

## Introduction

**Goal:** Minimum number of samples required to learn the graph from data.

### Examples

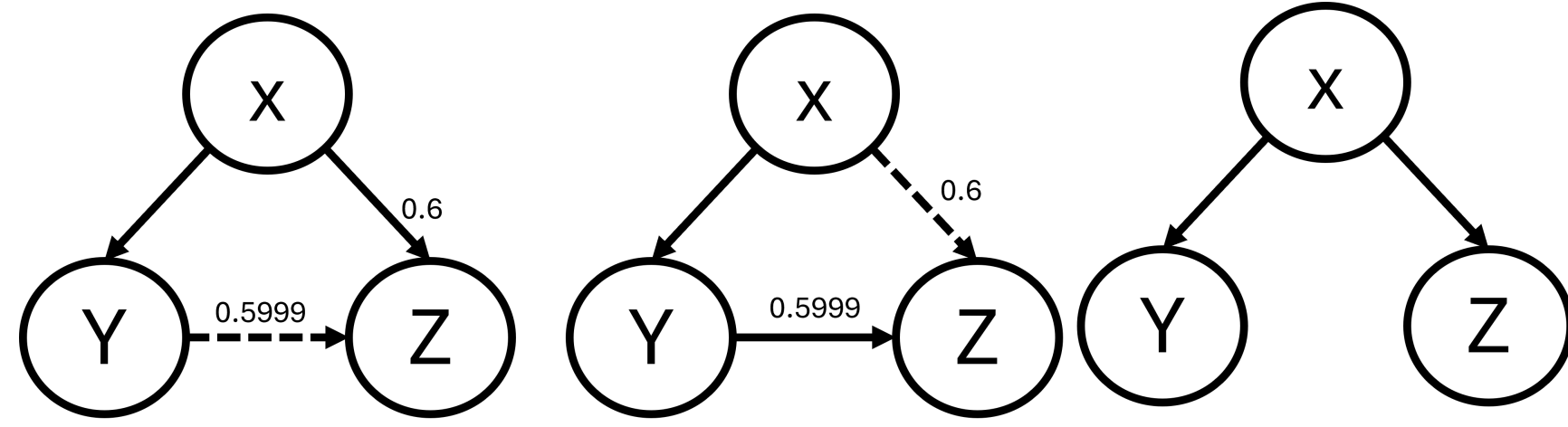


Figure 1:(a) and (b) Distribution learning and (c) Structure learning

## Questions

- (Non-realizable setting)**  $P$  might not representable by any tree, how many samples are required to learn a tree-structured distribution  $Q$ ?
- (Realizable setting)**  $P$  itself is tree-structured, how many samples are required to learn a tree-structured distribution  $Q$ ?
- (Faithful setting)**  $P$  is faithful to some tree  $T$ , how many samples are required to learn  $T$  up to Markov equivalence?

## Bayesian networks and Tree-Faithfulness

**Distribution learning:** For a distribution  $P$  and a directed tree  $T$ , let

$$P_T := \arg \min_{T\text{-structured distribution } Q} D_{\text{KL}}(P||Q),$$

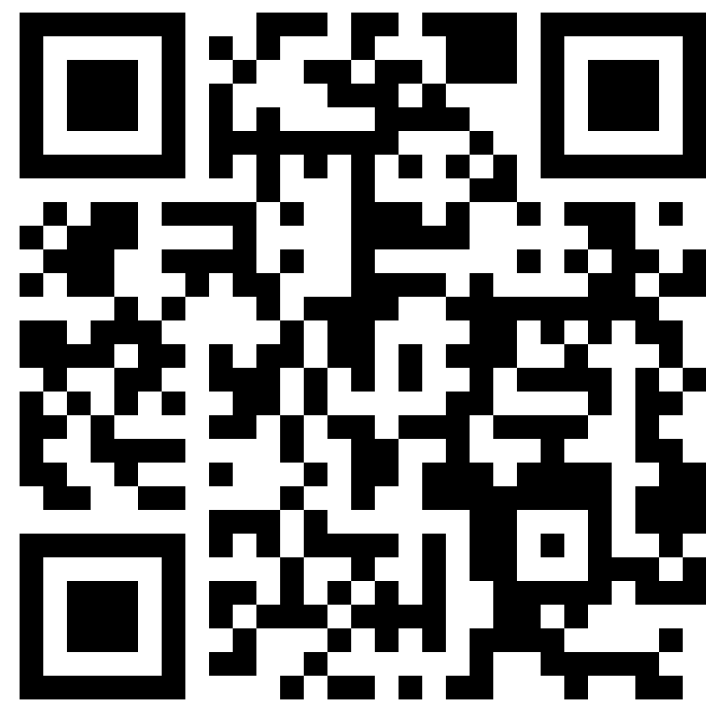
where  $D_{\text{KL}}\{\cdot||\cdot\}$  denotes the KL-divergence.

**Definition 1 : [Tree-faithfulness]** We say distribution  $P$  is tree-faithful to a polytree  $T$  if

- For any two nodes connected  $X_j - X_k$ , we have  $X_k \perp\!\!\!\perp X_j | X_\ell$  for all  $\ell \in V \cup \{\emptyset\} \setminus \{k, j\}$ ;
- For any  $v$ -structure  $X_k \rightarrow X_\ell \leftarrow X_j$ , we have  $X_k \perp\!\!\!\perp X_j | X_\ell$ .

**Definition 2 : [ $c$ -strong tree-faithfulness]** We say that  $P$  is  $c$ -strong tree-faithful to a polytree  $T$  if

- For any two nodes connected  $X_j - X_k$ , we have  $\rho(X_k, X_j | X_\ell) \geq c$  for  $\ell \in V \cup \{\emptyset\} \setminus \{k, j\}$ ;
- For any  $v$ -structure  $X_k \rightarrow X_\ell \leftarrow X_j$ , we have  $\rho(X_k, X_j | X_\ell) \geq c$ .



## Our Contribution

### Non-realizable Setting

Without making additional assumptions on  $P$ , we show that

$$n = \tilde{\Theta}\left(\frac{d^2}{\varepsilon^2}\right) \quad (1)$$

samples are necessary and sufficient to learn (with probability at least  $2/3$ ) a tree-structured distribution that is  $\varepsilon$ -close to the closest tree-structured distribution for  $P$ .

**Realizable Setting** When  $P$  itself is Markov to a tree  $T$  (i.e. it is *tree-structured*), then

$$n = \tilde{\Theta}\left(\frac{d}{\varepsilon}\right) \quad (2)$$

samples are necessary and sufficient to learn (with probability at least  $2/3$ ) a tree-structured distribution that is  $\varepsilon$ -close to  $P$  itself.

**Faithful Polytrees** Assuming that  $P$  is faithful to some *polytree*  $T$ , we show that the optimal sample complexity of learning  $\bar{T}$ , the CPDAG of  $T$ , is

$$n = \Theta\left(\frac{\log d}{c^2}\right), \quad (3)$$

where  $c$  is the strong faithfulness parameter.

## Learning Tree-Structured Gaussians

**Theorem 1 :** [Non-realizable setting] Let  $P$  be a Gaussian distribution. Given  $n$  i.i.d. samples from  $P$ , for any  $\varepsilon, \delta > 0$ , if  $n \gtrsim \frac{d^2}{\varepsilon^2} \log \frac{d}{\delta}$ , then  $\hat{T}$  returned by Algorithm 1 satisfies

$$D_{\text{KL}}(P||P_{\hat{T}}) \leq \min_{T \in \mathcal{T}} D_{\text{KL}}(P||P_T) + \varepsilon,$$

with probability at least  $1 - \delta$ . Besides, if  $n = o(d^2/\varepsilon^2)$ , no algorithm returns a directed tree  $\hat{T}$  such that

$$D_{\text{KL}}(P||P_{\hat{T}}) \leq \min_{T \in \mathcal{T}} D_{\text{KL}}(P||P_T) + \varepsilon$$

with probability at least  $2/3$ .

**Theorem 2 :** [Realizable setting] Let  $T^*$  be a directed tree and  $P_{T^*}$  be a  $T^*$ -structured Gaussian. Given  $n$  i.i.d. samples from  $P_{T^*}$ , for any  $\varepsilon, \delta > 0$ , if  $n \gtrsim \frac{d}{\varepsilon} \log \frac{d}{\delta}$ , then  $\hat{T}$  returned by Algorithm 1 satisfies

$$D_{\text{KL}}(P_{T^*}||P_{\hat{T}}) \leq \varepsilon,$$

with probability at least  $1 - \delta$ . Besides, suppose  $P$  is an unknown Gaussian distribution such that  $P = P_{T^*}$ . Given  $n$  i.i.d. samples drawn from  $P$ . For any small  $\varepsilon > 0$ , if  $n = o(d/\varepsilon)$ , no algorithm returns a directed tree  $\hat{T}$  such that

$$D_{\text{KL}}(P||P_{\hat{T}}) \leq \varepsilon$$

with probability at least  $2/3$ .

### Algorithm 1 Modified Chow-Liu algorithm

- Input:**  $n$  i.i.d. samples  $(X_1^{(i)}, \dots, X_d^{(i)})$
- For each  $j = 1, \dots, d$ :
  - $\hat{\sigma}_j^2 \leftarrow \frac{1}{n} \sum_{i=1}^n (X_j^{(i)})^2$
- For each pair  $(j, k), 1 \leq j < k \leq d$ :
  - $\hat{\rho}_{jk} \leftarrow \frac{1}{n} \sum_{i=1}^n X_j^{(i)} X_k^{(i)}$
- For each pair  $(j, k), 1 \leq j < k \leq d$ :
  - $\hat{T}(X_j; X_k) \leftarrow -\frac{1}{2} \log \left(1 - \frac{\hat{\rho}_{jk}^2}{\hat{\sigma}_j^2 \hat{\sigma}_k^2}\right)$  which is same as  $\frac{1}{2} \log \left(1 + \frac{\hat{\rho}_{jk}^2}{\hat{\sigma}_j^2 \hat{\sigma}_k^2}\right)$
- $G \leftarrow$  the weighted complete undirected graph on  $[d]$  whose edge weight for  $(j, k)$  is  $\hat{T}(X_j; X_k)$
- $\hat{S} \leftarrow$  the maximum weighted spanning tree of  $G$
- $\hat{T} \leftarrow$  any directed tree with skeleton to be  $\hat{S}$
- return** A directed tree  $\hat{T}$

## Optimal Faithful Tree Learning

### Algorithm 2 PC-Tree algorithm

- Input:**  $n$  i.i.d. samples  $(X_1^{(i)}, \dots, X_d^{(i)})$
- Let  $\hat{E} = \emptyset$ .
- For each pair  $(j, k), 0 \leq j < k \leq d$ :
  - For all  $\ell \in [d] \cup \{\emptyset\} \setminus \{j, k\}$ :
    - Test  $H_0 : X_j \perp\!\!\!\perp X_k | X_\ell$  vs.  $H_1 : X_j \not\perp\!\!\!\perp X_k | X_\ell$ , store the results.
  - If all tests reject, then  $\hat{E} \leftarrow \hat{E} \cup \{j - k\}$ .
  - Else (if some test accepts), let  $S(j, k) = \{\ell \in [d] \cup \{\emptyset\} \setminus \{j, k\} : X_j \perp\!\!\!\perp X_k | X_\ell\}$ .
- Return:**  $\hat{T} = ([d], \hat{E})$ , separation set  $S$

**Theorem 3 :** [Structure learning] For any  $T \in \tilde{\mathcal{T}}$ , assuming  $P$  is  $c$ -strong tree-faithful to  $T$ , applying Algorithm 2 with sample correlation for CI testing, if the sample size

$$n \gtrsim \frac{1}{c^2} \left( \log d + \log(1/\delta) \right),$$

then  $\Pr(\hat{T} = \text{sk}(T)) \geq 1 - \delta$ , and  $\Pr(\text{Orient}(\hat{T}, S) = \bar{T}) \geq 1 - \delta$ .

Besides, assuming  $c^2 \leq 1/5, d \geq 4$ , if the sample size is bounded as

$$n \leq \frac{1 - 2\delta}{8} \times \frac{\log d}{c^2},$$

then for any estimator  $\hat{T}$  for  $\bar{T}$ ,

$$\inf_{\hat{T}} \sup_{\substack{T \in \tilde{\mathcal{T}} \\ P \text{ is } c\text{-strong} \\ \text{tree-faithful to } T}} \Pr(\hat{T} \neq \bar{T}) \geq \delta - \frac{\log 2}{\log d}.$$

## Experiment results

- PC-Tree algorithm does perform the best, especially on PRR over the baselines.
- We have not analyzed the performance of Chow-Liu under the goal of structure learning, and we conjecture a similar sample complexity is shared with PC-Tree.

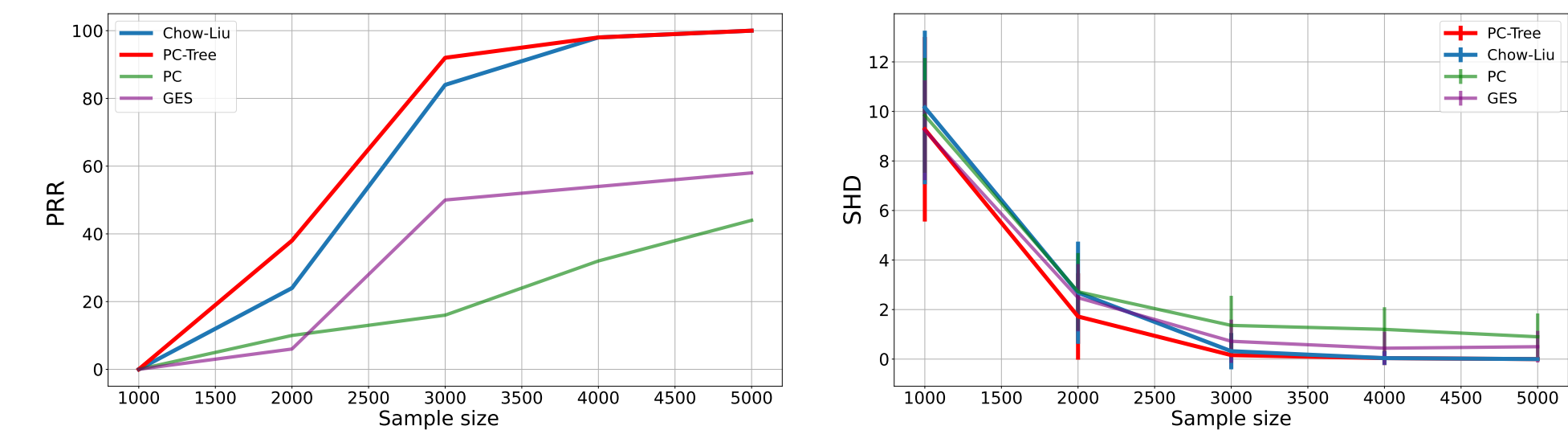


Figure 2: Performance comparison for PC-Tree, Chow-Liu, PC and GES algorithm evaluated on SHD and PRR. The red, blue, green, purple lines are for PC-Tree, Chow-Liu, PC and GES respectively.

## Conclusion and Future Work

We treat both problems in a unified setting, allowing for an explicit comparison of these problem:

- In regime  $\varepsilon \ll dc^2$ , distribution learning is harder (in terms of sample size needed);
- In regime  $c^2 \ll \varepsilon \ll dc^2$ , distribution learning does not automatically imply structure learning;
- Extending these results beyond the Gaussians we consider here (as well as finite alphabets as in previous work) is a promising direction for future research.