

## \* Structural Learning Recap.

### Approach 1: ML for DAG

$$\text{SCORE}(G) = \sum_{i=1}^n I_p(x_i; j | X_{\pi_i}) \quad \leftarrow \{P(x_i | X_{\pi_i})\}$$

[Chow-Liu Algorithm]  $\longleftrightarrow$  Max-weight Spanning  
If search over all trees, then optimization  
can be solved efficiently exactly.

### Approach 2. ML for Ising models

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

$\log Z_G(\theta)$   
 $\langle \hat{M}, \theta \rangle$   
 $\|\theta\|_{L_1}$

$\hat{M}_{11}, \hat{M}_{12}, \dots$   
 $\hat{M}_{11} = \frac{1}{N} \sum_{j=1}^N x_1^{(j)}$   
 $\hat{M}_{12} = \frac{1}{N} \sum_{j=1}^N x_1^{(j)} x_2^{(j)}$

$\theta_{11}, \theta_{12}, \dots$   
 $\theta_{11} = \frac{1}{N} \sum_{j=1}^N q_1^{(j)}$

Complex to evaluate

Q. Can we design an algorithm that is efficient & works well in Practice?

Binary random variables  $X = \{0, 1\}$

undirected graphical models

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

Strategy is to predict  $x_i$  using  $x_{-i} \triangleq \underbrace{\{x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n\}}_{\text{subset of this set}}$

Claim: A condition distribution of  $x_i$  given the rest is

$$(*) \quad P(x_i=1 | \overbrace{x_2 \dots x_n}^{x_{-i}}) = \frac{P(x_i=1 | x_2 \dots x_n)}{P(x_i=0 | x_2 \dots x_n)} = \exp \left\{ \theta_{ii} + \sum_{j \neq i} \theta_{ij} x_j \right\}$$

this is a logistic Regression problem,

to predict  $x_i$  from  $x_2^n$ , can be solved very efficiently:

$$(*) \rightarrow P(x_i=1 | x_{-i}) = \frac{e^{\theta_{ii} + \sum_j \theta_{ij} x_j}}{1 + e^{\theta_{ii} + \sum_j \theta_{ij} x_j}}$$

Proof (\*)

$$\begin{aligned} P(x_i=1 | x_2^n) &= \frac{P(x_i=1 \wedge x_2^n)}{P(x_2^n)} \\ &= \frac{P(x_i=1)}{P(x_2^n)} \cdot \exp \left\{ \theta_{ii} + \sum_{j=2}^n \theta_{ij} x_j + \underbrace{\sum_{j=2}^n \theta_{1j} x_j \cdot 1}_{\text{linear in } \theta_1} \right\} \\ &\quad \left. + \sum_{j \neq 1} \theta_{1j} x_j \right\} \end{aligned}$$

$$P(x_i=0 | x_2^n) =$$

$$P(X_1=1 | X_2^u) = \frac{\exp\{\theta_{01} + \sum_j X_j \theta_{1j}\}}{1 + \exp\{\theta_{01} + \sum_j X_j \theta_{1j}\}}$$

\* Review of L logistic regression

- We are given labelled examples of structures  $\mathcal{Z}^{(l)} = (\mathcal{Z}_1^{(l)}, \dots, \mathcal{Z}_L^{(l)})$

$$\text{Data} = (\mathcal{Z}^{(l)}, Y^{(l)}) \quad l \in \{1, \dots, N\}$$

label  $Y^{(l)} \in \{0, 1\}$

- Suppose a parametric model with  $w \in \mathbb{R}^L$

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

$$P(Y=0 | Z) = \frac{1}{1 + \exp(\cdot)}$$

- Maximum likelihood

$$\text{log-likelihood } L(\{(Z^{(l)}, Y^{(l)})\}_{l=1}^N, w \in \mathbb{R}^L) \triangleq \sum_{l=1}^N \log P_w(Y^{(l)} | Z^{(l)})$$

$$= \frac{1}{N} \sum_{l=1}^N \left\{ Y^{(l)} \cdot P_w(Y=1 | Z^{(l)}) + (1 - Y^{(l)}) \cdot P_w(Y=0 | Z^{(l)}) \right\}$$

$$= \frac{1}{N} \sum_{l=1}^N \left\{ Y^{(l)} \cdot \left( \prod_{k=1}^L w_k Z_k^{(l)} \right) - \log \left( 1 + \exp \left( \sum_{k=1}^L w_k Z_k^{(l)} \right) \right) \right\}$$

utility function over  $w$

\* this is a concave maximization over  $w \in \mathbb{R}^L$   
use gradient ascent to solve efficiently.

structural learning.

\* Logistic Regression for neighborhood selection

Logistic regression

$$P(X_1=1|X_2) = \frac{e^{\theta_{11} + \sum_{j \neq 1} \theta_{1j} X_j}}{1 + e^{\theta_{11} + \sum_{j \neq 1} \theta_{1j} X_j}}$$

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

for node 1, if we know  $G$ , then features are  $(1, X_{j1})$   
label is  $X_1$ .

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

this motivates following Sparse logistic regression.

$$\min_{\theta_{10}} -L(\underbrace{\{(1, x_1^{(1)}, \dots, x_n^{(1)})\}}_{\text{feature}}, \underbrace{\{x_1^{(1)}\}_{l=1}^N}_{\text{label}}, \theta_{10})$$

$$(0, 0, \dots, 0)$$

$$\text{s.t. } \|\theta_{10}\|_{L_1} \leq K \quad \text{or equivalently}$$

$$\sum_j |\theta_{1j}|$$

$$(\star\star) \quad \min_{\theta_{10}} -L(\cdot, \theta_{10}) + \lambda \cdot \|\theta_{10}\|_{L_1}$$

$$\Rightarrow \hat{\theta}_{10} \approx (0, \underbrace{\hat{\theta}_{13}, \hat{\theta}_{14}, \dots, \hat{\theta}_{118}}_{\text{Graph structure}}, \dots)$$

Theorem [Klivans, Meka 2017]

If max degree  $\leq k$

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

for some  $\varepsilon > 0$

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

Then (ex) achieve  $\|\hat{\theta} - \theta^*\|_\infty \leq \varepsilon$

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

and runtime  $O(n^2 \text{polylog } \frac{n}{\varepsilon})$

if all  $|\theta_{ij}| > 2\varepsilon$ , then threshold  $|\hat{\theta}_{ij}| \geq \varepsilon$

learn structure exactly

\*

$$P(x^n) = \frac{1}{Z(\theta)} \exp \left\{ \sum_j \theta_{jj} x_j + \sum_{(i,j) \in E} \theta_{ij} x_i x_j \right\}$$

\* This principle can be used for more general graphical models

- Gaussian Graphical Models.

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

$$J = \Sigma^{-1}$$

$$\text{u.e.} \Leftrightarrow h=0$$

$$\mathcal{X} = \mathbb{R}$$

① Compute conditional distribution

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

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

2 Neighbors of node 1

② Apply Maximum Likelihood principle to get  $x_i$  from  $x_i^{(l)}$

$$\begin{aligned} \min -\frac{1}{N} \sum_{l=1}^N \log P(x_i^{(l)} | x_2^{(l)}, \dots, x_n^{(l)}) \\ = \frac{\mathcal{J}_{11}}{2N} \cdot \sum_{l=1}^N \left( x_i^{(l)} - \sum_{j \neq 1} \underbrace{\frac{\mathcal{J}_{1j}}{-\mathcal{J}_{11}} \cdot x_j^{(l)}}_{w_j} \right)^2 \end{aligned}$$

↑  
label      ↑  
model parameter      ↓-less  
↑  
features.

$$-\frac{1}{2} \log \mathcal{J}_{11}$$

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

$$w \in \mathbb{R}^{n-1} \quad \frac{1}{N} \sum_{l=1}^N \left( x_i^{(l)} - \sum_{j=2}^n w_j x_j^{(l)} \right)^2 + \lambda \|w\|_1.$$

### \* Graphical Lasso.

Motivation: — Get  $\mathcal{J}$  directly, without separating each neighborhood.

Step 1. Compute the empirical Covariance matrix

$$S_{ij} = \frac{1}{N} \sum_{l=1}^N x_i^{(l)} x_j^{(l)}$$

Step 2. maximize likelihood over Information matrix  $\mathcal{J} \in \mathbb{R}^{n \times n}$

$$\hat{\mathcal{J}} \in \arg \max_{\mathcal{J}} \left\{ \underbrace{\log |\mathcal{J}|}_{\text{determinant}} - \underbrace{\text{Tr}(S \cdot \mathcal{J})}_{\text{Tr}(A) = \sum_i A_{ii}} - \lambda \|\mathcal{J}\|_F^2 \right\}$$

likelihood      regularizer

$$\text{s.t. } \mathcal{J} \succ 0$$

Others: linear regression

→ Semidefinite program.