

# Non-Parametric Structure Learning on Hidden Tree-Shaped Distributions

<sup>1</sup>Konstantinos E. Nikolakakis, <sup>2</sup>Dionysios S. Kalogieras,  
<sup>1</sup>Anand D. Sarwate

<sup>1</sup>Department of Electrical & Computer Engineering, Rutgers University

<sup>2</sup>Department of Electrical & Systems Engineering, University of Pennsylvania



# Why learn from noisy data?

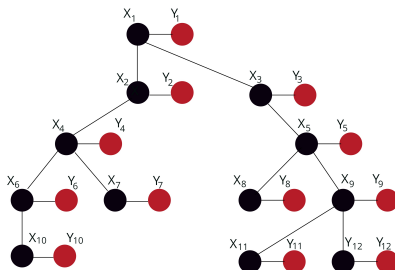
- Data acquisition devices or sensors introduce noise
- Local differential privacy
- Communication constrains and quantization error
- Adversarial attacks



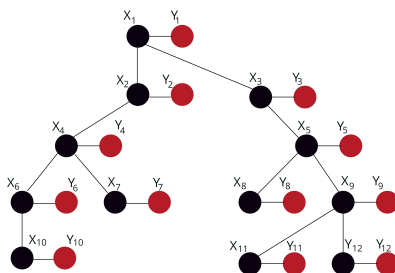
# Problem Statement (Learning Hidden Tree Structures)

Observe noisy values for each node of the unknown tree structure

- $X_1, X_2, \dots, X_p$  are hidden variables (black nodes)
- $Y_1, Y_2, \dots, Y_p$  are observable variables (red nodes)



# Learning a tree structure



## Assumptions:

- Distribution of  $\mathbf{X}$  is nondegenerate and factorizes according to a tree  $T$ .
- $T = (\mathcal{V}, \mathcal{E})$  is connected.
- $I(X_i; X_j) > 0$  for all  $i, j \in \mathcal{V}$ .



# Chow-Liu Algorithm

**Given:** Data set  $\mathcal{D} = \mathbf{Y} \in \mathcal{Y}^{|\mathcal{V}| \times n}$

- 1 Compute empirical distribution on each edge:

$$\hat{p}_{i,j}(\ell, m) = \frac{1}{n} \sum_{k=1}^n \mathbf{1}_{\{Y_{i,k}=\ell, Y_{j,k}=m\}} \quad \forall i, j \in \mathcal{V}$$

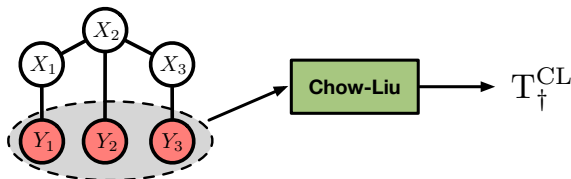
- 2 Find plug-in estimate of mutual information:

$$\hat{I}(Y_i; Y_j) = \sum_{\ell, m} \hat{p}_{i,j}(\ell, m) \log_2 \frac{\hat{p}_{i,j}(\ell, m)}{\hat{p}_i(\ell) \hat{p}_j(m)}$$

- 3 Output  $T_{\dagger}^{\text{CL}} = \text{MST} \left( \{ \hat{I}(Z_i; Z_j) : i, j \in \mathcal{V} \} \right)$



# Main questions



Given noise corrupted data:

- Is Chow-Liu consistent?
- How does noise affect the sample complexity?

Prior work: Finite sample complexity for Ising and Gaussian Models.

- Tan et al. (2011), Liu et al. (2011), Bresler & Karzand (2018)
- Hidden models: Our work (2019), Goel-Kane-Klivans (2019)



# A motivating example: 3-node hidden model

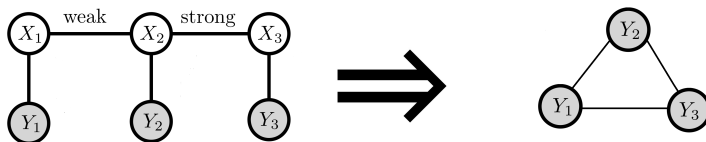
What can go wrong when we have noise?



# A motivating example: 3-node hidden model

What can go wrong when we have noise?

The MRF of the observable is a complete graph!

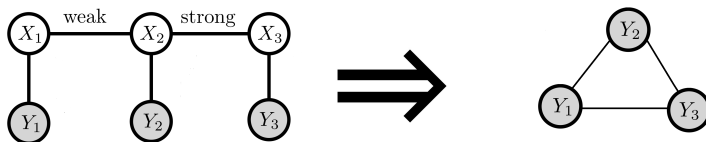




# A motivating example: 3-node hidden model

What can go wrong when we have noise?

The MRF of the observable is a complete graph!

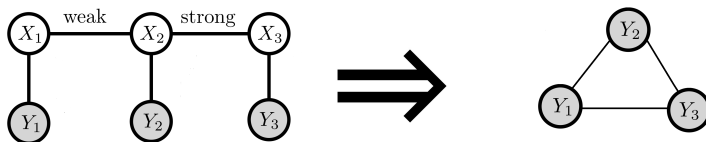


**Questions:**

# A motivating example: 3-node hidden model

What can go wrong when we have noise?

The MRF of the observable is a complete graph!



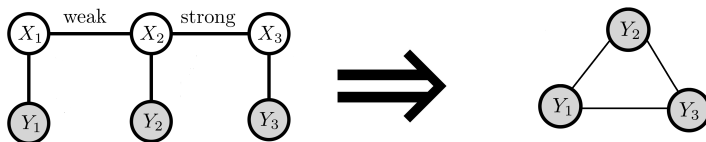
**Questions:**

- Is Chow-Liu consistent? **NO**

# A motivating example: 3-node hidden model

What can go wrong when we have noise?

The MRF of the observable is a complete graph!



## Questions:

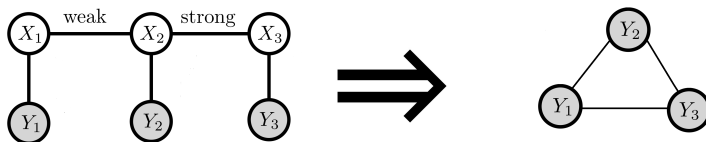
- Is Chow-Liu consistent? **NO**
- When does  $\lim_{n \rightarrow \infty} T^{\text{CL}} \rightarrow T$  w.p. 1? **A sufficient condition**



# A motivating example: 3-node hidden model

What can go wrong when we have noise?

The MRF of the observable is a complete graph!



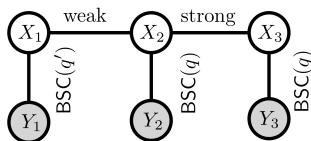
## Questions:

- Is Chow-Liu consistent? **NO**
- When does  $\lim_{n \rightarrow \infty} T^{\text{CL}} \rightarrow T$  w.p. 1? **A sufficient condition**
- Can we tweak Chow-Liu to fix it? **Sometimes**



# A closer look at the example

$$X_1, X_2, X_3 \in \{-1, +1\}, \quad 0 < \underbrace{|\mathbb{E}[X_1 X_2]|}_{\text{weak}} \leq \underbrace{|\mathbb{E}[X_2 X_3]|}_{\text{strong}} < 1$$



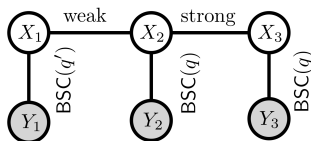
$$I(X_2; X_3) > I(X_1; X_2) \stackrel{\text{DPI}}{>} I(X_1; X_3)$$

$$\lim_{n \rightarrow \infty} \hat{I}(X_i; X_j) \rightarrow I(X_i; X_j) \text{ and } \lim_{n \rightarrow \infty} \mathcal{E}_{\text{TCL}} \rightarrow \{(1, 2), (2, 3)\} \equiv \mathcal{E}_{\text{T}}$$

- Does a similar condition hold for the observables?
- Could we have  $\lim_{n \rightarrow \infty} \mathcal{E}_{\text{T}^\dagger_{\text{CL}}} \neq \mathcal{E}_{\text{T}}$ ?



# Feasibility Threshold



- If  $I(Y_1; Y_2) > I(Y_1; Y_3) > I(Y_2; Y_3)$
- then  $\lim_{n \rightarrow \infty} \mathcal{E}_{T_{\dagger}^{cl}} \neq \mathcal{E}_T$



$$I(Y_1; Y_3) > I(Y_2; Y_3) \iff |\mathbb{E}[Y_1 Y_3]| > |\mathbb{E}[Y_2 Y_3]| \iff$$

$$|\mathbb{E}[X_1 X_2]| > \frac{1 - 2q}{1 - 2q'}, \quad q, q' \in [0, 1/2).$$



# Unprocessed vs Processed Data

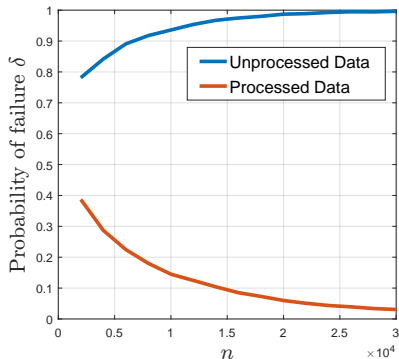
- What if  $|\mathbb{E}[X_1 X_2]| > (1 - 2q)/(1 - 2q')$ ?
- We have to pre-process

$$Z_1 \triangleq Y_1/(1 - 2q'), \quad Z_2 \triangleq Y_2/(1 - 2q), \quad Z_3 \triangleq Y_3/(1 - 2q)$$

- Correct order,  $I(Z_2; Z_3) > I(Z_1; Z_2) > I(Z_1; Z_3)$
- Then  $\lim_{n \rightarrow \infty} \mathcal{E}_{\text{T}^\dagger \text{CL}} = \mathcal{E}_{\text{T}}$  with probability 1.



# Unprocessed vs Processed Data



*Figure: Synthetic data,  $q' = 0.2$ ,  $q = 0.25$*



## Definition

**(The set  $\mathcal{EV}^2$ )** Let  $e \equiv (w, \bar{w}) \in \mathcal{E}_T$  be an edge and  $u, \bar{u} \in \mathcal{V}_T$  be a pair of nodes such that  $e \in \text{path}_T(u, \bar{u})$  and  $|\text{path}_T(u, \bar{u})| \geq 2$ . Then

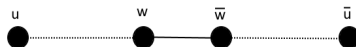
$$\mathcal{EV}^2 \triangleq \{(w, \bar{w}), u, \bar{u} \in \mathcal{E}_T \times \mathcal{V}_T \times \mathcal{V}_T : \\ (w, \bar{w}) \in \text{path}_T(u, \bar{u}) \text{ and } |\text{path}_T(u, \bar{u})| \geq 2\}.$$



## Definition

(**The set  $\mathcal{EV}^2$** ) Let  $e \equiv (w, \bar{w}) \in \mathcal{E}_T$  be an edge and  $u, \bar{u} \in \mathcal{V}_T$  be a pair of nodes such that  $e \in \text{path}_T(u, \bar{u})$  and  $|\text{path}_T(u, \bar{u})| \geq 2$ . Then

$$\mathcal{EV}^2 \triangleq \{(w, \bar{w}), u, \bar{u} \in \mathcal{E}_T \times \mathcal{V}_T \times \mathcal{V}_T : \\ (w, \bar{w}) \in \text{path}_T(u, \bar{u}) \text{ and } |\text{path}_T(u, \bar{u})| \geq 2\}.$$



Error Characterization of CL algorithm (Bresler & Karzand 2018):

$$\text{If } T_{\dagger}^{\text{CL}} \neq T \implies \exists ((w, \bar{w}), u, \bar{u}) \in \mathcal{EV}^2 : \hat{I}(Y_w; Y_{\bar{w}}) \leq \hat{I}(Y_u; Y_{\bar{u}})$$



# Sufficient Condition for Exact Recovery

Exact recovery:

$$\text{If } \forall ((w, \bar{w}), u, \bar{u}) \in \mathcal{EV}^2 : \hat{I}(Y_w; Y_{\bar{w}}) > \hat{I}(Y_u; Y_{\bar{u}}) \implies T_{\dagger}^{\text{CL}} = T$$

$$\begin{aligned} & \hat{I}(Y_w; Y_{\bar{w}}) > \hat{I}(Y_u; Y_{\bar{u}}) \iff \\ & I(Y_w; Y_{\bar{w}}) - I(Y_u; Y_{\bar{u}}) > \\ & \quad \left[ \hat{I}(Y_u; Y_{\bar{u}}) - I(Y_u; Y_{\bar{u}}) \right] - \left[ \hat{I}(Y_w; Y_{\bar{w}}) - I(Y_w; Y_{\bar{w}}) \right]. \end{aligned}$$



# Sufficient Condition for Exact Recovery

Exact recovery:

$$\forall ((w, \bar{w}), u, \bar{u}) \in \mathcal{EV}^2 : \hat{I}(Y_w; Y_{\bar{w}}) > \hat{I}(Y_u; Y_{\bar{u}}) \iff T_{\dagger}^{\text{CL}} = T$$

$$\begin{aligned} & \hat{I}(Y_w; Y_{\bar{w}}) > \hat{I}(Y_u; Y_{\bar{u}}) \iff \\ & I(Y_w; Y_{\bar{w}}) - I(Y_u; Y_{\bar{u}}) > \\ & \left[ \hat{I}(Y_u; Y_{\bar{u}}) - I(Y_u; Y_{\bar{u}}) \right] - \left[ \hat{I}(Y_w; Y_{\bar{w}}) - I(Y_w; Y_{\bar{w}}) \right]. \end{aligned}$$

## Sufficient Condition

$$\text{If } \left| \hat{I}(Y_{\ell}; Y_{\bar{\ell}}) - I(Y_{\ell}; Y_{\bar{\ell}}) \right| < \frac{1}{2} \min_{(e, u, \bar{u}) \in \mathcal{EV}^2} \{I(Y_w; Y_{\bar{w}}) - I(Y_u; Y_{\bar{u}})\}$$

for all  $\ell, \ell' \in \mathcal{V}$  then  $T_{\dagger}^{\text{CL}} = T$ .



## Definition

### (Information Thresholds $\mathbf{I}^o$ , $\mathbf{I}_\dagger^o$ )

$$\mathbf{I}^o \triangleq \frac{1}{2} \min_{((w, \bar{w}), u, \bar{u}) \in \mathcal{E}\mathcal{V}^2} [I(X_w; X_{\bar{w}}) - I(X_u; X_{\bar{u}})]$$

$$\mathbf{I}_\dagger^o \triangleq \frac{1}{2} \min_{((w, \bar{w}), u, \bar{u}) \in \mathcal{E}\mathcal{V}^2} [I(Y_w; Y_{\bar{w}}) - I(Y_u; Y_{\bar{u}})]$$

- Always  $\mathbf{I}^o \geq 0$ , DPI
- $\mathbf{I}_\dagger^o \leq 0$  generalizes the condition  $\frac{1-2q}{1-2q'} \leq |\mathbb{E}[X_1 X_2]|$  to non-parametric models and general channels
- $\mathbf{I}_\dagger^o < 0$  implies that structure learning is infeasible without post-processing



# Sample Complexity

- Sufficient condition  $\left| \hat{I}(Y_\ell; Y_{\bar{\ell}}) - I(Y_\ell; Y_{\bar{\ell}}) \right| < \mathbf{I}_\dagger^o$
- Concentration of measure of mutual information estimates
- Union bound over the pairs  $\ell, \ell' \in \mathcal{V}$

## Theorem

Fix  $\delta \in (0, 1)$ . There exist constants  $C > 0$  and  $c \in (1, 2]$  independent of  $\delta$  such that, if  $\mathbf{I}_\dagger^o > 0$  and

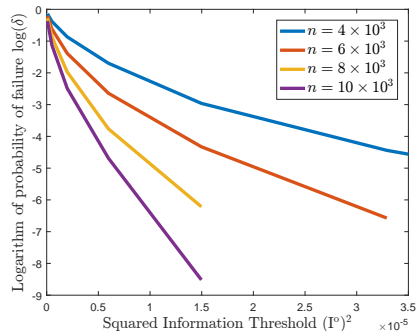
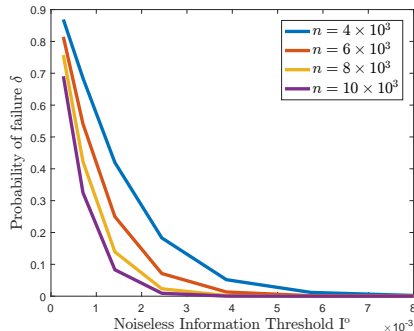
$$\frac{n}{\log_2^2 n} \geq \frac{72 \log\left(\frac{p}{\delta}\right)}{\left(\mathbf{I}_\dagger^o - Cn^{\frac{1-c}{c}}\right)^2} \quad \text{and} \quad \mathbf{I}_\dagger^o > Cn^{\frac{1-c}{c}},$$

then CL with input  $\mathcal{D} = \mathbf{Y}^{1:n}$  returns  $\mathbf{T}_\dagger^{CL} = \mathbf{T}$  w.p. at least  $1 - \delta$ .

Almost logarithmic order:  $\mathcal{O}(\log^{1+\zeta}(p/\delta))$ , for all  $\zeta > 0$



# Experiments: Noiseless Binary Data



*Figure: Left:  $\hat{\mathbb{P}}(T^{\text{CL}} \neq T)$  vs  $I^0$ , Right:  $\log \hat{\mathbb{P}}(T^{\text{CL}} \neq T)$  vs  $(I^0)^2$*



# Experiments: Noisy Binary Data (BSC)

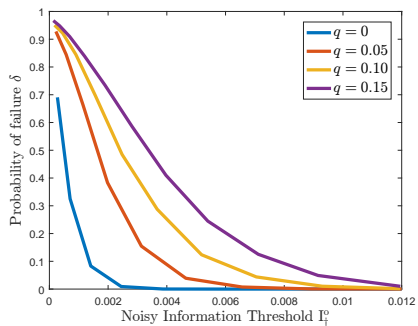


Figure:  $\hat{\mathbb{P}} \left( T_{\dagger}^{\text{CL}} \neq T \right)$  vs  $I_{\dagger}^o$



# Further Questions and Future Directions

- What is the relationship of  $\mathbf{I}^o$  and  $\mathbf{I}_{\dagger}^o$ ? Connection with SDPI
- How to estimate  $\mathbf{I}_{\dagger}^o$  from training data?
- How to preserve privacy while structure learning remains feasible?
- Find robust methods for pre-processing against adversarial attacks



# Thank you!

