

\* Structure Learning Recap.

Approach 1. ML for DAG

$$\text{SCORE}(G) = \sum_{i=1}^n I_{\hat{\rho}}(X_i; X_{\pi_i})$$

tractable for  $G \in \text{TREE}$ , using max-weight spanning tree  
[Chow-Liu algorithm]

Approach 2. ML for Ising Model

$$\min_{\theta \in \mathbb{R}^{n \times n}} \log Z_G(\theta) - \langle \hat{M}, \theta \rangle + \lambda \|\theta\|_{L_1}$$

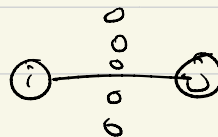
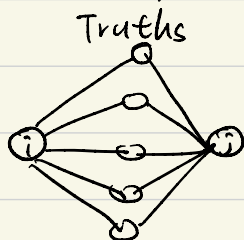
↑  
intractable for general graph  
even if it is sparse.

Approach 3. Local Independence Tests.

Enumerate all  $|S|=k$  neighborhood for each node  $i$   
and select one with best "Independence" score.

$$O(n^{k+1})$$

$O(n^2)$  → alternatively, one could check "Correlation" for each pair  
Algorithm output.



even if  $i, j$  are not directly connected  
but if there are many paths connecting  $i$  to  $j$ ,

\* None of the algorithms above work in large scale practical problems.  
 Q. Can we design an algorithm that is efficient & works well in practice?

Consider Binary Random variables  $X = \{0, 1\}$ , and an undirected G.M

$$P(X_1, X_2, \dots, X_n) = \frac{1}{Z(\theta)} \cdot \exp \left\{ \sum_{i=1}^n \theta_{ii} X_i + \sum_{(i,j) \in E} \theta_{ij} X_i X_j \right\}$$

claim: A conditional distribution of  $X_i$  given the rest is

$$\frac{P(X_i=1 | X_2, \dots, X_n)}{P(X_i=0 | X_2, \dots, X_n)} = \exp \left\{ \theta_{ii} + \sum_{j \in \partial_i} \theta_{ij} \cdot X_j \right\}$$

↑ this is a logistic Regression problem, which can be solved efficiently.

$$P(X_i=1 | X_2, \dots, X_n) = \frac{e^{\theta_{ii} + \sum_{j \in \partial_i} \theta_{ij} X_j}}{1 + e^{\theta_{ii} + \sum_{j \in \partial_i} \theta_{ij} X_j}}$$

= sigmoid  $\left( \theta_{ii} + \sum_{j \in \partial_i} \theta_{ij} X_j \right)$   
 ↑  
 Logistic function  
 ↑  
 linear in  $\theta_{ij}$   
 ↑  
 first row of  $\theta \in \mathbb{R}^{n \times n}$

proof:

$$P(X_i=1 | X_2^n) = \frac{P(X_i=1 \wedge X_2^n)}{P(X_2^n)}$$

$$= \frac{\frac{1}{Z(\theta)} \cdot e^{\left\{ \theta_{ii} + \sum_{j \in \partial_i} \theta_{ij} X_j + \sum_{\substack{(i,j) \in E \\ (i,j) \neq (i,i)}} \theta_{ij} X_i X_j \right\}}}{\frac{1}{Z(\theta)} \cdot e^{\left\{ \sum_{\substack{(i,j) \in E \\ (i,j) \neq (i,i)}} \theta_{ij} X_i X_j \right\}}}$$

$$P(X_i=0 | X_2^n) = \frac{\frac{1}{Z(\theta)} \cdot e^{\left\{ \sum_{\substack{(i,j) \in E \\ (i,j) \neq (i,i)}} \theta_{ij} X_i X_j \right\}}}{\frac{1}{Z(\theta)} \cdot e^{\left\{ \sum_{\substack{(i,j) \in E \\ (i,j) \neq (i,i)}} \theta_{ij} X_i X_j \right\}}}$$

\* Review of Logistic Regression.

We are given labelled samples of features  $\mathbf{z} = (z_1^{(d)}, \dots, z_L^{(d)}) \in \mathbb{R}^L$  and labels  $Y \in \{0, 1\}$

$$\text{Data: } \{(\mathbf{z}^{(d)}, Y^{(d)})\}_{d=1}^N$$

We want to find a model parameter  $W \in \mathbb{R}^L$  such that

$$P(Y=1 | \mathbf{z}) = \frac{\exp(\sum_{k=1}^L w_k z_k)}{1 + \exp(\sum_{k=1}^L w_k z_k)}$$

$$P(Y=0 | \mathbf{z}) = \frac{1}{1 + \exp(\sum_{k=1}^L w_k z_k)}$$

We will find the maximum likelihood estimator.

$$\text{log-likelihood } \mathcal{L}(\{\mathbf{z}^{(d)}, Y^{(d)}\}_{d=1}^N, w) = \frac{1}{N} \sum_{d=1}^N \log P_w(Y^{(d)} | \mathbf{z}^{(d)})$$

$$= \frac{1}{N} \sum_{d=1}^N \left\{ Y^{(d)} \log P(Y=1 | \mathbf{z}^{(d)}) + (1 - Y^{(d)}) \log P(Y=0 | \mathbf{z}^{(d)}) \right\}$$

$$= \frac{1}{N} \sum_{d=1}^N \left\{ Y^{(d)} \left( \sum_{k=1}^L w_k z_k^{(d)} \right) - \log \left( 1 + \exp \left( \sum_{k=1}^L w_k z_k^{(d)} \right) \right) \right\}$$

this is a concave function in  $W \in \mathbb{R}^L$ , which is maximised using Gradient Descent / Ascent.

Initialize  $W^{(0)} = \mathbf{1}$ .

$$\text{Repeat } W^{(t+1)} = W^{(t)} + \frac{1}{t} \cdot \nabla_w \mathcal{L}(\{\mathbf{z}^{(d)}, Y^{(d)}\}_{d=1}^N, W^{(t)})$$

$$\frac{1}{N} \sum_{d=1}^N \left\{ Y^{(d)} \cdot \mathbf{z}^{(d)} - \frac{\exp(\sum_{k=1}^L w_k z_k^{(d)})}{1 + \exp(\sum_{k=1}^L w_k z_k^{(d)})} \cdot \mathbf{z}^{(d)} \right\}$$

$$= W^{(t)} + \frac{1}{t} \cdot \frac{1}{N} \sum_{d=1}^N \mathbf{z}^{(d)} \cdot (Y^{(d)} - P_{W^{(t)}}(Y=1 | \mathbf{z}^{(d)}))$$

\* Logistic regression for neighborhood selection.

Structural Learning

$$P(X_i=1 | X_2, \dots, X_n) = \frac{e^{\theta_{i1} + \sum_{j \neq i} \theta_{ij} X_j}}{1 + e^{\theta_{i1} + \sum_{j \neq i} \theta_{ij} X_j}}$$

$$P(X_i=0 | X_2, \dots, X_n) = \frac{1}{1 + \exp\{\theta_{i1} + \sum_{j \neq i} \theta_{ij} X_j\}}$$

Logistic regression

$$P(Y=1 | Z) = \frac{\exp(\sum_{k=1}^L w_k z_k)}{1 + \exp(\sum_{k=1}^L w_k z_k)}$$

$$P(Y=0 | Z) = \frac{1}{1 + \exp(\sum_{k=1}^L w_k z_k)}$$

for node 1, if we know  $G$ , then we can find  $\theta_{1j}$ 's using Logistic Regression approaches with features  $(1, X_{\partial 1})$  and label  $X_1$

As we do not know  $G$ , suppose the ground truths  $|\theta_{ij}| \leq 1$  and degree  $\leq K$ .

This motivates the following formulation

$$\min_{\theta_{1 \cdot}} - \mathcal{L}(\{X_i^{(l)}, (1, X_2^{(l)}, \dots, X_n^{(l)})\}_{l=1}^N, \theta_{1 \cdot})$$

"  $(\theta_{11}, \theta_{12}, \dots, \theta_{1n}) \in \mathbb{R}^n$

s. t.  $\|\theta_{1 \cdot}\|_{L_1} \leq K.$

or equivalently, we can solve

$$(*) \min_{\theta_{1 \cdot}} - \mathcal{L}(\{X_i^{(l)}, (1, X_2^{(l)}, \dots, X_n^{(l)})\}_{l=1}^N, \theta_{1 \cdot}) + \lambda \cdot \|\theta_{1 \cdot}\|_{L_1}.$$



Theorem [Klivans, Meka, 2017]

If max degree  $\leq K$

$P(X_i=1 | X_{-i}) \in [\delta, 1-\delta]$  for some  $\delta > 0$

number of samples  $N \geq C \cdot \log n \cdot \frac{1}{\epsilon^2}$

Then (\*) achieves  $\|\hat{\theta} - \theta^*\|_\infty \leq \epsilon$

in run-time  $O(n^2 \text{polylog}(\frac{n}{\epsilon}))$

if all nonzero  $\theta_{ij}^* > 2\epsilon$ , then we can threshold <sup>at  $\theta_{ij} \geq \epsilon$</sup>  to recover the structure exactly.

\* This principle can be used for more general graphical models.  
For example, Gaussian Graphical Models.

$$P(X) = \frac{1}{Z} \exp\left\{-\frac{1}{2} X^T \theta X\right\}$$

w.l.o.g. suppose  $h=0$ .

$X \in \mathbb{R}$

We first compute the conditional.

$$\begin{aligned} P(x_1 | X_2^n) &= \frac{1}{Z_1} \cdot \exp\left\{-\frac{1}{2} \left(\theta_{11} x_1^2 + 2 \overbrace{\left(\sum_{j \in \partial 1} \theta_{1j} x_j\right)}^\beta x_1\right)\right\} \\ &= \frac{1}{\sqrt{2\pi/\theta_{11}}} \cdot \exp\left\{-\frac{1}{2} \frac{(x_1 + \beta/\theta_{11})^2}{1/\theta_{11}}\right\} \end{aligned}$$

Applying maximum likelihood to estimate  $x_1$  from  $X_2^n$ ,

$$\min \sum_{l=1}^N \log P(x_1^{(l)} | x_2^{(l)} \dots x_n^{(l)}) = \frac{\theta_{11}}{2} \sum_{l=1}^N \left(x_1^{(l)} - \sum_{j \in \partial 1} \frac{\theta_{1j}}{-\theta_{11}} \cdot x_j^{(l)}\right)^2$$

$$-\frac{1}{2} \log \theta_{11}$$

ensures we don't make  $\theta_{11} \downarrow$  small.

If we only care about the graph structure, for now, we can re-parametrize and solve for each node.

$$\min \sum_{i=1}^N \left( X_i^{(a)} - \sum_{j=2}^n W_j \cdot X_j^{(a)} \right)^2 + \lambda \|W\|_{L1} \quad [\text{HW4}]$$

and threshold the resulting  $w$  to recover the neighborhood.

\* An alternative (and also very popular) way to learn the structure of a Gaussian Graphical Model is "Graph Lasso".

step 1. Compute Covariance Matrix  $S_{ij} = \frac{1}{N} \sum_{k=1}^N X_i^{(k)} X_j^{(k)}$

step 2. maximize likelihood for information matrix  $J$ .

$$J^* = \arg \max_J \left\{ \underbrace{\log |J|}_{\substack{\text{log-likelihood} \\ \uparrow \\ \text{Determinant}}} - \underbrace{\text{Tr}(S \cdot J)}_{\substack{\uparrow \\ \text{Trace}(A) \\ \sum_{i=1}^n A_{ii}}} - \underbrace{\lambda \|J\|_{L1}}_{\substack{\uparrow \\ \text{sparsity} \\ \text{encouraging} \\ \text{regularizer} \\ \|J\|_{L1} = \sum_{ij} |J_{ij}|}} \right\}$$

\* This is a concave maximization, and solved with gradient ascent.

\* This is more accurate than node-wise neighborhood learning as all edges are learnt jointly.