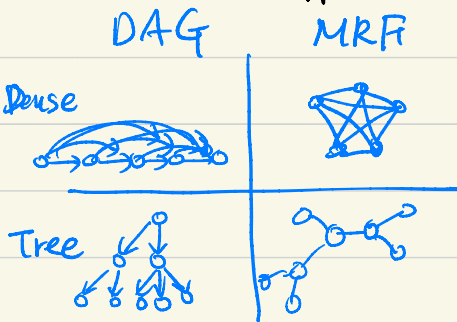


* Structure Learning Recap.

Approach 1. ML for DAG

Derived under B.N.

$$\text{SCORE}(G) = \sum_{i=1}^n I_{\beta}(X_i; X_{\pi_i}) \leftarrow \{P_{X_i|K(\pi_i)}\}$$



tractable for $G \in \text{Tree} \leftrightarrow \text{max weight spanning tree}$
 [Chow-Liu algorithm]

Approach 2. ML for Ising Models.

$$\min_{\theta \in \mathbb{R}^{\text{max}}} \underbrace{\log Z_G(\theta) - \langle \hat{M}, \theta \rangle}_{-\text{log-likelihood}} + \underbrace{\lambda \|\theta\|_{L_1}}_{\text{regularization}}$$

intractable for general graphs even for sparse graphs. $O(K^n)$

Approach 3. Local Independence Tests. with $\text{max-degree} \leq K$

Enumerate over all nodes.

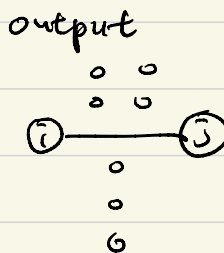
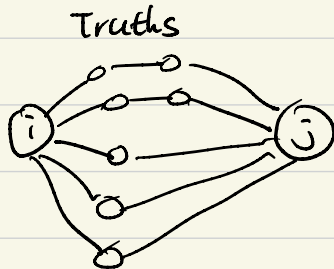
$$O(n^{K+1})$$

select a neighborhood S_i .

with best $\text{SCORE}(S_i) \cong \text{conditional independence}$

Approach 4. pairwise tests.

$$\text{SCORE}(i, j) = \frac{1}{N} \sum_{k=1}^N X_i^{(k)} X_j^{(k)} \quad : \text{correlation}$$



* 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 $X = \{0, 1\}$

undirected graphical models.

$$P(x_1, \dots, x_n) = \frac{1}{Z(\theta)} \exp \left\{ \sum_{i=1}^n \theta_{i,i} x_i + \sum_{(i,j) \in E} \theta_{i,j} x_i x_j \right\}$$

Strategy is to predict x_i using x_{-i} , but with a subset.

Claim: A conditional distribution of x_1 given the rest is

$$(*) \quad \frac{P(x_1=1 | x_2 \dots x_n)}{P(x_1=0 | x_2 \dots x_n)} = \exp \left\{ \theta_{1,1} + \sum_{j \in \partial 1} \theta_{1,j} x_j \right\}$$

this is a logistic regression problem, where we try to predict x_1 from x_2^n can be solved efficiently.

$$(*) \rightarrow P(x_1=1 | x_2 \dots x_n) = \frac{e^{\theta_{1,1} + \sum_j \theta_{1,j} x_j}}{1 + e^{\theta_{1,1} + \sum_j \theta_{1,j} x_j}}$$

$$= \text{sigmoid} \left(\theta_{1,1} + \sum_j \theta_{1,j} x_j \right)$$

↓
logistic function.

Linear function

θ_1 .

$(\theta_{1,1}, \theta_{1,2}, \dots, \theta_{1,n})$

proof of (*) >

$$P(x_1=1 | x_2^n) = \frac{P(x_1=1 \wedge x_2^n)}{P(x_2^n)} = \frac{\frac{1}{Z(\theta)} \cdot \exp \left\{ \theta_{1,1} + \sum_{j \in \partial 1} \theta_{1,j} x_j + \sum_{\substack{(i,j) \in E \\ i,j \neq 1}} \theta_{i,j} x_i x_j \right\}}{\frac{1}{Z(\theta)} \cdot \exp \left\{ \sum_{\substack{(i,j) \in E \\ i,j \neq 1}} \theta_{i,j} x_i x_j \right\}}$$

$$P(x_1=0 | x_2^n) = \frac{\frac{1}{Z(\theta)} \cdot \exp \left\{ \sum_{\substack{(i,j) \in E \\ i,j \neq 1}} \theta_{i,j} x_i x_j \right\}}{P(x_2^n)}$$

$$P(X_1=1|X_2^n) = \frac{\exp\{\theta_0 + \sum \theta_j X_j\}}{1 + \exp\{\theta_0 + \sum \theta_j X_j\}}$$

* Review of Logistic Regression.

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

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

We suppose a parametric model with $W \in \mathbb{R}^L$

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

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

Apply maximum likelihood.

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

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

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

this is a concave maximization in $W \in \mathbb{R}^L$.

use Gradient Ascent to find \hat{W} .

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

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

$$\frac{1}{N} \sum_{d=1}^N Y^{(d)} \cdot z^{(d)} - \frac{\exp(\langle W^{(t)}, z^{(d)} \rangle)}{1 + \exp(\langle W^{(t)}, z^{(d)} \rangle)} \cdot z^{(d)}$$

$$= W^{(t)} + \frac{1}{t} \cdot \frac{1}{N} \sum_{d=1}^N z^{(d)} \left(\underset{\substack{\uparrow \\ \text{sample} \\ \text{label}}}{Y^{(d)}} - \underset{\substack{\uparrow \\ \text{current prediction}}}{P_{W^{(t)}}(Y=1|z^{(d)})} \right)$$

* Logistic regression for neighborhood selection.

structural learning

$$P(X_i=1 | X_2, \dots, X_n) = \frac{e^{\theta_{i1} + \sum_{j \in \mathcal{N}_1} \theta_{ij} X_j}}{1 + e^{\theta_{i1} + \sum_{j \in \mathcal{N}_1} \theta_{ij} X_j}}$$

\uparrow
 G

logistic regression

$$P(Y=1 | Z) = \frac{e^{\sum W_k Z_k}}{1 + e^{\sum W_k Z_k}}$$

for node 1, if we know G , then features are $(1, X_{\mathcal{N}_1})$
label is X_1 .

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

this motivates the following algorithm.

$$\begin{aligned} \min_{\theta_1} & - \mathcal{L}(\{X_1^{(1)}, (1, X_2^{(1)}, \dots, X_n^{(1)})\}_{i=1}^N, \theta_1) \\ & \text{feature.} \quad \theta_{11}, \theta_{12}, \dots, \theta_{1n} \\ & \text{feature 1.} \quad \theta_{11} + \theta_{12}X_2 + \theta_{13}X_3 + \dots \\ \text{s.t.} & \|\theta_1\|_{L_1} \leq K \end{aligned}$$

OR equivalently,

$$(**) \min_{\theta_1} - \mathcal{L}(\cdot, \theta_1) + \lambda \|\theta_1\|_{L_1}$$

$\exists \lambda^*(K)$ gives same solution.

$$\hat{\theta}_1 \rightarrow (0 \ \hat{\theta}_{13} \ 0 \ \hat{\theta}_{17} \ 0 \ 0 \ 0)$$

$\uparrow \quad \quad \quad \uparrow$
 structure
 values \cong factors

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 $N \geq C \cdot \log n \cdot \frac{1}{\delta^2} \times K \times C_\delta$

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

$$\max_{i,j} |\hat{\theta}_{ij} - \theta_{ij}^*| \leq \epsilon$$

run-time $O(N^2 \text{ poly } \log \frac{n}{\epsilon})$.

if all $|\theta_{ij}^*| > 2\epsilon$, then threshold $|\hat{\theta}_{ij}| \gtrsim \epsilon$
learn structure exactly.

* This principle can be used for more general graphical models
example of Gaussian Graphical models.

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

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

$X = \mathbb{R}$

① Compute conditional distribution.

$$P(x_1 | x_2^n) = \frac{1}{Z_1} \cdot \exp\left\{-\frac{1}{2} \left(J_{11} x_1^2 + 2 \left(\sum_{j \in \mathcal{B}} J_{1j} \cdot x_j \right) \cdot x_1 \right)\right\}$$

$$= \frac{1}{\sqrt{2\pi \cdot \frac{1}{J_{11}}}} \cdot \exp\left\{-\frac{1}{2} \frac{\left(x_1 + \frac{\mathcal{B}}{J_{11}}\right)^2}{\frac{1}{J_{11}}}\right\}$$

② Apply maximum likelihood to estimate x_1 from x_2^n

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

Annotations:
 - J_{11}, J_{1j} small
 - J_{11} can be *very small*
 - $\frac{J_{1j}}{-J_{11}}$ ratio
 - $x_1^{(l)}$ label
 - $x_j^{(l)}$ features
 - $-\frac{1}{2} \log J_{11}$ model

as we care about structure,
we re-parametrize it by $w \in \mathbb{R}^{n-1}$

[hw4]

$$\min_{w \in \mathbb{R}^{n-1}} \frac{1}{N} \sum_{l=1}^N \left(X_i^{(l)} - \sum_{j=2}^n w_j \cdot X_j^{(l)} \right)^2 + \lambda \|w\|_1$$

* Motivation: $\left\{ \begin{array}{l} \text{Get } J \text{ directly} \\ \text{without separating the neighborhoods} \end{array} \right.$

Graphical Lasso.

Step 1: Compute the Covariance Matrix $S_{ij} = \frac{1}{N} \sum_{l=1}^N X_i^{(l)} X_j^{(l)}$

Step 2: Maximize likelihood for information matrix J .

$$\hat{J} = \arg \max_J \left\{ \underbrace{\log |J|}_{\substack{\uparrow \\ \text{determinant}}} - \underbrace{\text{Tr}(S \cdot J)}_{\substack{\uparrow \\ \text{Tr}(A) \\ = \sum_i A_{ii}}} - \underbrace{\lambda \|J\|_{L_1}}_{\substack{\uparrow \\ \text{regularize} \\ \sum_{i,j} |J_{ij}|}} \right\}$$

s.t. $J \succ 0$

* This concave maximization