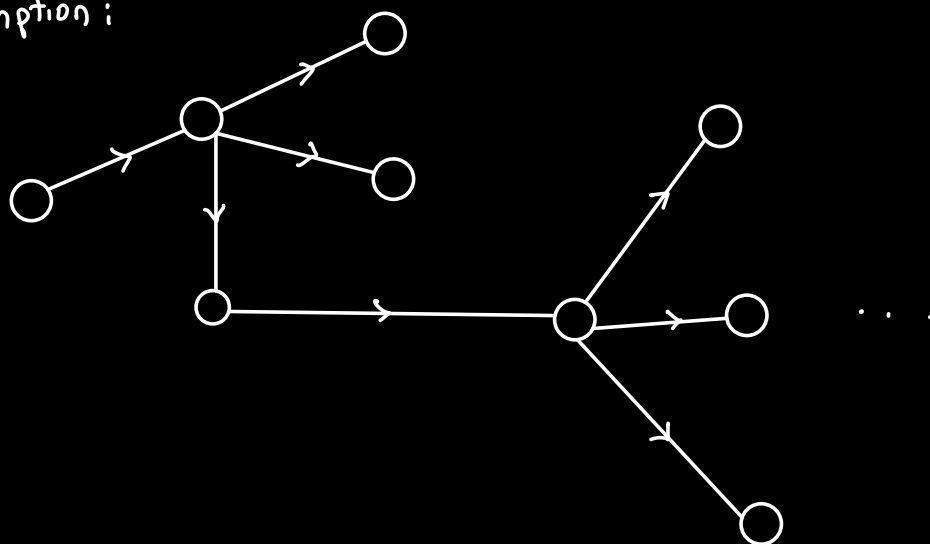


Learning Bayesian networks defined by trees.

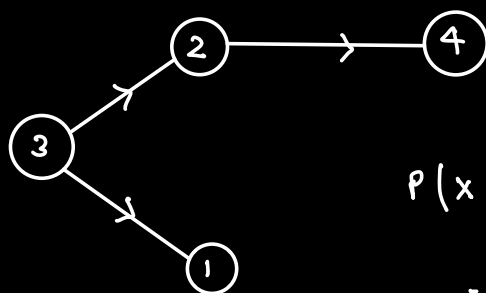
Distribution: $(x_1, x_2, \dots, x_d) \in \mathcal{X}^d \quad (\{0,1\}^d)$

Assumption:



Rooted tree; Each node (except root) has exactly one parent.

Example:



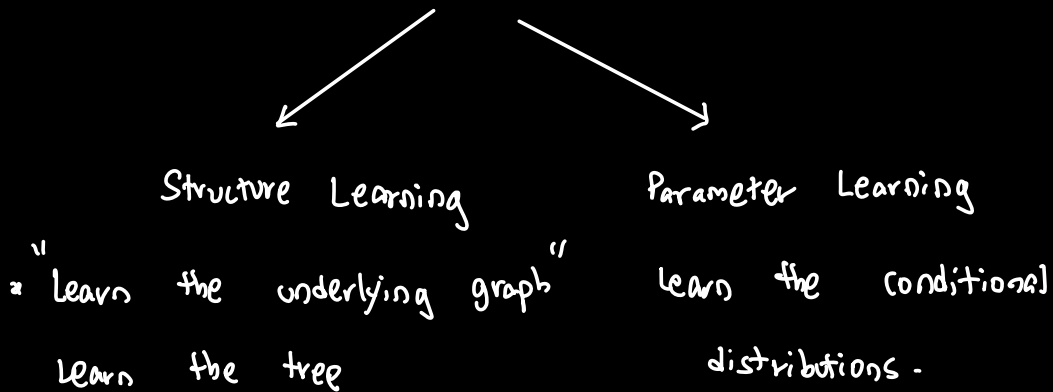
$$P(x = x_1, x_2, x_3, x_4)$$

$$= P(x_3 = x_3) \cdot \Pr(x_2 = x_2 | x_3 = x_3)$$

$$\cdot \Pr(x_1 = x_1 | x_3 = x_3)$$

$$\cdot \Pr(x_4 = x_4 | x_2 = x_2)$$

Suppose we see samples from a Bayes net as above, can you learn the distribution.



Goal:

Let us try to learn a distribution D' such that $\text{dist}(D, D')$ is small.

KL - DIVERGENCE:

"measures distance between distributions"

We have two distributions p, q over some space Ω .

(cross entropy)

$$KL(p \parallel q) = \sum_{s \in \Omega} p(s) \log \frac{p(s)}{q(s)}$$

↓
Probability s
happens under p

↘ Probability s happens
under q .

EXAMPLE:

$$\Omega = \{0, 1\}$$

$$p(0) = 1/2, \quad p(1) = 1/2$$

$$q(0) = 3/4, \quad q(1) = 1/4$$

$$K_L(p|q) = p(0) \cdot \log \frac{p(0)}{q(0)} + p(1) \cdot \log \frac{p(1)}{q(1)}$$

$$= \frac{1}{2} \cdot \log\left(\frac{2}{3}\right) + \frac{1}{2} \cdot \log(2)$$

$$= \frac{1}{2} \cdot \log\left(\frac{4}{3}\right)$$

What is $K_L(p|p)$?

$$\rightarrow 0$$

Property:

$$K_L(p|q) \geq 0 \quad \text{and} \quad 0 \Leftrightarrow p = q.$$

NOT SYMMETRIC!

INPUT : Samples from a distribution P on \mathcal{Z}^d ($\mathcal{Z} = \{0,1\}$)
 (being generated by some unknown Bayes net which
 is a tree : T^*)

OUTPUT : Find a tree T and a corresponding Bayes net
 P_T such that $KL(P|P_T) \leq \epsilon$.

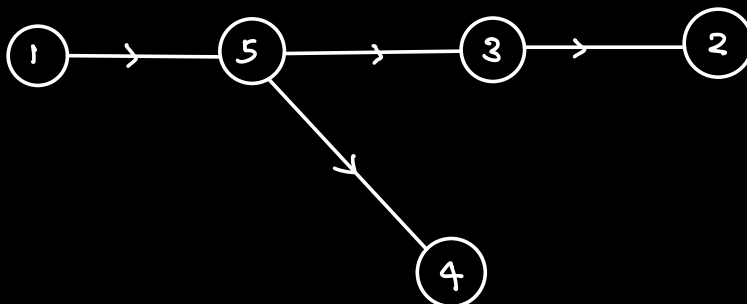
EXAMPLE :



$(x_1, x_2, x_3, x_4, x_5)$

If I give a tree T and ask find me the

Bayes net P_T that minimizes $KL(P|P_T)$.



Idea: well, use the "empirical" conditional probabilities.

1. Estimate the conditional probabilities along the tree T .

(e.g.: Estimate $\Pr[x_1=1] \cdot \Pr[x_5=1|x_1=0]$, $\Pr[x_5=1, x_1=1]$
:

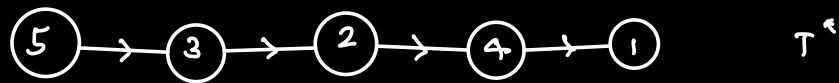
INPUT: Samples from P and a tree T .

Algorithm: For each edge (i, j) in T , estimate $\Pr[x_j|x_i]$ and use these to define P_T .

CHOW - LIU ALGORITHM:

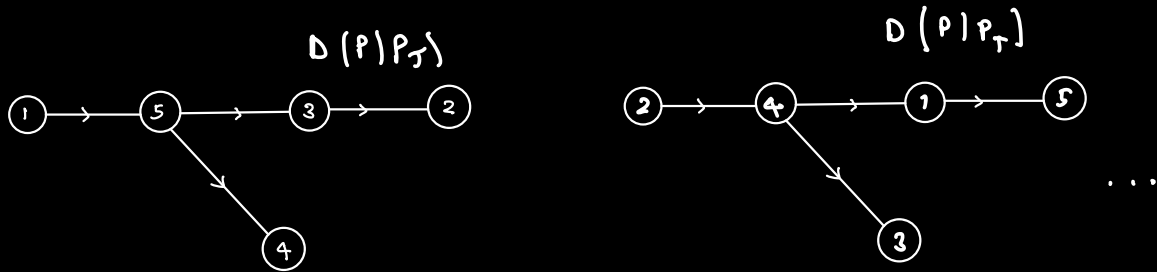
→ For any tree T ,

$$\begin{aligned} KL(P|P_T) &= J_P - \sum_{(i,j) \text{ is an edge in } T} I(x_i; x_j) \\ &\downarrow \\ &\text{Some number depending} \\ &\text{on } P \end{aligned}$$



$(x_1, x_2, x_3, x_4, x_5)$

Possible trees:



Given a tree T , we can find P_T .

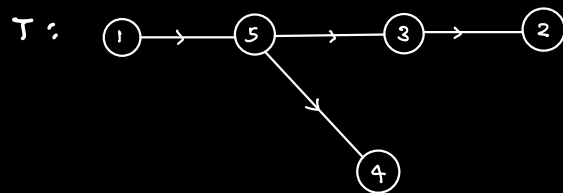
$I(x_i; x_j) \equiv$ "Measures how much information about x_j does x_i have".

$$I(x_i; x_j) = \sum_{a,b} \Pr[x_i = a, x_j = b] \cdot \log \frac{\Pr[x_i = a] \Pr[x_j = b]}{\Pr[x_i = a, x_j = b]}$$

we can estimate $I(x_i, x_j)$ from sample.



$(x_1, x_2, x_3, x_4, x_5)$



Chow-Liu Bound (1968):

$$KL(p|p_T) = J_p - I(x_1; x_5) - I(x_5; x_3) - I(x_5; x_4) \\ - I(x_3; x_2).$$

Summary: Find T that minimizes $KL(p|p_T)$ or find

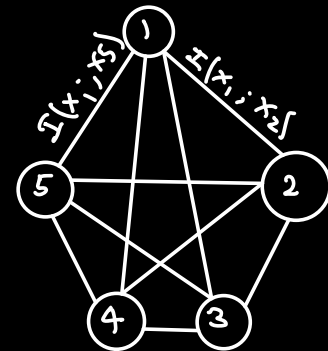
$$\arg \min_T J_p - \sum_{(i,j) \in T} I(x_i; x_j)$$

$$\equiv \arg \min_T \sum_{(i,j) \in T} -I(x_i; x_j)$$

$$\equiv \arg \max_T \sum_{(i,j) \in T} I(x_i; x_j)$$

Idea:

We have samples from $P \rightsquigarrow$



Maximum Spanning Tree.

CHOW - LIU Algorithm:

- Use samples to estimate $I(x_i; x_j)$ for all i, j .
- Form a weighted graph where weights are exactly $I(x_i; x_j)$.
- Compute the maximum spanning tree T' in G .
- Output $P_{T'}$.

Maximum Spanning Tree: Given a weighted graph G on d vertices, find tree T that has maximum total weight. - We can do this very fast. We can do this in linear time.

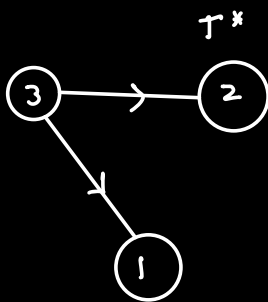
$(x_1, x_2, \dots, x_5) :$ $\square \rightarrow \text{compute } \mathcal{I}(x_1; x_2)$

$\square \rightarrow \text{compute } \mathcal{I}(x_1; x_3)$

$\square \rightarrow \text{compute } \mathcal{I}(x_1; x_4)$

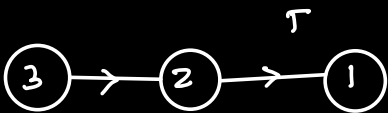
\vdots

EXAMPLE:



True distribution, P

$$P[x = x_1, x_2, x_3] = P[x_3 = x_3] \cdot P[x_2 = x_2 | x_3 = x_3] \\ \cdot P[x_1 = x_1 | x_3 = x_3]$$



$$P_T[x = x_1, x_2, x_3] = P[x_3 = x_3] \cdot P[x_2 = x_2 | x_3 = x_3] \cdot P[x_1 = x_1 | x_2 = x_2]$$

\downarrow
 P_T

SUMMARY: "DIRECTED GRAPHICAL MODELS"

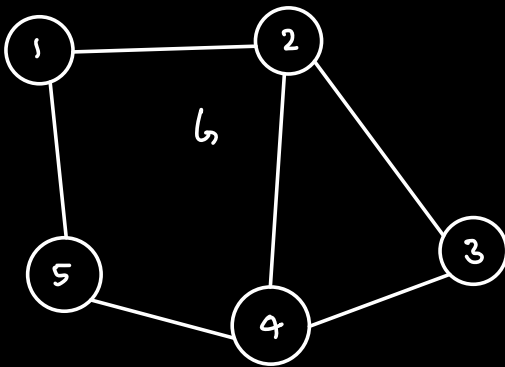
- Can learn tree-structured Bayes networks
- Active research: what are other structural assumptions
Can we learn Bayes networks.

UNDIRECTED GRAPHICAL MODELS:

- "Markov Random Fields" (MRFs).

Distribution: $D \equiv (x_1, x_2, \dots, x_d)$

Dependency graph G for D .

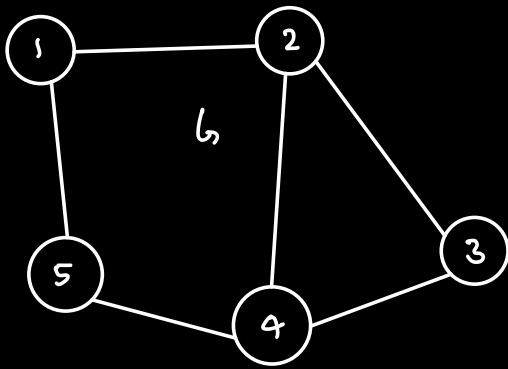


eg: $x_1 \perp x_4 \mid x_2, x_3, x_5$

D satisfies "Pairwise Markov Property" with respect to G if

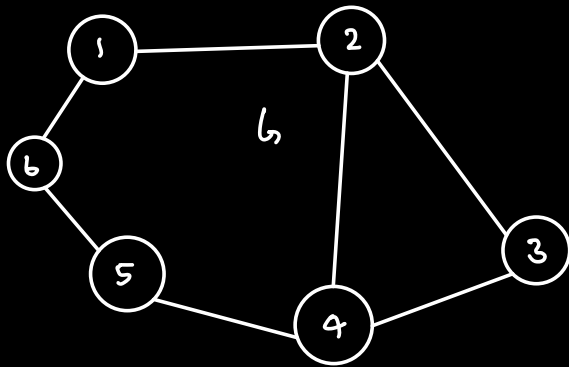
if i, j have no edge,

then $x_i \perp x_j \mid x_{\text{rest}}$



eg: $x_1 \perp x_4 \mid x_2, x_5$

Δ satisfies "Local Markov Property" with respect to b if
 If i, j have no edge,
 then $x_i \perp x_j \mid x_{\{\text{neighbors of } i\}}$



eg: $x_1 \perp x_4 \mid x_2, x_5$

Δ satisfies "Global Markov Property" with respect to b if
 If i, j have no edge,
 then $x_i \perp x_j \mid x_{\{\text{any separating set}\}}$

subset of vertices removing which disconnects i and j in b .

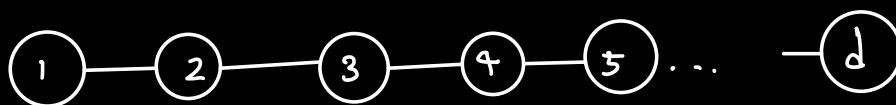
Global MP \Rightarrow Local MP \Rightarrow Pairwise MP

For "most reasonable" distributions all three are equivalent.

Example:

Markov Chain

$$x_1, x_2, \dots, x_d \quad x_{i+1} \perp x_{i-1} \mid x_i$$



Remark: We will say D has dependency graph G if it satisfies Markov property with respect to G .

MAIN LEARNING CHALLENGES:

Structure Learning: Given samples from D learn its dependency graph G .

Parametric Learning: Given samples from D learn the full distribution.

Inference: You know dependency graph want to find most likely value of $x_i \mid x_{\{\text{partial assignment}\}}$

MAP

Goal: Structure Learning

INPUT: Samples $x^1, x^2, \dots, x^n \sim D$

OUTPUT: Dependency graph of D .

ILL-POSED: Complete graph is a dependency graph
for all distributions!

"Minimal Dependency graph?"

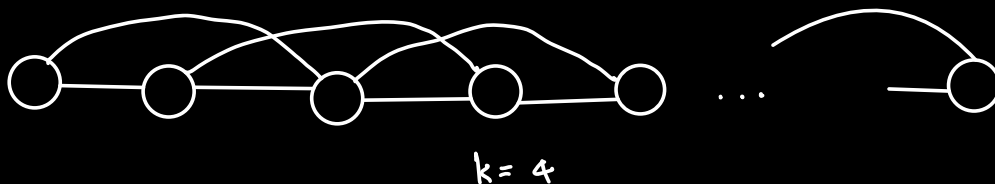
ASSUMPTION:

Dependency graph of D is sparse.

- Each vertex has a bounded degree k .

INPUT: Samples $x^1, x^2, \dots, x^n \sim D$

OUTPUT: A dependency graph for D where each
vertex has degree $\leq k$.



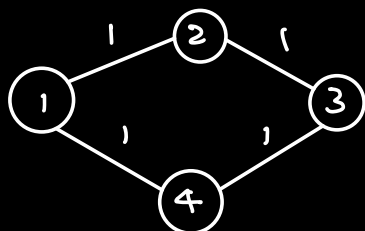
→ If we know the structure, the general learning problem becomes easier.

LEARNING BOLTZMANN MACHINES



D on $\{1, -1\}^d$

$$Pr[x = x] \propto \exp\left(\sum_{\{i,j\} \in G} w_{ij} x_i x_j\right)$$



$$Pr[x = x] \propto \exp(x_1 x_2 + x_2 x_3 + x_3 x_4 + x_4 x_1)$$

Distributions as defined above,
they satisfy Markov property
with respect to G .

GAUSSIAN GRAPHICAL MODELS



D on \mathbb{R}^d

Distribution D is
 $N(0, \Sigma)$

Σ is the covariance
matrix.

$$x \sim N(0, \Sigma)$$

$$\Sigma_{ij} = E[x_i x_j]$$

$\Sigma_{ij} = 0 \Rightarrow x_i, x_j$ are
independent.

(Dempsey 1972):

Precision Matrix

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

THEOREM: Gaussian distribution
has dependency graph

with support (Θ) .

