

Returning the Favour:
When Regression Benefits from Probabilistic Causal Knowledge



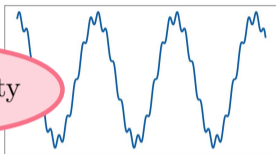
Shahine Bouabid^{*1} Jake Fawkes^{*1} Dino Sejdinovic²

* Equal contribution

¹ Department of Statistics, University of Oxford

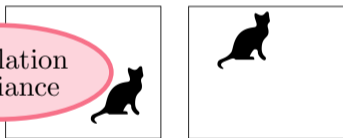
² School of CMS & AIML, University of Adelaide, Adelaide

Periodicity

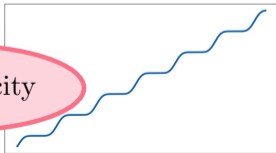


x

Translation invariance

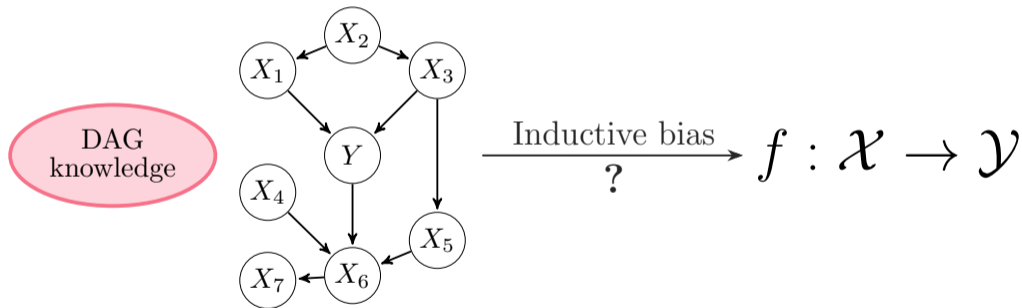


Monotonicity



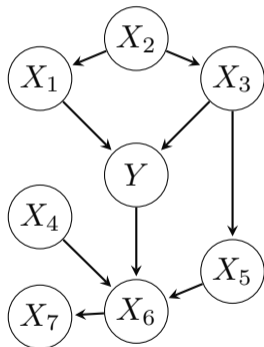
x

Inductive bias $\rightarrow f : \mathcal{X} \rightarrow \mathcal{Y}$



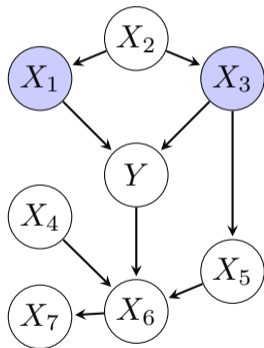
What does a DAG tell us about $\mathbb{P}(Y|X)$?

3



What does a DAG tell us about $\mathbb{P}(Y|X)$?

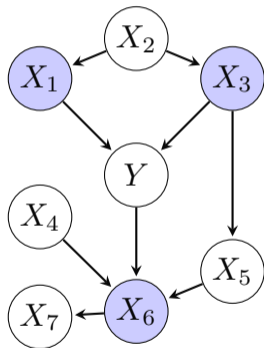
3



Parents
of Y

What does a DAG tell us about $\mathbb{P}(Y|X)$?

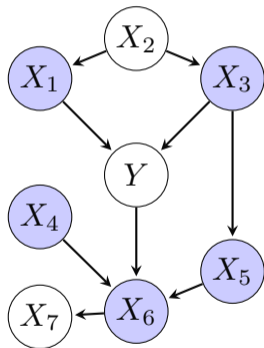
3



Parents of Y \cup Children of Y

What does a DAG tell us about $\mathbb{P}(Y|X)$?

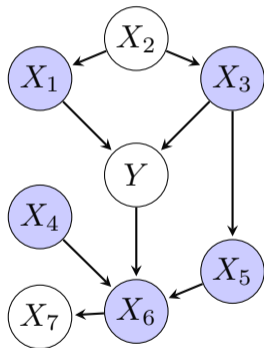
3



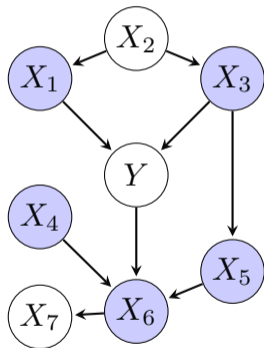
Parents of Y \cup Children of Y \cup Spouses of Y

What does a DAG tell us about $\mathbb{P}(Y|X)$?

3

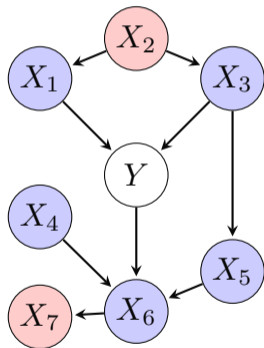


Parents of Y \cup Children of Y \cup Spouses of Y
= Markov Boundary of Y



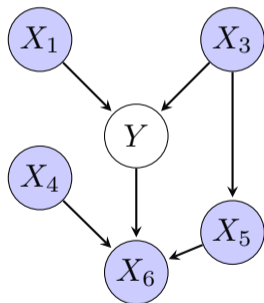
$$\begin{aligned} & \text{Parents of } Y \cup \text{Children of } Y \cup \text{Spouses of } Y \\ & = \text{Markov Boundary of } Y \end{aligned}$$

- The Markov Boundary contains all relevant information for $\mathbb{P}(Y|X)$



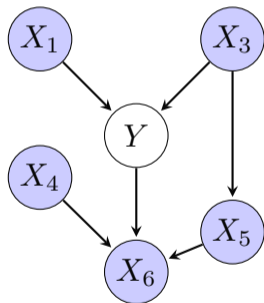
$$\begin{aligned} & \text{Parents of } Y \cup \text{Children of } Y \cup \text{Spouses of } Y \\ & = \text{Markov Boundary of } Y \end{aligned}$$

- ▶ The Markov Boundary contains all relevant information for $\mathbb{P}(Y|X)$
- ▶ X_2 and X_7 do not inform $\mathbb{P}(Y|X)$

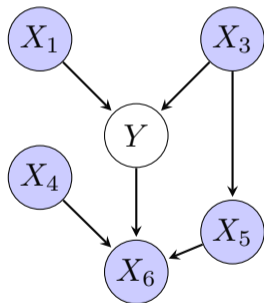


$$\begin{aligned} & \text{Parents of } Y \cup \text{Children of } Y \cup \text{Spouses of } Y \\ & = \text{Markov Boundary of } Y \end{aligned}$$

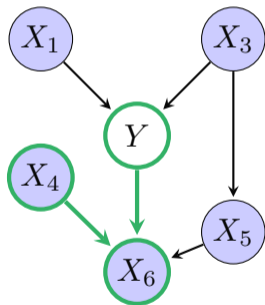
- ▶ The Markov Boundary contains all relevant information for $\mathbb{P}(Y|X)$
- ▶ X_2 and X_7 do not inform $\mathbb{P}(Y|X)$



► Presence of **collider** X_6

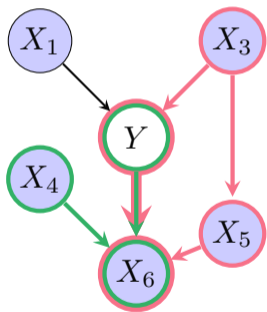


► Presence of **collider** X_6



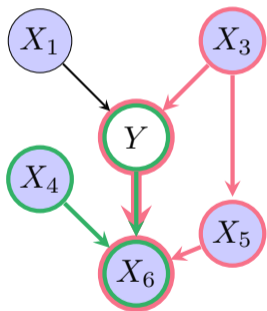
$$Y \perp\!\!\!\perp X_4$$

► Presence of **collider** X_6



$$Y \perp\!\!\!\perp X_4$$

$$Y \perp\!\!\!\perp X_5 \mid X_3$$



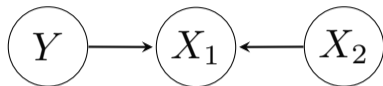
- Presence of **collider** X_6

$$Y \perp\!\!\!\perp X_4$$

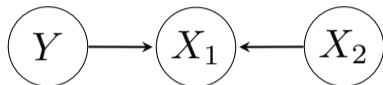
$$Y \perp\!\!\!\perp X_5 \mid X_3$$

- Colliders provide additional information on $\mathbb{P}(Y|X)$ which is unused

Simple collider

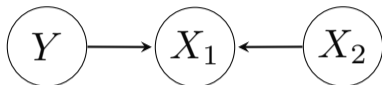


Simple collider



$$Y \perp\!\!\!\perp X_2, \quad Y \not\perp\!\!\!\perp X_2 \mid X_1$$

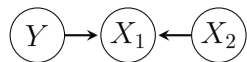
Simple collider



$$Y \perp\!\!\!\perp X_2, \quad Y \not\perp\!\!\!\perp X_2 \mid X_1$$

Squared-loss regression problem

$$\arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_{1i}, x_{2i}))^2 + \lambda \Omega(f)$$





Optimal regressor: $f^*(x_1, x_2) = \mathbb{E}[Y | X_1 = x_1, X_2 = x_2]$



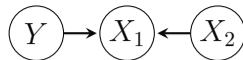
Optimal regressor: $f^*(x_1, x_2) = \mathbb{E}[Y | X_1 = x_1, X_2 = x_2]$

$$\mathbb{E}[f^*(X_1, X_2) | X_2]$$



Optimal regressor: $f^*(x_1, x_2) = \mathbb{E}[Y | X_1 = x_1, X_2 = x_2]$

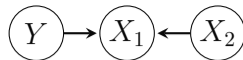
$$\mathbb{E}[f^*(X_1, X_2) | X_2] = \mathbb{E}[\mathbb{E}[Y | X_1, X_2] | X_2]$$



Optimal regressor: $f^*(x_1, x_2) = \mathbb{E}[Y | X_1 = x_1, X_2 = x_2]$

$$\begin{aligned}\mathbb{E}[f^*(X_1, X_2) | X_2] &= \mathbb{E}[\mathbb{E}[Y | X_1, X_2] | X_2] \\ &= \mathbb{E}[Y | X_2]\end{aligned}$$

(Tower property)



Optimal regressor: $f^*(x_1, x_2) = \mathbb{E}[Y | X_1 = x_1, X_2 = x_2]$

$$\begin{aligned} \mathbb{E}[f^*(X_1, X_2) | X_2] &= \mathbb{E}[\mathbb{E}[Y | X_1, X_2] | X_2] \\ &= \mathbb{E}[Y | X_2] && \text{(Tower property)} \\ &= \mathbb{E}[Y] && (Y \perp\!\!\!\perp X_2) \end{aligned}$$



Optimal regressor: $f^*(x_1, x_2) = \mathbb{E}[Y | X_1 = x_1, X_2 = x_2]$

$$\begin{aligned} \mathbb{E}[f^*(X_1, X_2) | X_2] &= \mathbb{E}[\mathbb{E}[Y | X_1, X_2] | X_2] \\ &= \mathbb{E}[Y | X_2] && \text{(Tower property)} \\ &= \mathbb{E}[Y] && (Y \perp\!\!\!\perp X_2) \\ &= 0 && \text{(wlog).} \end{aligned}$$



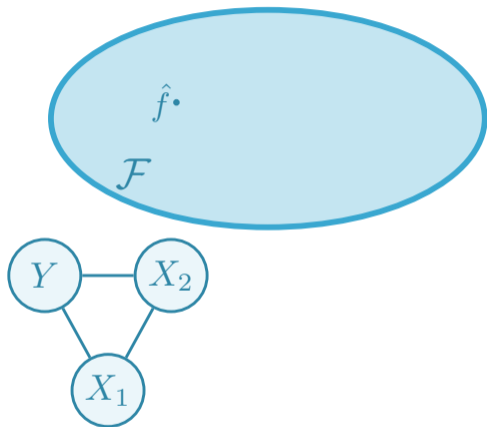
Optimal regressor: $f^*(x_1, x_2) = \mathbb{E}[Y | X_1 = x_1, X_2 = x_2]$

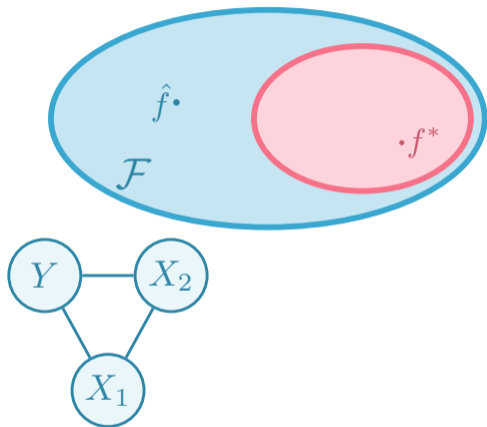
$$\begin{aligned} \mathbb{E}[f^*(X_1, X_2) | X_2] &= \mathbb{E}[\mathbb{E}[Y | X_1, X_2] | X_2] \\ &= \mathbb{E}[Y | X_2] && \text{(Tower property)} \\ &= \mathbb{E}[Y] && (Y \perp\!\!\!\perp X_2) \\ &= 0 && \text{(wlog).} \end{aligned}$$

Objective

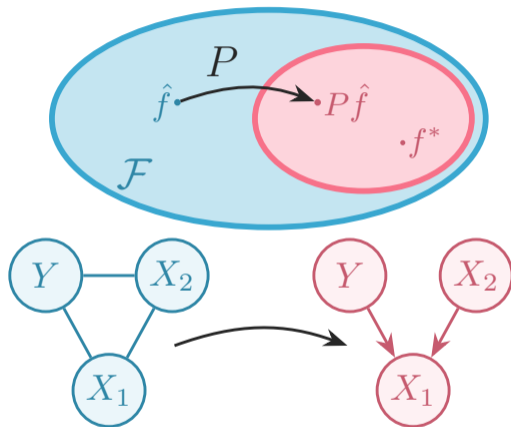
Find a regressor $\hat{f} : \mathcal{X} \rightarrow \mathcal{Y}$ in hypothesis space \mathcal{F} that satisfies

$$\hat{f} \in \{f \in \mathcal{F} \mid \mathbb{E}[f(X_1, X_2) | X_2] = 0\}.$$



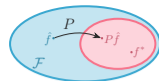


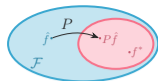
$$\{f \in \mathcal{F} \mid \mathbb{E}[f(X_1, X_2) \mid X_2] = 0\}$$



$$\{f \in \mathcal{F} \mid \mathbb{E}[f(X_1, X_2) \mid X_2] = 0\} = \text{Range}(P)$$

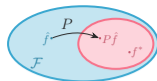
Collider regression in $\mathcal{F} = L^2(X)$





Take $P : L^2(X) \rightarrow L^2(X)$,

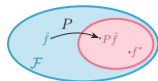
$$Pf(x_1, x_2) = f(x_1, x_2) - \mathbb{E}[f(X_1, X_2) | X_2 = x_2].$$



Take $P : L^2(X) \rightarrow L^2(X)$,

$$Pf(x_1, x_2) = f(x_1, x_2) - \mathbb{E}[f(X_1, X_2) | X_2 = x_2].$$

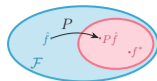
► P is an orthogonal projection



Take $P : L^2(X) \rightarrow L^2(X)$,

$$Pf(x_1, x_2) = f(x_1, x_2) - \mathbb{E}[f(X_1, X_2) | X_2 = x_2].$$

- ▶ P is an orthogonal projection
- ▶ $f^* \in \text{Range}(P) = \{f \in \mathcal{F} \mid \mathbb{E}[f(X_1, X_2) \mid X_2] = 0\}$



Take $P : L^2(X) \rightarrow L^2(X)$,

$$Pf(x_1, x_2) = f(x_1, x_2) - \mathbb{E}[f(X_1, X_2) | X_2 = x_2].$$

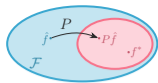
- ▶ P is an orthogonal projection
- ▶ $f^* \in \text{Range}(P) = \{f \in \mathcal{F} \mid \mathbb{E}[f(X_1, X_2) \mid X_2] = 0\}$

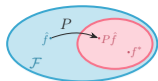
Proposition

Let $h \in L^2(X)$, then we have

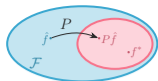
$$\Delta(h, Ph) = \mathbb{E}[(Y - h(X))^2] - \mathbb{E}[(Y - Ph(X))^2] = \|(\text{Id} - P)h\|_{L^2(X)}^2.$$

Collider regression in a RKHS $\mathcal{F} = \mathcal{H}$

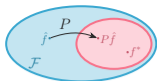




- ▶ Reproducing kernel Hilbert space with kernel $k(x, x')$



- ▶ Reproducing kernel Hilbert space with kernel $k(x, x')$
- ▶ Convenient mathematical properties and dense in $L^2(X)$ under mild assumptions



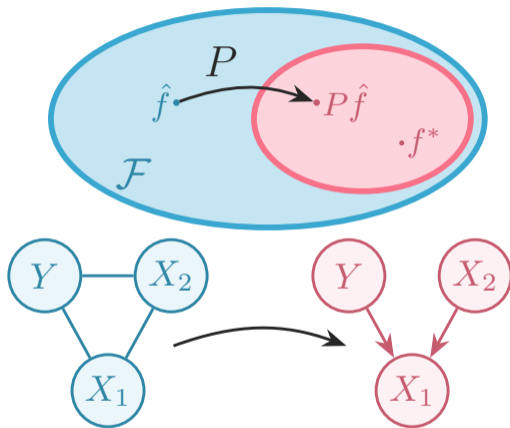
- ▶ Reproducing kernel Hilbert space with kernel $k(x, x')$
- ▶ Convenient mathematical properties and dense in $L^2(X)$ under mild assumptions

Theorem

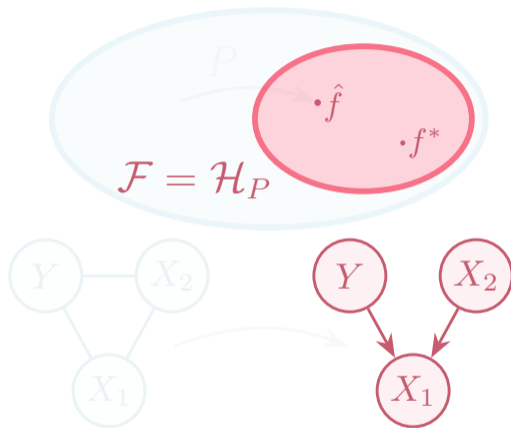
Assume $M = \sup_{x \in \mathcal{X}} k(x, x) < \infty$ and $\text{Var}(Y|X) \geq \eta > 0$. Then, the generalisation gap between \hat{f} and $P\hat{f}$ satisfies

$$\mathbb{E}[\Delta(\hat{f}, P\hat{f})] \geq \frac{\eta \mathbb{E}[\|\mu_{X|X_2}(X)\|_{L^2(X)}^2]}{(\sqrt{n}M + \lambda/\sqrt{n})^2} \quad (1)$$

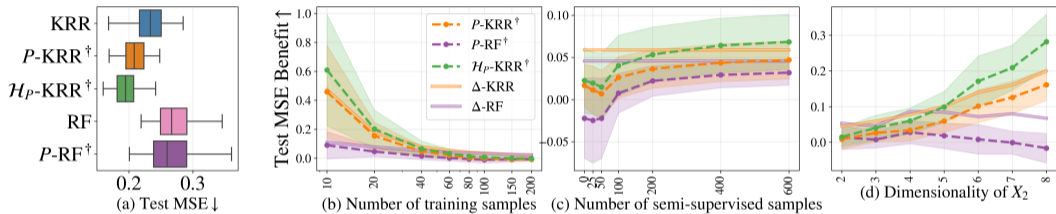
where $\mu_{X|X_2} = \mathbb{E}[k(X, \cdot)|X_2]$ is a RKHS representation of $\mathbb{P}(X|X_2)$.



$\mathcal{F} = \mathcal{H}$ with kernel $k(x, x') = \langle k(x, \cdot), k(x', \cdot) \rangle$



$$\mathcal{F} = \mathcal{H}_P \text{ with kernel } k_P(x, x') = \langle P^*k(x, \cdot), P^*k(x', \cdot) \rangle$$



(a) : Test MSEs for the simulation experiment ; dataset is generated using $d_1 = 3$, $d_2 = 3$, $n = 50$ and 100 semi-supervised samples ; experiments is run for 100 datasets generated with different seeds ; statistical significance is confirmed via Wilcoxon signed-rank ; (b, c, d) : Ablation study on the number of training samples, number of semi-supervised samples and dimensionality of X_2 .

Table 1: MSE, signal-to-noise ratio (SNR) and correlation on test data for the aerosol radiative forcing experiment ; $n = 50$ and 200 semi-supervised samples ; statistical significance is confirmed via Wilcoxon signed-rank ; experiments are run for 100 datasets generated by FaIR with different seeds ; \uparrow/\downarrow indicates higher/lower is better; we report 1 standard deviation; \dagger indicates our proposed methods.

	MSE \downarrow	SNR \uparrow	Correlation \uparrow
RF	0.90 \pm 0.04	0.44 \pm 0.19	0.32 \pm 0.08
P -RF \dagger	0.89\pm0.03	0.49\pm0.15	0.34\pm0.07
KRR	0.88 \pm 0.04	0.56 \pm 0.21	0.37 \pm 0.05
P -KRR \dagger	0.86\pm0.03	0.65\pm0.14	0.40\pm0.02
\mathcal{H}_P -KRR \dagger	0.86\pm0.03	0.64\pm0.14	0.39 \pm 0.03

- ▶ Collider structures within causal graphs constitute a useful form of inductive bias within supervised learning

**To find out more come to our poster:
Poster Session 4, Wednesday 2:00-3:30, Exhibit hall 1.**

- ▶ Collider structures within causal graphs constitute a useful form of inductive bias within supervised learning
- ▶ Provable benefit of semi-supervised learning arising from causal structure

**To find out more come to our poster:
Poster Session 4, Wednesday 2:00-3:30, Exhibit hall 1.**

- ▶ Collider structures within causal graphs constitute a useful form of inductive bias within supervised learning
- ▶ Provable benefit of semi-supervised learning arising from causal structure
- ▶ Collider regression for more general DAGs in the paper

**To find out more come to our poster:
Poster Session 4, Wednesday 2:00-3:30, Exhibit hall 1.**

- ▶ Collider structures within causal graphs constitute a useful form of inductive bias within supervised learning
- ▶ Provable benefit of semi-supervised learning arising from causal structure
- ▶ Collider regression for more general DAGs in the paper
- ▶ Similar reasoning can be applied to other forms of inference about $\mathbb{P}(Y|X)$, e.g. classification or quantile regression

**To find out more come to our poster:
Poster Session 4, Wednesday 2:00-3:30, Exhibit hall 1.**