

# Incorporating Causal Graphical Prior Knowledge into Predictive Modeling via Simple Data Augmentation

Takeshi Teshima, Masashi Sugiyama

The University of Tokyo, RIKEN

Presentors: Tong Zhang, Zhongqi Yang

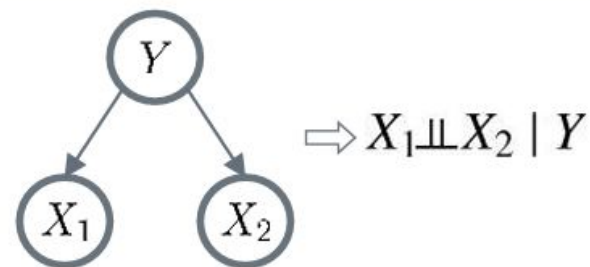
# Agenda

- Background & Problem Definition
- General Idea of ADMGs Data Augmentation
- Practical Algorithm
- Results on Real World Data

## Causal Graphs (CGs) (Pearl, 2009)

Representation of our knowledge of data generating processes.

CGs imply **conditional independence (CI)** relations (Pearl, 2009) (Richardson, 2003) .



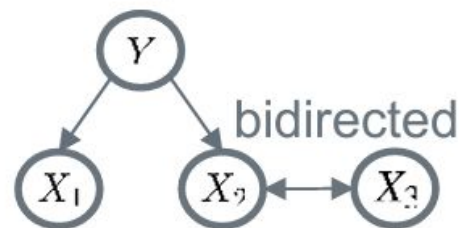
## Acyclic Directed Mixed Graphs (ADMGs) (Richardson, 2003) (Richardson et al., 2017)

Directed acyclic graphs (possibly) with bidirected edges.

Used for causal models with latent variables

(semi-Markov models; cf. Latent projection (Tian et al., 2002)).

$$\mathcal{G} = ([D], \mathcal{E}, \mathcal{B})$$

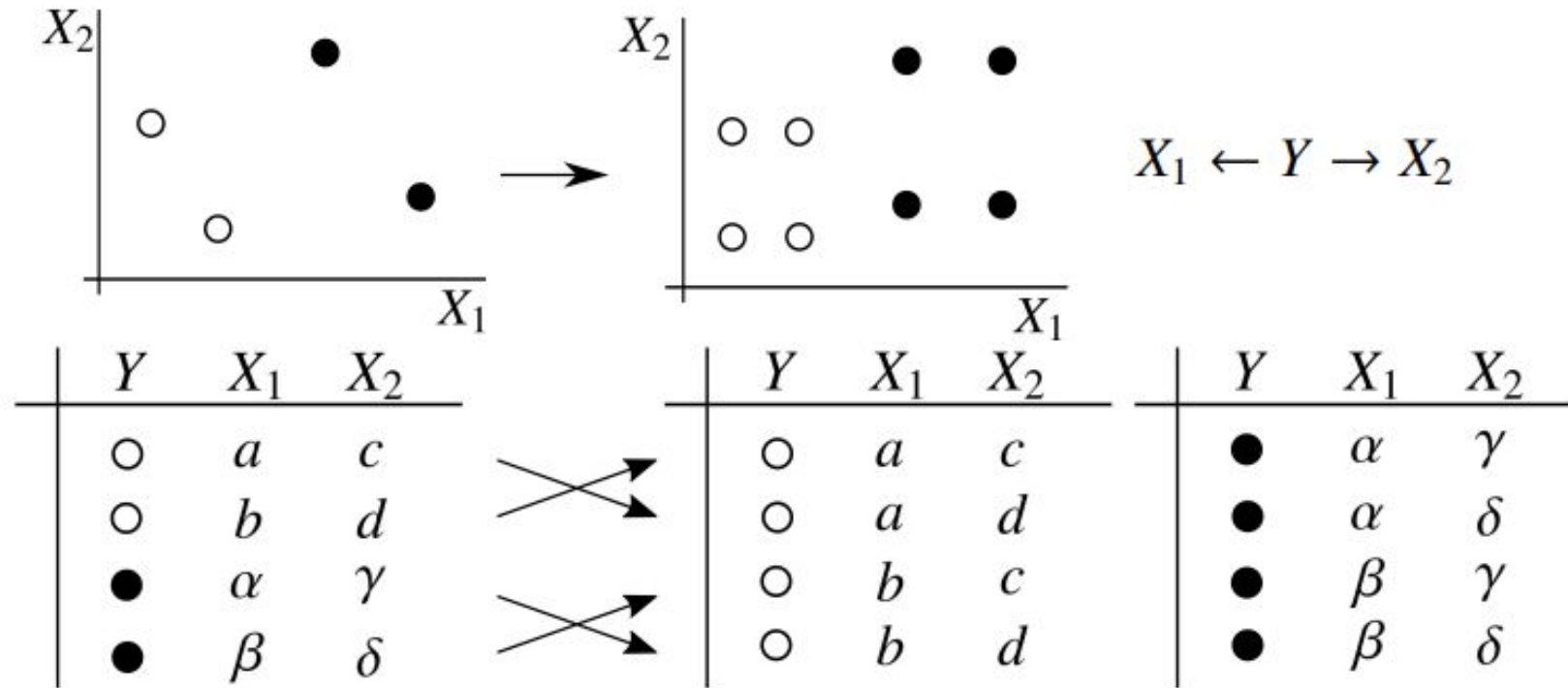


## Topological ADMG Factorization (Tian et al., 2002) (Bhattacharya et al., 2020)

Given a semi-Markov model,  $p(\mathbf{Z}) = \prod_{j=1}^D p_{j|\text{mp}(j)}(\mathbf{Z}^j | \mathbf{Z}^{\text{mp}(j)})$  holds.

$\text{mp}(j)$ : “Markov pillow” of variable  $\mathbf{Z}^j$  (Generalization of “parents” in ADMGs.)

# Augmentation with a CG



# Problem Definition

$\mathbf{Z} = (Z^1, \dots, Z^D) \sim p$  : joint data of  $X$  and  $Y$ .

(each  $Z^j$  may be continuous or discrete)

## Main Assumption

- $p(\mathbf{Z})$  satisfies the topological ADMG factorization w.r.t.  $\mathcal{G}$   
(Bhattacharya *et al.*, 2020)

## We are given:

- Labeled data  $\mathcal{D} = \{\mathbf{Z}_i\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} p$ .
- Estimator  $\hat{\mathcal{G}}$  of the underlying ADMG  $\mathcal{G}$ .

## Goal

Find a predictor  $f : X \mapsto Y$  with small  $R(f) = \mathbb{E}[\ell(f, \mathbf{Z})]$ .

$$R(f) = \int_{\mathbf{Z}} \ell(f, \mathbf{Z}) \prod_{j=1}^D \underbrace{p_{j|\text{mp}(j)}(Z^j | \mathbf{Z}^{\text{mp}(j)})}_{(*)} d\mathbf{Z}.$$

# Proposed Method for ADMG

- Recall topological ADMG factorization:  $p(\mathbf{Z}) = \prod_{j=1}^D p_{j|\text{mp}(j)}(Z^j | \mathbf{Z}^{\text{mp}(j)})$ .
- Approximate each conditional by kernel-based estimator.  
Let  $K^j: \bar{\mathbf{Z}}^{\text{mp}(j)} \rightarrow \mathbb{R}_{\geq 0}$  and

$$p(\mathbf{Z}) \simeq \prod_{j=1}^D \hat{p}_{j|\text{mp}(j)}(Z^j | \mathbf{Z}^{\text{mp}(j)}) := \frac{\sum_{i=1}^n \delta_{Z_i^j}(Z^j) K^j(\mathbf{Z}^{\text{mp}(j)} - \mathbf{Z}_i^{\text{mp}(j)})}{\sum_{k=1}^n K^j(\mathbf{Z}^{\text{mp}(j)} - \mathbf{Z}_k^{\text{mp}(j)})}$$

Empirical conditional density

- Plug-in risk estimator

$$\hat{R}_{\text{aug}}(f) = \int_{\mathbf{Z}} \ell(f, \mathbf{Z}) \prod_{j=1}^D \hat{p}_{j|\text{mp}(j)}(Z^j | \mathbf{Z}^{\text{mp}(j)}) d\mathbf{Z} = \sum_{\mathbf{i} \in [n]^D} \hat{w}_{\mathbf{i}} \cdot \ell(f, \mathbf{Z}_{\mathbf{i}})$$

Augmented data + instance weights

- Considering all the possible resampling candidates
- instance-weighted data augmentation procedure:

$$\hat{R}_{\text{aug}}(f) = \sum_{\mathbf{i} \in [n]^D} \hat{w}_{\mathbf{i}} \cdot \ell(f, \mathbf{Z}_{\mathbf{i}}),$$

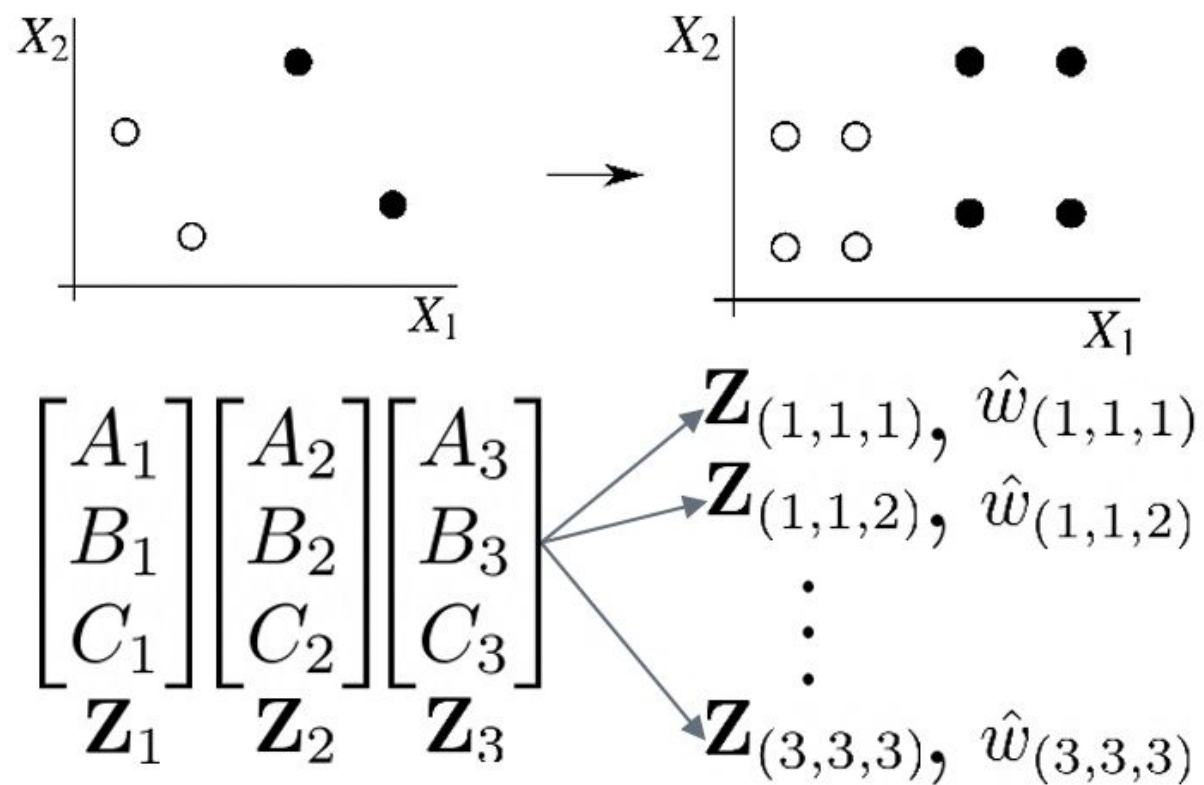
$$\hat{w}_{\mathbf{i}} = \prod_{j=1}^D \frac{K^j(\mathbf{Z}_{\mathbf{i}_{1:j-1}}^{\text{mp}(j)} - \mathbf{Z}_{\mathbf{i}}^{\text{mp}(j)})}{\sum_{k=1}^n K^j(\mathbf{Z}_{\mathbf{i}_{1:j-1}}^{\text{mp}(j)} - \mathbf{Z}_k^{\text{mp}(j)})},$$

$$\mathbf{Z}_{\mathbf{i}} = (Z_{i_1}^1, \dots, Z_{i_D}^D), \quad \mathbf{Z}_{\mathbf{i}_{1:j-1}} = (Z_{i_1}^1, \dots, Z_{i_{j-1}}^{j-1}),$$

$$\hat{f} \in \arg \min_{f \in \mathcal{F}} \{(1 - \lambda) \hat{R}_{\text{emp}}(f) + \lambda \hat{R}_{\text{aug}}(f) + \Omega(f)\}$$



$$\begin{bmatrix} A_1 \\ B_1 \\ C_1 \\ \mathbf{Z}_1 \end{bmatrix} \begin{bmatrix} A_2 \\ B_2 \\ C_2 \\ \mathbf{Z}_2 \end{bmatrix} \begin{bmatrix} A_3 \\ B_3 \\ C_3 \\ \mathbf{Z}_3 \end{bmatrix} \rightarrow \mathbf{Z}_{(1,1,2)} = \begin{bmatrix} A_1 \\ B_1 \\ C_2 \end{bmatrix}$$





# In Practice: Compute the weights

- Construct the probability tree

$$\hat{w}_{\mathbf{i}_{1:0}} = 1, \quad \hat{w}_{\mathbf{i}_{1:j}} = \hat{w}_{i_j|\mathbf{i}_{1:j-1}} \cdot \hat{w}_{\mathbf{i}_{1:j-1}} \quad (j \in [D], \mathbf{i}_{1:j-1} \in [n]^{j-1}),$$

$$\hat{w}_{i_j|\mathbf{i}_{1:j-1}} := \frac{K^j(\mathbf{Z}_{\mathbf{i}_{1:j-1}}^{\text{mp}(j)} - \mathbf{Z}_i^{\text{mp}(j)})}{\sum_{i=1}^n K^j(\mathbf{Z}_{\mathbf{i}_{1:j-1}}^{\text{mp}(j)} - \mathbf{Z}_i^{\text{mp}(j)})},$$

Kernel Function:

$$K^j(\boldsymbol{x} - \boldsymbol{y}) := \prod_{j' \in \text{mp}(j)} \frac{1}{h^{j'}} K_{j'}^j \left( \frac{\boldsymbol{x}^{j'} - \boldsymbol{y}^{j'}}{h^{j'}} \right)$$

Continuous:  $K_{j'}^j(x - y) := (2\pi)^{-1/2} \exp\left(-\frac{(x-y)^2}{2}\right).$

Discrete:  $K_{j'}^j(x - y) := \mathbb{1}[x = y]$

# Algorithm

---

**Algorithm 1** Proposed method: ADMG data augmentation

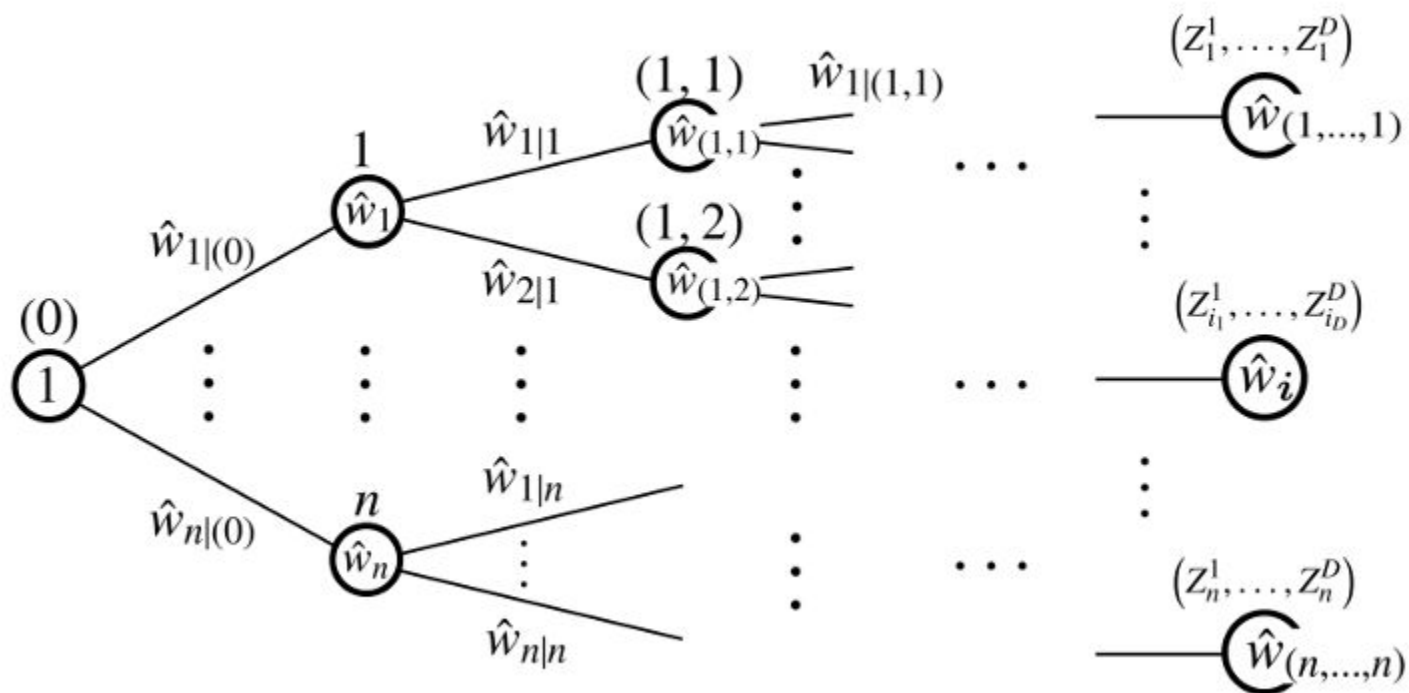
---

**Input:** Training data  $\mathcal{D}$ , ADMG  $\hat{\mathcal{G}}$ , coefficient  $\lambda \in [0, 1]$ , regularization functional  $\Omega$ , pruning threshold  $\theta \in [0, 1)$ , hypothesis class  $\mathcal{F}$ , kernel functions  $\{K^j\}_{j=1}^D$ , loss function  $\ell$ .

```
1: function FILLPROBTREE( $\mathcal{D}, \hat{\mathcal{G}}, \theta, \{K^j\}_{j=1}^D$ )                                ▶ see Fig. 2
2:   for  $j \in [D]$                                                                 ▶ for each variable  $j$ 
3:     for  $i_{1:j-1} \in [n]^{j-1}$                                                 ▶ current node (depth  $j$ )
4:       for  $i_j \in [n]$                                                         ▶ next node (depth  $j + 1$ )
5:          $\hat{w}_{i_{1:j-1}} \leftarrow \hat{w}_{i_{1:j-1}} \mathbb{1}[\hat{w}_{i_{1:j-1}} \geq \theta]$                                 ▶ pruning
6:          $\hat{w}_{i_{1:j}} \leftarrow \hat{w}_{i_j|i_{1:j-1}} \cdot \hat{w}_{i_{1:j-1}}$ 
7:   return  $\mathcal{W}_{\text{aug}} := \{\hat{w}_i\}_{i \in [n]^D}$ 
8: Let  $\mathcal{W}_{\text{aug}} = \text{FILLPROBTREE}(\mathcal{D}, \hat{\mathcal{G}}, \theta, \{K^j\}_{j=1}^D)$ .
9: Let  $\hat{R}_{\text{aug}}(f) := \sum_{i \in [n]^D} \hat{w}_i \cdot \ell(f, \mathbf{Z}_i)$ .
10: Let  $\tilde{R}_\lambda(f) := (1 - \lambda)\hat{R}_{\text{emp}}(f) + \lambda\hat{R}_{\text{aug}}(f) + \Omega(f)$ .
Output:  $\hat{f} \in \arg \min_{f \in \mathcal{F}} \tilde{R}_\lambda(f)$ : the predictor.
```

---

# Algorithm



# Real World Data Experiment

- Goal: Confirm that the proposed method contributes to the performance of the trained predictor, especially in the small-data regime.

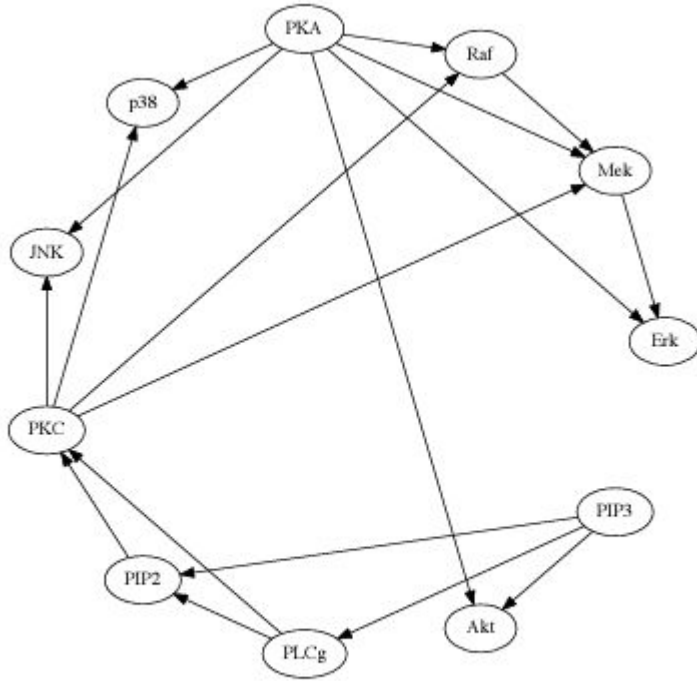
• **Compared:**  $\hat{f} \in \arg \min_{f \in \mathcal{F}} \{\hat{R}_{\text{emp}}(f) + \Omega(f)\}$

$$\hat{f} \in \arg \min_{f \in \mathcal{F}} \{(1 - \lambda)\hat{R}_{\text{emp}}(f) + \lambda\hat{R}_{\text{aug}}(f) + \Omega(f)\}$$

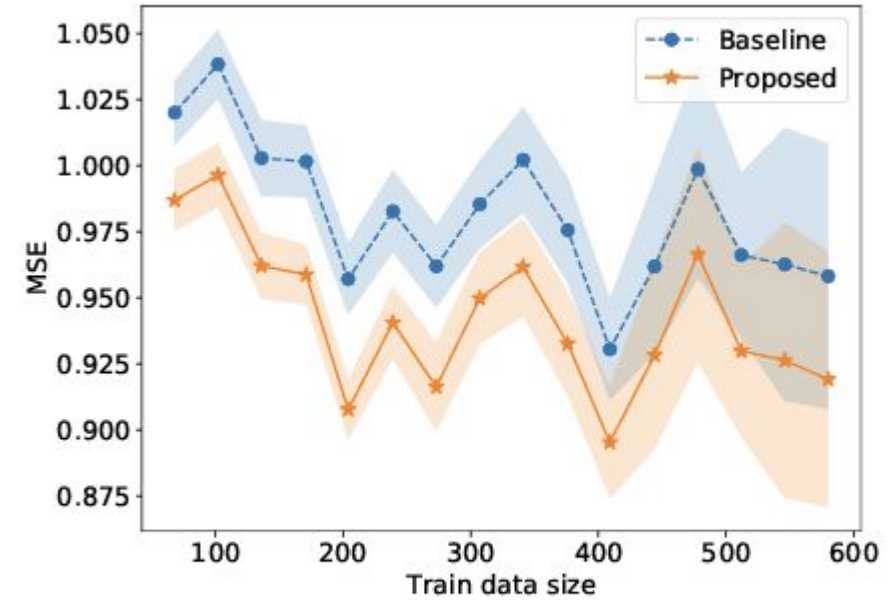
NAME	#VAR	#OBS	GRAPH
<i>Sachs</i>	11	853	Consensus
<i>GSS</i>	6	1380	Domain
<i>Boston Housing</i>	14	506	LiNGAM
<i>Auto MPG</i>	7	392	LiNGAM
<i>White Wine</i>	12	4898	LiNGAM
<i>Red Wine</i>	12	1599	LiNGAM

# Real World Data Experiment

Sachs : Continuous, flow cytometry of proteins and phospholipids in human immune system cells



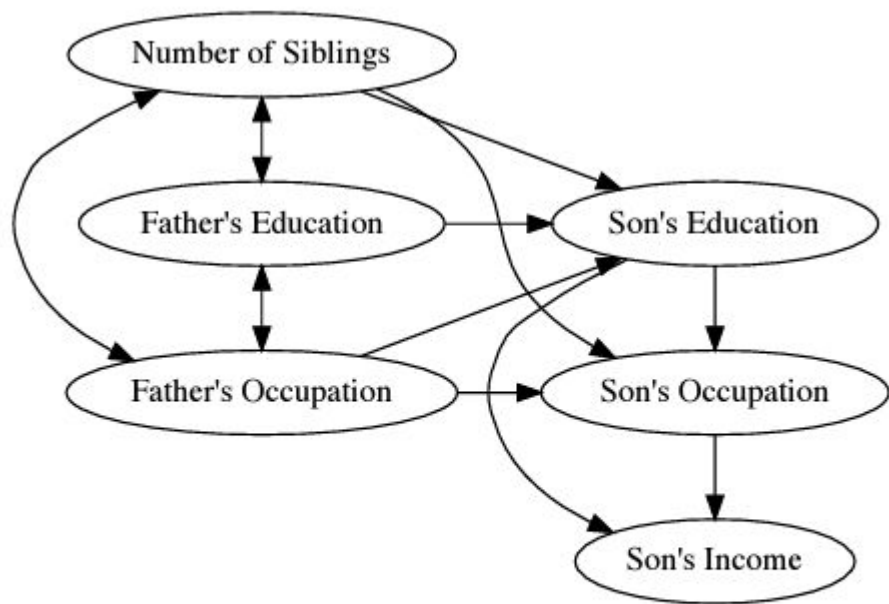
(a) Reference graph for Sachs data.



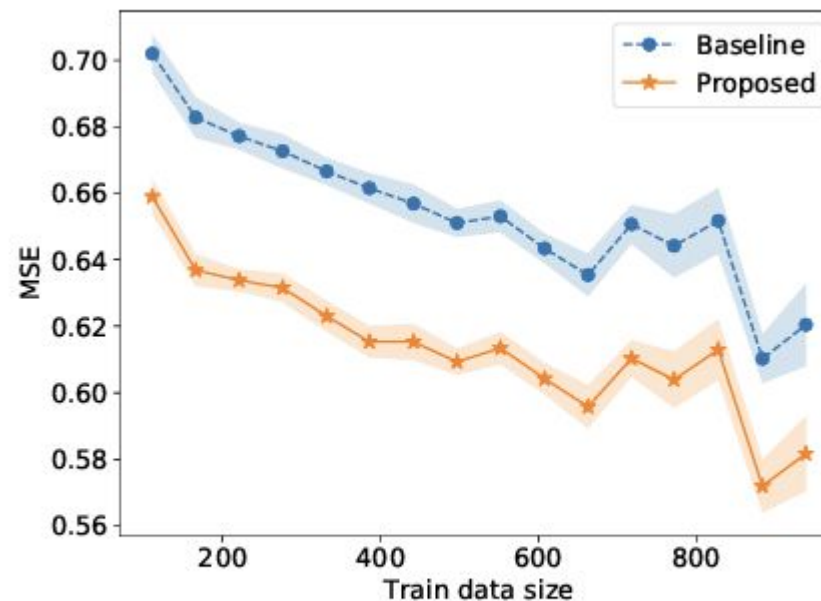
(a) Sachs data.

# Real World Data Experiment

GSS: Part of General Social Survey



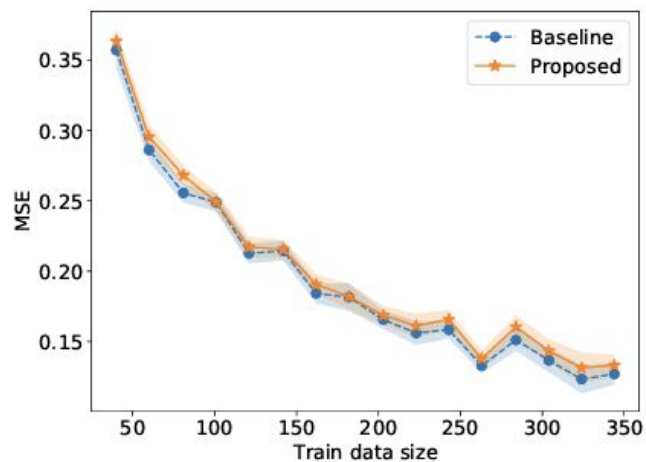
(b) Reference graph for GSS data.



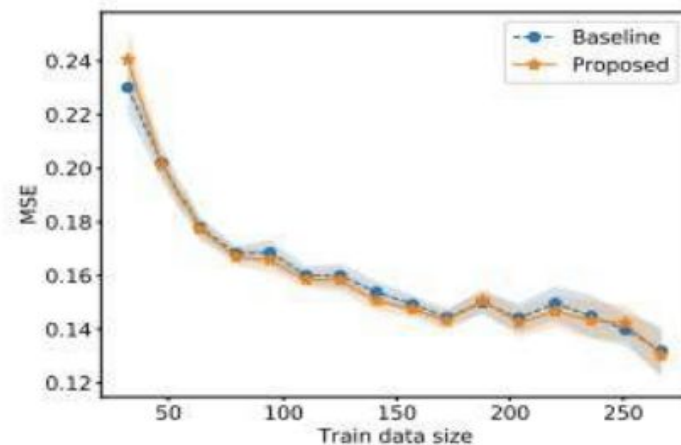
(b) GSS data.



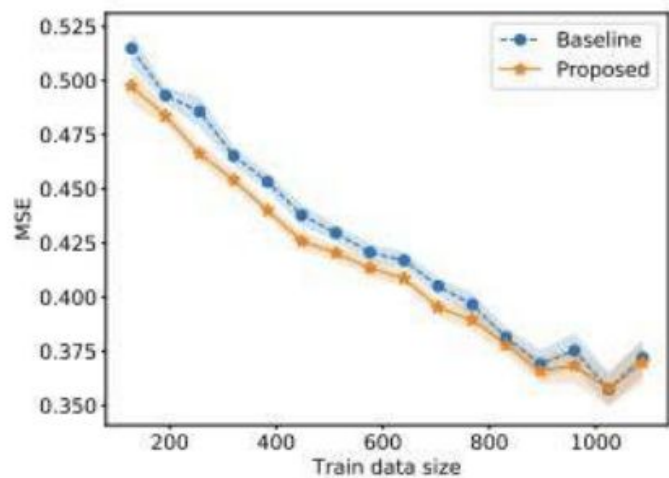
# Real World Data Experiment



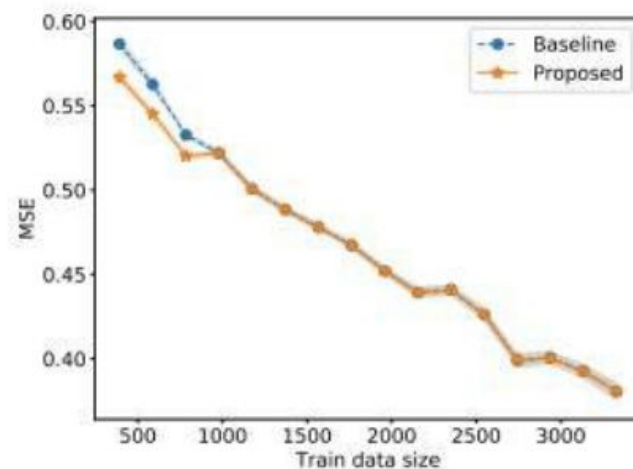
(c) Boston Housing data.



(d) Auto MPG data.



(e) Red Wine data.



(f) White Wine data.

# Conclusion

- Proposed a data augmentation method to use the causal graphical prior knowledge in predictive modeling.
- Experimentally, the benefit may be worth the extra complexity and bias in small-data regime when domain knowledge is available.

# Question

- What is the difference between **markov pillow** and **markov blanket**?

Theorem (Excess Risk Bound; informal)  $\hat{f} \in \arg \min_{f \in \mathcal{F}} \{\hat{R}_{\text{aug}}(f)\}$ ,  $f^* \in \arg \min_{f \in \mathcal{F}} \{R(f)\}$

$$R(\hat{f}) - R(f^*) \leq \underbrace{C_1 R_{\mathbf{H}} + C_p}_{\text{Kernel Bias}} + \underbrace{C_2 R_K + C_3 R_{\mathcal{F}, K}}_{\text{Complexity terms}} + \underbrace{C_4 \sqrt{\frac{\log(4D/\delta)}{2n}}}_{\text{Uncertainty}}$$

w/ high probability.

- The complexity terms have a better sample-size dependency than the usual Rademacher complexity, **implying mitigated overfitting**. (Intuition: Synthesized data  $\Rightarrow$  Reduced possibility of overfitting.)
- But the **bias due to the kernel approximation** is introduced.