\* Overview

**Lec 1-4**  · Graphical Models & Markov Properties

· Inference Problems : given $P_G(x)$ find
$$\begin{cases} P(x_i) \\ \arg\max_x P(x) \\ \text{compute } Z. \end{cases}$$

**Lec 5-11**  — Belief Propagation

**Lec 12-14**  — Variational methods

**Lec 15-16**  — Gibbs sampling

**Lec 17-19**  · Learning Graphical models : given samples $x^{(1)}, x^{(2)}, \cdots, x^{(N)} \in \mathcal{X}^n$

— Learn the structure of the graph

\* Structural learning.

- $X \in \mathbb{R}^n$ 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 $X^{(1)}, \cdots, X^{(N)}$ 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 $P = \{P(X_i \mid X_{\pi_i})\}_{j=1}^{n}$ are deterministic but unknown | · Assume $(G, P)$ is randomly drawn from some known prior distribution $Q_{G,P}$ |
| · Maximum Likelihood (ML) estimation finds $(G, P)$ that maximizes log likelihood | · Maximum a Posteriori (MAP) estimation, which |
| $\max_{G, P} \underbrace{\sum_{j=1}^{N} \log P_{G,P}(X^{(j)})}_{\text{SCORE}(G, P)}$ | $\max_{G, P} P(G, P \mid X^{(1)}, X^{(2)}, \cdots, X^{(N)})$ |

\* 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 $n=2$ , $x=(x_1, x_2) \in \{0,1\}^2$, $\hat{P}_{12}(x_1,x_2)=$

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

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

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

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

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

$\underbrace{\quad}_{\mathbb{I}(x_1^{(j)}=0)}$

definition of empirical distribution $\longrightarrow$ $= \max_{P_1=[-]} \boxed{\hat{P}_1(0) \log P_1(0)} + \hat{P}_1(1) \log P_1(1)$

$$= \max_{P_1} \sum_{x_1} \hat{P}_1(x) \log P_1(x)$$

adding/subtracting $H(\hat{P}_1)$ $\longrightarrow$ $= \max_{P_1} \underbrace{\sum_{x_1} \hat{P}_1(x) \log \hat{P}_1(x)}_{-H_{\hat{P}_1}(X_1)} + \underbrace{\sum_{x_1} \hat{P}_1(x_1) \log \frac{P_1(x_1)}{\hat{P}_1(x_1)}}_{-D_{KL}(\hat{P}_1 \| P_1)}$

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

for $G_1, \textcircled{1} \textcircled{2}$ $\max_{P_1 \times P_2} \frac{1}{N}\sum_{j=1}^{N} \log P_1(x_1^{(j)}) P_2(x_2^{(j)})$ $\leftarrow$ max achieved with $\begin{cases} P_1 = \hat{P}_1 \\ P_2 = \hat{P}_2 \end{cases}$

$\text{SCORE}(G_1) = \underbrace{\sum_{x_1} \hat{P}_1(x_1) \cdot \log \hat{P}_1(x)}_{-H(\hat{P}_1)} + \underbrace{\sum_{x_2} \hat{P}_2(x_2) \log \hat{P}_2(x_2)}_{-H(\hat{P}_2)}$

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

$$\max_{P_{12}(x_1,x_2)} \frac{1}{N}\sum_{j=1}^{N} \log P_{12}(x^{(j)}) = -H(\hat{P}_{12}) - \underbrace{D_{KL}(\hat{P}_{12} \| P_{12})}_{}$$

maximum achieved with $\hat{P}_{12} = P_{12}$

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

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

$$\mathcal{L}(① \rightarrow ②) = -H(\hat{P}_{12}) \qquad H(\hat{P}_{12}) \leqq H(\hat{P}_1 \cdot \hat{P}_2)$$

<span style="color:blue">independt</span>

· Remark 1: $\qquad -H(\hat{P}_{12}) \geqq -H(\hat{P}_1) - H(\hat{P}_2)$

<span style="color:blue">$G_2 > G_1$</span>

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 $\beta$, decision is made by

$$\mathcal{L}(① \rightarrow ②) - \mathcal{L}(① \quad ②) = H(\hat{P}_{12}) - H(\hat{P}_1) - H(\hat{P}_2)$$

<span style="color:blue">↑ depart $X_1, X_2$</span>     <span style="color:blue">↑ $X_1 \perp X_2$</span>    $\triangleq \boxed{I_{\hat{P}_{12}}(X_1; X_2)}$

output $① \rightarrow ②$ if $\quad I_{\hat{P}_{12}}(X_1; X_2) > \dfrac{t_\beta}{N}$

$① \quad ②$    otherwise.

<div style="border:2px solid red; padding:8px; color:blue">
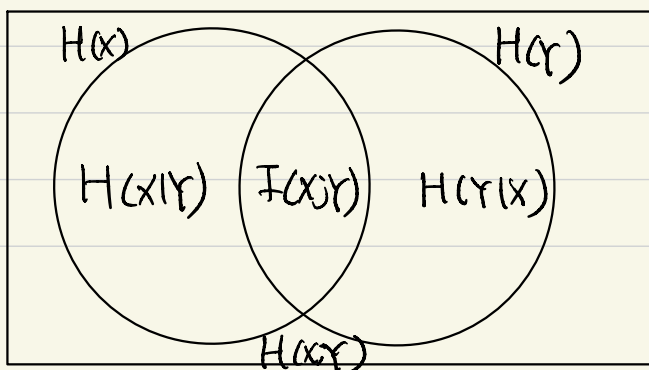* We need to restrict the model class OR control false discovery rate
</div>

Refresh notations:

$$I(X;Y) \triangleq \sum_{X,Y} P(X,Y) \log \frac{P(X,Y)}{P(X) \cdot P(Y)}$$

$$H(X) \triangleq \sum_X -P(X) \log P(X)$$

$$H(Y|X) \triangleq \sum_{X,Y} -P(X,Y) \log P(Y|X)$$



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

**\*Maximum Likelihood Approach for a DAG.**

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

<span style="color:blue">the maximum is achieved at</span>

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

<span style="color:blue">↑ the empirical distribution</span>

<span style="color:blue">SCORE(G)</span>

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

<span style="color:blue">Sufficient statistics.</span>    <span style="color:blue">choice G</span>
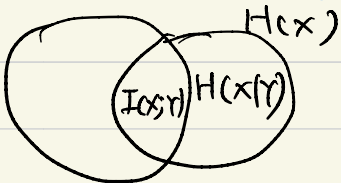
$$= \sum_{i=1}^{n} \frac{1}{N}\sum_{j=1}^{N} \log \hat{P}(X_i|X_{\pi_i})$$

<span style="color:blue">def. empirical distribution →</span>
$$= \sum_{i=1}^{n} \left\{ \sum_{X_i, X_{\pi_i}} \hat{P}_i(X_i, X_{\pi_i}) \cdot \log \hat{P}_i(X_i|X_{\pi_i}) \right\}$$

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

<span style="color:blue">H(X)</span>
<span style="color:blue">=H(X|Y)+I(X;Y)</span>

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

<span style="color:blue">find G from a family of graphs that maximize this term</span>    <span style="color:blue">Does not depend on G.</span>

H(x)

(I(x;y)) H(x|y)

---

**\*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.

- $I_{\hat{P}}(X_i; X_{\pi_i}) \geq I_{\hat{P}}(X_i; X_{\pi_i'})$   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 undirected graph over $V = \{1, \dots, n\}$ with edge weights $I_{\hat{p}}(X_i, X_j) = w_{ij}$
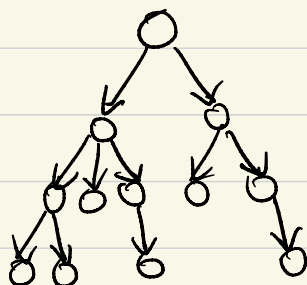
Step 2: Use Kruskal's algorithm, for example, to find the <u>max-weight spanning tree</u>.

Claim: Chow-Liu algorithm finds the <u>optimal</u> tree that maximizes

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

$\underset{\text{single node.}}{\overset{\shortparallel}{}}$

proof: each node only has one parent.



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

\* Another impractical approach for <u>Undirected graph</u> learning.

Consider learning an <u>Ising Model</u> $(G, \theta)$, $\mathcal{X} = \{\pm 1\}$
from samples $\{X^{(1)}, X^{(2)}, \dots, X^{(N)}\} = D$ $\quad \{\theta_i\}_{i \in V}$
$\mathcal{X}^n$ $\qquad \{\theta_{ij}\}_{(i,j) \in E}$

the <u>likelihood</u> is

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

$$= \prod_{\ell=1}^{N} \frac{1}{Z_G(\theta)} \prod_{i \in V} e^{\theta_i \cdot X_i^{(\ell)}} \cdot \prod_{(i,j) \in E} e^{\theta_{ij} X_i^{(\ell)} X_j^{(\ell)}}$$

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

$$\underset{\frac{1}{N}\sum_{\ell=1}^{N} X_i^{(\ell)}}{\overset{\shortparallel}{\phantom{=}}} \qquad \underset{\frac{1}{N}\sum_{\ell=1}^{N} X_i^{(\ell)} X_j^{(\ell)}}{\overset{\shortparallel}{\underset{G}{\phantom{=}}}}$$

the <u>log-likelihood</u> is

min $\quad L(G, \theta, D) = -\frac{1}{N} \log P_{(G,\theta)}(D)$ $\qquad\qquad \langle M, \theta \rangle = \sum_{ij} M_{ij} \theta_{ij}$

$$= \Phi(\theta) - \langle \hat{M}, \theta \rangle$$

$\log Z_G(\theta)$ $\quad$ log-partition function

$$n \left\{ \begin{bmatrix} \hat{M}_{11} & \hat{M}_{12} & \cdots \\ & \hat{M}_{22} & \\ & & \ddots \end{bmatrix} \right. \qquad \begin{bmatrix} \theta_1 & \theta_{12} & 0 \cdots 0 \\ & \theta_2 & 0 \\ 0 & 0 & \ddots \\ 0 & 0 & \end{bmatrix}$$

$\underbrace{\phantom{xxxxxx}}_{n}$ $\qquad\qquad$ non-zero for
$\qquad\qquad\qquad\qquad (i,j) \in E \, \& \, (i,i)$

Remark: this is <u>strictly convex</u> in $\theta$,
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 $|E| \leq m$.

$$\underset{\theta}{\text{minimize}} \ \mathcal{L}(K_n, \theta, D) \qquad , \text{ where } K_n \text{ is the}$$
$$\text{complete graph}$$
$$\text{s. t.} \quad \|\theta\|_0 \leq m \qquad \|\theta\|_{L_0} = \sum_{i,j} \mathbb{I}(\theta_{ij} \neq 0)$$

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

$$\underset{\theta}{\text{minimize}} \ \underbrace{\Phi(\theta) - \langle \hat{M}, \theta \rangle}_{\mathcal{L}(K_n, \theta, D)} + \lambda \cdot \|\theta\|_1$$
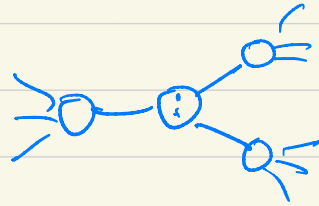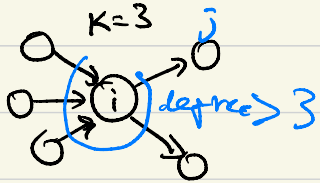
where $\|\theta\|_1 = \sum_{i,j} |\theta_{ij}|$

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 $O(n^K)$.

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

- $E = \emptyset$
- For each $i \in V$
-    For each $S \subseteq V \setminus \{i\}$   s.t. $|S| \leq K$    $O(\binom{n}{K}) = O(n^K)$
-       Compute $SCORE(S,i) = H_{\hat{P}}(X_i | X_S)$
-    Set $S^* \leftarrow \arg\min_{S} SCORE(S,i)$
-    $E \leftarrow E \cup \{(i,j)\}_{j \in S^*}$

- Prune the resulting graph.



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

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

**Remark 2:** Other scores have been proposed.

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

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



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

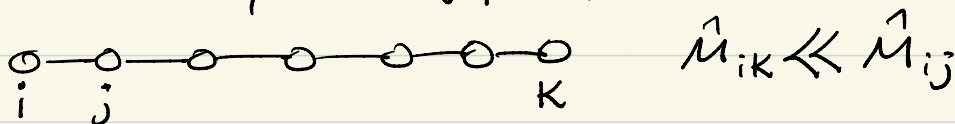These approaches are still of more theoretical interest, as the run-time is $O(n^{k+1}) \ll O(|\mathcal{X}|^n)$

Here is a practical algorithm.

Alg 2: Thresholding ( samples $\{X^{(l)}\}_{l=1}^{N}$, threshold $\tau$ )
- Compute the empirical correlation $\{\hat{M}_{ij}\}_{(i,j) \in V \times V}$
- For each $(i,j) \in V \times V$    $O(n^2)$
- If $|\hat{M}_{ij}| \geq \tau$,   set $(i,j) \in E$

where   $\hat{M}_{ij} = \frac{1}{N} \sum_{l=1}^{N} \left( X_i^{(l)} - \overline{X_i} \right) \left( X_j^{(l)} - \overline{X_j} \right)$ ,    $\overline{X_i} = \frac{1}{N} \sum_{l=1}^{N} X_i^{(l)}$

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



$\hat{M}_{ik} \ll \hat{M}_{ij}$

in general, this can fail if.

True graph      learned graph