

Non-Parametric Structure Learning on Hidden Tree-Shaped Distributions

¹Konstantinos E. Nikolakakis, ²Dionysios S. Kalogerias,
¹Anand D. Sarwate

¹Department of Electrical & Computer Engineering, Rutgers University

²Department 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

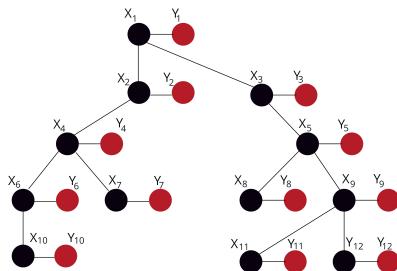


Problem Statement (Learning Hidden Tree Structures)

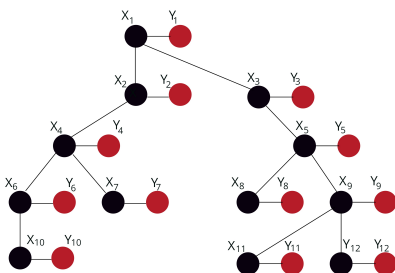
Observe noisy values for each node of the unknown tree structure

$X_1; X_2; \dots; X_p$ 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 .

$T = (V; E)$ is connected.

$I(X_i; X_j) > 0$ for all $i; j \in V$.



Chow-Liu Algorithm

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

- 1 Compute empirical distribution on each edge:

$$\hat{p}_{i,j}(\cdot; m) = \frac{1}{n} \sum_{k=1}^n \mathbf{1}_{f_{Y_i;k}=\cdot; Y_j;k=m} \quad \forall i,j \in V$$

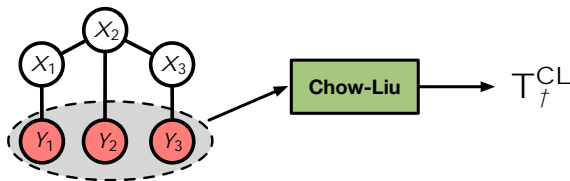
- 2 Find plug-in estimate of mutual information:

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

- 3 Output $T_y^{\text{CL}} = \text{MST} \ f\hat{I}(Z_i; Z_j) : i,j \in V$



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

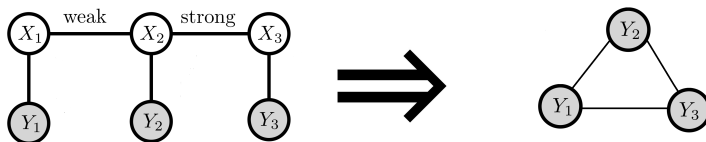
What can go wrong when we have noise?



A motivating example: 3-node hidden model

What can go wrong when we have noise?

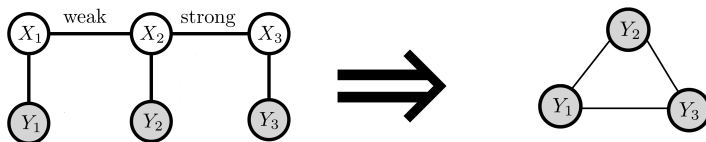
The MRF of the observable is a complete graph!



A motivating example: 3-node hidden model

What can go wrong when we have noise?

The MRF of the observable is a complete graph!

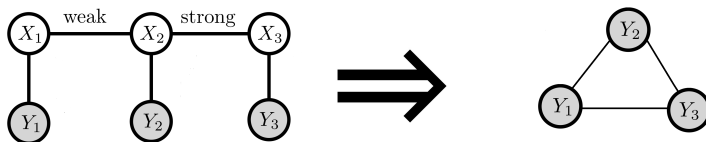


Questions:

A motivating example: 3-node hidden model

What can go wrong when we have noise?

The MRF of the observable is a complete graph!



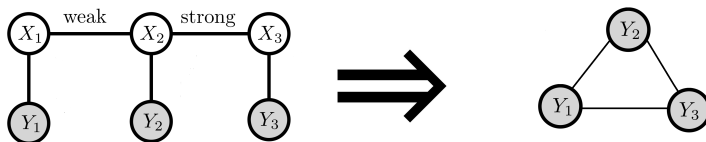
Questions:

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**

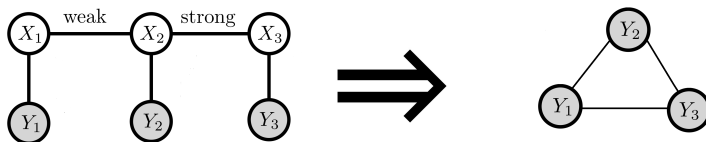
When does $\lim_{n \rightarrow \infty} T^{\text{CL}} = T$ w.p. 1? **A sufficient condition**



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**

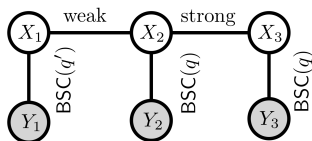
When does $\lim_{n \rightarrow \infty} T^{\text{CL}} = T$ w.p. 1? **A sufficient condition**

Can we tweak Chow-Liu to fix it? **Sometimes**



A closer look at the example

$$X_1; X_2; X_3 \stackrel{2}{\neq} f(1; +1g); \quad 0 < \underbrace{\sqrt{\frac{E[X_1 X_2]}{Z}}}_{\text{weak}} < \underbrace{\sqrt{\frac{E[X_2 X_3]}{Z}}}_{\text{strong}} < 1$$



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

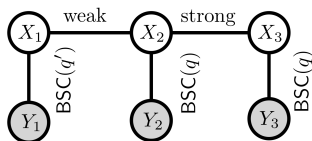
$$\lim_{n \rightarrow \infty} \hat{I}(X_i; X_j) \neq I(X_i; X_j) \text{ and } \lim_{n \rightarrow \infty} E_{\text{TCL}} \neq f(1;2);(2;3)g \quad E_{\text{T}}$$

Does a similar condition hold for the observables?

Could we have $\lim_{n \rightarrow \infty} E_{\text{TCL}} \not\subseteq E_{\text{T}}$?



Feasibility Threshold



If $I(Y_1; Y_2) > I(Y_1; Y_3) > I(Y_2; Y_3)$

then $\lim_{n \rightarrow \infty} E_{T_y^{CL}} \notin E_T$



$I(Y_1; Y_3) > I(Y_2; Y_3) \iff jE[Y_1 Y_3]j > jE[Y_2 Y_3]j \iff$

$$jE[X_1 X_2]j > \frac{1 - 2q}{1 - 2q^0}; \quad q; q^0 \in [0; 1=2):$$



Unprocessed vs Processed Data

What if $jE[X_1 X_2]j > (1 - 2q) = (1 - 2q^0)$?

We have to pre-process

$$Z_1, Y_1 = (1 - 2q^0); Z_2, Y_2 = (1 - 2q); Z_3, Y_3 = (1 - 2q)$$

Correct order, $I(Z_2; Z_3) > I(Z_1; Z_2) > I(Z_1; Z_3)$

Then $\lim_{n \rightarrow \infty} E_{T_y^{CL}} = E_T$ with probability 1.



Unprocessed vs Processed Data

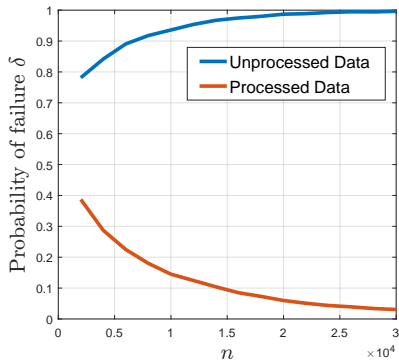


Figure: Synthetic data, $q^0 = 0:2$, $q = 0:25$



Definition

(The set \mathcal{EV}^2) Let $e = (w; \bar{w}) \in E_T$ be an edge and $u; \bar{u} \in V_T$ be a pair of nodes such that $e \in \text{path}_T(u; \bar{u})$ and $|\text{path}_T(u; \bar{u})| = 2$. Then

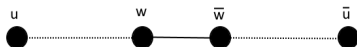
$$EV^2 = \{ (w; \bar{w}); u; \bar{u} \in E_T \times V_T \times V_T : \\ (w; \bar{w}) \in \text{path}_T(u; \bar{u}) \text{ and } |\text{path}_T(u; \bar{u})| = 2 \}$$



Definition

(The set \mathcal{EV}^2) Let $e = (w; \bar{w}) \in E_T$ be an edge and $u; \bar{u} \in V_T$ 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

$$\mathcal{EV}^2 = \{ (w; \bar{w}); u; \bar{u} \in E_T \times V_T \times V_T : \\ (w; \bar{w}) \in \text{path}_T(u; \bar{u}) \text{ and } |\text{path}_T(u; \bar{u})| \geq 2 \}$$



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

$$\text{If } T_y^{\text{CL}} \notin T \Rightarrow \exists ((w; \bar{w}); u; \bar{u}) \in \mathcal{EV}^2 : \hat{T}(Y_w; Y_{\bar{w}}) \neq \hat{T}(Y_u; Y_{\bar{u}})$$



Sufficient Condition for Exact Recovery

Exact recovery:

$$\text{If } \delta((w; w); u; u) \geq EV^2 : \hat{I}(Y_w; Y_{\bar{w}}) > \hat{I}(Y_u; Y_{\bar{u}}) \Rightarrow T_y^{\text{CL}} = T$$

$$\hat{I}(Y_w; Y_{\bar{w}}) > \hat{I}(Y_u; Y_{\bar{u}}) \quad ()$$

$$I(Y_w; Y_{\bar{w}}) \underset{h}{>} I(Y_u; Y_{\bar{u}}) > \underset{i}{I(Y_u; Y_{\bar{u}})} \underset{h}{I(Y_w; Y_{\bar{w}})} \underset{i}{I(Y_w; Y_{\bar{w}})} :$$



Sufficient Condition for Exact Recovery

Exact recovery:

$$\mathcal{S}((w; \bar{w}); (u; \bar{u})) \supseteq EV^2 : \hat{I}(Y_w; Y_{\bar{w}}) > \hat{I}(Y_u; Y_{\bar{u}}) \quad () \quad T_y^{\text{CL}} = T$$

$$\hat{I}(Y_w; Y_{\bar{w}}) > \hat{I}(Y_u; Y_{\bar{u}}) \quad ()$$

$$I(Y_w; Y_{\bar{w}}) - I(Y_u; Y_{\bar{u}}) > \frac{1}{2} \min_{(e; u; \bar{u}) \in EV^2} \left(I(Y_w; Y_{\bar{w}}) - I(Y_u; Y_{\bar{u}}) \right)$$

Sufficient Condition

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

for all $\epsilon > 0 \supseteq V$ then $T_y^{\text{CL}} = T$:



Definition

(Information Thresholds I^o , I_y^o)

$$I^o = \frac{1}{2} \min_{((w;\bar{w});u;\bar{u}) \in \mathcal{E}V^2} [I(X_{w_i}; X_{\bar{w}}) - I(X_{u_i}; X_{\bar{u}})]$$

$$I_y^o = \frac{1}{2} \min_{((w;\bar{w});u;\bar{u}) \in \mathcal{E}V^2} [I(Y_{w_i}; Y_{\bar{w}}) - I(Y_{u_i}; Y_{\bar{u}})]$$

Always $I^o \geq 0$, DPI

$I_y^o \geq 0$ generalizes the condition $\frac{1}{1-2q^0} \geq j \in [X_1 X_2] j$ to non-parametric models and general channels

$I_y^o < 0$ implies that structure learning is infeasible without post-processing



Sample Complexity

Sufficient condition $\hat{I}(Y; Y^c) - I(Y; Y^c) < \mathbf{I}_y^o$

Concentration of measure of mutual information estimates

Union bound over the pairs $y; y^c \in \mathcal{V}$

Theorem

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

$$\frac{n}{\log_2^2 n} \geq \frac{72 \log \epsilon}{\mathbf{I}_y^o C n^{\frac{1-c}{c}} \epsilon^2} \text{ and } \mathbf{I}_y^o > C n^{\frac{1-c}{c}} \epsilon;$$

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

Almost logarithmic order: $O(\log^{1+}(\frac{1}{\epsilon}))$, for all $\epsilon > 0$



Experiments: Noiseless Binary Data

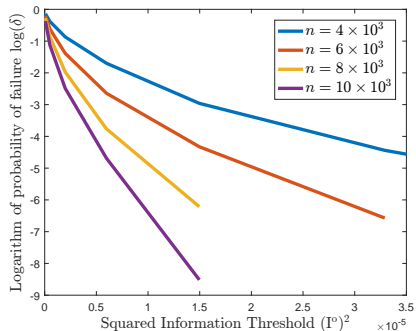
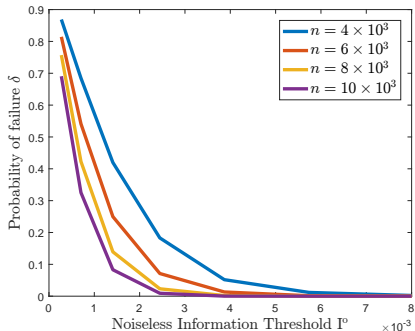


Figure: Left: $\hat{P} T^{\text{CL}} \notin T$ vs I^0 , Right: $\log \hat{P} T^{\text{CL}} \notin T$ vs $(I^0)^2$



Experiments: Noisy Binary Data (BSC)

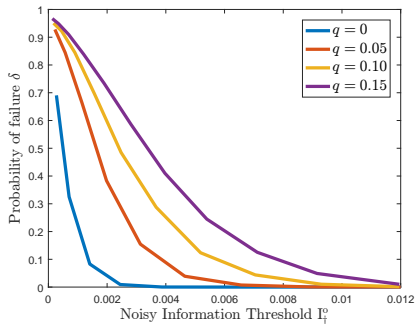


Figure: $\hat{P} T_y^{\text{CL}} \notin T$ vs I_+^o



Further Questions and Future Directions

What is the relationship of I^o and I_y^o ? Connection with SDPI

How to estimate I_y^o from training data?

How to preserve privacy while structure learning remains feasible?

Find robust methods for pre-processing against adversarial attacks



Thank you!

