# GRAPHICAL MODELS:

## CONDITION INDEPENDENCE (CI):

$$\boxed{x \perp y \mid z}$$

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

Structure: Rooted Tree.

## KL - DIVERGENCE / CROSS - ENTROPY:

Distance between 2 distributions over some space, $\Omega \rightarrow d_1, d_2$

$$\boxed{KL(d_1 \| d_2) = \sum_{s \in \Omega} d_1(s) \log \frac{d_2(s)}{d_1(s)}}$$

$d_1(s)$: Probability s happens under $d_i$.

Goal: Given distribution $P$ on $\Sigma^d$ generated from unknown Bayes net (Tree: $T^*$), find the $T$ and Bayes Net $P_T$ such that

$$KL(P, P_T) \leq \varepsilon.$$

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

# Chow-Liu Bound:

For any tree T,

$$k_L(P \| P_T) = J_P - \sum_{(i,j) \text{ is an edge in } T} I(x_i, x_j)$$

$\underbrace{\phantom{J_P}}$ function of P

## Mutual Information:

$I(x_i; x_j) \equiv$ Measures how much information $x_i$ has about $x_j$.

$\Downarrow$

can be estimated from samples [Independent of T].

## 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 max. spanning tree $T'$ in G

→ Output $P_{T'}$.

# UNDIRECTED GRAPHICAL MODELS:

## MARKOV RANDOM FIELDS:

Given $D = (x_1, x_2 \ldots x_d) \quad \{0,1\}^d$

$G:$ Dependency Graph for $D$.

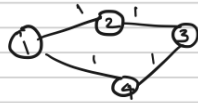| D satisfies | Pairwise Markov Property | Local | Global |
|---|---|---|---|
| with respect to $G$ <br><br> If $i,j$ have no edge, then | $x_i \perp x_j \mid x_{rest}$ | $x_i \perp x_j \mid$ $x_{\{neighbors\ of\ i\}}$ | $x_i \perp x_j \mid$ $x_{\{any\ separating\ set\}}$ $\downarrow$ vertices removing which $i,j$ disconnected. |

Global $\Rightarrow$ Local $\Rightarrow$ Pairwise.

**Goal:** Find $D$ such that each vertex has degree $\leq k$.

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

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



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

Distributions as defined above, They satisfy Markov property wrt $G$.

---

Gaussian Graphical Models

$D$ on $\mathbb{R}^d$

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

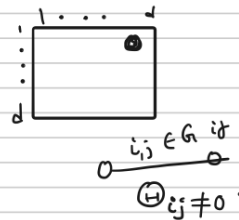$\Sigma$ is the covariance matrix.  $X \sim N(0, \Sigma)$

$\Sigma_{ij} = \mathbb{E}[X_i X_j]$

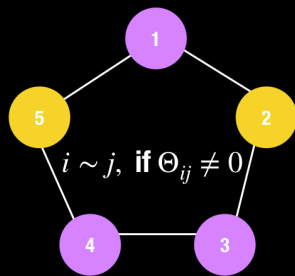$\Sigma_{ij} = 0 \implies X_i, X_j$ are independent.

(Dempsey 1972):

Precision Matrix:

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

Thm: Gaussian dist has dependency graph with $Supp(\Theta)$.



$i, j \in G$ if $\Theta_{ij} \neq 0$.

**Example:** $(X_1 \,|\, X_2, X_5)$ independent of $(X_3 \,|\, X_2, X_5)$

$i \sim j$, if $\Theta_{ij} \neq 0$

**Markov property:** $\Theta_{ij} = 0 \Rightarrow X_i, X_j$ are independent conditioned on neighbors of $i$.
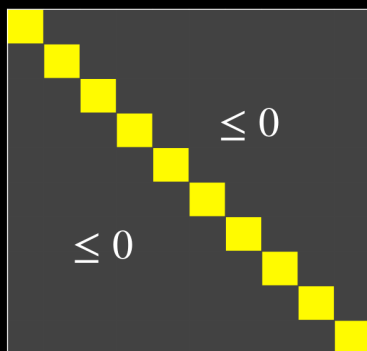
# Structure Learning for GGMs

Given samples $X^1, X^2, \ldots, X^n$ from a GGM of degree $d \ll p$, can we efficiently find the dependency graph with $n \ll p$?

(**Think:** $n = O_d(\log p)$.)

# Attractive GGMs

GGM is attractive if all covariances are non-negative.

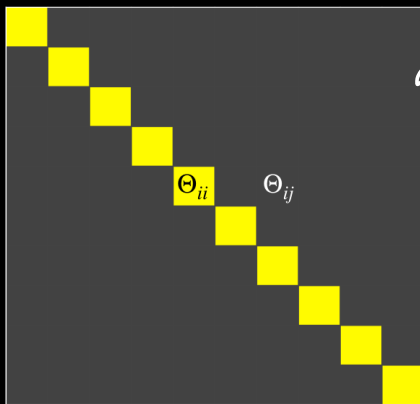(Equivalently, $\Theta$ has non-positive off-diagonals.)



$\leq 0$

$\leq 0$

$\Theta$ : Precision Matrix

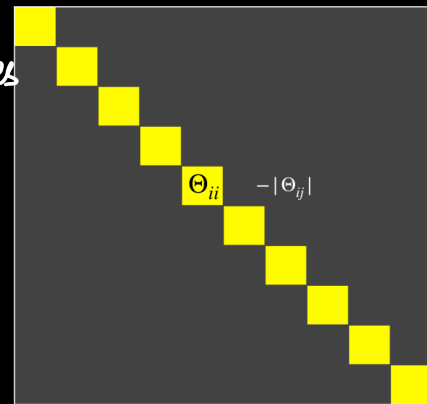## Ex: Gaussian Free Fields

- Many applications via Gaussian processes

# Walk-Summable GGMs

GGM walk-summable if making off-diagonals of precision matrix negative preserves positive semi-definiteness.

all eigenvalues > 0

$\Theta_{ii}$     $\Theta_{ij}$

$\Theta_{ii}$     $-|\Theta_{ij}|$

$\Theta$ : Precision Matrix

Offdiagonals negative $\geq 0$

---

## GREEDYPRUNE

1. **Recover neighborhood of each vertex in parallel.**

2. **Grow a candidate neighborhood.**

3. **Prune out some vertices.**

---

## GREEDY-GROWING

1. Set $S \leftarrow \emptyset$

2. While S is small enough:

   1. Find **j** to minimize estimate of $Var(X_1 | X_{S \cup j})$.

   2. $S \leftarrow S \cup \{j\}$.

Intuition: Add vertex that gives maximum decrease in conditional variance.

$\longrightarrow$ That $j$ and $x_1$ are clearly dependent as $j$ minimizes $x_1$'s variance.

$\hookrightarrow$ For each $var(X_1 | X_j)$ we need $\frac{1}{\epsilon^2}$ samples

to get $(1-\epsilon)$ accuracy.

If $var(X_1 | X_2, X_3, X_4) \rightarrow \frac{3}{\epsilon^2}$

---

**GREEDY-PRUNING**

1. **For each j in S:**

   1. **If** $Var(X_1 | X_{S \setminus \{j\}}) < (1 + \tau)Var(X_1 | X_S)$, **drop j from S.**

---

Intuition: If dropping a vertex, does not hurt
too much, drop it.

Can learn Attractive and Walk-Summable GGMs

with $O\left(d^2 \log p / \kappa^b\right)$ samples and quadratic

run-time.

d - degree (max)

p - total number of parameters

$$\kappa(\Theta) = \min_{i,j:\Theta_{ij} \neq 0} \frac{|\Theta_{ij}|}{\sqrt{\Theta_{ii}\Theta_{jj}}}$$

# Summary:

1. KL - Divergence:

$$KL(d_1 \| d_2) = \sum_{s \in \Omega} d_1(s) \log \frac{d_2(s)}{d_1(s)}$$

2. Chow-Liu Bound:

$$KL(P \| P_T) = J_P - \sum_{(i,j) \text{ is an edge in } T} I(x_i, x_j)$$

3. Chow-Liu Algorithm:
   - Find $I(x_i, x_j)$ - graph
   - Max Spanning Tree $\to T'$

4. Markov Random Fields:

   D (Markov Property) with respect to $G$: if no edge $i,j$

   $$x_i \perp x_j \mid x_{\{neighbors \text{ of } i\}}$$

5. Learning Boltzmann Machines: D on $\{1,-1\}^d$

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

6. Gaussian Graphical models:

   $$D: N(0, \Sigma) \qquad \Sigma_{ij} = E[x_i x_j]$$

   $$\Sigma_{ij} = 0 \implies \text{independent}$$

   $$\textcircled{H} = \Sigma^{-1}$$

   $$\text{support}(\textcircled{H}) \to G \qquad i,j \in G \text{ if } \textcircled{H}_{ij} \neq 0$$

7. Attractive: $\Sigma_{ij} > 0$, $(H)_{ij} \leq 0$ $\forall i \neq j$

8. Walk Summable: $(H)_{ij} \rightarrow -(H)_{ij}$ $\forall i \neq j$ => Still positive Semi Definite.

9. Greedy Prune

For each vertex $i \rightarrow$ find neighborhood in parallel

$\rightarrow \arg\min_j var(x_i \mid x_{S \cup j})$

$\rightarrow var(x_i \mid x_{S - \{j\}}) < (1 + \tau) \, var(x_i \mid x_S)$

$\rightarrow$ learns 2 types $O(d^2 \log p / k^b)$ samples, quadratic time

$d$ - max degree

$p$ - number of parameters

$$k(H) = \min_{i,j \, : \, (H)_{ij} \neq 0} \frac{|\Theta_{ij}|}{\sqrt{(H)_{ii} (H)_{jj}}}$$