

## \*Overview

- Graphical Models & Markov Properties
- Inference Problems: given  $P_G(x)$  find
  - $P(x_i)$
  - $\max_x P_G(x)$
  - Conj. Z

- Belief Propagation
- Variational methods
- Gibbs Sampling

- Learning Graphical model

given samples :  $x^{(1)}, x^{(2)}, \dots, x^{(N)} \in \mathcal{X}^n$

- ↳ Learn the structure of the graph  $G$
- ↳ Learn the parameters of factor.

## \* Structural Learning

- $X \in \mathbb{R}^n$  Random Vector
- Directed Graphical Model :  $|E| = \binom{n}{2} = \frac{n(n-1)}{2}$

# of possible graphs  $\leq 3^{|E|} \approx 3^{\frac{n(n-1)}{2}}$

- Given  $x^{(1)}, \dots, x^{(N)}$  indep. samples from unknown  $P(x)$ 
  - How do we "score" each graph?
  - How do we find the graph with highest score?
- There are 2 ways to approach such statistical problems.

### Frequentist

- Assume graph  $G$  and  $P = \{P(x_i | X_{\pi_i})\}_{i=1}^n$  are deterministic but unknown
- Maximum likelihood (ML) estimate  
find  $(G, P)$  that maximizes:

$$\max_{G, P} \sum_{j=1}^N \underbrace{\log P_{G, P}(x^{(j)})}_{\text{SCORE}(G, P)}$$

### Bayesian

- Assume  $(G, P)$  is randomly drawn from known dist.

$$Q_{G, P}$$

- Maximum A Posteriori (MAP) estimate, maximizes:

$$\max_{G, P} P(G, P | x^{(1)}, \dots, x^{(N)})$$

## \* Frequentist : Structural Learning

$$\hat{G} = \underset{G}{\operatorname{argmax}} \left\{ \max_P \frac{1}{N} \sum_{j=1}^N \log P_{G, P}(x^{(j)}) \right\}$$

$\text{SCORE}(G)$

\* Simple case with  $n=2$ ,  $X = (X_1, X_2) \in \{0, 1\}^2$

Samples  $(x_1, x_2)$ :  $(0,0), (0,1), (1,0), (1,1)$

$x_1$	$x_2=0$	$x_2=1$
$x_1=0$	.1	.3
$x_1=1$	.2	.4

Case 1:  $G_1 = (\textcircled{1} \quad \textcircled{2})$  independent  $\rightarrow P_{12}(x_1, x_2) = P_1(x_1) \cdot P_2(x_2)$

ML estimate:

$$\max_{P_1} \frac{1}{N} \sum_{j=1}^N \log P_1(x_1^{(j)}) + \max_{P_2} \frac{1}{N} \sum_{j=1}^N \log P_2(x_2^{(j)})$$

def. of empirical distribution

$$= \max_{P_1 = (P_1(0), P_1(1))} \hat{P}_1(0) \cdot \log P_1(0) + \hat{P}_1(1) \cdot \log P_1(1) + \dots$$

$$= \max_{P_1} \sum_{x_1} \hat{P}_1(x_1) \cdot \log P_1(x_1) + \dots$$

$$= \max_{P_1} \sum_{x_1} \hat{P}_1(x_1) \log \underline{\hat{P}_1(x_1)} + \sum_{x_1} \hat{P}_1(x_1) \log \frac{P_1(x_1)}{\hat{P}_1(x_1)} + \dots$$

$$-H(\hat{P}_1)$$

↳ Does not depend  $P_1$

$$-D_{KL}(\hat{P}_1 \parallel P_1) \leq 0$$

max achieved when

$$P_1 = \hat{P}_1, P_2 = \hat{P}_2$$

for  $G_1(\textcircled{1} \textcircled{2})$   $\max_{P_1, P_2} \frac{1}{N} \sum_{j=1}^N \log P_1(x_1^{(j)}) \cdot P_2(x_2^{(j)})$

$$SCORE(G_1) = -H(\hat{P}_1) - H(\hat{P}_2)$$

Case 2:  $G_2 = (\textcircled{1} \rightarrow \textcircled{2})$

$$\max_{P_{12}(x_1, x_2)} \frac{1}{N} \sum_{j=1}^N P_{12}(x_{12}^{(j)}) = -H(\hat{P}_{12}) - \underbrace{D_{KL}(\hat{P}_{12} \parallel P_{12})}_{\text{max achieved at}}$$

$$SCORE(G_2) = -H(\hat{P}_{12})$$

$$P_{12} = \hat{P}_{12}$$

$$\chi(①②) = -H(\beta_1) - H(\beta_2)$$

$$\lambda(\beta \rightarrow \gamma) = -H(\beta_{12})$$

follows from Independence makes entropy ↑

• Remark 1:  $-H(\vec{P}_{12}) \leq -H(\vec{P}_1) - H(\vec{P}_2)$

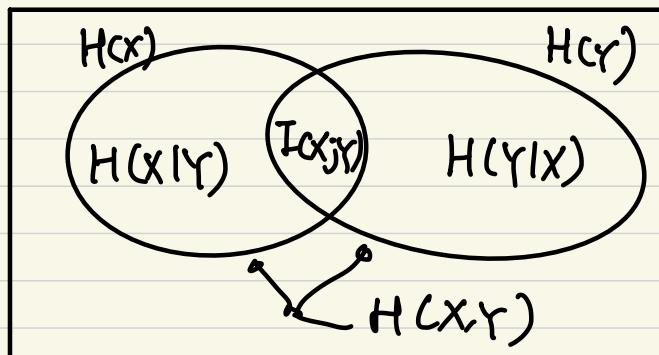
$$H(P_{12}) \leq H(P_1) + H(P_2)$$

Refresher on H, I.  $X, Y$  random vectors from joint  $P_{X,Y}(x,y)$

$$I(X;Y) \triangleq \sum_{x,y} P(x,y) \cdot \text{I}_{\{x,y\}} \frac{P(x,y)}{P(x) \cdot P(y)}$$

$$H(x) \triangleq \sum_x -p(x) \log p(x)$$

$$H(y|x) \triangleq \sum_{x,y} -P(x,y) \log P(y|x)$$



$$H(Y) = H(Y|X) + I(X;Y)$$

$$H(x) + H(y) = H(xy) + I(xy)$$

$$L(1 \rightarrow 2) \geq L(1) + L(2)$$

$\Downarrow$

$$-H(\vec{p}_1) - H(\vec{p}_2) \geq -H(\vec{p}_1) - H(\vec{p}_2)$$

always get the densest graph = Overfitting

Even if the truth is ① ②,  
you choose ① → ②

• Remark 2. depending on  $N$  and target false positive rate  $\beta$  decision is made by:

$$L(①\rightarrow②) - L(①②) = -H(\vec{p}_{12}) + H(\vec{p}_1) + H(\vec{p}_2)$$

$$= + I \beta_{12} (x_1; x_2)$$

output  $\begin{cases} \textcircled{1} \textcircled{2} & \text{if } I_{P_{12}}(x_1, x_2) > \frac{t\beta}{N} \\ \textcircled{1} \rightarrow \textcircled{2} & \text{else} \end{cases}$

\* Maximum Likelihood Approach for a DAG

$$G^* \in \operatorname{arg\max}_G \left\{ \max_{\{\pi_i\}} \frac{1}{N} \sum_{j=1}^N \log \prod_{i=1}^n P_i(x_i^{(j)} | x_{\pi_i}^{(j)}) \right\}$$

$$P_i(x_i | X_{\pi_i}) = \hat{P}_i(x_i | X_{\pi_i})$$

SCORE(G)

$$\Rightarrow \frac{1}{N} \sum_{j=1}^N \log \left( \prod_{i=1}^n \hat{P}_i(x_i^{(j)} | X_{\pi_i}^{(j)}) \right)$$

$$= \sum_{i=1}^n \frac{1}{N} \sum_{j=1}^N \log \hat{P}_i(x_i^{(j)} | X_{\pi_i}^{(j)})$$

*def empirical distribution*

$$= \sum_{i=1}^n \sum_{x_i, X_{\pi_i}} \hat{P}_i(x_i, X_{\pi_i}) \cdot \log \hat{P}_i(x_i | X_{\pi_i})$$

$$= \sum_{i=1}^n \{-H_{\hat{P}}(X_i | X_{\pi_i})\}$$

$H(x|y) = H(x) - I(x:y)$

$$= \sum_{i=1}^n I_{\hat{P}}(X_i; j X_{\pi_i})$$

*depends on Graph*      *does not depend on G.*

\* Remark: this gives a "SCORE" = likelihood for any given DAG G.

[ Search over all DAGs. ]

$$I_{\hat{P}}(X_i; j X_{\pi_i}) \geq I_{\hat{P}}(X_i; j X_{\pi'_i}) \quad \text{if } \pi_i \supset \pi'_i$$

$\Rightarrow$  we need to restrict the family of graphs we are searching over. = set of trees

\* Chow-Liu algorithm: Searches over all Trees, efficiently, exactly.

Step 1: Create a complete graph (Undirected) over  $V = \{1, 2, \dots, n\}$

fix edge weight  $w_{ij} = I_p(x_i; x_j)$

$$W = \begin{bmatrix} w_{ij} \end{bmatrix}$$

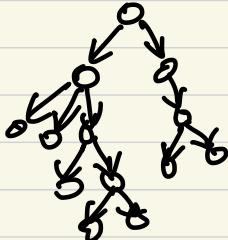


Step 2: find the max-weight Spanning Tree, for example, using Kruskal's algorithm.

Claim: Chow-Liu finds the optimal tree that maximizes

$$\max_{\text{GTrees}} \sum_{i=1}^n I_p(x_i; x_{\pi_i})$$

↑  
single node.



\* Impractical approach for undirected graphical model learning.

Consider learning an Ising Model  $(G, \{\theta_i\}_{i=1}^n, \{\theta_{ij}\}_{i,j=1}^n)$

$$x \in \{-1, +1\}$$

$$P_{G,\theta}(x) = \frac{1}{Z_\theta} \cdot \prod_{i=1}^N e^{\theta_i x_i} \cdot \prod_{(i,j) \in E} e^{\theta_{ij} x_i x_j}$$

$$\theta = \begin{bmatrix} \theta_1 & \theta_2 & \dots \\ & \ddots & \ddots \end{bmatrix}$$

Samples  $\{x^{(1)}, x^{(2)}, \dots, x^{(N)}\} = D$

The likelihood

$$P_{G,\theta}(D) = \prod_{l=1}^N P_{G,\theta}(x^{(l)})$$

$$= \prod_{l=1}^N \frac{1}{Z_\theta} \prod_{i \in V} e^{\theta_i x_i^{(l)}} \cdot \prod_{(i,j) \in E} e^{\theta_{ij} x_i^{(l)} x_j^{(l)}}$$

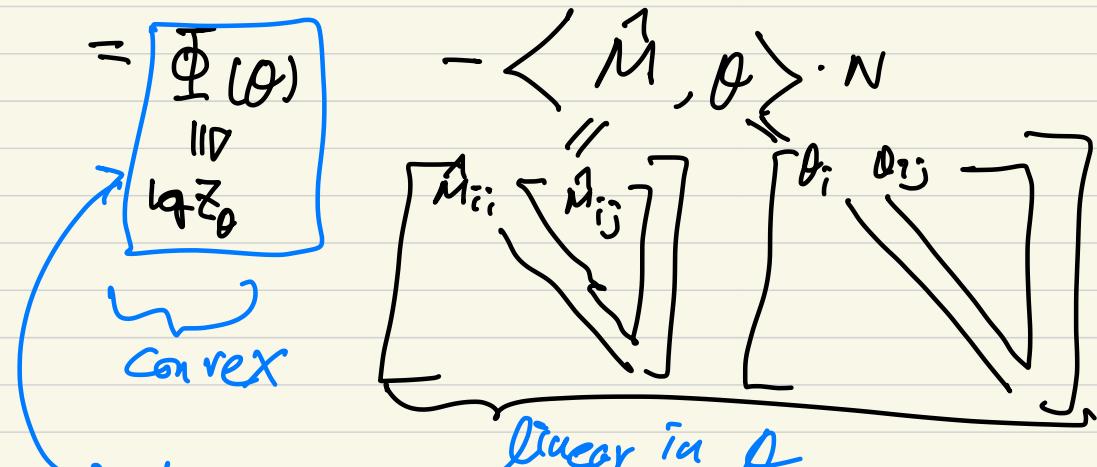
$$= \exp \left\{ -N \lg Z_\theta + \sum_{i \in V} N \cdot \hat{M}_{ii} \cdot \theta_i + \sum_{(i,j) \in E} N \cdot \hat{M}_{ij} \cdot \theta_{ij} \right\}$$

$$\hat{M}_{ii} \equiv \frac{1}{N} \sum_{l=1}^N x_i^{(l)}$$

$$\hat{M}_{ij} \equiv \frac{1}{N} \sum_{l=1}^N x_i^{(l)} x_j^{(l)}$$

the log-likelihood is

$$\min \quad L(G, \theta, D) = -\frac{1}{N} \lg P_{G,\theta}(D)$$



Remark:

involves exp-time

\* Learning is easier if inference is easier.