

* Overview

- Graphical Models & Markov Properties
- Inference Problems : given $P_G(x)$ find $\left\{ \begin{array}{l} P(x_i) \\ \arg \max_x P(x) \\ \text{compute } Z. \end{array} \right.$

- Belief Propagation
- Variational methods
- Gibbs sampling

- Learning Graphical models : given samples $x^{(1)}, x^{(2)}, \dots, x^{(n)} \in \mathcal{X}^n$
 - Learn the structure of the graph
 - Learn the parameters of the factors.

* How to find Z from a black-box that gives $P(x_i)$ from $P(x)$.

* Structural learning.

• $X \in \mathcal{X}^n$ Random Vector.

• directed graphical model: $|E| = \binom{n}{2} = \frac{n(n-1)}{2}$

of possible DAGs $\leq 3^{|E|} = 3^{\frac{n(n-1)}{2}}$

• Given $X^{(1)}, \dots, X^{(N)}$ independent 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

Bayesian

• Assume graph G and conditionals $P = \{P(x_i | x_{\pi_i})\}_{i=1}^n$ are deterministic but unknown

• Assume graph G and conditionals P are drawn from some known prior distribution $P_{G,P}(G,P)$

• Maximum Likelihood (ML) estimation finds (G, P) that maximizes log likelihood

• Maximum a Posteriori (MAP) estimation finds (G, P) that maximizes the posterior distribution

$$\max_{G, P} \sum_{j=1}^N \log P_{G, P}(X^{(j)})$$

$$\max_{G, P} P(G, P | X^{(1)}, X^{(2)}, \dots, X^{(N)})$$

* Frequentist's approach to structural learning

$$\hat{G} = \arg \max_G \max_P \frac{1}{N} \sum_{i=1}^N \log P_{G, P}(X^{(i)})$$

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

samples $(0,0), (0,1), (1,1), (0,0)$

empirical distribution: $\hat{P}_1(x_1) = \begin{cases} \frac{3}{4}, & x_1=0 \\ \frac{1}{4}, & x_1=1 \end{cases}$, $\hat{P}_2(x_2) = \begin{cases} \frac{1}{2}, & x_2=0 \\ \frac{1}{2}, & x_2=1 \end{cases}$

case 1: for $G_1 = \textcircled{1} \textcircled{2}$

the maximum likelihood estimate of $P_1(x_1), P_2(x_2)$ are

$$\begin{aligned} & \max_{P_1} \frac{1}{N} \sum_{j=1}^N \log P_1(X_1^{(j)}) \\ &= \max_{P_1} \hat{P}_1(0) \cdot \log P_1(0) + \hat{P}_1(1) \cdot \log P_1(1) \\ &= \max_{P_1} \underbrace{\sum_{x_1} \hat{P}_1(x_1) \log \hat{P}_1(x_1)}_{-H(\hat{P}_1)} + \underbrace{\sum_{x_1} \hat{P}_1(x_1) \log \frac{\hat{P}_1(x_1)}{P_1(x_1)}}_{-D_{KL}(\hat{P}_1 \| P_1)} \end{aligned}$$

is maximized when $P_1 = \hat{P}_1$ *empirical distribution is the maximum likelihood.

for G_1 ,

$$\begin{aligned} & \max_{P_1 \times P_2} \frac{1}{N} \sum_{j=1}^N \log P_1(x_1^{(j)}) P_2(x_2^{(j)}) \\ &= \underbrace{\sum_{x_1} \hat{P}_1(x_1) \cdot \log \hat{P}_1(x_1)}_{-H(\hat{P}_1)} + \underbrace{\sum_{x_2} \hat{P}_2(x_2) \log \hat{P}_2(x_2)}_{-H(\hat{P}_2)} \end{aligned}$$

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

$$\begin{aligned} & \max_{P_2(x_1, x_2)} \frac{1}{N} \sum_{j=1}^N \log P(x^{(j)}) = -H(\hat{P}_{12}) - \underbrace{D_{KL}(\hat{P}_{12} \| P_{12})}_{\substack{\text{maximum} \\ \text{achieved} \\ \text{with } \hat{P}_{12} = P_{12}}} \\ &= -H(\hat{P}_{12}) \end{aligned}$$

$$L(\textcircled{1} \textcircled{2}) = -H(\hat{P}_1) - H(\hat{P}_2)$$

$$L(\textcircled{1} \rightarrow \textcircled{2}) = -H(\hat{P}_{12})$$

• Remark 1: $-H(\hat{P}_{12}) \underline{\underline{>}} -H(\hat{P}_1) - H(\hat{P}_2)$

hence, blindly choosing a more likely model results in overfitting.

↑ even if the true distribution was independent, we always choose dependent model.

• Remark 2: depending on the sample size N and the target false positive rate β , decision is made by

$$L(\textcircled{1} \rightarrow \textcircled{2}) - L(\textcircled{1} \textcircled{2}) = H(\hat{P}_{12}) - H(\hat{P}_1) - H(\hat{P}_2) \\ \hat{=} I_{\hat{P}_{12}}(x_1, x_2)$$

output $\textcircled{1} \rightarrow \textcircled{2}$ if $I_{\hat{P}_{12}}(x_1, x_2) > \frac{t_\beta}{n}$
 $\textcircled{1} \textcircled{2}$ otherwise.

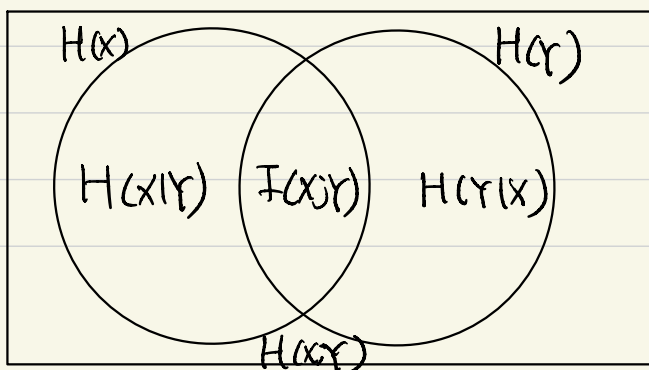
* We need to restrict the model class OR control false discovery rate

Refresh notations:

$$I(X;Y) \hat{=} \sum_{x,y} P(x,y) \log \frac{P(x,y)}{P(x)P(y)}$$

$$H(X) \hat{=} \sum_x -P(x) \log P(x)$$

$$H(Y|X) \hat{=} \sum_{x,y} -P(x,y) \log P(y|x)$$



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

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

* Maximum Likelihood Approach for a DAG.

$$G^* = \arg \max_G \max_{\{P_i(X_i | X_{\pi_i})\}} \frac{1}{N} \sum_{i=1}^N \log \prod_{i=1}^n P_i(X_i | X_{\pi_i})$$

the maximum is achieved at

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

↑ the empirical distribution

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

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

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

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

$$= \sum_{i=1}^n \left\{ \underset{\uparrow}{I_{\hat{P}_i}(X_i; X_{\pi_i})} - \underbrace{H_{\hat{P}_i}(X_i)} \right\}$$

↑ find G from a family of graphs that maximize this term

Does not depend on G .

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

we can now search over a class of graphs to find the best one

$I_{\hat{P}_i}(X_i; X_{\pi_i}) \geq I_{\hat{P}_i}(X_i; X_{\pi'_i})$ if $\pi_i \supset \pi'_i$
and hence denser graphs are preferred (and overfitted)

so we need an appropriate class of graphs to search over.

* Chow-Liu algorithm: searches over all trees., efficiently.

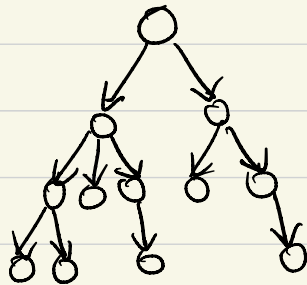
Step 1: Create a complete graph over $V = \{1, \dots, n\}$
with edge weights $I_p(X_i, X_j) = w_{ij}$

Step 2: Use Kruskal's algorithm, for example, to
find the max-weight spanning tree.

Claim: Chow-Liu algorithm finds the optimal tree that maximizes

$$\max_{G \in \text{Tree}} \sum_{i=1}^n I_{\hat{p}}(X_i; X_{\pi_i})$$

proof: each node only has one parent.



for general graphs, $n!$ orderings make it intractable.

*Another impractical approach for Undirected graph learning.

Consider learning an Ising Model (G, θ) , $\mathcal{X} = \{\pm 1\}$
 from samples $\{x^{(1)}, x^{(2)}, \dots, x^{(N)}\} = \mathcal{D}$

the likelihood is

$$\begin{aligned}
 P_{(G, \theta)}(\mathcal{D}) &= \prod_{l=1}^N P_{(G, \theta)}(x^{(l)}) \\
 &= \prod_{l=1}^N \frac{1}{Z_G(\theta)} \prod_{(i,j) \in E} e^{x_i^{(l)} x_j^{(l)} \theta_{ij}} \prod_{i \in V} e^{x_i^{(l)} \cdot \theta_i} \\
 &= \exp \left\{ -N \cdot \log Z_G(\theta) + \sum_{(i,j) \in E} N \cdot \hat{M}_{ij} + \sum_{i \in V} N \cdot \hat{M}_{ii} \right\} \\
 &\quad \quad \quad \parallel \quad \quad \quad \parallel \\
 &\quad \quad \quad \frac{1}{N} \sum_{l=1}^N x_i^{(l)} x_j^{(l)} \quad \quad \quad \frac{1}{N} \sum_{l=1}^N x_i^{(l)}
 \end{aligned}$$

the log-likelihood is

$$L(G, \theta, \mathcal{D}) = -\frac{1}{N} \log P_{(G, \theta)}(\mathcal{D})$$

$$\begin{aligned}
 &= \Phi(\theta) - \langle \hat{M}, \theta \rangle \\
 &\quad \uparrow \\
 &\quad \text{log-partition function} \\
 &\quad \quad \quad \left[\hat{m}_{11} \hat{m}_{12} \dots \right] \quad \quad \quad \left[\theta_1 \theta_2 \dots \right] \\
 &\quad \quad \quad \quad \quad \quad \quad \quad \quad \quad \quad \quad \uparrow \\
 &\quad \quad \quad \quad \quad \quad \quad \quad \quad \quad \quad \quad \text{non-zero for } (i,j) \in E \text{ \& } (i,i)
 \end{aligned}$$

Remark: this is strictly convex in θ ,
 but log-partition function requires inference.

* learning is easier when inference is easier.

but in general this is computationally intractable, even if
 the graph is sparse, requiring $O(|\mathcal{X}|^n)$ computations.

but continuing our (theoretical) investigation,
we want to apply this method to learn the structure
of the graph as follows

$$\underset{G}{\text{minimize}} \underset{\theta}{\text{minimize}} \mathcal{L}(G, \theta, D)$$

As we did previously, we need to restrict our search
to a class of "simple" graphs, as otherwise dense graphs
always win. A natural condition is $|E| \leq M$.

$$\underset{\theta}{\text{minimize}} \mathcal{L}(K_n, \theta, D) \quad , \text{ where } K_n \text{ is the complete graph}$$

s. t. $\|\theta\|_0 \leq M$

As $\|\theta\|_0$ constraint is intractable, people have proposed

$$\underset{\theta}{\text{minimize}} \Phi(\theta) - \langle \hat{M}, \theta \rangle + \lambda \cdot \|\theta\|_1$$

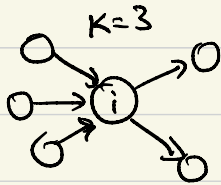
$$\text{where } \|\theta\|_1 = \sum_{i,j} |\theta_{ij}|$$

A different approach: Local Independence Test for undirected graphical models.

tries to take advantage of sparsity of the learned graph. to make learning faster than $O(n^4)$.

Alg 1: Local Independence Test (samples $\{X^{(t)}\}_{t=1}^N$, neighborhood size K)

- $E = \emptyset$
- For each $i \in V$
- For each $S \subseteq V \setminus \{i\}$ s.t. $|S| \leq K$
- Compute $\text{SCORE}(S, i) = H_{\hat{p}}(X_i | X_S)$
- Set $S^* \leftarrow \arg \min_S \text{SCORE}(S, i)$
- $E \leftarrow E \cup \{(i, j)\}_{j \in S^*}$
- Prune the resulting graph.



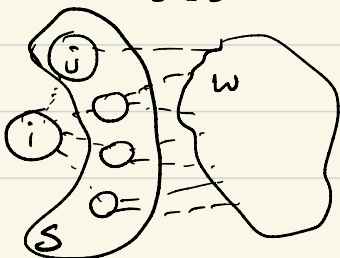
Remark 1: $X_i \perp\!\!\!\perp X_{\text{rest}} | X_S \iff$

$$H(X_i | X_S \cup X_{\text{rest}}) = H(X_i | X_S) \leq H(X_i | X_T) \quad \forall T \subseteq S.$$

Remark 2: Other scores have been proposed.

$$S^* = \arg \max \{ |S| : \epsilon < \text{SCORE}(S, i) \}$$

$$\text{SCORE}(S, i) = \min_{\substack{W \subseteq V \setminus S \\ j \in S}} \max_{\substack{x_i, x_w \\ x_s, x_j = a}} \left| \hat{p}\{X_i = x_i | X_w = x_w, X_S = x_s\} - \hat{p}\{X_i = x_i | X_w = x_w, X_{S \setminus j} = x_{S \setminus j}, X_j = a\} \right|$$



Remark: $\text{SCORE}(S, i)$ will be small if it includes any non-significant node.

These approaches are ^{still} of more theoretical interest, as the run-time is $O(N^{k+1})$.

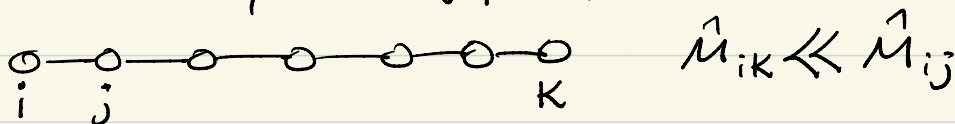
Here is a practical algorithm.

Alg 2: Thresholding (samples $\{x^{(l)}\}_{l=1}^N$, threshold τ)

- Compute the empirical correlation $\{\hat{M}_{ij}\}_{(i,j) \in V \times V}$
- For each $(i,j) \in V \times V$
- If $\hat{M}_{ij} \geq \tau$, set $(i,j) \in E$

where $\hat{M}_{ij} = \frac{1}{N} \sum_{l=1}^N (x_i^{(l)} - \bar{x}_i)(x_j^{(l)} - \bar{x}_j)$, $\bar{x}_i = \frac{1}{N} \sum_{l=1}^N x_i^{(l)}$

Remark: a heuristic based on the fact that two nodes faraway in the graph might be less correlated.



in general, this can fail if.

True graph learned graph

