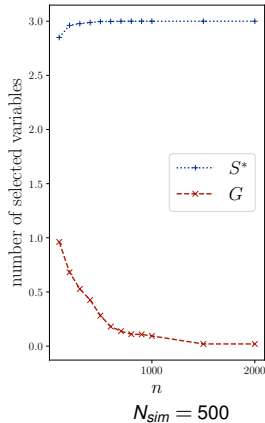


# Simulation Results ( $\gamma = 20$ )

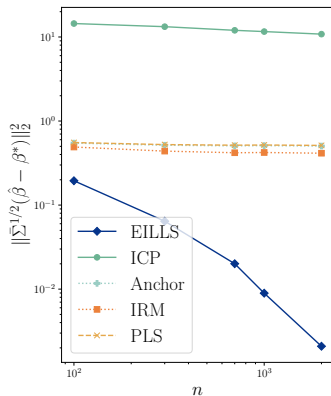
## $\ell_2$ error rate



## Variable selection



## Comparisons



# Non-asymptotic Result for EILLS

**EILLS estimator:**  $\hat{\beta}_Q = \operatorname{argmin}_{\beta} \hat{Q}(\beta; \gamma).$

**Theorem 2.** Under Cond 1-3 & IDF, if  $\gamma \geq C\gamma^*$  and  $p\gamma = o(n)$ , then

- (1) **Sure screening:**  $S^* \subseteq \operatorname{supp}(\hat{\beta}_Q) \subseteq \mathcal{G}^c$  holds w.h.p. for large  $n$ .
- (2)  **$\ell_2$ -error.** With high probability,

$$\|\hat{\beta}_Q - \beta^*\|_2 \leq C\gamma \left\{ \sqrt{\frac{|\mathcal{G}^c|}{n \cdot K}} + \frac{|\mathcal{G}^c|}{n} \right\};$$

**Endogenous spurious:**  $\mathcal{G} = \left\{ j : \sum_{e=1}^K \mathbb{E}[X_j^{(e)} \varepsilon^{(e)}] \neq 0 \right\}.$

# Non-asymptotic Result for EILLS

**EILLS estimator:**  $\hat{\beta}_Q = \operatorname{argmin}_{\beta} \hat{Q}(\beta; \gamma).$

**Theorem 2.** Under Cond 1-3 & IDF, if  $\gamma \geq C\gamma^*$  and  $p\gamma = o(n)$ , then

- (1) **Sure screening:**  $S^* \subseteq \operatorname{supp}(\hat{\beta}_Q) \subseteq \mathcal{G}^c$  holds w.h.p. for large  $n$ .
- (2)  **$l_2$ -error.** With high probability,

$$\|\hat{\beta}_Q - \beta^*\|_2 \leq C\gamma \left\{ \sqrt{\frac{|\mathcal{G}^c|}{n \cdot K}} + \frac{|\mathcal{G}^c|}{n} \right\};$$

## Selection consistency?

# Non-asymptotic Result for EILLS

EILLS estimator:  $\hat{\beta}_Q = \operatorname{argmin}_{\beta} \hat{Q}(\beta; \gamma)$ .

Theorem 2. Under Cond 1-3 & IDF, if  $\gamma \geq C\gamma^*$  and  $p\gamma = o(n)$ , then

- (1) Sure screening:  $S^* \subseteq \operatorname{supp}(\hat{\beta}_Q) \subseteq \mathcal{G}^c$  holds w.h.p. for large  $n$ .
- (2)  $l_2$ -error. With high probability,

$$\|\hat{\beta}_Q - \beta^*\|_2 \leq C\gamma \left\{ \sqrt{\frac{|\mathcal{G}^c|}{n \cdot K}} + \frac{|\mathcal{G}^c|}{n} \right\};$$

**Endogenous Spurious**  $\times$  **by  $J(\beta)$**

# Non-asymptotic Result for EILLS

EILLS estimator:  $\hat{\beta}_Q = \operatorname{argmin}_{\beta} \hat{Q}(\beta; \gamma)$ .

Theorem 2. Under Cond 1-3 & IDF, if  $\gamma \geq C\gamma^*$  and  $p\gamma = o(n)$ , then

- (1) Sure screening:  $S^* \subseteq \operatorname{supp}(\hat{\beta}_Q) \subseteq \mathcal{G}^c$  holds w.h.p. for large  $n$ .
- (2)  $l_2$ -error. With high probability,

$$\|\hat{\beta}_Q - \beta^*\|_2 \leq C\gamma \left\{ \sqrt{\frac{|\mathcal{G}^c|}{n \cdot K}} + \frac{|\mathcal{G}^c|}{n} \right\};$$

Endogenous Spurious  $\times$  by  $J(\beta)$

Exogenous Spurious  $\times$  by  $\ell_0(\beta)$

# Variable Selection Consistency in High-dims

Regularized EILLS:  $\hat{\beta}_L = \operatorname{argmin}_{\beta} \hat{Q}(\beta; \gamma) + \lambda \|\beta\|_0.$

**Theorem 3.** Under Conditions 1-3 & IDF, if  $\gamma \geq C\gamma^*$ , for sufficiently large  $n$  and proper choice of  $\lambda$ , we have

$$\mathbb{P}[\operatorname{supp}(\hat{\beta}_L) = S^*] \geq 1 - p^{-10}.$$

★ When  $|S^*| + \gamma = O(1)$ , choose  $\{K^{-1} + \sqrt{\frac{\log p}{n}}\} \frac{\log p}{n} \ll \lambda \ll \beta_{\min}^2.$

# Neural Causal Learning

- ★ Gu, Y., Fang, C., Buehlmann, P., and Fan, J. (2024). Causality Pursuit from Heterogeneous Environments via Neural Adversarial Invariance Learning. *arxiv.org*

# Nonparametric Causality Pursuit

- ★ Collect  $n$  data from  $K$  **heterogeneous** environment with dist  $\mu^{(e)}$ . For  $e \in [K]$ ,

$$Y^{(e)} = m^*(X_{S^*}^{(e)}) + \varepsilon^{(e)} \quad \text{with} \quad \mathbb{E}[\varepsilon^{(e)} | X_{S^*}^{(e)}] = 0$$

—  $S^*$  **unknown** variable set;

—  $m^*$  **invariant** assoc.

◆ **Much** weaker than standard reg:  $\mathbb{E}[\varepsilon | X] = 0$ .

- ★ Goal: estimate  $S^*$  and  $m^*$  using  $n \cdot K$  data.



# Focused Adversarial Invariance Regularizer (FAIR)

## ★ Endogeneity (FAIR) Penalty:

♣ delete endogenous variables

$$\max_{\mathbf{f} \in \mathbf{S}_g} \left\{ \sum_{e \in [K]} \mathbb{E}_{\mu^{(e)}} \left[ \{\mathbf{Y} - \mathbf{g}(\mathbf{X})\} \mathbf{f}_e(\mathbf{X}) \right] \right\} \text{ with } \mathbb{E}_{\mu^{(e)}} f_e^2(\mathbf{X}) = 1.$$

—  $\mathbf{S}_g$  is the support of function  $g$

★ When  $\text{supp}(g) = S$ , maximizing all  $f_e(X_S)$  gives

$$\frac{1}{2} \sum_{e \in [K]} \mathbb{E}_{\mu^{(e)}} \left[ |\mathbb{E}[Y|X_S] - g(X_S)|^2 \right].$$

# Focused Adversarial Invariance Regularizer (FAIR)

## ★ Endogeneity (FAIR) Penalty:

♣ delete endogenous variables

$$\max_{\mathbf{f} \in \mathbf{S}_g} \left\{ \sum_{\mathbf{e} \in [K]} \mathbb{E}_{\mu^{(\mathbf{e})}} \left[ \{ \mathbf{Y} - \mathbf{g}(\mathbf{X}) \} \mathbf{f}_{\mathbf{e}}(\mathbf{X}) - \lambda f_{\mathbf{e}}^2(\mathbf{X}) \right] \right\}.$$

—  $\mathbf{S}_g$  is the support of function  $g$

★ When  $\text{supp}(g) = S$ , maximizing all  $f_{\mathbf{e}}(X_S)$  gives

$$\frac{1}{2} \sum_{\mathbf{e} \in [K]} \mathbb{E}_{\mu^{(\mathbf{e})}} \left[ |\mathbb{E}[Y|X_S] - g(X_S)|^2 \right].$$

# Focused Adversarial Invariance Regularizer (FAIR)

## ★ Endogeneity (FAIR) Penalty:

♣ delete endogenous variables

$$J(g) = \max_{f \in S_g} \left\{ \sum_{e \in [K]} \mathbb{E}_{\mu^{(e)}} \left[ \{Y - g(X)\} f_e(X) - \frac{1}{2} f_e^2(X) \right] \right\}.$$

—  $S_g$  is the support of function  $g$

★ When  $\text{supp}(g) = S$ , maximizing all  $f_e(X_S)$  gives

$$\frac{1}{2} \sum_{e \in [K]} \mathbb{E}_{\mu^{(e)}} \left[ |\mathbb{E}[Y|X_S] - g(X_S)|^2 \right].$$

# Focused Adversarial Invariance Regularizer (FAIR)

## ★ Endogeneity (FAIR) Penalty:

♣ delete endogenous variables

$$J(g) = \max_{f \in S_g} \left\{ \sum_{e \in [K]} \mathbb{E}_{\mu^{(e)}} \left[ \{Y - g(X)\} f_e(X) - \frac{1}{2} f_e^2(X) \right] \right\}.$$

—  $S_g$  is the support of function  $g$

★ When  $\text{supp}(g) = S$ , maximizing all  $f_e(X_S)$  gives

$$\frac{1}{2} \sum_{e \in [K]} \mathbb{E}_{\mu^{(e)}} \left[ |\mathbb{E}[Y|X_S] - g(X_S)|^2 \right].$$

# FAIR Estimation Method

■ Predictor class  $\mathcal{G}$ ,      Discriminator class  $\mathcal{F}$ .

★ Population-level Objective Function:

$$Q(g, f; \gamma) = \sum_{e \in [K]} \mathbb{E}_{\mu^{(e)}} [\ell(g(X), Y)] + \gamma J(g)$$

★ Empirical FAIR Estimator:  $\mathbb{E} \rightsquigarrow \hat{\mathbb{E}}$

$$\hat{g} \in \operatorname{argmin}_{g \in \mathcal{G}} \max_{f \in \mathcal{F}_{Sg}} \hat{Q}(g, f; \gamma)$$

# FAIR Estimation Method

■ Predictor class  $\mathcal{G}$ ,      Discriminator class  $\mathcal{F}$ .

★ Population-level Objective Function:

$$Q(g, f; \gamma) = \sum_{e \in [K]} \mathbb{E}_{\mu^{(e)}} [\ell(g(X), Y)] + \gamma J(g)$$

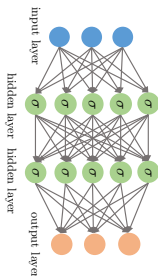
★ Empirical FAIR Estimator:  $\mathbb{E} \rightsquigarrow \hat{\mathbb{E}}$

$$\hat{g} \in \operatorname{argmin}_{g \in \mathcal{G}} \max_{f \in \mathcal{F}_{Sg}} \hat{Q}(g, f; \gamma)$$

# Causal Adversarial Networks

## FAIR-NN:

- ★  $\mathcal{G}$ : ReLU network with width  $N$  and depth  $L$ .
- ★  $\mathcal{F}$ : ReLU network with width  $2N$  and depth  $L + 2$ .



Identifiability: IDF-A

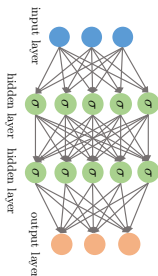
$$\star m^{(e,S)}(x) = \mathbb{E}[Y^{(e)} | X_S = x_S]$$

$$\forall S \text{ if } \bar{m}^{(SUS^*)} \neq m^* \implies \exists e, e' \in [K], \text{ s.t. } m^{(e,S)} \neq m^{(e',S)}$$

# Causal Adversarial Networks

## FAIR-NN:

- ★  $\mathcal{G}$ : ReLU network with width  $N$  and depth  $L$ .
- ★  $\mathcal{F}$ : ReLU network with width  $2N$  and depth  $L + 2$ .



Identifiability: IDF-A

$$\star m^{(e,S)}(x) = \mathbb{E}[Y^{(e)} | X_S = x_S]$$

$$\forall S \text{ if } \bar{m}^{(S, S^*)} \neq m^* \implies \exists e, e' \in [K], \text{ s.t. } m^{(e,S)} \neq m^{(e',S)}$$



# Properties for FAIR-NN

$$\gamma^* = \sup_{S: m^* \neq \bar{m}(SUS^*)} \frac{\|\bar{m}(SUS^*) - m^*\|_2^2}{\frac{1}{|\mathcal{E}|} \|\mathbf{m}^{(e,S)} - \bar{m}(S)\|_{2,e}^2}$$

← Bias of LS w/ all data  
← Variance of biases

## Theorem 4. (Oracle-type of Inequality)

END

Under Conditions IDF-A, if  $\gamma \geq 8\gamma^*$ , for large enough  $n$ ,

$$\|\hat{g} - m^*\|_2 \leq \tilde{C} \left\{ \inf_{g \in \mathcal{G}_{S^*}} \|g - m^*\|_2 + \frac{NL}{\sqrt{n}} \right\}.$$

① Rates depends on approx. errors of  $m^*$

Adaptive Learning

② For HCM  $m^* = f_1 \circ \dots \circ f_q$ , rate is  $n^{-\frac{\alpha^*}{2\alpha^*+1}}$ , with  $\alpha^* = \min(\beta_j/d_j)$ .

e.g.  $m^* = f_1(x_1) + \dots + f_p(x_p)$

$$m^* = f_1(x_1, f_2(x_2, x_3)) + f_4(x_2, x_9)$$

# Properties for FAIR-NN

$$\gamma^* = \sup_{S: m^* \neq \bar{m}(SUS^*)} \frac{\|\bar{m}(SUS^*) - m^*\|_2^2}{\frac{1}{|\mathcal{E}|} \|\mathbf{m}^{(e,S)} - \bar{m}(S)\|_{2,e}^2}$$

← Bias of LS w/ all data  
← Variance of biases

## Theorem 4. (Oracle-type of Inequality)

END

Under Conditions IDF-A, if  $\gamma \geq 8\gamma^*$ , for large enough  $n$ ,

$$\|\hat{g} - m^*\|_2 \leq \tilde{C} \left\{ \inf_{g \in \mathcal{G}_{S^*}} \|g - m^*\|_2 + \frac{NL}{\sqrt{n}} \right\}.$$

① Rates depends on approx. errors of  $m^*$

**Adaptive Learning**

② For HCM  $m^* = f_1 \circ \dots \circ f_q$ , rate is  $n^{-\frac{\alpha^*}{2\alpha^*+1}}$ , with  $\alpha^* = \min(\beta_j/d_j)$ .

e.g.  $m^* = f_1(x_1) + \dots + f_p(x_p)$

$$m^* = f_1(x_1, f_2(x_2, x_3)) + f_4(x_2, x_9)$$

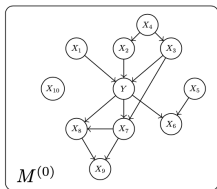
# Causality under SCM

$S^*$  is direct causes under non-degenerate interventions

# Structural Causal Model with Intervention

**SCM Model:** For each env  $e \in \mathcal{E}$ ,  $(X^{(e)}, Y^{(e)}) = (Z_1^{(e)}, \dots, Z_d^{(e)}, Z_{d+1}^{(e)})$

$$X_j^{(e)} \leftarrow f_j^{(e)}(Z_{\text{pa}(j)}^{(e)}, U_j) \quad \forall j \in [d], \quad Y^{(e)} \leftarrow f_{d+1}(X_{\text{pa}(d+1)}^{(e)}, U_{d+1})$$



**Intervention:** Some  $X_j$  intervened: SCM  $\tilde{M}$  of  $(X, Y, E)$  is  $E \leftarrow \text{Unif}([K])$

$$X_j \leftarrow \begin{cases} f_j(Z_{\text{pa}(j)}, U_j, \mathbf{E}) & \forall j \in I \\ f_j(Z_{\text{pa}(j)}, U_j) & \forall j \in [d] \setminus I \end{cases} \quad Y \leftarrow f_{d+1}(X_{\text{pa}(d+1)}, U_{d+1})$$

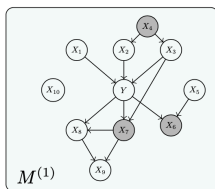
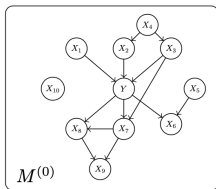
★ DAG induced graph

★ Unknown interventions, not on  $Y$ .

# Structural Causal Model with Intervention

**SCM Model:** For each env  $e \in \mathcal{E}$ ,  $(X^{(e)}, Y^{(e)}) = (Z_1^{(e)}, \dots, Z_d^{(e)}, Z_{d+1}^{(e)})$

$$X_j^{(e)} \leftarrow f_j^{(e)}(Z_{\text{pa}(j)}^{(e)}, U_j) \quad \forall j \in [d], \quad Y^{(e)} \leftarrow f_{d+1}(X_{\text{pa}(d+1)}^{(e)}, U_{d+1})$$



$$I = \{4, 6, 7\}$$

**Intervention:** Some  $X_i$  intervened: SCM  $\tilde{M}$  of  $(X, Y, E)$  is  $E \leftarrow \text{Unif}([K])$

$$X_j \leftarrow \begin{cases} f_j(Z_{\text{pa}(j)}, U_j, \mathbf{E}) & \forall j \in I \\ f_j(Z_{\text{pa}(j)}, U_j) & \forall j \in [d] \setminus I \end{cases} \quad Y \leftarrow f_{d+1}(X_{\text{pa}(d+1)}, U_{d+1})$$

★ DAG induced graph

★ Unknown interventions, not on  $Y$ .

# Characterizing $S^*$ under Intervention

## Theorem 5. Existence of Maximum Invariant Set

Under nondegenerate interventions, Condition IDF-A holds with

$$S^* = \text{pa}(d+1) \cup A(I) \cup \bigcup_{j \in A(I)} (\text{pa}(j) \setminus \{d+1\})$$

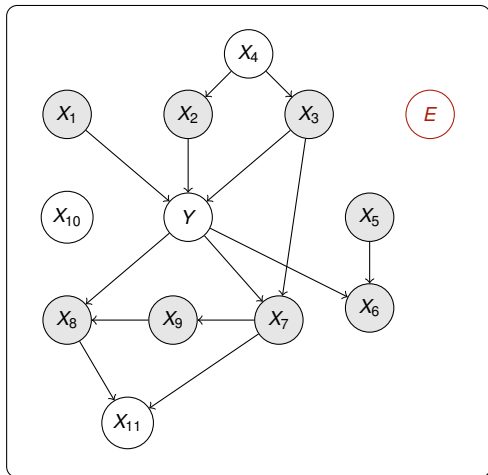
where  $A(I) = \{j : j \in \text{ch}(d+1), \text{at}(j) \cap \text{ch}(d+1) \cap I = \emptyset\}$

Invariant variables: ★parents of  $Y$ ;

★uninterviewed children of  $Y$ ;

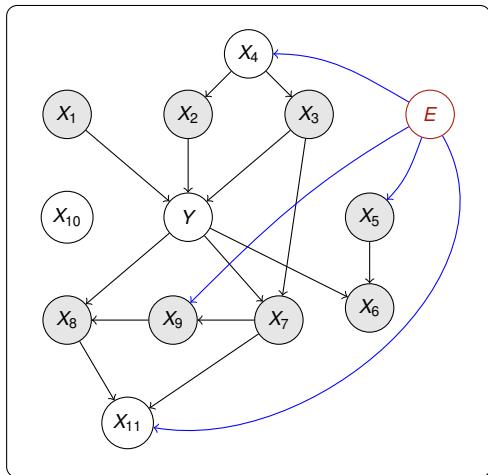
★parents of uninterviewed children of  $Y$ .

# An Illustration: $k = 1$



$$0 \leftrightarrow 0, S^* = \{1, 2, 3, 5, 6, 7, 8, 9\}$$

# An Illustration: $k = 2$

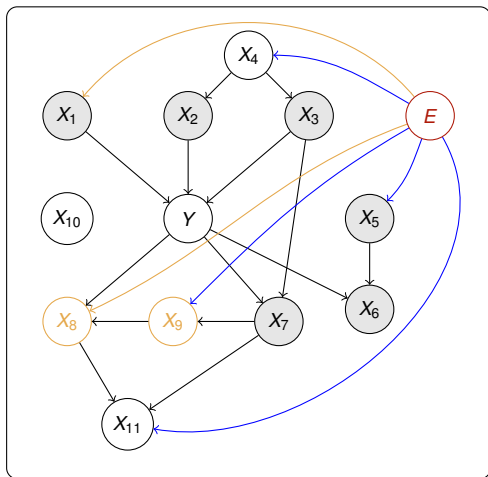


$0 \leftrightarrow 0, S^* = \{1, 2, 3, 5, 6, 7, 8, 9\}$

$0 \leftrightarrow \mathbf{1}, S^* = \{1, 2, 3, 5, 6, 7, 8, 9\}$



# An Illustration: $k = 3$

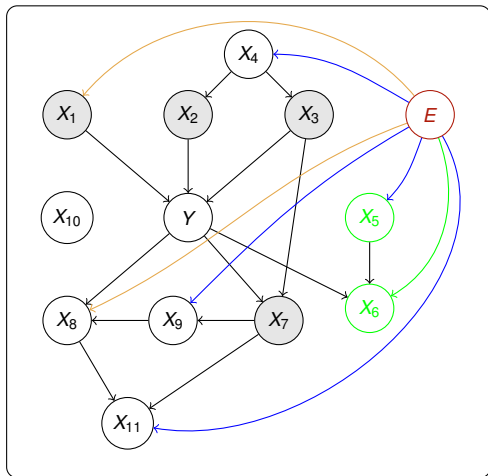


$$0 \leftrightarrow 0, S^* = \{1, 2, 3, 5, 6, 7, 8, 9\}$$

$$0 \leftrightarrow 1, S^* = \{1, 2, 3, 5, 6, 7, 8, 9\}$$

$$0 \leftrightarrow 2, S^* = \{1, 2, 3, 5, 6, 7\}$$

# An Illustration: $k = 4$



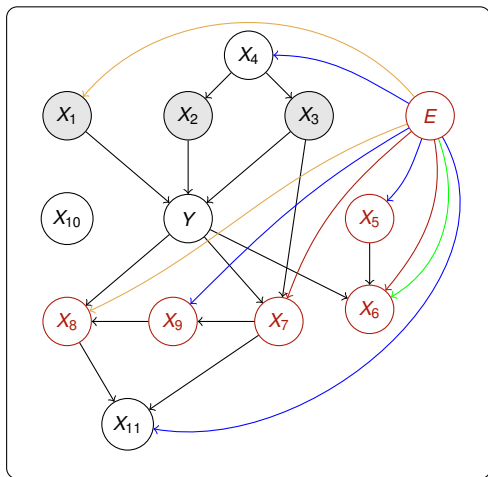
$$0 \leftrightarrow 0, S^* = \{1, 2, 3, 5, 6, 7, 8, 9\}$$

$$0 \leftrightarrow 1, S^* = \{1, 2, 3, 5, 6, 7, 8, 9\}$$

$$0 \leftrightarrow 2, S^* = \{1, 2, 3, 5, 6, 7\}$$

$$0 \leftrightarrow 3, S^* = \{1, 2, 3, 7\}$$

# An Illustration: $k = 5$



$$0 \leftrightarrow 0, S^* = \{1, 2, 3, 5, 6, 7, 8, 9\}$$

$$0 \leftrightarrow 1, S^* = \{1, 2, 3, 5, 6, 7, 8, 9\}$$

$$0 \leftrightarrow 2, S^* = \{1, 2, 3, 5, 6, 7\}$$

$$0 \leftrightarrow 3, S^* = \{1, 2, 3, 7\}$$

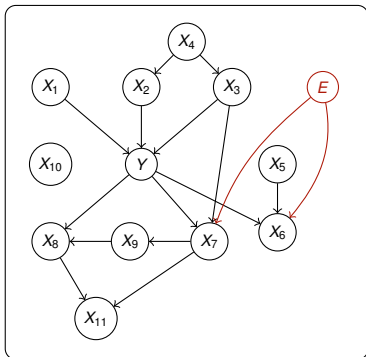
$$0 \leftrightarrow 4, S^* = \{1, 2, 3\}$$

# Exact Direct Causal Recovery

## Proposition 1. Sufficient and Necessary Condition for Causal Discovery

When all root-children are intervened (★),  $S^* = \text{pa}(d+1)$ . The condition is also necessary, if  $Y$  does not have degenerate children.

★  $I \supseteq I^*$ , where  $I^* = \{j : j \in \text{ch}(d+1), \text{pa}(j) \cap \text{ch}(d+1) = \emptyset\}$ .



$$0 \leftrightarrow 4$$

$$S^* = \{1, 2, 3\}$$

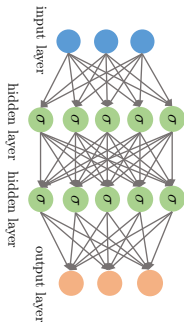
# Implementation and Simulations

# Challenge of Implementation

Parameterization:  $g_\theta, f_{e, \phi_e}$  with  $e \in [K]$

Objective:  $\hat{g} \in \operatorname{argmin}_{g \in \mathcal{G}} \max_{\{f_e \in \mathcal{F}_{Sg}\}_{e \in [K]}} \hat{Q}(g, f^{[K]}; \gamma)$

- ★ min-max optimization.  $\rightsquigarrow$  gradient descent ascent
- ★  $f$  has the same  $X$ -variables as  $g$ .



$$\min_{\theta, a \in \{0,1\}^d} \max_{\phi_1, \dots, \phi_k} \mathcal{L}(g_\theta(a \odot x), \{f_{e, \phi_e}(a \odot x)\}_{e=1}^K)$$

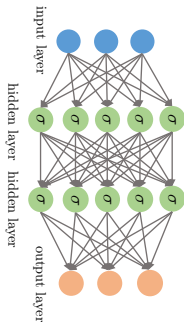
$\rightsquigarrow$  Enumerate all  $a \in \{0,1\}^d$ !

# Challenge of Implementation

Parameterization:  $g_\theta, f_{e, \phi_e}$  with  $e \in [K]$

Objective:  $\hat{g} \in \operatorname{argmin}_{g \in \mathcal{G}} \max_{\{f_e \in \mathcal{F}_{S_g}\}_{e \in [K]}} \hat{Q}(g, f^{[K]}; \gamma)$

- ★ min-max optimization.  $\rightsquigarrow$  gradient descent ascent
- ★  $f$  has the same  $X$ -variables as  $g$ .



$$\min_{\theta, a \in \{0, 1\}^d} \max_{\phi_1, \dots, \phi_k} \mathcal{L}(g_\theta(a \odot x), \{f_{e, \phi_e}(a \odot x)\}_{e=1}^K)$$

$\rightsquigarrow$  Enumerate all  $a \in \{0, 1\}^d$ !

# Gradient Method with Gumbel Approximation

Equivalence Problem:  $\sigma(u) = 1/(1 + e^{-u})$

$$\min_{\theta, w} \max_{\phi_1, \dots, \phi_k} \mathbb{E}_{A \sim \text{Bern}(\sigma(w))} \mathcal{L}(g_{\theta}(A \odot x), \{f_{e, \phi_e}(A \odot x)\}_{e=1}^K)$$



# Gradient Method with Gumbel Approximation

$$\min_{\theta, w} \max_{\phi_1, \dots, \phi_k} \mathbb{E}_{A \sim \text{Bern}(\sigma(w))} \mathcal{L}(g_{\theta}(A \odot x), \{f_{e, \phi_e}(A \odot x)\}_{e=1}^k)$$

Gumbel Approx:  $\text{Bern}(\sigma(w)) = I(U - \sigma(w) < 0)$

# Gradient Method with Gumbel Approximation

$$\min_{\theta, w} \max_{\phi_1, \dots, \phi_k} \mathbb{E}_{A \sim \text{Bern}(\sigma(w))} \mathcal{L}(g_{\theta}(A \odot x), \{f_{e, \phi_e}(A \odot x)\}_{e=1}^k)$$

Gumbel Approx:  $\text{Bern}(\sigma(w)) = I(U - \sigma(w) < 0) = I(\text{logit}(U) - w < 0)$

# Gradient Method with Gumbel Approximation

$$\min_{\theta, w} \max_{\phi_1, \dots, \phi_k} \mathbb{E}_{A \sim \text{Bern}(\sigma(w))} \mathcal{L}(g_{\theta}(A \odot x), \{f_{e, \phi_e}(A \odot x)\}_{e=1}^K)$$

Gumbel Approx:  $\text{Bern}(\sigma(w)) = I(U - \sigma(w) < 0) \approx \frac{1}{1 + \exp((\text{logit}(U) - w)/\tau)}$ , as  $\tau \rightarrow 0$

# Gradient Method with Gumbel Approximation

$$\min_{\theta, w} \max_{\phi_1, \dots, \phi_k} \mathbb{E}_{A \sim \text{Bern}(\sigma(w))} \mathcal{L}(g_{\theta}(A \odot x), \{f_{e, \phi_e}(A \odot x)\}_{e=1}^K)$$

Gumbel Approx:  $\text{Bern}(\sigma(w)) = I(U - \sigma(w) < 0) \approx \frac{1}{1 + \exp((V_2 - V_1 - w)/\tau)} \equiv B_{\tau}(V, w)$  as  $\tau \rightarrow 0$

# Gradient Method with Gumbel Approximation

$$\min_{\theta, w} \max_{\phi_1, \dots, \phi_k} \mathbb{E}_{A \sim \text{Bern}(\sigma(w))} \mathcal{L}(g_{\theta}(A \odot x), \{f_{e, \phi_e}(A \odot x)\}_{e=1}^K)$$

Gumbel Approx:  $\text{Bern}(\sigma(w)) = I(U - \sigma(w) < 0) \approx \frac{1}{1 + \exp((V_2 - V_1 - w)/\tau)} \equiv B_{\tau}(V, w)$  as  $\tau \rightarrow 0$

$$\min_{\theta, w} \max_{\phi_1, \dots, \phi_k} \mathbb{E}_{V \sim \text{Gum}} \mathcal{L}(g_{\theta}(B_{\tau}(V, w) \odot x), \{f_{e, \phi_e}(B_{\tau}(V, w) \odot x)\}_{e=1}^K)$$

# Gradient Method with Gumbel Approximation

$$\min_{\theta, w} \max_{\phi_1, \dots, \phi_k} \mathbb{E}_{A \sim \text{Bern}(\sigma(w))} \mathcal{L}(g_{\theta}(A \odot x), \{f_{e, \phi_e}(A \odot x)\}_{e=1}^K)$$

Gumbel Approx:  $\text{Bern}(\sigma(w)) = I(U - \sigma(w) < 0) \approx \frac{1}{1 + \exp((V_2 - V_1 - w)/\tau)} \equiv B_{\tau}(V, w)$  as  $\tau \rightarrow 0$

$$\min_{\theta, w} \max_{\phi_1, \dots, \phi_k} \mathbb{E}_{V \sim \text{Gum}} \mathcal{L}(g_{\theta}(B_{\tau}(V, w) \odot x), \{f_{e, \phi_e}(B_{\tau}(V, w) \odot x)\}_{e=1}^K)$$

## Algorithm:

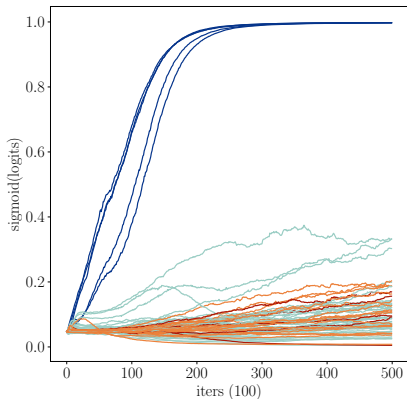
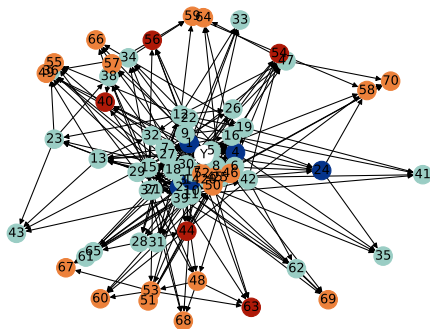
- 1 Sample  $V$ , batch of  $(X^{(e)}, Y^{(e)})$  from each environment.
- 2 Gradient ascent update for  $\phi_1, \dots, \phi_k$ .
- 3 Gradient descent update for  $\theta, w$ ,

with decreasing temperature  $\tau$ .

# Linear Model, $\mathcal{G} = \mathcal{F}$ linear

- $k = 2, d = 70$ ,
- Random generated SCM sharing same cause-effect relationship.
- All  $X$  are intervened (randomly).

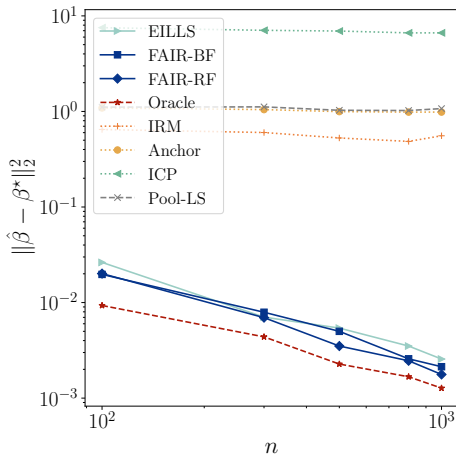
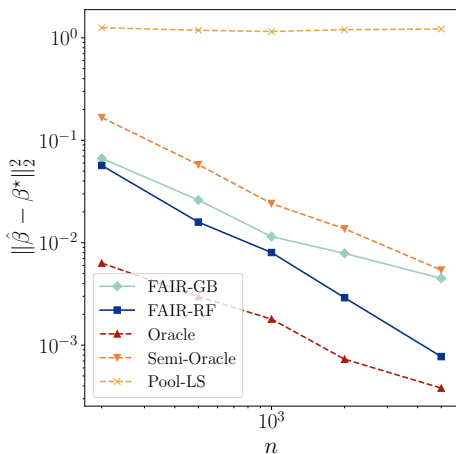
brute force search is impossible



Relations with Y: blue = parent, red = child, orange = offspring, lightblue = other

# Performance of FAIR-Linear

- FAIR-GB: implementation using Gumbel approximation.
- FAIR-RF: refitting after running FAIR-GB.



- ★ log-log plot of med(MSE) based on  $N_{sim} = 50$  for ● (a)  $p = 70$ ,  $n \in \{200, 500, 1000, 2000, 5000\}$  and ●  $p = 15$  and  $n \in \{100, 200, 500, 800, 1000\}$

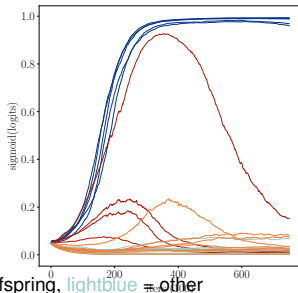
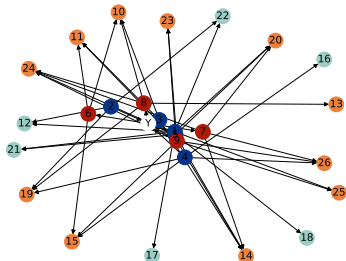


# Simulations for FAIR-NN

$d = 26, k = 2$

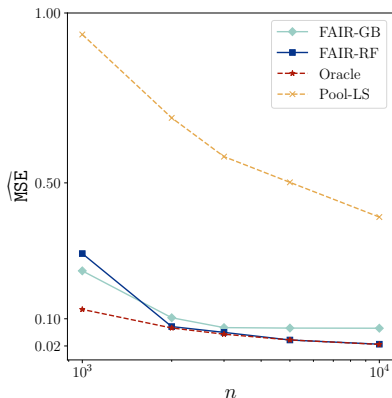
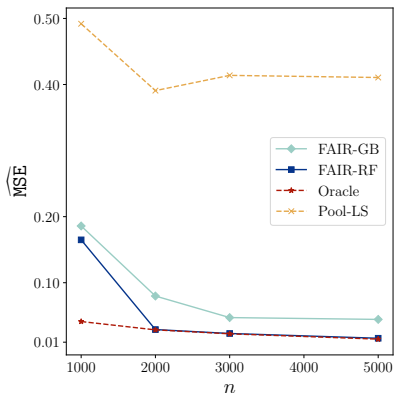
$$X_i^{(e)} \leftarrow \begin{cases} \varepsilon_i^{(e)} & i \leq 5 \\ f_{i,0}^{(e)}(Y^{(e)}) + \varepsilon_i^{(e)} & 6 \leq i \leq 9 \\ \sum_{j \in \text{pa}(i) \subseteq [8]} f_{i,j}^{(e)}(X_j^{(e)}) + \varepsilon_i^{(e)} & 10 \leq i \leq 26 \end{cases}$$
$$Y^{(e)} \leftarrow m_k^*(X_1^{(e)}, \dots, X_5^{(e)}) + \varepsilon_0,$$

$m_1^*(x)$  additive,  $m_2^*(x) = x_1 x_2^3 + \log(1 + e^{\tanh(x_3)} + e^{x_4}) + \sin(x_5)$  HCM.



Relations with Y: blue = parent, red = child, orange = offspring, lightblue = other

# Performance of FAIR-NN



★ MSE over  $N_{sim} = 50$  over 30K x-values. ● (a) additive  $m_1^*$  and  $n \in \{1000, 2000, 3000, 5000\}$

● (b)  $m_2^*$  and  $n \in \{1000, 2000, 3000, 5000, 10000\}$ .

# Application I: Transfer Learning

## Waterbird Classification

- ◆  $Y = 1$  (**water bird**) and  $Y = 0$  (**land bird**)
- ◆  $X \in \mathbb{R}^{500}$  extracted from ResNet pre-trained on ImageNet.

## Data

- ★ Training data with spurious background (n=50k).
  - ◆  $\mathcal{D}_1$ : 95% **water birds** on **water**, 90% **land birds** on **land**.
  - ◆  $\mathcal{D}_2$ : 75% **water birds** on **water**, 70% **land birds** on **land**.
- ★ Test data with reverse spurious background (n=30k).
  - ◆  $\mathcal{D}_3$ : 98% **water birds** on **land**, 98% **land birds** on **water**.



—Build a linear classifier using  $(\mathcal{D}_1, \mathcal{D}_2)$ .

# Application I: Transfer Learning

## Waterbird Classification

- ◆  $Y = 1$  (**water bird**) and  $Y = 0$  (**land bird**)
- ◆  $X \in \mathbb{R}^{500}$  extracted from ResNet pre-trained on ImageNet.

## Data

- ★ Training data with spurious background (n=50k).
  - ◆  $\mathcal{D}_1$ : 95% **water birds** on **water**, 90% **land birds** on **land**.
  - ◆  $\mathcal{D}_2$ : 75% **water birds** on **water**, 70% **land birds** on **land**.
- ★ Test data with reverse spurious background (n=30k).
  - ◆  $\mathcal{D}_3$ : 98% **water birds** on **land**, 98% **land birds** on **water**.



—Build a linear classifier using  $(\mathcal{D}_1, \mathcal{D}_2)$ .

# Application I: Transfer Learning

## Waterbird Classification

- ◆  $Y = 1$  (**water bird**) and  $Y = 0$  (**land bird**)
- ◆  $X \in \mathbb{R}^{500}$  extracted from ResNet pre-trained on ImageNet.

## Data

- ★ Training data with spurious background (n=50k).
  - ◆  $\mathcal{D}_1$ : 95% **water birds** on **water**, 90% **land birds** on **land**.
  - ◆  $\mathcal{D}_2$ : 75% **water birds** on **water**, 70% **land birds** on **land**.
- ★ Test data with reverse spurious background (n=30k).
  - ◆  $\mathcal{D}_3$ : 98% **water birds** on **land**, 98% **land birds** on **water**.



—Build a linear classifier using  $(\mathcal{D}_1, \mathcal{D}_2)$ .

# Application I: Transfer Learning

## Waterbird Classification

- ◆  $Y = 1$  (**water bird**) and  $Y = 0$  (**land bird**)
- ◆  $X \in \mathbb{R}^{500}$  extracted from ResNet pre-trained on ImageNet.

## Data

- ★ Training data with spurious background (n=50k).
  - ◆  $\mathcal{D}_1$ : 95% **water birds** on **water**, 90% **land birds** on **land**.
  - ◆  $\mathcal{D}_2$ : 75% **water birds** on **water**, 70% **land birds** on **land**.
- ★ Test data with reverse spurious background (n=30k).
  - ◆  $\mathcal{D}_3$ : 98% **water birds** on **land**, 98% **land birds** on **water**.



—Build a linear classifier using  $(\mathcal{D}_1, \mathcal{D}_2)$ .

# Bias Reduction from Two Biased Samples

## Methods and Results:

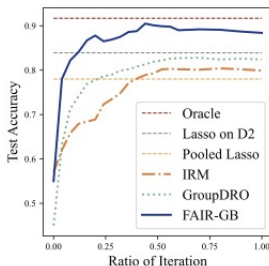
- ★ **FAIR-GB** FAIR estimator with linear  $(\mathcal{G}, \mathcal{F})$ , cross-entropy loss and Gumbel approx.
- ★ **PooledLasso** on  $\mathcal{D}_1 \cup \mathcal{D}_2$ ; **Lasso on D2** Lasso on  $\mathcal{D}_2$ .
- ★ **Oracle**: Lasso on  $\mathcal{D}_4$  where label/background independent.
- ★ **IRM** invariant risk minimization; **GroupDRO** group distributionally robust optimization.

# Bias Reduction from Two Biased Samples

## Methods and Results:

- ★ **FAIR-GB** FAIR estimator with linear  $(\mathcal{G}, \mathcal{F})$ , cross-entropy loss and Gumbel approx.
- ★ **PooledLasso** on  $\mathcal{D}_1 \cup \mathcal{D}_2$ ; **Lasso on D2** Lasso on  $\mathcal{D}_2$ .
- ★ **Oracle**: Lasso on  $\mathcal{D}_4$  where label/background independent.
- ★ **IRM** invariant risk minimization; **GroupDRO** group distributionally robust optimization.

Method	Test Accuracy
<b>Oracle</b>	<b>91.06 <math>\pm</math> 0.24 %</b>
Lasso on D2	84.57 $\pm$ 0.71 %
Pooled Lasso	79.08 $\pm$ 0.54 %
IRM	80.32 $\pm$ 0.67 %
GroupDRO	82.75 $\pm$ 1.10 %
<b>FAIR-GB</b>	<b>89.56 <math>\pm</math> 0.53 %</b>



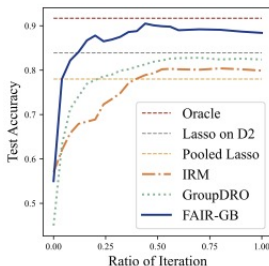


# Bias Reduction from Two Biased Samples

## Methods and Results:

- ★ **FAIR-GB** FAIR estimator with linear  $(\mathcal{G}, \mathcal{F})$ , cross-entropy loss and Gumbel approx.
- ★ **PooledLasso** on  $\mathcal{D}_1 \cup \mathcal{D}_2$ ; **Lasso on D2** Lasso on  $\mathcal{D}_2$ .
- ★ **Oracle**: Lasso on  $\mathcal{D}_4$  where label/background independent.
- ★ **IRM** invariant risk minimization; **GroupDRO** group distributionally robust optimization.

Method	Test Accuracy
<b>Oracle</b>	<b>91.06 ± 0.24 %</b>
Lasso on D2	84.57 ± 0.71 %
Pooled Lasso	79.08 ± 0.54 %
IRM	80.32 ± 0.67 %
GroupDRO	82.75 ± 1.10 %
<b>FAIR-GB</b>	<b>89.56 ± 0.53 %</b>



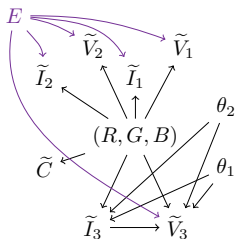
★ **FAIR can correct bias from two biased samples!**



# Data and Ground Truth

**Data**  $\check{D}_1, \check{D}_2$  sub-sample of  $D_1, D_2$  with equal size  $n$ .

## Augmented SCM graph



★ Direct Causes  $S^* = (R, G, B, \theta_1, \theta_2)$ .

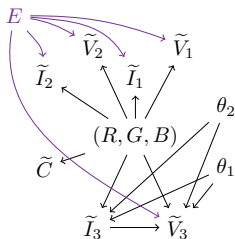
★ Challenges

- ◆ weak & nonlinear signal  $\tilde{I}_3 \propto \cos^2(\theta_1 - \theta_2)$ .
- ◆ strong spurious association  $\tilde{I}_3 \leftrightarrow \tilde{V}_3$ .
- ◆ strong explained  $R^2$  for  $\tilde{I}_2, \tilde{I}_1$  ( $\geq 0.9$ ).

# Data and Ground Truth

**Data**  $\check{D}_1, \check{D}_2$  sub-sample of  $D_1, D_2$  with equal size  $n$ .

## Augmented SCM graph



★ Direct Causes  $S^* = (R, G, B, \theta_1, \theta_2)$ .

★ Challenges

- ◆ weak & nonlinear signal  $\tilde{I}_3 \propto \cos^2(\theta_1 - \theta_2)$ .
- ◆ strong spurious association  $\tilde{I}_3 \leftrightarrow \tilde{V}_3$ .
- ◆ strong explained  $R^2$  for  $\tilde{I}_2, \tilde{I}_1$  ( $\geq 0.9$ ).

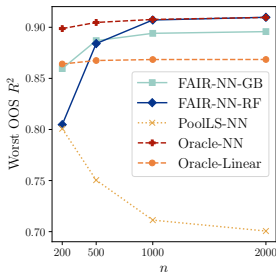
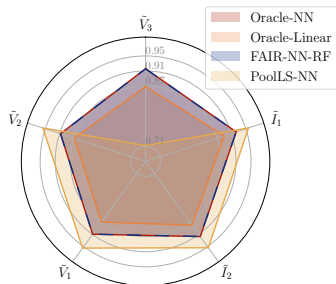
# Prediction Performance Evaluation

## Methods

- ★ **FAIR-NN-GB**: Gumbel implented FAIR-NN, **FAIR-NN-RF** refitted estimator
- ★ **Oracle- $M$** : Regress  $Y$  on  $X_{S^*}$  using  $M \in \{\text{Linear}, \text{NN}\}$ .
- ★ **PoolLS-NN**: Regress  $Y$  on  $X_{S^*}$  using all the data and NN.

## Evaluate the Dependency on Variables Other Than $X_{S^*}$

- ★ Out-of-sample(OOS)- $R^2$  on  $\mathcal{D}_{3,Z}$  with  $Z \in \{\tilde{V}_j\}_{j=1}^3 \cup \{\tilde{I}_j\}_{j=1}^2$  where  $Z$  is **strongly** intervened.



- n=1000** ♦ Remove strong spurious var  $\tilde{V}_3$  (otherwise  $R^2 \downarrow 0.2$ ) ♦ Detect weak signals  $(\theta_1, \theta_2)$ :  $R^2 \uparrow 0.04$  as Linear  $\rightarrow$  NN.

~ Near oracle performance.

# Attaining Variable Selection Consistency

## Methods

- ★ **FAIR-M**: Gumbel implented FAIR-M method  $M \in \{\text{Linear}, \text{NN}\}$ ,  $\hat{S} = \{j : \sigma(w_j) > 0.9\}$ .
- ★ **ForestVarSel**: Select by importance measure using RandomForest
- ★ **NonlinearICP**: Previous invariance learning estimators.

**Results** (blue=parent, red=child, orange=neither ancestor nor descendants.)

		$R$	$G$	$B$	$\theta_1$	$\theta_2$	$\tilde{V}_1$	$\tilde{V}_2$	$\tilde{V}_3$	$\tilde{I}_1$	$\tilde{I}_2$	$\tilde{C}$
$n = 1000$	FAIR-Linear	1.0	1.0	1.0	0.09	0.06	0.0	0.0	0.0	0.0	0.0	0.01
	FAIR-NN	1.0	1.0	1.0	0.99	0.99	0.0	0.0	0.0	0.03	0.01	0.0
	ForestVarSel	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	1.0	0.0
	NonlinearICP	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
FAIR-NN	$n = 200$	1.0	0.99	0.74	0.61	0.65	0.04	0.02	0.29	0.04	0.04	0.01
	$n = 500$	1.0	1.0	0.95	0.86	0.89	0.03	0.03	0.03	0.01	0.02	0.0
	$n = 2000$	1.0	1.0	1.0	1.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0

★ Variable Selection Consistency ★ NN detect nonlinear Malus's law  $\tilde{I}_3 \propto \cos^2(\theta_1 - \theta_2)$ .

# Attaining Variable Selection Consistency

## Methods

- ★ **FAIR-M**: Gumbel implented FAIR-M method  $M \in \{\text{Linear}, \text{NN}\}$ ,  $\hat{S} = \{j : \sigma(w_j) > 0.9\}$ .
- ★ **ForestVarSel**: Select by importance measure using RandomForest
- ★ **NonlinearICP**: Previous invariance learning estimators.

## Results (blue=parent, red=child, orange=neither ancestor nor descendants.)

		$R$	$G$	$B$	$\theta_1$	$\theta_2$	$\tilde{V}_1$	$\tilde{V}_2$	$\tilde{V}_3$	$\tilde{I}_1$	$\tilde{I}_2$	$\tilde{C}$
$n = 1000$	FAIR-Linear	1.0	1.0	1.0	0.09	0.06	0.0	0.0	0.0	0.0	0.0	0.01
	FAIR-NN	1.0	1.0	1.0	0.99	0.99	0.0	0.0	0.0	0.03	0.01	0.0
	ForestVarSel	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	1.0	0.0
	NonlinearICP	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
FAIR-NN	$n = 200$	1.0	0.99	0.74	0.61	0.65	0.04	0.02	0.29	0.04	0.04	0.01
	$n = 500$	1.0	1.0	0.95	0.86	0.89	0.03	0.03	0.03	0.01	0.02	0.0
	$n = 2000$	1.0	1.0	1.0	1.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0

- ★ Variable Selection Consistency
- ★ NN detect nonlinear Malus's law  $\tilde{I}_3 \propto \cos^2(\theta_1 - \theta_2)$ .

# Summary

1 Introduce a FAIR-NN method to learn causal predictors w/o knowledge of

★cause-effect      ★function structure

◆ Neural network  $\leadsto$  learn feature representation from data

◆ Invariance  $\leadsto$  distinguish causal/non-causal via **FAIR-penalty**  $J_0$

2 Establish sample efficiency in different aspects.

- ▶ Minimal identification condition.
- ▶ Convergence rate depends on  $m^*$ , adapt to low-dimension structures.
- ▶ Regularization hyper-parameter minor impact.

3 Give an efficient implementation via Gumble Approx using SGD.



# Summary

- 1 Introduce a FAIR-NN method to learn causal predictors w/o knowledge of
  - ★cause-effect      ★function structure
  - ◆ Neural network  $\leadsto$  learn feature representation from data
  - ◆ Invariance  $\leadsto$  distinguish causal/non-causal via **FAIR-penalty**  $J_0$
- 2 Establish sample efficiency in different aspects.
  - ▶ Minimal identification condition.
  - ▶ Convergence rate depends on  $m^*$ , adapt to low-dimension structures.
  - ▶ Regularization hyper-parameter minor impact.
- 3 Give an efficient implementation via Gumble Approx using SGD.

# Summary

- 1 Introduce a FAIR-NN method to learn causal predictors w/o knowledge of
  - ★ cause-effect      ★ function structure
  - ◆ Neural network  $\leadsto$  learn feature representation from data
  - ◆ Invariance  $\leadsto$  distinguish causal/non-causal via **FAIR-penalty**  $J_0$
- 2 Establish sample efficiency in different aspects.
  - ▶ Minimal identification condition.
  - ▶ Convergence rate depends on  $m^*$ , adapt to low-dimension structures.
  - ▶ Regularization hyper-parameter minor impact.
- 3 Give an efficient implementation via Gumble Approx using SGD.

# The End

*Thank*



*You*

- ★ Fan, J., Fang, C., Gu, Y., and Zhang, T. (2024+). Environment Invariant Linear Least Squares. AOS
- ★ Gu, Y., Fang, C., Buehlmann, P., and Fan, J. (2024). Causality Pursuit from Heterogeneous Environments via Neural Adversarial Invariance Learning. *arxiv.org*