

Distribution-free learning of Bayesian network structure

Xiaohai Sun

Max Planck Institute for Biological Cybernetics, Tübingen, Germany

September 16, 2008



MAX-PLANCK-GESELLSCHAFT



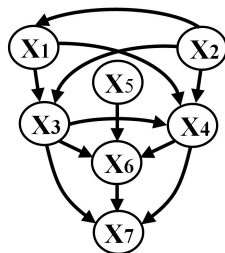
BIOLOGISCHE KYBERNETIK

Learning causal Bayesian network structure

- ▶ Purely **observational** (non-experimental) data on a number of random variables are given
- ▶ Estimate **direct cause-effect** relations between these variables, represented by a directed acyclic graph (causal BN structure)

	X_1	X_2	\dots	X_N
(1)	0	45.2	...	12
(2)	1	1.7	...	62
(3)	1	-3.1	...	90
...
(n)	0	19.1	...	22

\Rightarrow



Constraint-based learning

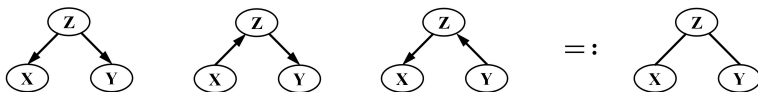
- ▶ Based on **independence** relations

- ▶ Weak commitments as to the nature of causal relationships
 1. **Markov** assumption states “given all its parents, every variable is independent of all its non-descendants”.

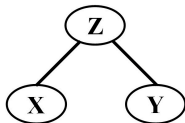
 2. **Faithfulness/stability** assumption (Spirtes et al. 1993; Pearl 2000) states “*only* the independence relations are true which are implied by the Markov assumption”.

From independence constraints to causal BN structure

Identifying \wedge -structure by independence constraints



$(X \perp\!\!\!\perp Y)$ and $(X \perp\!\!\!\perp Y | Z) \Rightarrow \wedge$ -structure



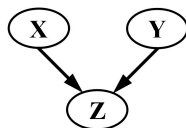
Markov equivalence class of BN structures

\hookrightarrow Learning **absence** of edges in causal BN structure

V-structure identification

Markov and faithfulness assumptions lead to a unique graph.

$$(X \perp\!\!\!\perp Y) \text{ and } (X \not\perp\!\!\!\perp Y \mid Z) \Rightarrow \text{V-structure}$$



Identification of V-structure (collider on Z)

\hookrightarrow Learning **orientation** of edges in causal BN structure

Inductive causation

- ▶ Finding conditional independence relations
- ▶ Taking Markov and faithfulness assumptions
 1. Learning absence of edges \leftrightarrow skeleton of BN structure
 2. Learning orientation of edges

\leftrightarrow Inductive causation (IC) algorithm (Pearl 2000)

Refinement: PC algorithm (Spirtes et al. 1993)

Using **correlation** analysis (assumption of normal distribution)

- ▶ Our goal:
non-parametric test of independence on arbitrary domains

Embedding of distributions in RKHS

- ▶ $\mathcal{H}_{\mathcal{X}}$: Hilbert space on measurable space \mathcal{X} , spanned by functions $k_{\mathcal{X}}(x, \cdot)$ ($x \in \mathcal{X}$) with $\langle k_{\mathcal{X}}(x, \cdot), k_{\mathcal{X}}(x', \cdot) \rangle = k_{\mathcal{X}}(x, x') \forall x, x' \in \mathcal{X}$.
 X : random variable on \mathcal{X} .
- ▶ Mean element in RKHS:
 $\mathfrak{M}_{\mathcal{X}} = \mathbb{E}[k_{\mathcal{X}}(X, \cdot)]$ and $\mathfrak{M}_{\mathcal{X}\mathcal{Y}} = \mathbb{E}[k_{\mathcal{X}}(X, \cdot)k_{\mathcal{Y}}(Y, \cdot)]$
- ▶ **Conditional** mean element in RKHS:
 $\mathfrak{M}_{\mathcal{X}|Y} = \mathbb{E}[k_{\mathcal{X}}(X, \cdot)|Y]$ and $\mathfrak{M}_{\mathcal{X}\mathcal{Y}|Z} = \mathbb{E}[k_{\mathcal{X}}(X, \cdot)k_{\mathcal{Y}}(Y, \cdot)|Z]$
- ▶ **Product** of mean elements in RKHS:
 $\mathfrak{M}_{\mathcal{X}}\mathfrak{M}_{\mathcal{Y}} = \mathfrak{M}_{\mathcal{X}} \otimes \mathfrak{M}_{\mathcal{Y}} = \mathbb{E}[k_{\mathcal{X}}(X, \cdot)]\mathbb{E}[k_{\mathcal{Y}}(Y, \cdot)]$
- ▶ **Product** of **conditional** mean elements in RKHS:
 $\mathfrak{M}_{\mathcal{X}|Z}\mathfrak{M}_{\mathcal{Y}|Z} = \mathfrak{M}_{\mathcal{X}|Z} \otimes \mathfrak{M}_{\mathcal{Y}|Z} = \mathbb{E}[k_{\mathcal{X}}(X, \cdot)|Z]\mathbb{E}[k_{\mathcal{Y}}(Y, \cdot)|Z]$

Cross-covariance operator

- Cross-covariance operator in RKHS:

$$\begin{aligned}
 \langle g, \Sigma_{YX} f \rangle_{\mathcal{H}_Y} &:= \langle \mathfrak{M}_{XY} - \mathfrak{M}_X \mathfrak{M}_Y, f \otimes g \rangle_{\mathcal{H}_X \otimes \mathcal{H}_Y} \\
 &= \mathbb{E}[f(X)g(Y)] - \mathbb{E}[f(X)]\mathbb{E}[g(Y)] \\
 &= \text{Cov}[f(X), g(Y)] \quad \forall f \in \mathcal{H}_X, g \in \mathcal{H}_Y
 \end{aligned}$$

- **Conditional** cross-covariance operator in RKHS:

$$\begin{aligned}
 \langle g, \Sigma_{YX|Z} f \rangle_{\mathcal{H}_Y} &:= \langle \mathfrak{M}_{XY} - \mathbb{E}_Z[\mathfrak{M}_{X|Z} \mathfrak{M}_{Y|Z}], f \otimes g \rangle_{\mathcal{H}_X \otimes \mathcal{H}_Y} \\
 &= \mathbb{E}_{XY}[f(X)g(Y)] - \mathbb{E}_Z[\mathbb{E}[f(X)|Z]\mathbb{E}[g(Y)|Z]] \\
 &= \mathbb{E}_Z[\text{Cov}[f(X), g(Y) | Z]] \quad \forall f \in \mathcal{H}_X, g \in \mathcal{H}_Y
 \end{aligned}$$

HS norm of operator and MMD

- Hilbert-Schmidt (HS) norm of operator $\Sigma: \mathcal{H}_X \rightarrow \mathcal{H}_Y$:

$$\|\Sigma\|_{\text{HS}}^2 = \text{Tr}(\Sigma^T \Sigma) = \sum_{i,j=1}^{\infty} \langle \varphi_j, \Sigma \phi_i \rangle_{\mathcal{H}_Y}^2,$$

$\{\phi_i\}_{i=1}^{\infty}, \{\varphi_j\}_{j=1}^{\infty}$: complete orthonormal systems of $\mathcal{H}_X, \mathcal{H}_Y$.

- kernel Maximum Mean Discrepancy (MMD)

(Borgwardt et al. Bioinformatics 2006)

$$\mathbb{D}_{\mathcal{H}}(\mathcal{P}, \mathcal{Q}) = \sup_{f \in \mathcal{F}} \mathbb{E}_{x \sim \mathcal{P}}[f(x)] - \mathbb{E}_{y \sim \mathcal{Q}}[f(y)].$$

\mathcal{P}, \mathcal{Q} : probability measures. \mathcal{F} : unit ball in RKHS \mathcal{H} .

$\hookrightarrow \mathcal{P} = \mathcal{Q} \iff \mathbb{D}_{\mathcal{H}}(\mathcal{P}, \mathcal{Q}) = 0$ (\mathcal{H} : characteristic RKHS).

(Fukumizu et al. NIPS 2007, 2008; Sriperumbudur et al. COLT 2008)

Unconditional independence with kernel

- ▶ HS norm of cross-covariance operator Σ_{YX} corresponds to MMD between P_{xy} and $P_x P_y$

$$\begin{aligned} \|\Sigma_{YX}\|_{\text{HS}}^2 &= \langle \mathfrak{M}_{XY} - \mathfrak{M}_X \mathfrak{M}_Y, \mathfrak{M}_{XY} - \mathfrak{M}_X \mathfrak{M}_Y \rangle \\ &= \|\mathfrak{M}_{XY} - \mathfrak{M}_X \mathfrak{M}_Y\|_{\mathcal{H}_X \otimes \mathcal{H}_Y}^2 \\ &= \mathbb{D}_{\mathcal{H}_X \otimes \mathcal{H}_Y}^2(P_{xy}, P_x P_y). \end{aligned}$$

- ▶ Given characteristic RKHS

$$\begin{aligned} \Sigma_{YX} = 0 &\iff \|\Sigma_{YX}\|_{\text{HS}}^2 = 0 \\ &\iff \mathfrak{M}_{XY} = \mathfrak{M}_X \mathfrak{M}_Y \\ &\iff P_{xy} = P_x P_y \\ &\iff X \perp\!\!\!\perp Y. \end{aligned}$$

Conditional cross-covariance operator and MMD

- ▶ HS norm of conditional cross-covariance operator $\Sigma_{YX|Z}$ corresponds to MMD between P_{xy} and $E_Z[P_{x|z}P_{y|z}]$

$$\begin{aligned} \|\Sigma_{YX|Z}\|_{\text{HS}}^2 &= \langle \mathfrak{M}_{XY} - E_Z[\mathfrak{M}_{X|Z}\mathfrak{M}_{Y|Z}], \mathfrak{M}_{XY} - E_Z[\mathfrak{M}_{X|Z}\mathfrak{M}_{Y|Z}] \rangle \\ &= \|\mathfrak{M}_{XY} - E_Z[\mathfrak{M}_{X|Z}\mathfrak{M}_{Y|Z}]\|_{\mathcal{H}_X \otimes \mathcal{H}_Y}^2 \\ &= \mathbb{D}_{\mathcal{H}_X \otimes \mathcal{H}_Y}^2(P_{xy}, E_Z[P_{x|z}P_{y|z}]) . \end{aligned}$$

- ▶ Given characteristic RKHS

$$\begin{aligned} \Sigma_{YX|Z} = O &\iff \|\Sigma_{YX|Z}\|_{\text{HS}}^2 = 0 \\ &\iff \mathfrak{M}_{XY} = E_Z[\mathfrak{M}_{X|Z}\mathfrak{M}_{Y|Z}] \\ &\iff P_{xy} = E_Z[P_{x|z}P_{y|z}] \\ &\iff X \perp\!\!\!\perp Y | Z . \end{aligned}$$

Conditional independence with kernel

Define $\dot{X} := (X, Z)$, $\dot{Y} := (Y, Z)$

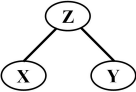
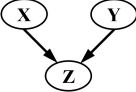
- ▶ HS norm / MMD

$$\begin{aligned} \|\Sigma_{\dot{Y}\dot{X}|Z}\|_{\text{HS}}^2 &= \mathbb{E}_Z \left[\|\mathfrak{M}_{\dot{X}\dot{Y}|Z} - \mathfrak{M}_{\dot{X}|Z} \mathfrak{M}_{\dot{Y}|Z}\|_{\mathcal{H}_{\dot{X}} \otimes \mathcal{H}_{\dot{Y}}}^2 \right] \\ &= \mathbb{E}_Z \left[\mathbb{D}_{\mathcal{H}_{\dot{X}} \otimes \mathcal{H}_{\dot{Y}}}^2 (P_{\dot{X}\dot{Y}|Z}, P_{\dot{X}|Z} P_{\dot{Y}|Z}) \right]. \end{aligned}$$

- ▶ Given characteristic RKHS

$$\begin{aligned} \Sigma_{\dot{Y}\dot{X}|Z} = O &\iff \|\Sigma_{\dot{Y}\dot{X}|Z}\|_{\text{HS}}^2 = 0 \\ &\iff \mathfrak{M}_{XY|Z} = \mathbb{E}_Z[\mathfrak{M}_{X|Z} \mathfrak{M}_{Y|Z}] \quad \forall Z \\ &\iff P_{xy|z} = \mathbb{E}_z[P_{x|z} P_{y|z}] \quad \forall z \\ &\iff X \perp\!\!\!\perp Y | Z. \end{aligned}$$

Constraint-based learning of BN structure

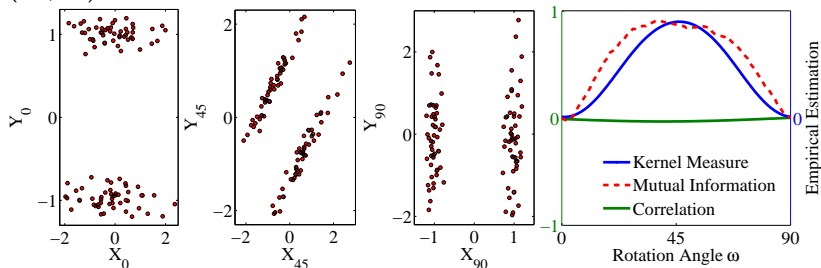
Constraints	\wedge -Structure 	\vee -Structure 
Independence relations	$X \not\perp\!\!\!\perp Y$ and $X \perp\!\!\!\perp Y Z$	$X \perp\!\!\!\perp Y$ and $X \not\perp\!\!\!\perp Y Z$
Joint distributions	$P_{xy} \neq P_x P_y$ $P_{xy z} = P_{x z} P_{y z}$	$P_{xy} = P_x P_y$ $P_{xy z} \neq P_{x z} P_{y z}$
Mean element in RKHS	$\mathfrak{M}_{XY} \neq \mathfrak{M}_X \mathfrak{M}_Y$ $\mathfrak{M}_{\dot{X}\dot{Y}} = \mathbb{E}_Z [\mathfrak{M}_{\dot{X} Z} \mathfrak{M}_{\dot{Y} Z}]$	$\mathfrak{M}_{XY} = \mathfrak{M}_X \mathfrak{M}_Y$ $\mathfrak{M}_{\dot{X}\dot{Y}} \neq \mathbb{E}_Z [\mathfrak{M}_{\dot{X} Z} \mathfrak{M}_{\dot{Y} Z}]$
HS norm of operators (MMD in RKHS)	$\ \Sigma_{XY}\ _{\text{HS}}^2 > 0$ $\ \Sigma_{\dot{X}\dot{Y} Z}\ _{\text{HS}}^2 = 0$	$\ \Sigma_{XY}\ _{\text{HS}}^2 = 0$ $\ \Sigma_{\dot{X}\dot{Y} Z}\ _{\text{HS}}^2 > 0$

Alternative measures of dependences

- ▶ Correlation coefficient (assumption of normal distribution)
- ▶ Mutual information (Kraskov et al. 2004)
(based on entropy estimates from k -nearest neighbor distances)

$$Y_0 \propto \frac{1}{2}\mathcal{N}(1, 0.01) + \frac{1}{2}\mathcal{N}(-1, 0.01), P(X_0|Y_0 < 0) \propto \mathcal{N}(0, 1), P(X_0|Y_0 \geq 0) \propto \mathcal{N}(0, 1)$$

(X_ω, Y_ω) : transformed data with rotation ω in an anticlockwise direction



Summary

- ▶ Kernel test of independence for constraint-based learning of causal BN structure

Further issues:

- ▶ Connection to mutual information? (Gretton et al. 2005)
- ▶ Useful MMD with higher-order tensors?
e.g., vanishing of

$$\|\mathfrak{M}_{XYZ} - \mathfrak{M}_X \mathfrak{M}_Y \mathfrak{M}_Z\|_{\mathcal{H}_X \otimes \mathcal{H}_Y \otimes \mathcal{H}_Z}^2 = \mathbb{D}_{\mathcal{H}_X \otimes \mathcal{H}_Y}^2(P_{xyz}, P_X P_Y P_Z),$$

indicates mutual independence (more than pairwise independence).

Thanks for your attention!