

Today

- What are graphical Models ?
- Bayesian Networks : Definition, Examples, Learning challenges, Chow-Liu Algorithm.

1. The class on 05/25/2022 will be remote on zoom.
2. Final exam:
 - Covers all topics until the exam
 - Easier than assignments .
 - In-person exam as per university schedule
 - MSOL students online.

Unsupervised Learning; Given a dataset, we want to build a "model" for the dataset.

PROBABILISTIC MODELS OF DATA

- Graphical models are probabilistic models for generating data
- Precursors for modern "deep" generative models

Cove issues : ① Succinct representations of distributions

② Modeling dependencies.

Applications : ① "Generate" new examples.

② "In-Painting"

③ Applied Sciences useful in identifying relations.

Independence: (X, Y) random variables (joint dist.)

$$X \perp Y \Leftrightarrow P_{\text{r}}[X=x \wedge Y=y] = P_{\text{r}}[X=x] \cdot P_{\text{r}}[Y=y]$$
$$\Leftrightarrow P_{\text{r}}[Y=y | X=x] = P_{\text{r}}[Y=y].$$

Weather today is independent of Stock market.

	Heart Problems	Wakes up	Age
P ₁	NO	11AM	19
P ₂	YES	6:30AM	45
P ₃	YES	7:AN	40
P ₄	NO	12PM	24
P ₅	NO	10AM	25
	:	:	
P ₁₀	YES	6:AM	59
P ₁₁	YES	6:30	62
	:	:	

Comparing features directly:

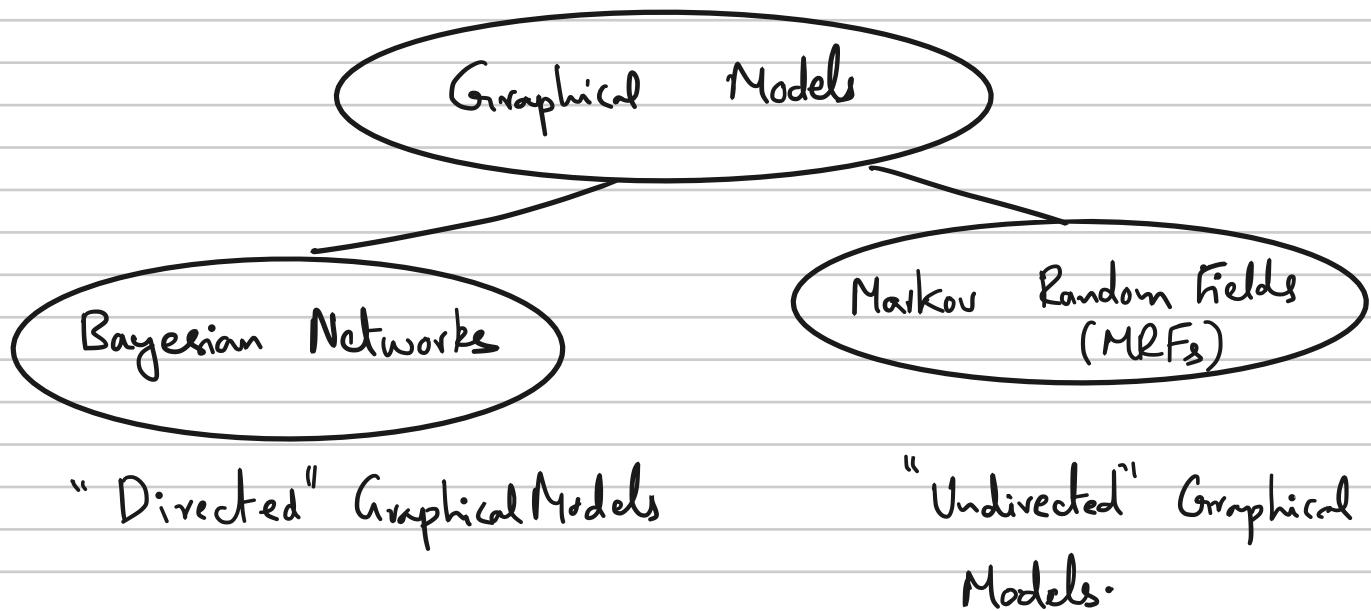
Waking up early
is bad for heart

Condition of Independence: $\forall i \forall j$ random variables

conditional independence: X, Y, Z random variables

$$\begin{aligned} X \perp Y \mid Z &\Leftrightarrow \Pr[X=x, Y=y \mid Z=z] \\ &= \Pr[X=x \mid Z=z] \cdot \Pr[Y=y \mid Z=z] \\ &\Leftrightarrow \Pr[Y=y \mid X=x, Z=z] = \Pr[Y=y \mid Z=z]. \end{aligned}$$

Graphical Models are distributions with
"constrained" conditional independence ("(I)"
relations.)



Bayesian Networks

$$X \in \Sigma^d \quad (x_1, x_2, \dots, x_d)$$

A Directed Acyclic Graph (DAG) on
[d] vertices: G

A distribution D is a Bayes Net with
graph G

$$P(x_1, x_2, \dots, x_d) = \prod_{i=1}^d P(x_i \mid x_{\pi(i)})$$

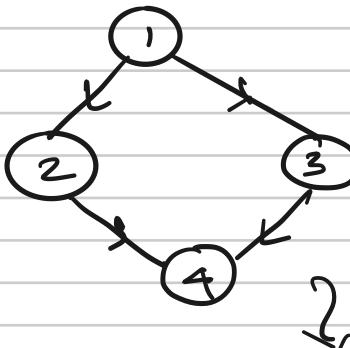
$$P(X=x_1, x_2, \dots, x_d) = \prod_{i=1}^d P(x_i | \text{Pa}(i))$$

$$\text{Pa}(1) = \emptyset$$

$$\text{Pa}(2) = \{1\}$$

$$\text{Pa}(3) = \{1\}$$

$$\text{Pa}(4) = \{2, 3\}$$

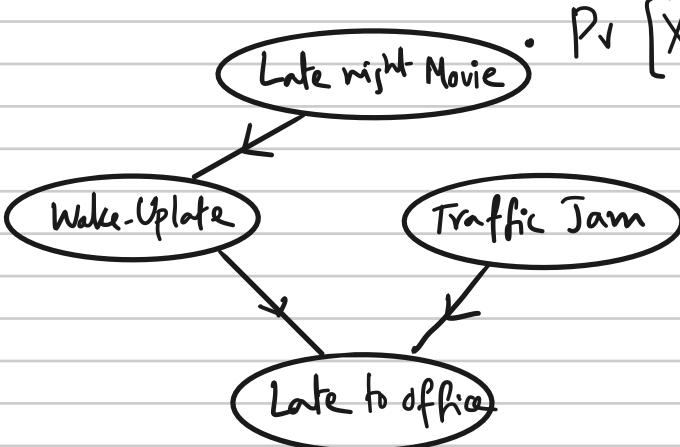


$\text{Pa}(i)$ = the parents of i .

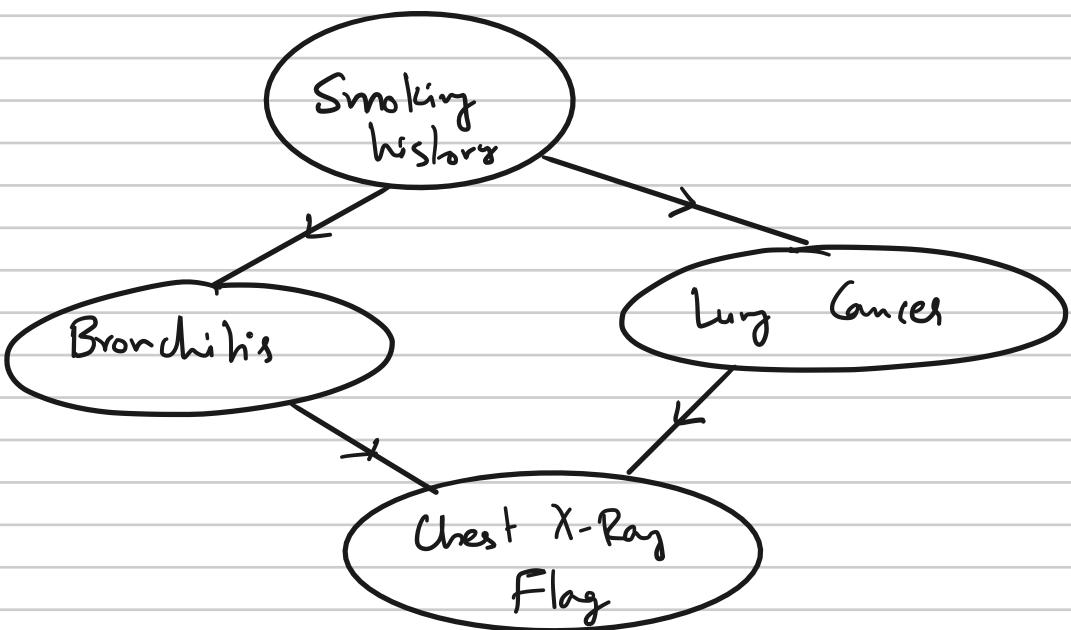
Example: $P(X=x_1, x_2, x_3, x_4)$

$$= P[X_1=x_1] \cdot P[X_2=x_2 | X_1=x_1] \cdot P[X_3=x_3 | X_1=x_1] \cdot P[X_4=x_4 | X_2=x_2, X_3=x_3]$$

Example:



Example:



Learn the Bayes Net from

Samples

→ Given n samples x^1, x^2, \dots, x^n

→ Learn the underlying directed dependency graph

Ex:



"Markov Chain":

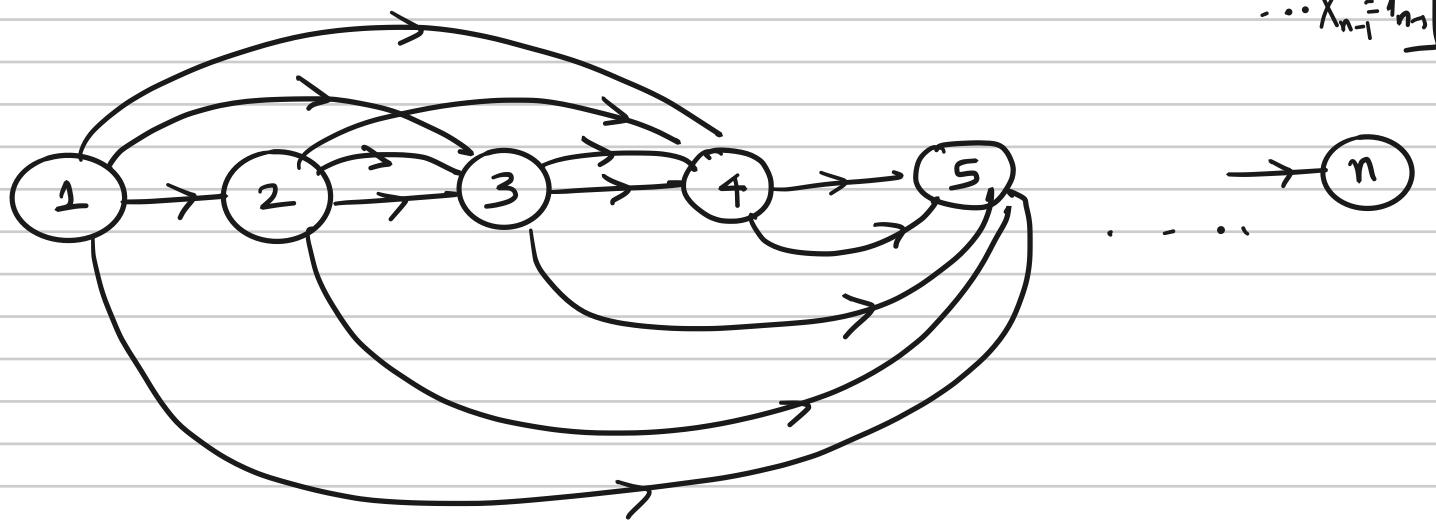
$$x_i \perp x_j \mid x_{i-1} \quad \text{for } j < i-1$$

$$x_1 \perp x_2 \mid x_0$$

$$P_r[x = x_1, x_2, \dots, x_n]$$

$$= P_r[x_1 = x_1] \cdot P_r[x_2 = x_2 \mid x_1 = x_1]$$

$$\cdot P_r[x_3 = x_3 \mid x_1 = x_1, x_2 = x_2] \dots P_r[x_n = x_n \mid x_1 = x_1, \dots, x_{n-1} = x_{n-1}]$$



Make an assumption about the true "graph"

(being "simple" in some sense).

→ Degree is bounded ?

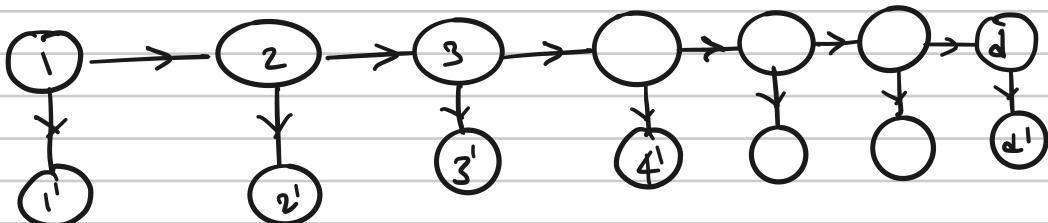
→ The "unknown" graph is a tree/path ...

→ Learn the structure graph under assumptions on the graph

→ Learn the distribution assuming you even know the graph.

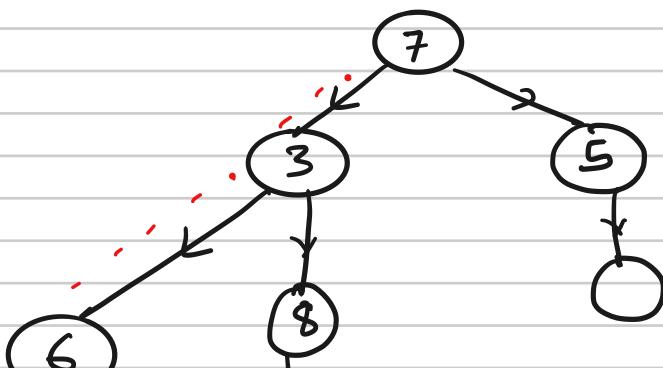
→ Some features are missing?

(Hidden Markov Models)



Suppose we get samples from a distribution generated by a tree Bayes network.

→ Can you learn the distribution from samples?



Each node has a single parent.

2

Chow-Liu Algorithm 1968:

→ We can learn tree-shaped Bayesian networks.

Last class

- Defined Bayesian Nets
- Structural assumptions
- Learning Bayes nets

Today

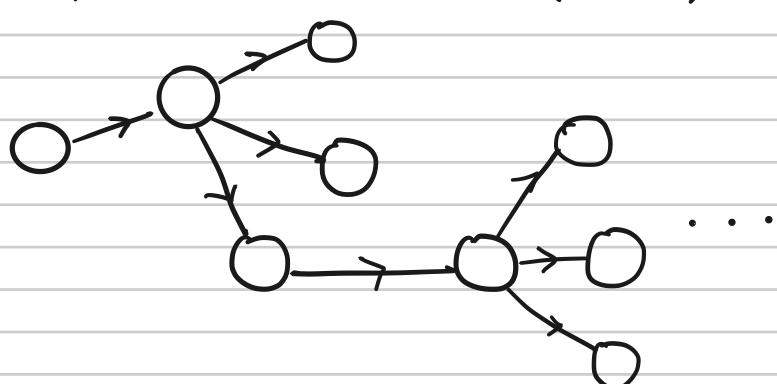
- Chow-Liu algorithm
- Undirected graphical Models

- Assignment 4 is due today @ 10PM
- Office hours on Wednesday 06/02 @ 4-5 PM.

Learning Bayesian Networks defined by trees

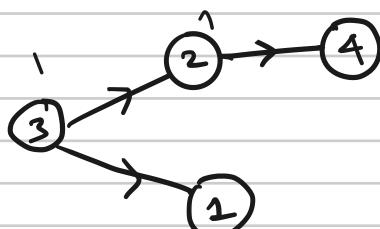
Distribution $(x_1, x_2, \dots, x_d) \in \Sigma^n (\{0,1\}^d)$
 D:

Assumption:



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

Example:



$$\begin{aligned} P(X = x_1, x_2, x_3, x_4) \\ = P(x_3=x_3) \cdot P(x_2=x_2 | X_3=x_3) \\ \cdot P(x_1=x_1 | X_3=x_3) \end{aligned}$$

$$\cdot \Pr[x_4 = x_4 \mid x_2 = x_2]$$

- Suppose we see samples from a Bayes net as above

- Can you learn the distribution?



Structure Learning

Parameter learning

• Learn the underlying graph
learn the tree

Learn the Conditional distributions.

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 dist $\underline{P}, \underline{q}$ 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: $\Omega = \{0, 1\}$

$$P(0) = \frac{1}{2}, \quad P(1) = \frac{1}{2}$$

$$q(0) = \frac{3}{4}, \quad q(1) = \frac{1}{4}$$

$$\begin{aligned}
 \rightarrow KL(P|q) &= P(0) \cdot \log \frac{P(0)}{q(0)} + P(1) \cdot \log \frac{P(1)}{q(1)} \\
 &= \frac{1}{2} \cdot \log \left(\frac{2}{3} \right) + \frac{1}{2} \cdot \log \left(2 \right) \\
 &= \frac{1}{2} \cdot \log \left(\frac{4}{3} \right).
 \end{aligned}$$

Question: What is $KL(P|P)$?

Answer: 0.

Property: $KL(P|q) \geq 0$ and is 0 $\Leftrightarrow P=q$.

INPUT: Samples from a distribution P on Σ^d ($\Sigma = \{0,1\}$)
 (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 $KL(P|P_T) \leq \varepsilon$.

Example:

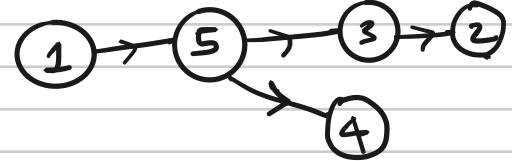


T^* .

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

If I give a tree T and ask
 find me the Bayes net P_T that minimises

$$KL(P|P_T)$$



$\rightarrow (x_1, x_2, x_3, x_4, x_5) \quad T^*$

Idea: Well, use the "empirical" conditional probabilities

1. Estimate the conditional probabilities along the tree T .

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

⋮

INPUT: Samples from P & 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 Bound:

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

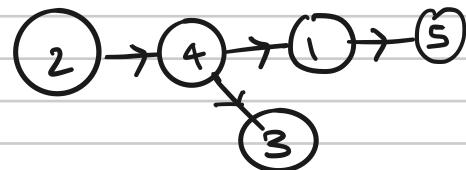
Some number depending
on P





$$(x_1, x_2, x_3, x_4, x_5) \quad T^*$$

$$D(P|P_T)$$



$$D(P|P_T)$$

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

$$\sum_{a,b} I(x_i; x_j) = P_{\{X_i=a, X_j=b\}} \cdot \log \frac{P_{\{X_i=a\}} \cdot P_{\{X_j=b\}}}{P_{\{X_i=a, X_j=b\}}}$$

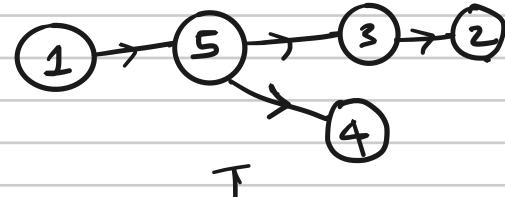


We can estimate $I(x_i; x_j)$ from Samples.

$$P\{X_i=x_i,$$



$$(x_1, x_2, x_3, x_4, x_5) \quad T^*$$



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

$$\underset{T}{\operatorname{argmin}} \quad J_P - \sum_{(i,j) \in T} I(x_i; x_j)$$

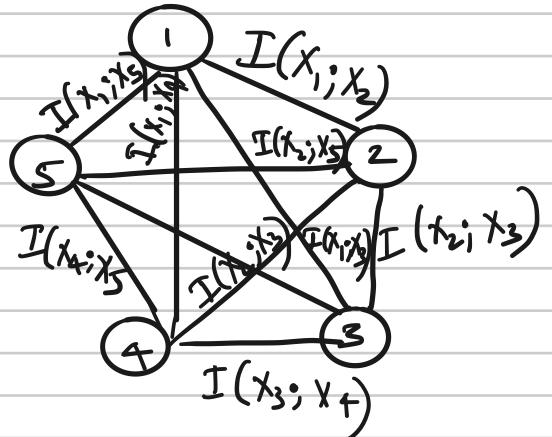
$$\equiv \underset{T}{\operatorname{argmin}} \sum_{(i,j) \in T} -I(x_i; x_j)$$

$$\equiv \underset{T}{\operatorname{argmax}} \sum_{(i,j) \in T} I(x_i; x_j)$$

Idea:

We have samples from

P



HOW-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 $P_{T'}$.

We can do this very fast.

Given a weighted graph G on d vertices, find tree T that has maximum total weight.

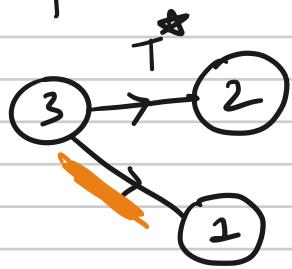
We can do this in linear time.

$(x_1, x_2, \dots, x_d) :$ → Compute $I(x_1; x_2)$

→ Compute $I(x_1; x_2)$

→ Compute $I(X_i; X_f)$

Example:



$$P[X_0 = x_1, x_2, x_3] = P[X_3 = x_3] \cdot P[X_2 = x_2 | X_3 = x_3]$$

$$\cdot P[X_1 = x_1 | X_3 = x_3]$$

"true distribution".

$$\begin{aligned} & T \\ & \text{---} \\ & \text{---} \\ & P_T[X = x_1, x_2, x_3] \\ & = P[X_3 = x_3] \cdot P[X_2 = x_2 | X_3 = x_3] \\ & \quad \cdot P[X_1 = x_1 | X_2 = x_2] \end{aligned}$$

Summary:

→ Can learn tree-structured Bayes networks

→ Active research: What are other structural assumptions
on we learn Bayes networks.

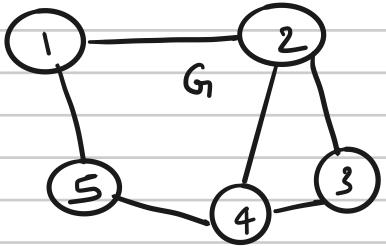
"Directed Graphical Models"

Undirected Graphical Models

- "Markov Random Fields". (MRF_S).

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

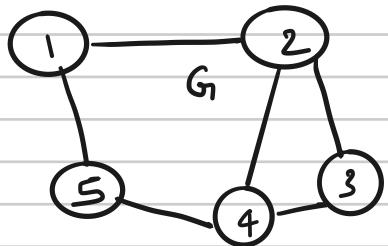
Dependency graph G for D .



e.g.: $X_1 \perp X_4 \mid X_2, X_3, X_5$

D satisfies "Pairwise Markov Property" with respect to G_i if
if i, j have no edge, then

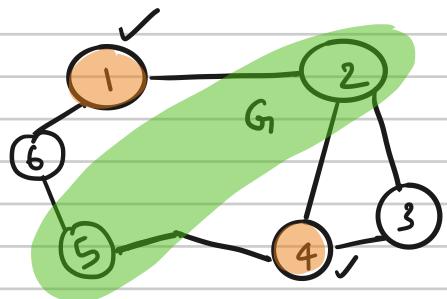
$$X_i \perp X_j \mid X_{\text{rest}}$$



e.g.: $X_1 \perp X_4 \mid X_2, X_5$

D satisfies "Local Markov Property" with respect to G_i if
if i, j have no edge, then

$$X_i \perp X_j \mid X_{\{\text{neighbors of } i\}}$$



e.g.: $X_1 \perp X_4 \mid X_2, X_5$

D satisfies "Global Markov Property" with respect to G_i if
if i, j have no edge, then

$$X_i \perp X_j \mid X_{\{\text{any separating set}\}}$$

$$X_1 \perp X_4 \mid X_2, X_5$$

Subset of vertices
removing which disconnects i
and j in G_i .

$$\text{Global MP} \Rightarrow \text{Local MP} \Rightarrow \text{Pairwise MP}$$

For "most reasonable" distributions all three are equivalent.

Example: Markov Chain

x_1, x_2, \dots, x_t

$x_{t+1} \perp x_{t-1} \mid x_t$



Remark: We'll say D has dependency graph G_D if it satisfies Markov property with respect to G_D .

Main Learning Challenges:

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

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

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

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 dep. graph for all distributions!

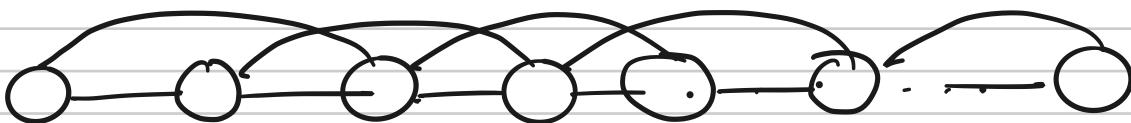
"Minimal dependency graph?

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 $\leq k$.



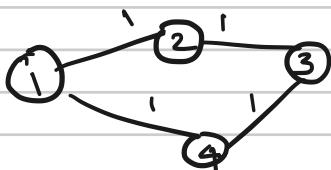
$$k = 4$$

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

Learning Boltzmann Machines

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

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



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

Gaussian Graphical Models

D on \mathbb{R}^d

Dist Dis $N(0, \Sigma)$

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

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

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

independent.

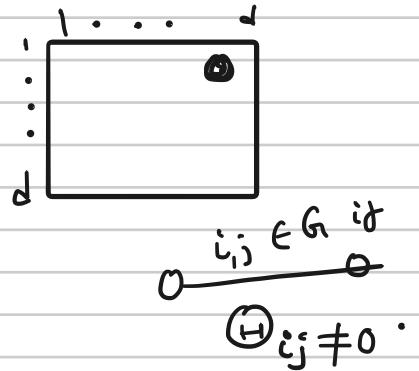
Distributions as defined above,
They satisfy Markov property
wrt Gr.

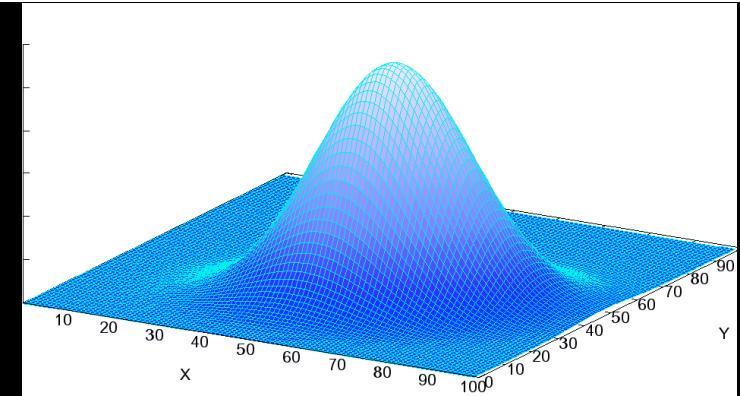
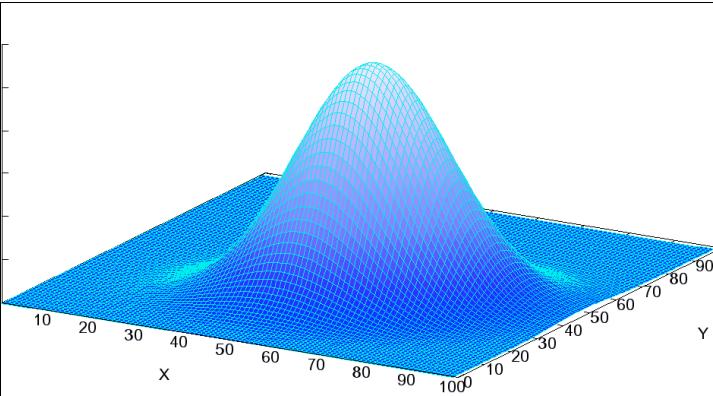
(Dempsey 1972):

Precision Matrix:

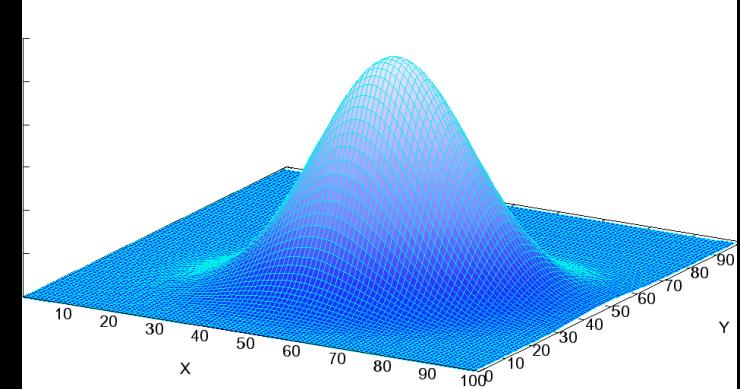
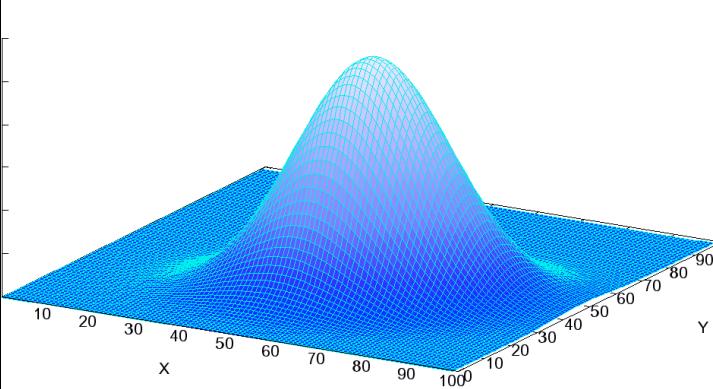
$$(\Theta) = \sum^{-1} .$$

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

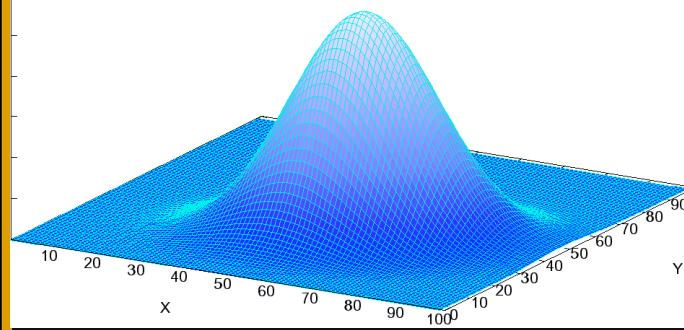




Learning Gaussian Graphical Models w/o condition number bounds



GGMs



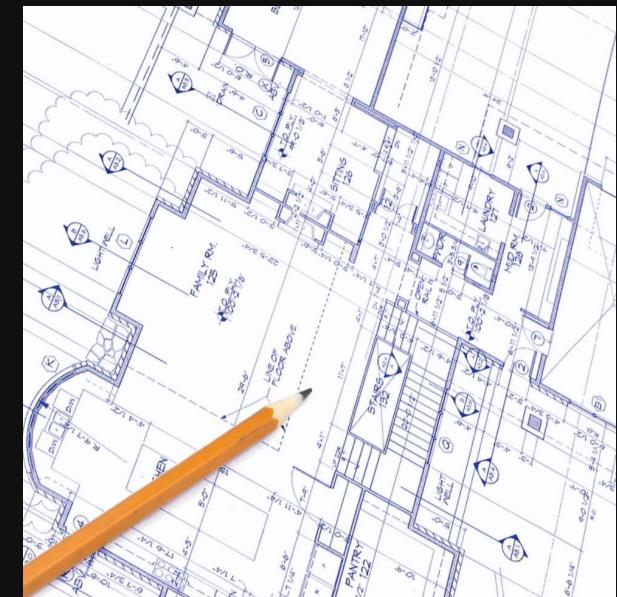
Motivation

Algorithm

```
While  $|S| < t$  :  
Add  $\arg \min_j \text{Var}(X_i | X_{S \cup j})$ .
```

Special Models

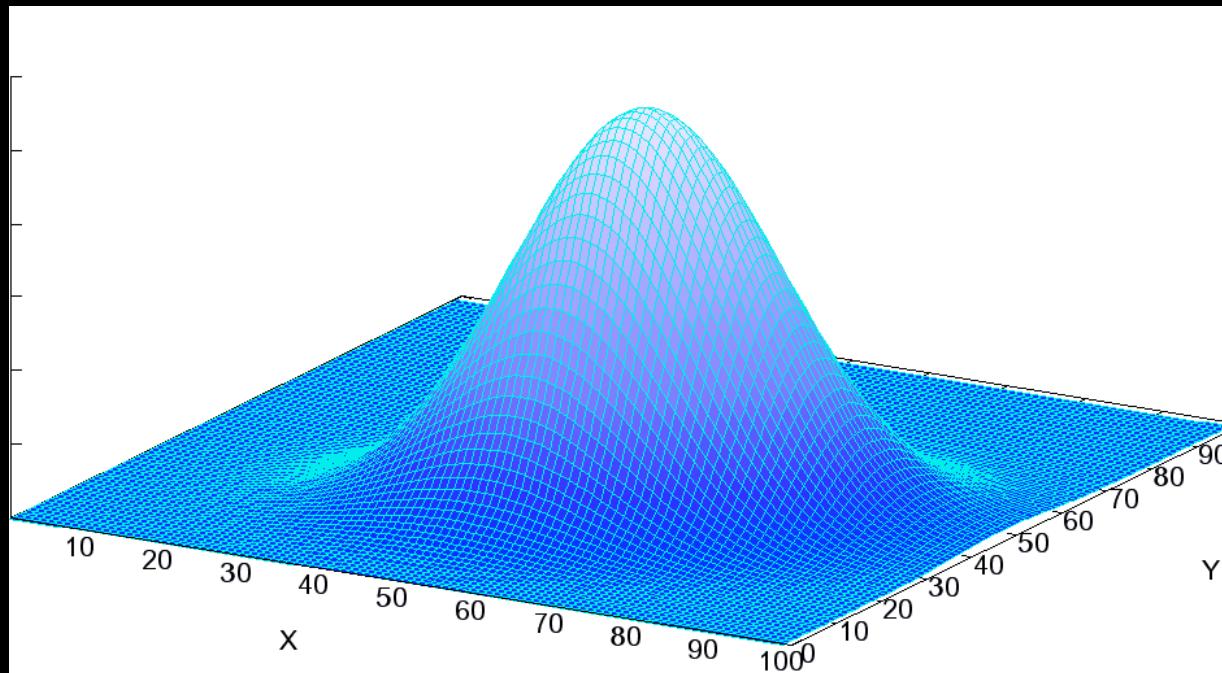
Analysis



Attractive. SDD

Gaussian Graphical Models (GGMs)

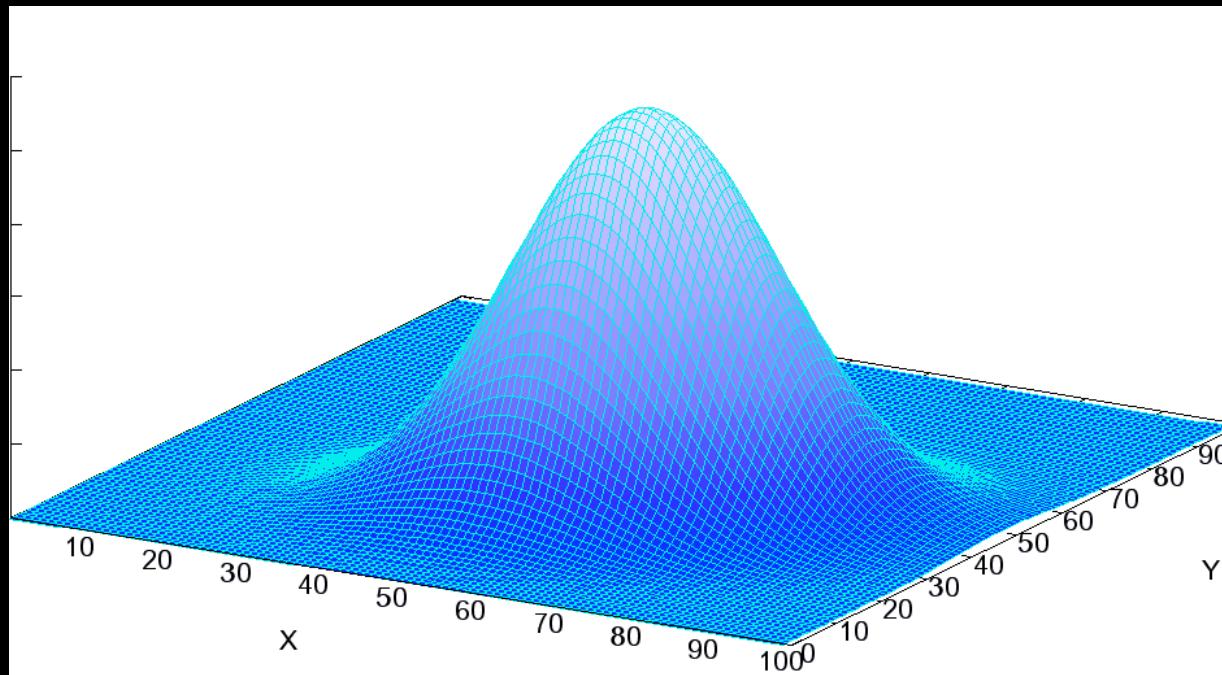
$X \sim N(0, \Sigma)$. $\Sigma \in R^{p \times p}$ covariance matrix.



$$Pr[X = x] = \frac{1}{\sqrt{(2\pi)^p \det(\Sigma)}} \exp(-x^T \Sigma^{-1} x / 2)$$

Gaussian Graphical Models (GGMs)

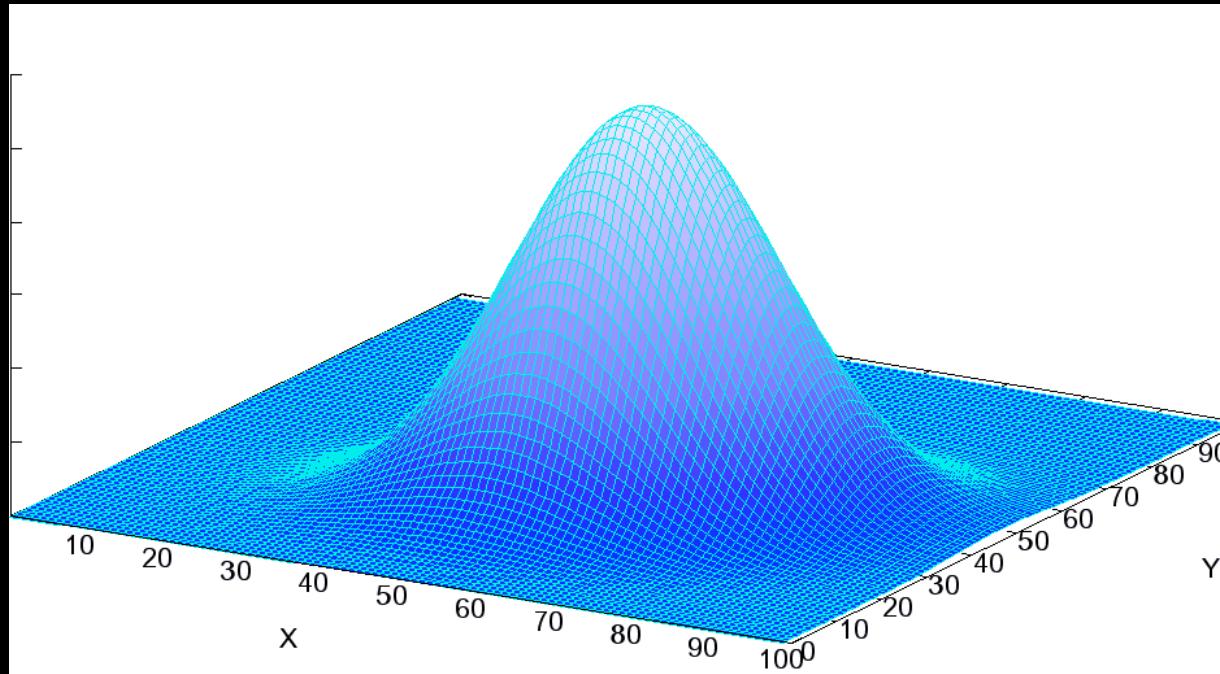
$X \sim N(0, \Sigma)$. $\Sigma \in R^{p \times p}$ covariance matrix.



Today's Focus: Precision Matrix $\Theta = \Sigma^{-1}$.

Gaussian Graphical Models (GGMs)

$X \sim N(0, \Sigma) . \Sigma \in R^{p \times p}$ covariance matrix.



Today's Focus: Precision Matrix $\Theta = \Sigma^{-1}$.

Dempster 72: Encodes conditional independence structure

Gaussian Graphical Models (GGMs)

$X \sim N(0, \Sigma)$. $\Sigma \in R^{p \times p}$ covariance matrix.

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

Gaussian Graphical Models (GGMs)

$X \sim N(0, \Sigma)$. $\Sigma \in R^{p \times p}$ covariance matrix.

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

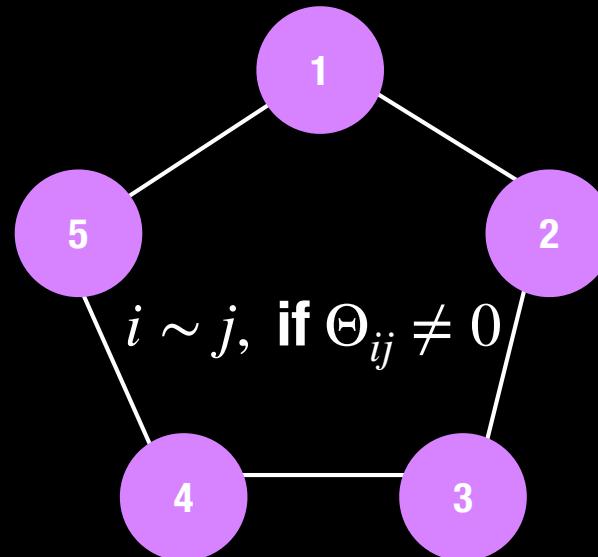
Dependency graph $G = Supp(\Theta)$.

Gaussian Graphical Models (GGMs)

$X \sim N(0, \Sigma)$. $\Sigma \in R^{p \times p}$ covariance matrix.

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

Dependency graph $G = Supp(\Theta)$.

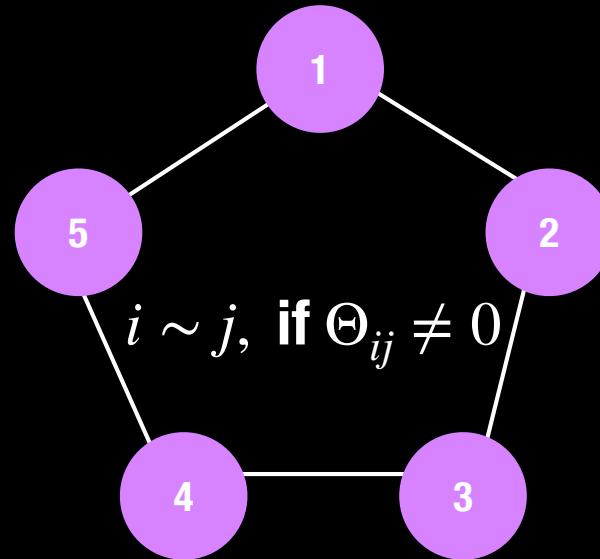


Gaussian Graphical Models (GGMs)

$X \sim N(0, \Sigma)$. $\Sigma \in R^{p \times p}$ covariance matrix.

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

Dependency graph $G = \text{Supp}(\Theta)$.



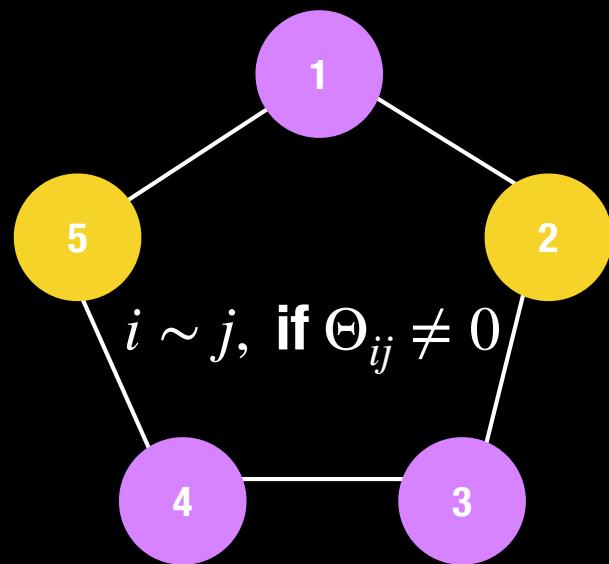
Markov property: $\Theta_{ij} = 0 \Rightarrow X_i, X_j \text{ are independent conditioned on neighbors of } i.$

Gaussian Graphical Models (GGMs)

$X \sim N(0, \Sigma)$. $\Sigma \in R^{p \times p}$ covariance matrix.

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

Dependency graph $G = \text{Supp}(\Theta)$.



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

Markov property: $\Theta_{ij} = 0 \Rightarrow X_i, X_j$ are independent
conditioned on neighbors of i .

Gaussian Graphical Models (GGMs)

$X \sim N(0, \Sigma)$. $\Sigma \in R^{p \times p}$ covariance matrix.

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

Dependency graph $G = Supp(\Theta)$.

Complexity of GGMs: max-degree of G $d \ll p$.

Gaussian Graphical Models (GGMs)

$X \sim N(0, \Sigma)$. $\Sigma \in R^{p \times p}$ covariance matrix.

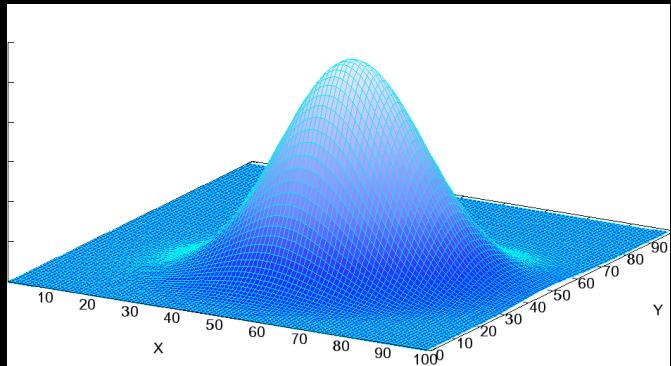
Precision Matrix $\Theta = \Sigma^{-1}$.

Dependency graph $G = Supp(\Theta)$.

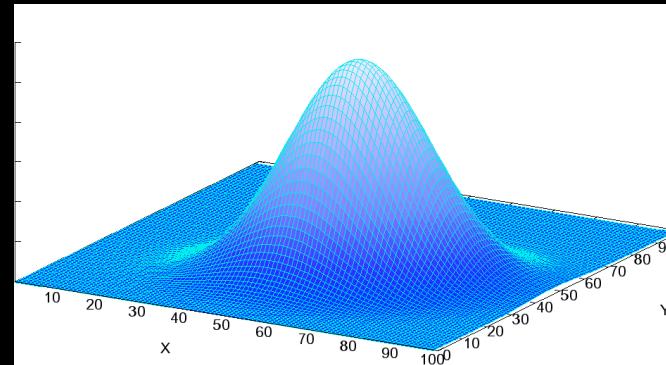
Complexity of GGMs: max-degree of G $d \ll p$.

The assumption that makes the formalism non-trivial
and useful ...

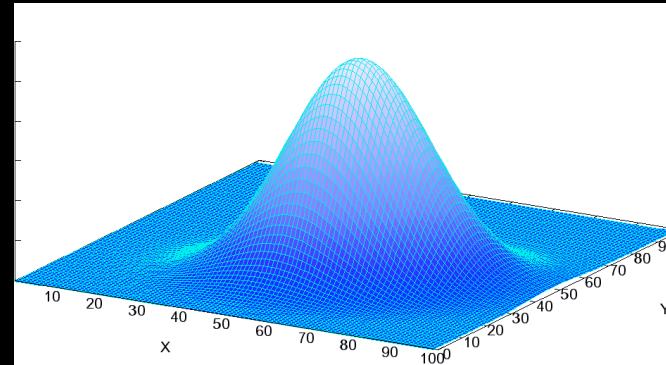
Why?



Fundamental model for modeling statistical relations between variables



Fundamental model for modeling statistical relations between variables



Many applications in machine learning, sciences:
Gene-regulatory networks
Brain connection networks from fMRI

...

Bigger Picture

Can we learn sparse dependency graphs from few samples?

(aka learning Markov random fields, undirected graphical models)

Bigger Picture

Can we learn sparse dependency graphs from few samples?

(aka learning Markov random fields, undirected graphical models)

$$X \sim \{1, -1\}^P.$$

Dependency graph of X : $i \not\sim j \Rightarrow X_i, X_j$ are independent conditioned on neighbors of i .

Example: “Random Walk” Model

$$X_i = X_{i-1} + Z_i, \text{ where } Z_1, \dots, Z_p \text{ i.i.d } N(0, 1)$$

1	1	1	1	1	1	1	1	1	1	1
1	2	2	2	2	2	2	2	2	2	2
1	2	3	3	3	3	3	3	3	3	3
1	2	3	4	4	4	4	4	4	4	4
1	2	3	4	5	5	5	5	5	5	5
1	2	3	4	5	6	6	6	6	6	6
1	2	3	4	5	6	7	7	7	7	7
1	2	3	4	5	6	7	8	8	8	8
1	2	3	4	5	6	7	8	9	9	9
1	2	3	4	5	6	7	8	9	10	

Σ : Covariance Matrix

Example: “Random Walk” Model

$$X_i = X_{i-1} + Z_i, \text{ where } Z_1, \dots, Z_p \text{ i.i.d } N(0, 1)$$

1	1	1	1	1	1	1	1	1	1	1
1	2	2	2	2	2	2	2	2	2	2
1	2	3	3	3	3	3	3	3	3	3
1	2	3	4	4	4	4	4	4	4	4
1	2	3	4	5	5	5	5	5	5	5
1	2	3	4	5	6	6	6	6	6	6
1	2	3	4	5	6	7	7	7	7	7
1	2	3	4	5	6	7	8	8	8	8
1	2	3	4	5	6	7	8	9	9	9
1	2	3	4	5	6	7	8	9	10	

Σ : Covariance Matrix

2	-1	0	0	0	0	0	0	0	0	0
-1	2	-1	0	0	0	0	0	0	0	0
0	-1	2	-1	0	0	0	0	0	0	0
0	0	-1	2	-1	0	0	0	0	0	0
0	0	0	-1	2	-1	0	0	0	0	0
0	0	0	0	-1	2	-1	0	0	0	0
0	0	0	0	0	-1	2	-1	0	0	0
0	0	0	0	0	0	-1	2	-1	0	0
0	0	0	0	0	0	0	-1	2	-1	0
0	0	0	0	0	0	0	0	-1	2	-1

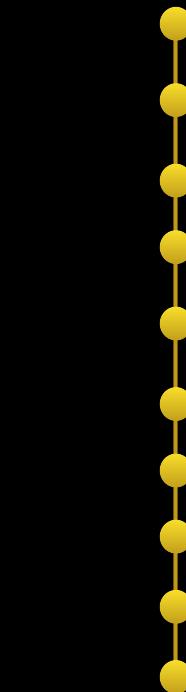
Θ : Precision Matrix

Example: “Random Walk” Model

$X_i = X_{i-1} + Z_i$, where Z_1, \dots, Z_p i.i.d $N(0, 1)$

1	1	1	1	1	1	1	1	1	1	1
1	2	2	2	2	2	2	2	2	2	2
1	2	3	3	3	3	3	3	3	3	3
1	2	3	4	4	4	4	4	4	4	4
1	2	3	4	5	5	5	5	5	5	5
1	2	3	4	5	6	6	6	6	6	6
1	2	3	4	5	6	7	7	7	7	7
1	2	3	4	5	6	7	8	8	8	8
1	2	3	4	5	6	7	8	9	9	9
1	2	3	4	5	6	7	8	9	9	10

Σ : Covariance Matrix



Dependency
graph

2	-1	0	0	0	0	0	0	0	0	0
-1	2	-1	0	0	0	0	0	0	0	0
0	-1	2	-1	0	0	0	0	0	0	0
0	0	-1	2	-1	0	0	0	0	0	0
0	0	0	-1	2	-1	0	0	0	0	0
0	0	0	0	-1	2	-1	0	0	0	0
0	0	0	0	0	-1	2	-1	0	0	0
0	0	0	0	0	0	-1	2	-1	0	0
0	0	0	0	0	0	0	-1	2	-1	0
0	0	0	0	0	0	0	0	-1	2	-1

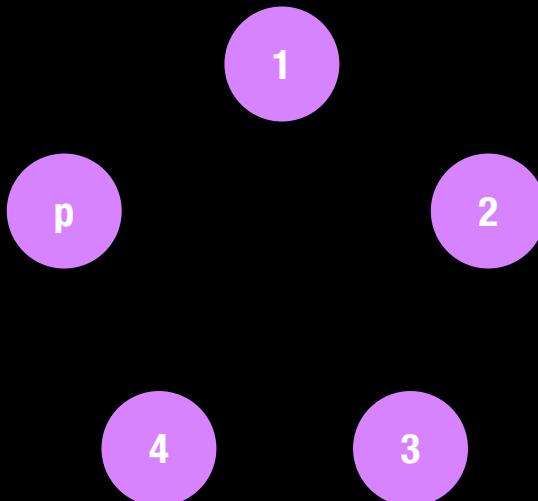
Θ : Precision Matrix

Learning Sparse GGMs

- **Structure learning - learn dependency graph.**
- **Parameter learning - learn the matrix.**

Structure Learning for GGMs

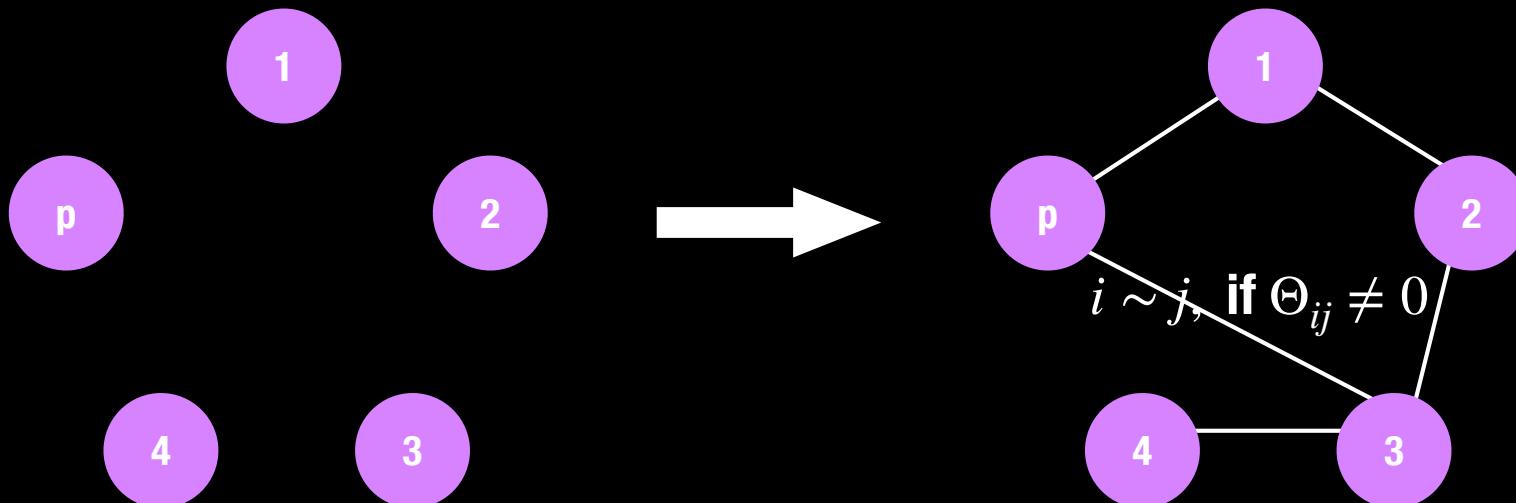
- Structure learning - learn dependency graph.



Input: Given samples X^1, X^2, \dots, X^n from a GGM
Output: Dependency graph of the GGM.

Structure Learning for GGMs

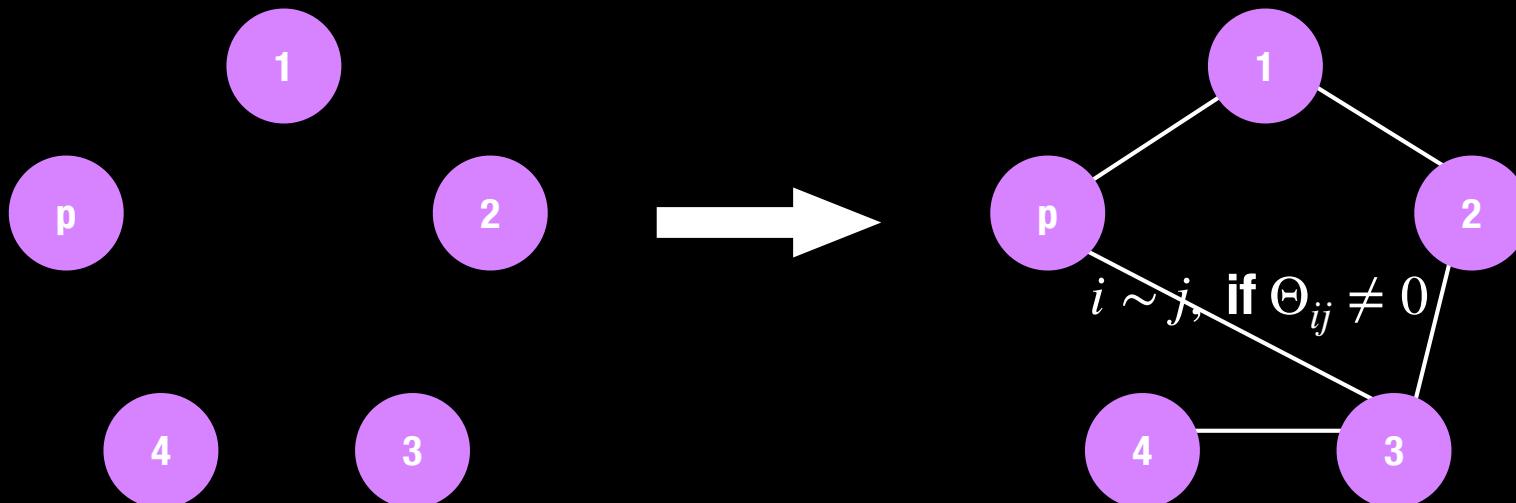
- Structure learning - learn dependency graph.



Input: Given samples X^1, X^2, \dots, X^n from a GGM
Output: Dependency graph of the GGM.

Structure Learning for GGMs

- Structure learning - learn dependency graph.

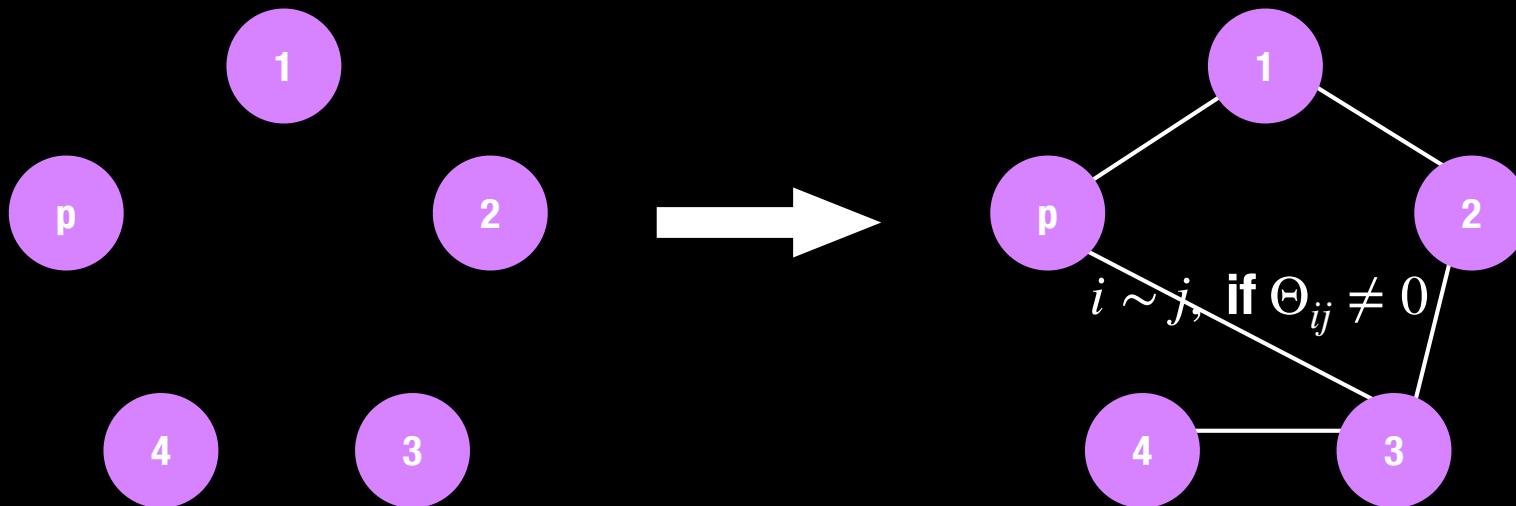


Input: Given samples X^1, X^2, \dots, X^n from a GGM
Output: Dependency graph of the GGM.

Challenge: Often $n \ll p$.

Structure Learning for GGMs

- Structure learning - learn dependency graph.



Assumption: Unknown dependency graph is sparse - each vertex has at most **d** edges.

Challenge: Often $n \ll p$.

Structure Learning for GGMs

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

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

Structure Learning for GGMs

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

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

Core of problem. Can
learn parameters easily
afterward.

Structure Learning for GGMs

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

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

Ideal: Practical algorithms with provable guarantees.

Example: Unknown order Random Walk

$X_{\pi(i)} = X_{\pi(i-1)} + Z_i$, where Z_1, \dots, Z_p i.i.d $N(0,1)$

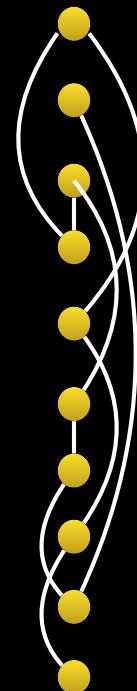
7	1	5	6	7	4	3	7	2	7
1	1	1	1	1	1	1	1	1	1
5	1	5	5	5	4	3	5	2	5
6	1	5	6	6	4	3	6	2	6
7	1	5	6	8	4	3	8	2	8
4	1	4	4	4	4	3	4	2	4
3	1	3	3	3	3	3	3	2	3
7	1	5	6	8	4	3	9	2	9
2	1	2	2	2	2	2	2	2	2
7	1	5	6	8	4	3	9	2	10

Σ : Covariance Matrix

Example: Unknown order Random Walk

$X_{\pi(i)} = X_{\pi(i-1)} + Z_i$, where Z_1, \dots, Z_p i.i.d $N(0, 1)$

7	1	5	6	7	4	3	7	2	7
1	1	1	1	1	1	1	1	1	1
5	1	5	5	5	4	3	5	2	5
6	1	5	6	6	4	3	6	2	6
7	1	5	6	8	4	3	8	2	8
4	1	4	4	4	4	3	4	2	4
3	1	3	3	3	3	3	3	2	3
7	1	5	6	8	4	3	9	2	9
2	1	2	2	2	2	2	2	2	2
7	1	5	6	8	4	3	9	2	10



Σ : Covariance Matrix

Dependency
graph?

Θ : Precision Matrix

2	0	0	-1	-1	0	0	0	0	0	0
0	2	0	0	0	0	0	0	0	-1	0
0	0	2	-1	0	-1	0	0	0	0	0
-1	0	-1	2	0	0	0	0	0	0	0
-1	0	0	0	2	0	0	-1	0	0	0
0	0	-1	0	0	2	-1	0	0	0	0
0	0	0	0	0	-1	2	0	-1	0	0
0	0	0	0	-1	0	0	2	0	-1	0
0	-1	0	0	0	0	-1	0	2	0	0
0	0	0	0	0	0	0	-1	0	2	0

Example: Unknown order Random Walk

$X_{\pi(i)} = X_{\pi(i-1)} + Z_i$, where Z_1, \dots, Z_p i.i.d $N(0, 1)$

7	1	5	6	7	4	3	7	2	7
1	1	1	1	1	1	1	1	1	1
5	1	5	5	5	4	3	5	2	5
6	1	5	6	6	4	3	6	2	6
7	1	5	6	8	4	3	8	2	8
4	1	4	4	4	4	3	4	2	4
3	1	3	3	3	3	3	3	2	3
7	1	5	6	8	4	3	9	2	9
2	1	2	2	2	2	2	2	2	2
7	1	5	6	8	4	3	9	2	10



2	0	0	-1	-1	0	0	0	0	0	0
0	2	0	0	0	0	0	0	0	-1	0
0	0	2	-1	0	-1	0	0	0	0	0
-1	0	-1	2	0	0	0	0	0	0	0
-1	0	0	0	2	0	0	-1	0	0	0
0	0	-1	0	0	2	-1	0	0	0	0
0	0	0	0	0	-1	2	0	-1	0	0
0	0	0	0	0	-1	0	0	2	0	-1
0	-1	0	0	0	0	-1	0	2	0	0
0	0	0	0	0	0	0	-1	0	2	0

Σ : Covariance Matrix

Dependency
graph?

Θ : Precision Matrix

How many samples of X to find the hidden permutation?

Previous Work

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

Well studied with several popular software packages:
GLASSO, CLIME, ACLIME

Previous Work

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

Previous Work

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

- GLASSO: Friedman, Hastie, Tibshirani 08
Can learn with $O(d^2 \log p)$ samples if precision matrix is incoherent.

Previous Work

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

Strong assumption: Violated by random walk ...

- GLASSO: Friedman, Hastie, Tibshirani 08
Can learn with $O(d^2 \log p)$ samples if precision matrix is incoherent.

Previous Work

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

- GLASSO: Friedman, Hastie, Tibshirani 08
Can learn with $O(d^2 \log p)$ samples if precision matrix is incoherent.
- CLIME: Cai, Liu, Luo 2011
Can learn with $O(M^4 \log p)$ samples where $M \sim \text{maximum ratio of entries in } \Sigma, \Theta$.

Previous Work

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

- GLASSO: Friedman, Hastie, Tibshirani 08
Can learn with $O(d^2 \log p)$ samples if precision
Strong assumption: Violated by random walk ...
- CLIME: Cai, Liu, Lub 2011
Can learn with $O(M^4 \log p)$ samples where
 $M \sim \text{maximum ratio of entries in } \Sigma, \Theta$.

Previous Work

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

- GLASSO: Friedman, Hastie, Tibshirani 08
Can learn with $O(d^2 \log p)$ samples if precision matrix is incoherent.
- CLIME: Cai, Liu, Luo 2011
Can learn with $O(M^4 \log p)$ samples where $M \sim \text{maximum ratio of entries in } \Sigma, \Theta$.

Previous Work

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

Summary: GLASSO, CLIME, ACLIME need ‘well-conditioned’ precision matrix.

Previous Work

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

Summary: GLASSO, CLIME, ACLIME need ‘well-conditioned’ precision matrix.

- Violated by simple models
- Not scale invariant
- Not just analysis ... fail empirically

Previous Work

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

What is possible information theoretically?

Previous Work: MVL18

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

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

Θ_{ii}		Θ_{ij}		
Θ_{ij}		Θ_{jj}		

Θ : Precision Matrix

Previous Work: MVL18

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

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

Θ_{ii}		Θ_{ij}		
Θ_{ij}		Θ_{jj}		

Measures signal in 2x2
submatrices ...
not whole matrix.

Θ : Precision Matrix

Previous Work: MVL18

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

$$\text{MVL18: } \kappa(\Theta) = \min_{i,j: \Theta_{ij} \neq 0} \frac{|\Theta_{ij}|}{\sqrt{\Theta_{ii}\Theta_{jj}}}$$

Dependency graph identifiable with $n = O(d \log p / \kappa^2)$ samples!

Previous Work: MVL18

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

$$\text{MVL18: } \kappa(\Theta) = \min_{i,j: \Theta_{ij} \neq 0} \frac{|\Theta_{ij}|}{\sqrt{\Theta_{ii}\Theta_{jj}}}$$

Dependency graph identifiable with $n = O(d \log p / \kappa^2)$ samples!

Example: Random walk model - $\kappa = 1/2$.

Information-Theoretic Limits: MVL18

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

$$\text{MVL18: } \kappa(\Theta) = \min_{i,j: \Theta_{ij} \neq 0} \frac{|\Theta_{ij}|}{\sqrt{\Theta_{ii}\Theta_{jj}}}$$

Dependency graph identifiable with $n = O(d \log p / \kappa^2)$ samples!

Wang, Wainwright, Ramachandran 10: Need $n = \Omega(\log p / \kappa^2)$.

Information-Theoretic Limits: MVL18

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

$$\text{MVL18: } \kappa(\Theta) = \min_{i,j: \Theta_{ij} \neq 0} \frac{|\Theta_{ij}|}{\sqrt{\Theta_{ii}\Theta_{jj}}}$$

Dependency graph identifiable with $n = O(d \log p / \kappa^2)$ samples!

Run-time of algorithm: $p^{O(d)}$.

Information-Theoretic Limits: MVL18

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

$$\text{MVL18: } \kappa(\Theta) = \min_{i,j: \Theta_{ij} \neq 0} \frac{|\Theta_{ij}|}{\sqrt{\Theta_{ii}\Theta_{jj}}}$$

Dependency graph identifiable with $n = O(d \log p / \kappa^2)$ samples!

Run-time of algorithm: $p^{O(d)}$.

Problematic for even moderate sized instances
... Can we do better?

GGMS: Main Learning Challenge

Given n samples from a GGM of degree $d \ll p$,
can we find the dependency graph with $n \approx d \log p / \kappa^2$, and
run-time fixed polynomial in p ? (Or even $p^{o(d)}$.)

GGMS: Main Learning Challenge

Given n samples from a GGM of degree $d \ll p$,
can we find the dependency graph with $n \approx d \log p / \kappa^2$, and
run-time fixed polynomial in p ? (Or even $p^{o(d)}$.)

This work — YES for a large class of models:
Attractive, SDD, ... more generally, Walk-Summable

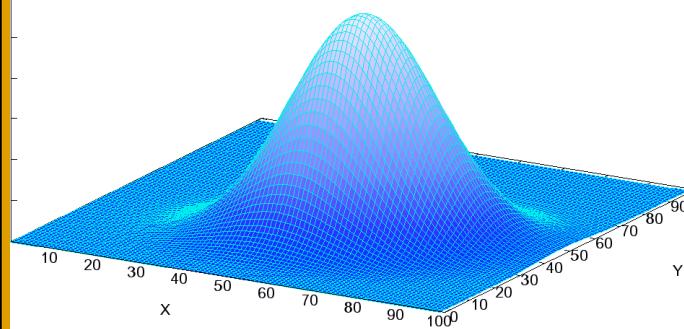
GGMS: Main Learning Challenge

Given n samples from a GGM of degree $d \ll p$,
can we find the dependency graph with $n \approx d \log p / \kappa^2$, and
run-time fixed polynomial in p ? (Or even $p^{o(d)}$.)

This work — YES for a large class of models:
Attractive, SDD, ... more generally, Walk-Summable

Main: A simple greedy algorithm solves above ill-conditioned cases
(and recovers guarantees of GLASSO, CLIME, ...)

GGMs



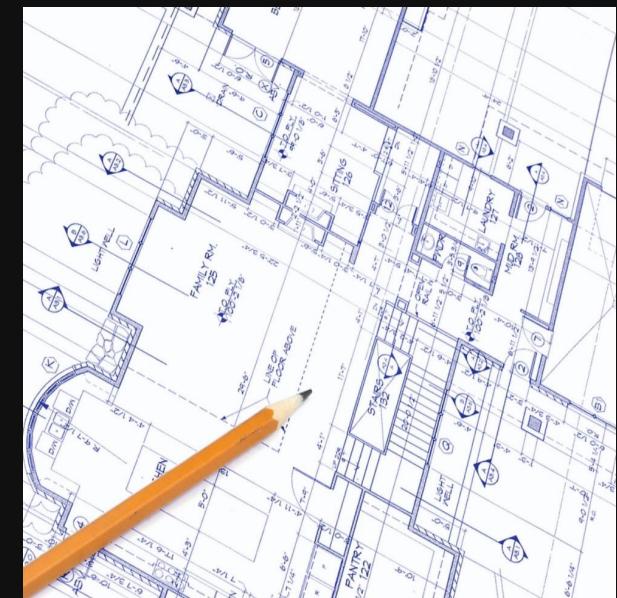
Motivation

Algorithm

```
While  $|S| < t$  :  
Add  $\arg \min_j \text{Var}(X_i | X_{S \cup j})$ .
```

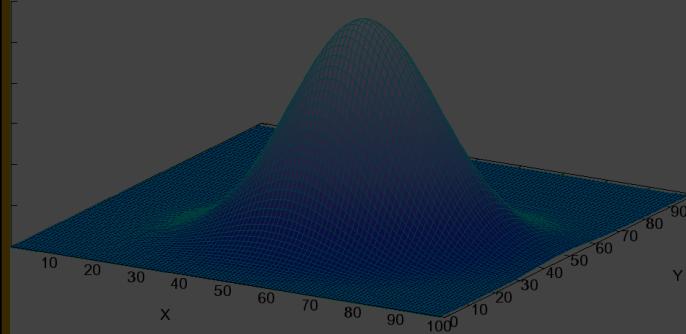
Special Models

Analysis



Attractive. SDD.

GGMs



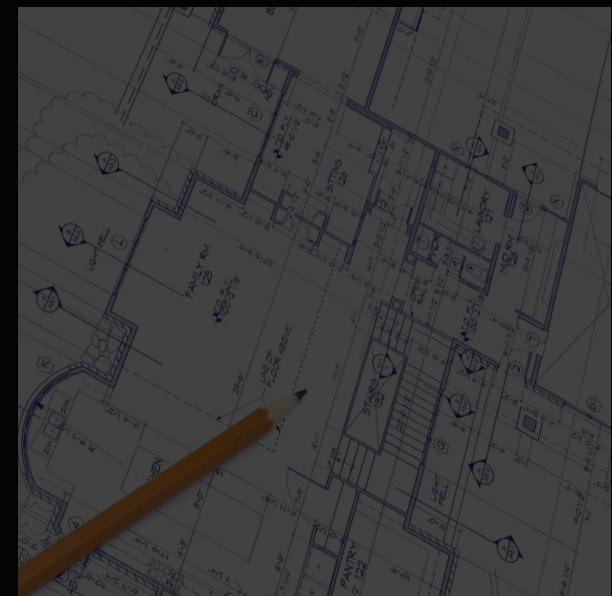
Motivation

Algorithm

```
While  $|S| < t$  :  
Add  $\arg \min_j Var(X_i | X_{S \cup j})$ .
```

Special Models

Analysis

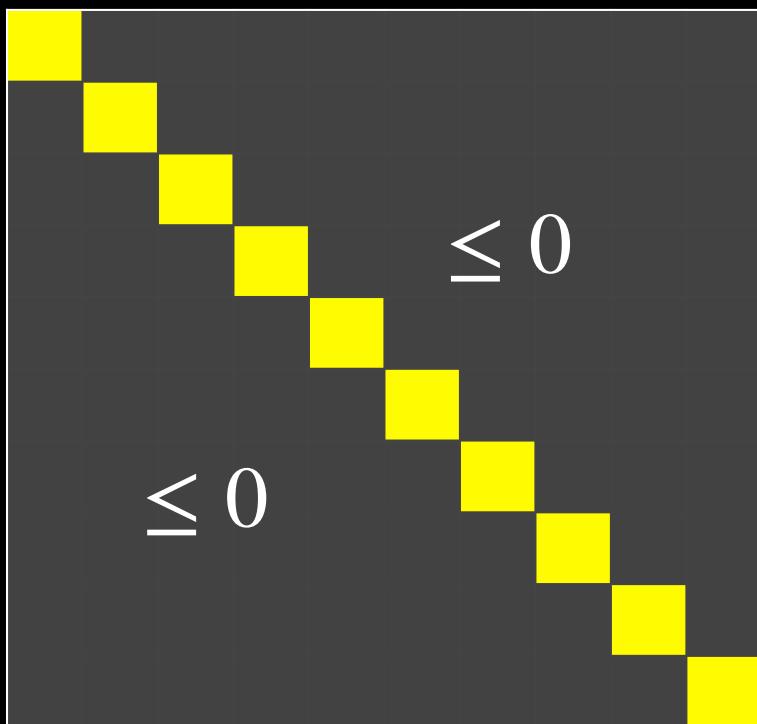


Attractive. SDD.

Attractive GGMs

GGM is attractive if all covariances are non-negative.

(Equivalently, Θ has non-positive off-diagonals.)

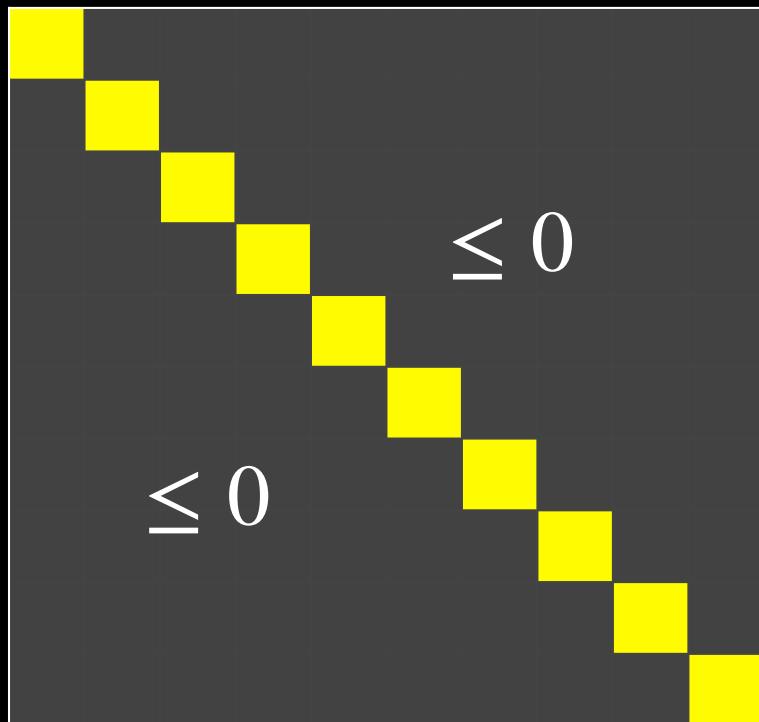


Θ : Precision Matrix

Attractive GGMs

GGM is attractive if all covariances are non-negative.

(Equivalently, Θ has non-positive off-diagonals.)



Θ : Precision Matrix

Ex: Gaussian Free Fields

- Many applications via Gaussian processes

Attractive GGMs

GGM is attractive if all covariances are non-negative.

(Equivalently, Θ has non-positive off-diagonals.)

2	-1	0	0	0	0	0	0	0	0
-1	2	-1	0	0	0	0	0	0	0
0	-1	2	-1	0	0	0	0	0	0
0	0	-1	2	-1	0	0	0	0	0
0	0	0	-1	2	-1	0	0	0	0
0	0	0	0	-1	2	-1	0	0	0
0	0	0	0	0	-1	2	-1	0	0
0	0	0	0	0	0	-1	2	-1	0
0	0	0	0	0	0	0	-1	2	-1
0	0	0	0	0	0	0	0	-1	2

Ex: Gaussian Free Fields

- Many applications via Gaussian processes
- Ex: Random walk model

Θ : Precision Matrix

Attractive GGMs

GGM is attractive if all covariances are non-negative.

(Equivalently, Θ has non-positive off-diagonals.)

2	-1	0	0	0	0	0	0	0	0
-1	2	-1	0	0	0	0	0	0	0
0	-1	2	-1	0	0	0	0	0	0
0	0	-1	2	-1	0	0	0	0	0
0	0	0	-1	2	-1	0	0	0	0
0	0	0	0	-1	2	-1	0	0	0
0	0	0	0	0	-1	2	-1	0	0
0	0	0	0	0	0	-1	2	-1	0
0	0	0	0	0	0	0	-1	2	-1
0	0	0	0	0	0	0	0	-1	2

Ex: Gaussian Free Fields

- Many applications via Gaussian processes
- Ex: Random walk model
- Ill-conditioned if ‘long paths’

Θ : Precision Matrix

Attractive GGMs

GGM is attractive if all covariances are non-negative.

Previous: No efficient algorithms with $O_d(\log p)$ sample complexity known.

Attractive GGMs

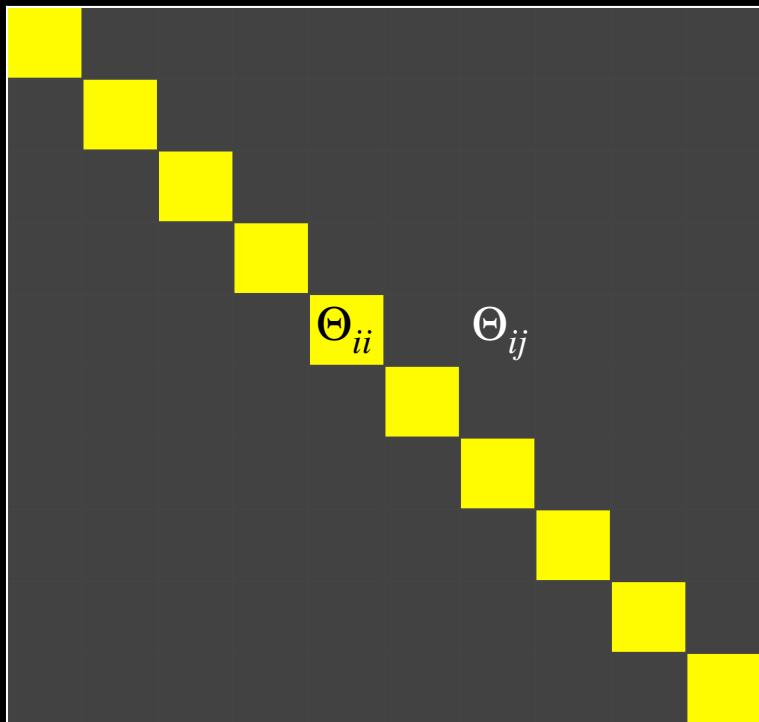
GGM is attractive if all covariances are non-negative.

KKMM: **GreedyPrune** learns attractive models with $\tilde{O}(d \log p / \kappa^2)$ samples and quadratic run-time.

Previous: No efficient algorithms with $O_d(\log p)$ sample complexity known.

Walk-Summable GGMs

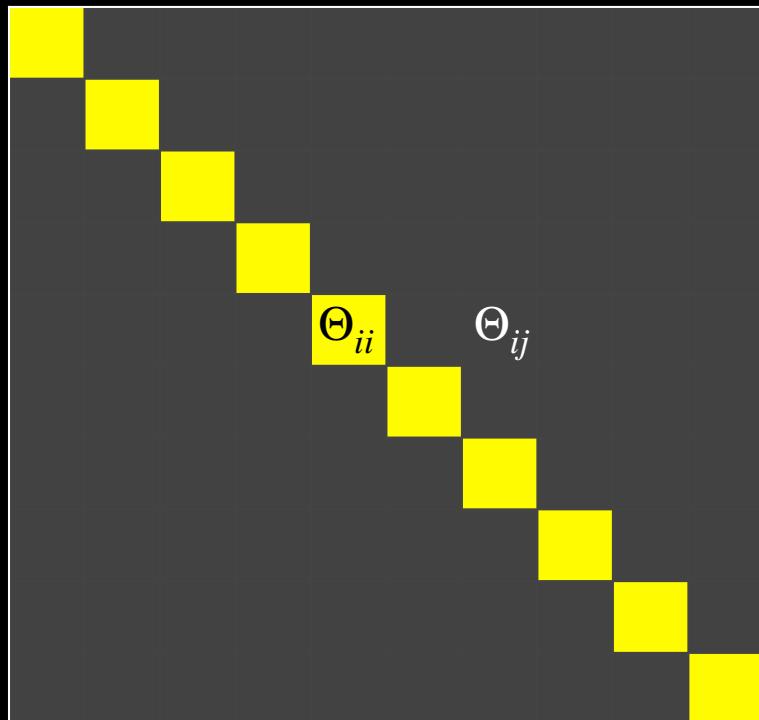
GGM walk-summable if making off-diagonals of precision matrix negative preserves positive semi-definiteness.



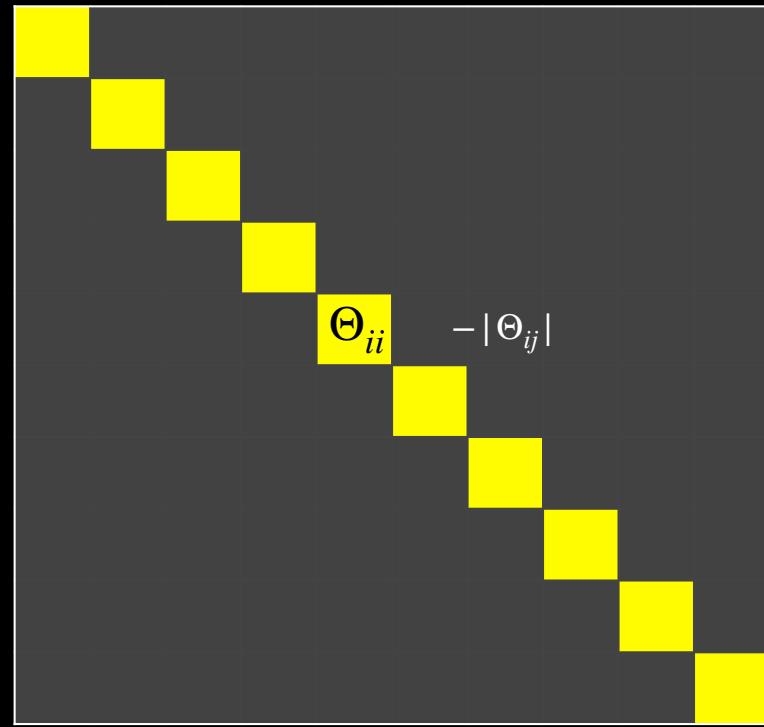
Θ : Precision Matrix

Walk-Summable GGMs

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



Θ : Precision Matrix

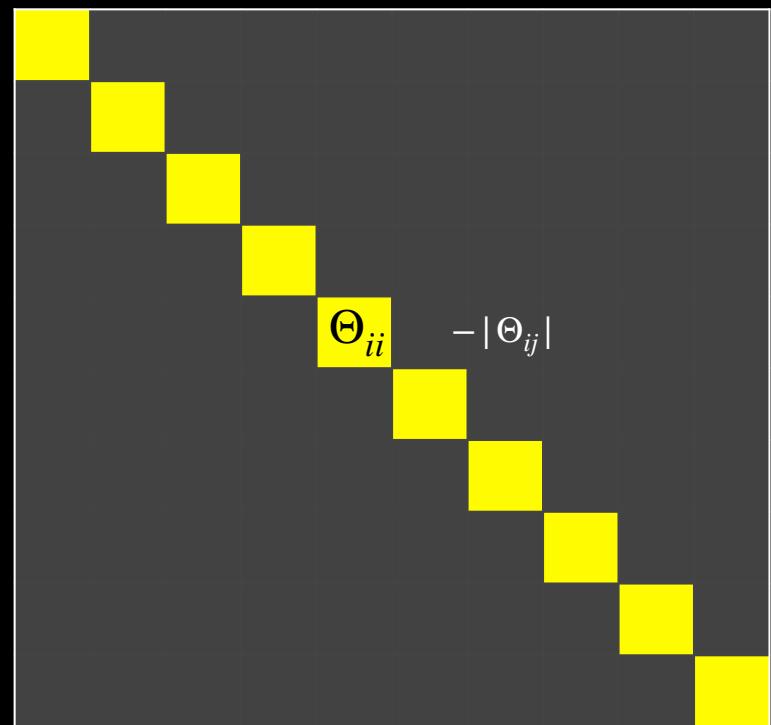


Offdiagonals negative ≥ 0

Walk-Summable GGMs

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

- **Introduced by Malioutov, Johnson, Willsky 2006**
- **Generalize many classes**
 - Attractive
 - Pairwise normalizable
 - Non-frustrated
 - Symmetric diagonally-dominant



Offdiagonals negative ≥ 0

Walk-Summable GGMs

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

Previous: ATHW12, MVL18 - $p^{O(d)}$ run-time algorithm with logarithmic samples.

Walk-Summable GGMs

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

KKMM: **GreedyPrune** learns walk-summable models with $O(d^2 \log p / \kappa^6)$ samples and quadratic run-time.

Previous: ATHW12, MVL18 - $p^{O(d)}$ run-time algorithm with logarithmic samples.

Walk-Summable GGMs

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

KKMM: Hybrid (Greedy+Lasso) learns walk-summable models with $O(d \log p / \kappa^4)$ samples.

Previous: ATHW12, MVL18 - $p^{O(d)}$ run-time algorithm with logarithmic samples.

Learning GGMs Greedily

Input: Samples from a sparse GGM $\sim \mathbf{X}$.
Output: Dependency graph of \mathbf{X} .

Similar greedy approaches for discrete GMs: [Bresler10],
[HKM17], [BKM19]

Learning GGMs Greedily

Input: Samples from a sparse GGM $\sim \mathbf{X}$.
Output: Dependency graph of \mathbf{X} .

GREEDYPRUNE

- 1. Recover neighborhood of each vertex in parallel.**
- 2. Grow a candidate neighborhood.**
- 3. Prune out some vertices.**

Similar greedy approaches for discrete GMs: [Bresler10],
[HKM17], [BKM19]

Phase 1: Growing a neighborhood

Input: Samples from a sparse GGM $\sim \mathbf{X}$.

Goal: Neighborhood of vertex 1.

Phase 1: Growing a neighborhood

Input: Samples from a sparse GGM $\sim \mathbf{X}$.

Goal: Neighborhood of vertex 1.

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.

Phase 2: Pruning a neighborhood

Input: Samples from a sparse GGM $\sim \mathbf{X}$.

Goal: Neighborhood of vertex 1.

Phase 2: Pruning a neighborhood

Input: Samples from a sparse GGM $\sim \mathbf{X}$.

Goal: Neighborhood of vertex 1.

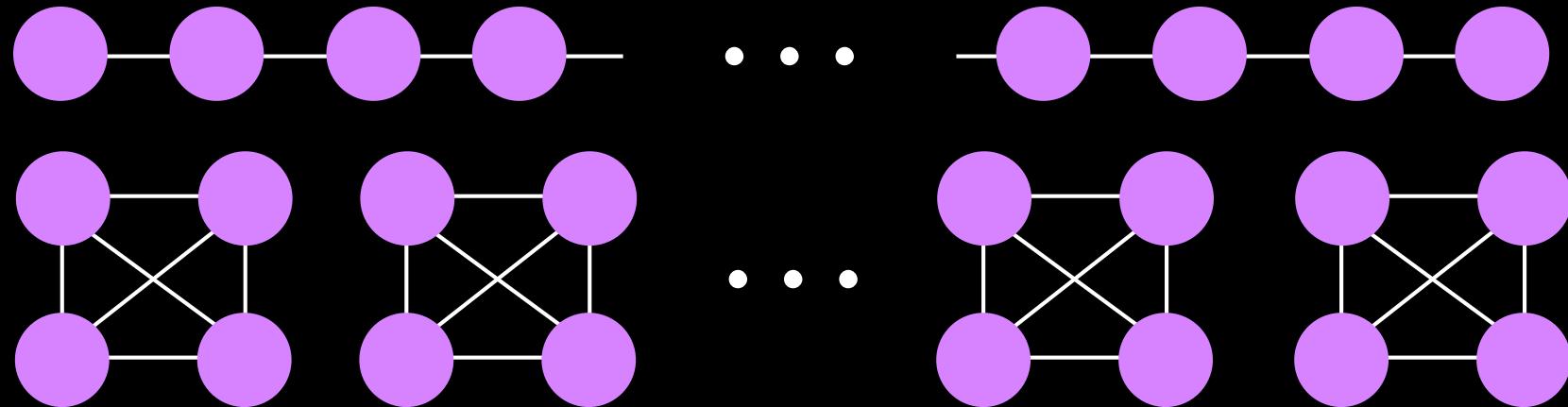
GREEDY-PRUNING

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.

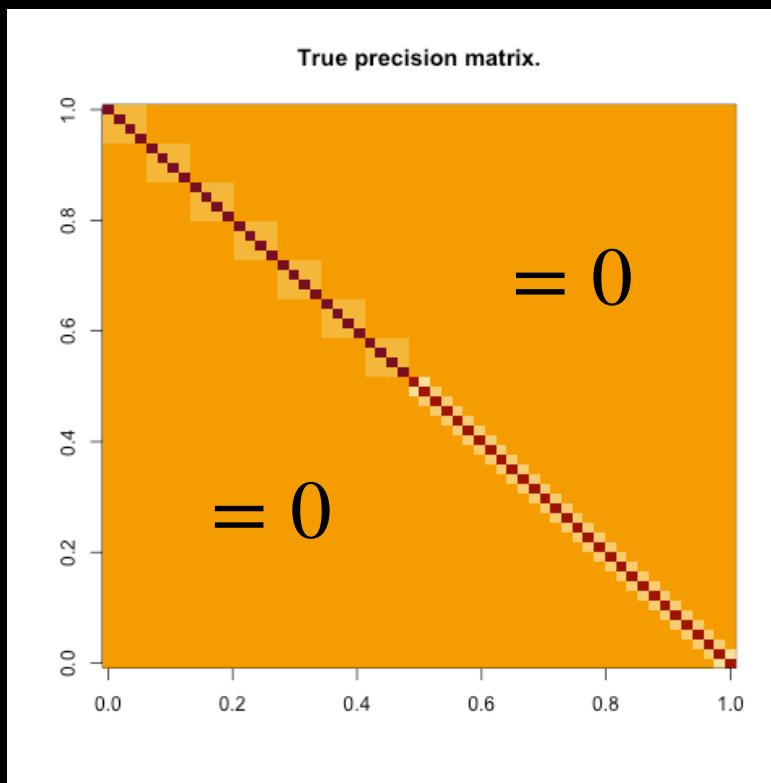
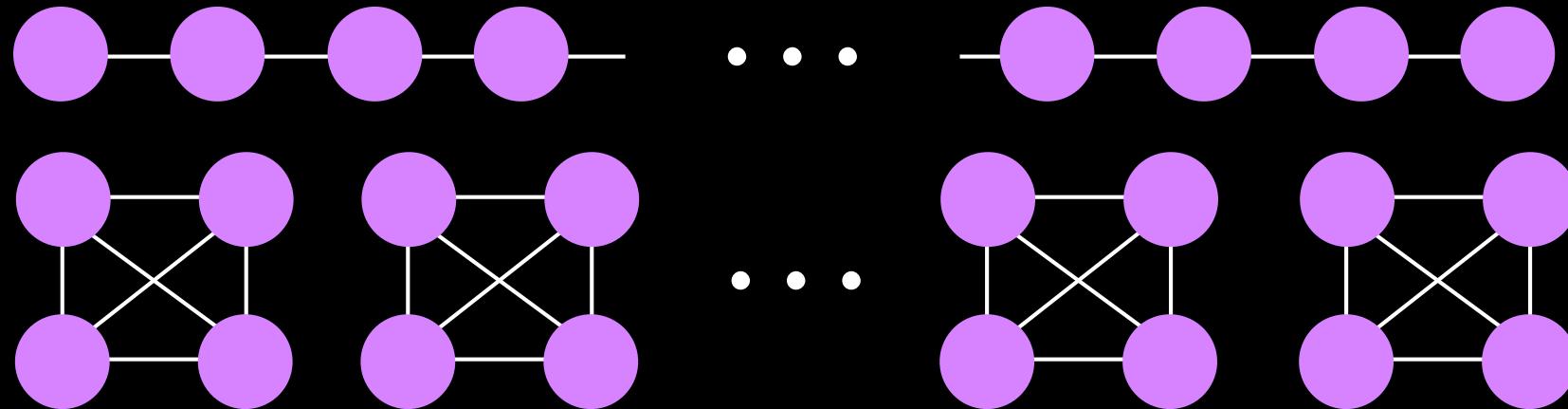
Experiments: A Simple Challenge



Precision: Path + Cliques

Max-degree = 3
Ill-conditioned

Experiments: A Simple Challenge

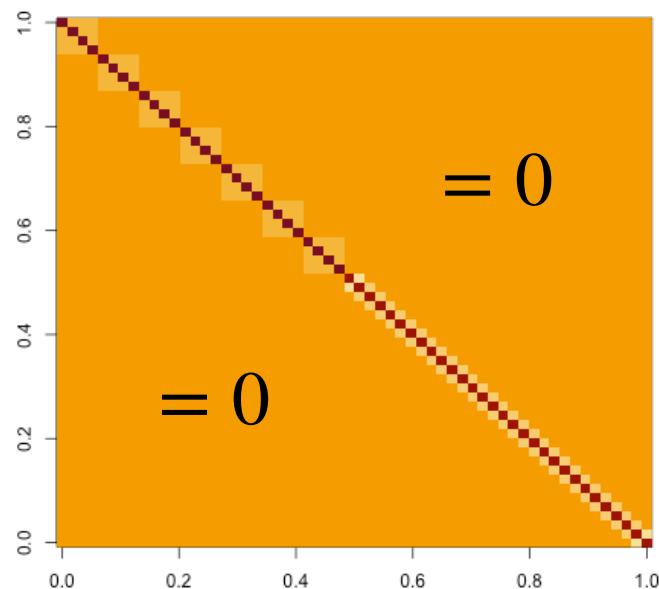


Precision: Path + Cliques

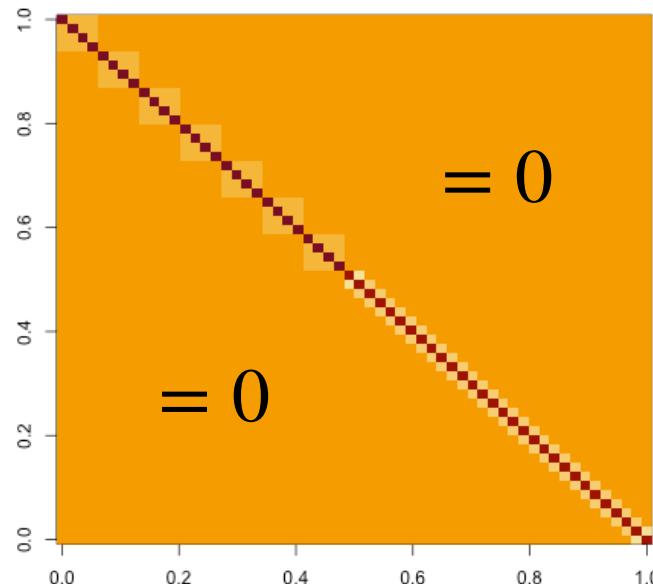
Max-degree = 3

III-conditioned

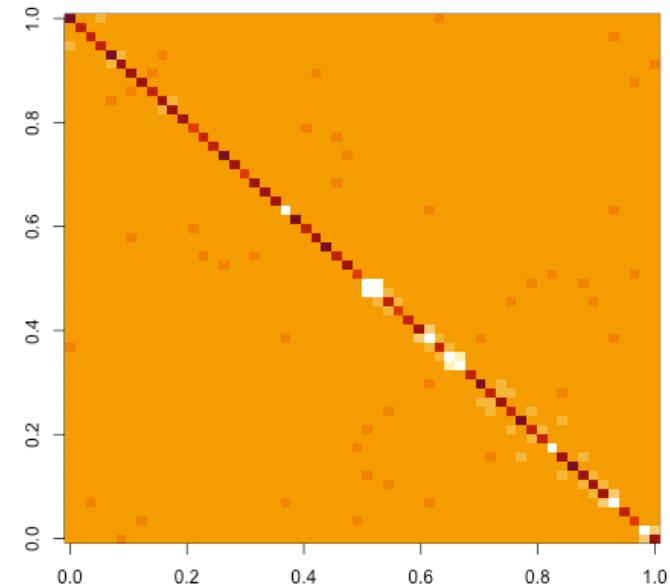
True precision matrix.



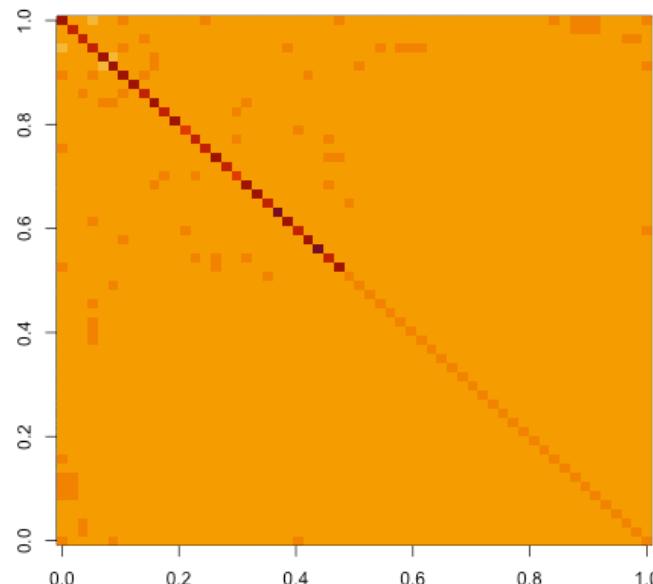
True precision matrix.



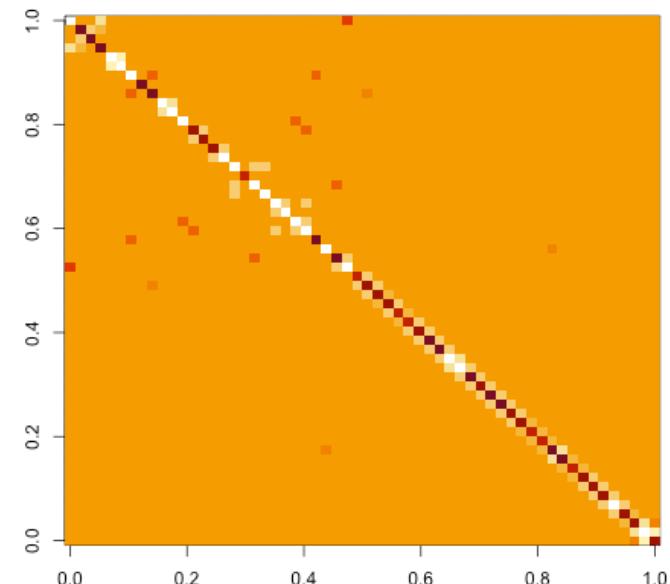
clime: 50



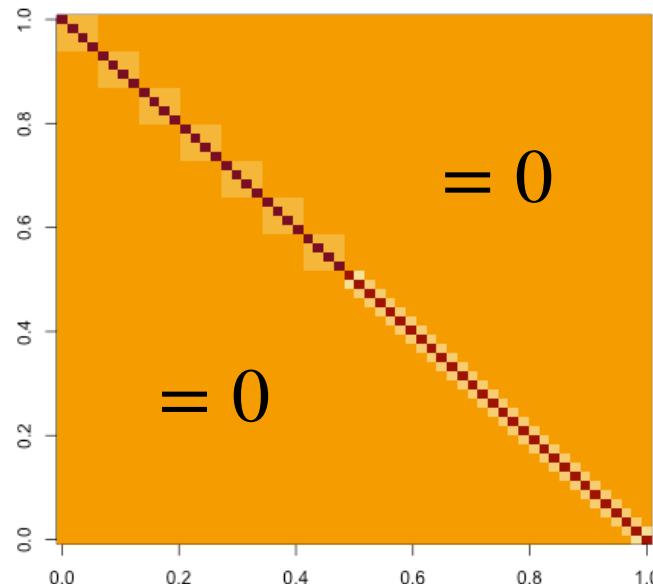
glasso: 50



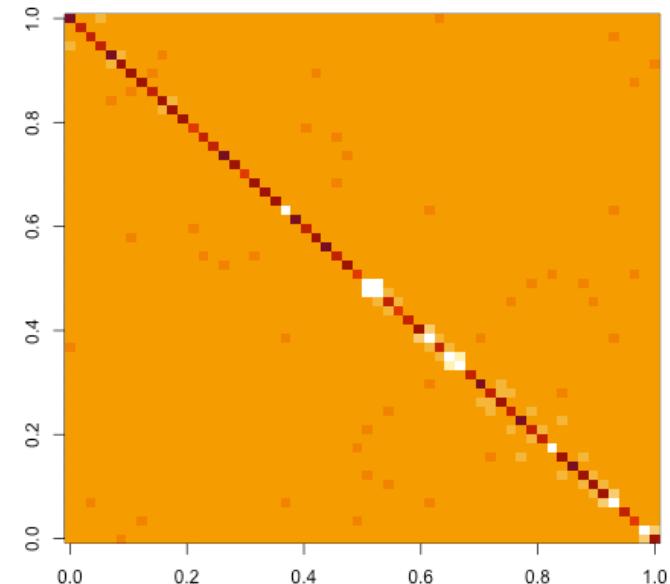
greedy: 50



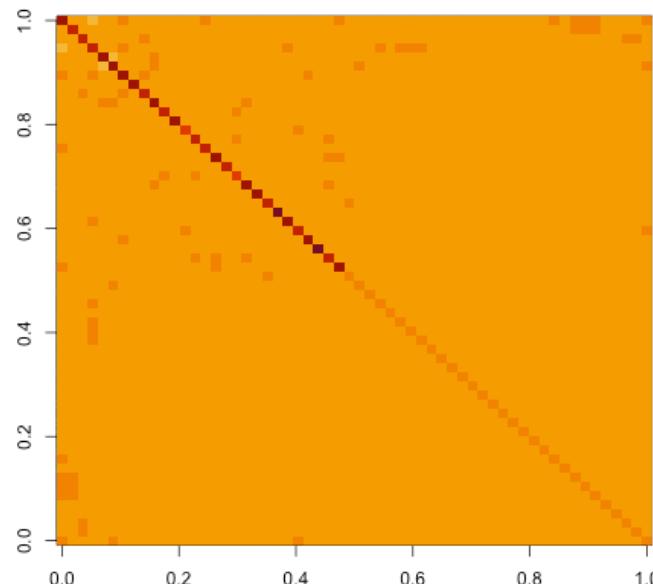
True precision matrix.



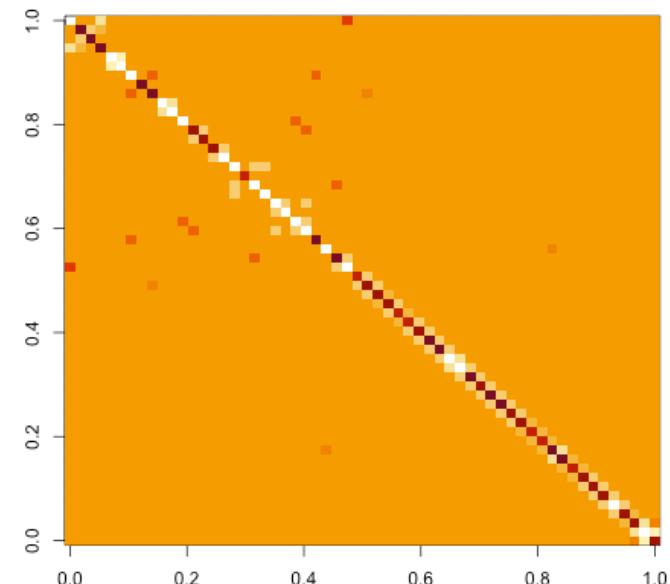
clime: 50



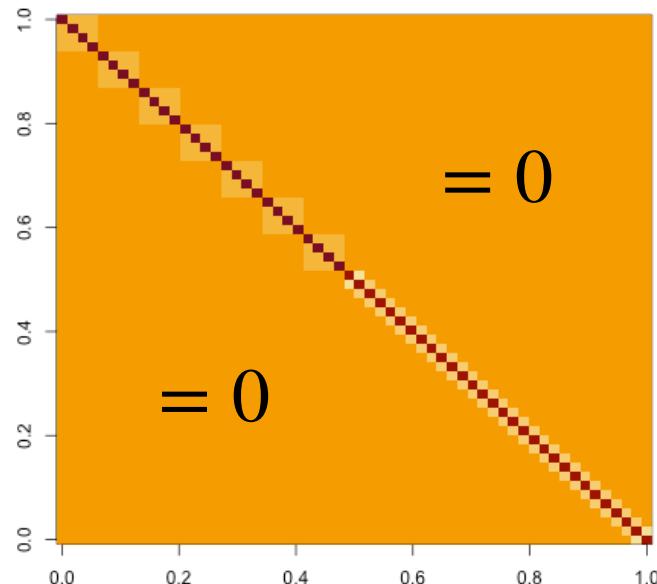
glasso: 50



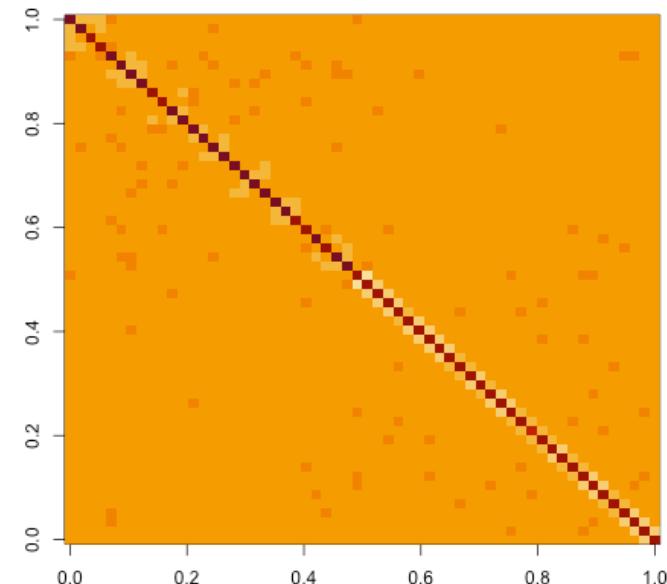
greedy: 50



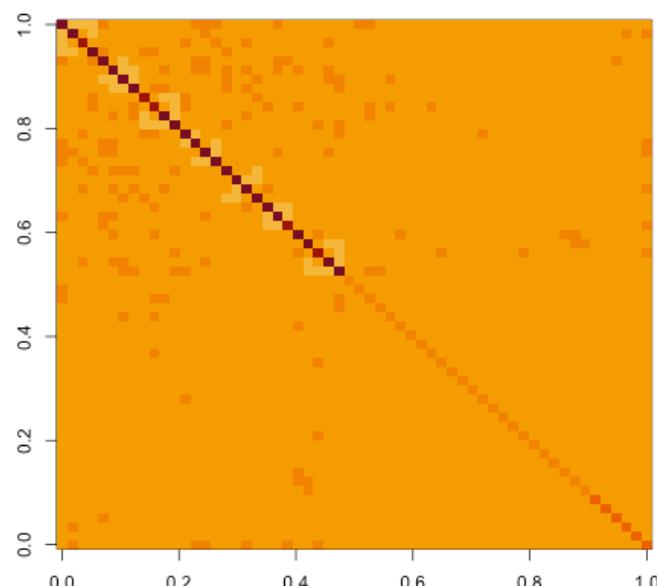
True precision matrix.



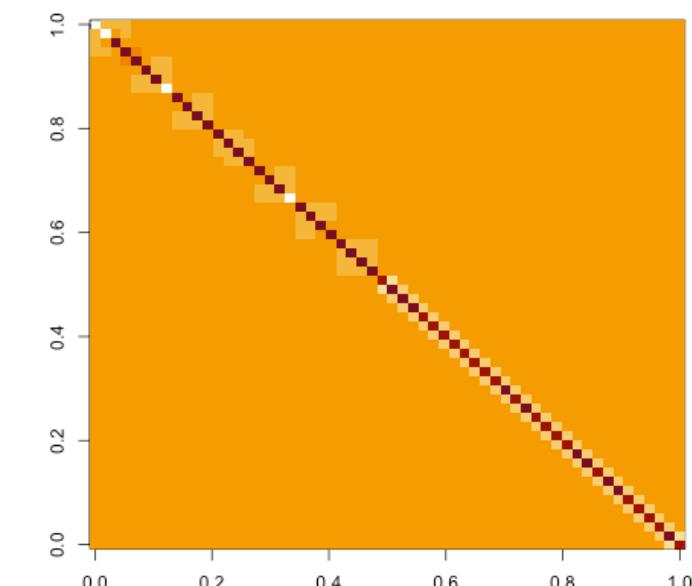
clime: 600



glasso: 600

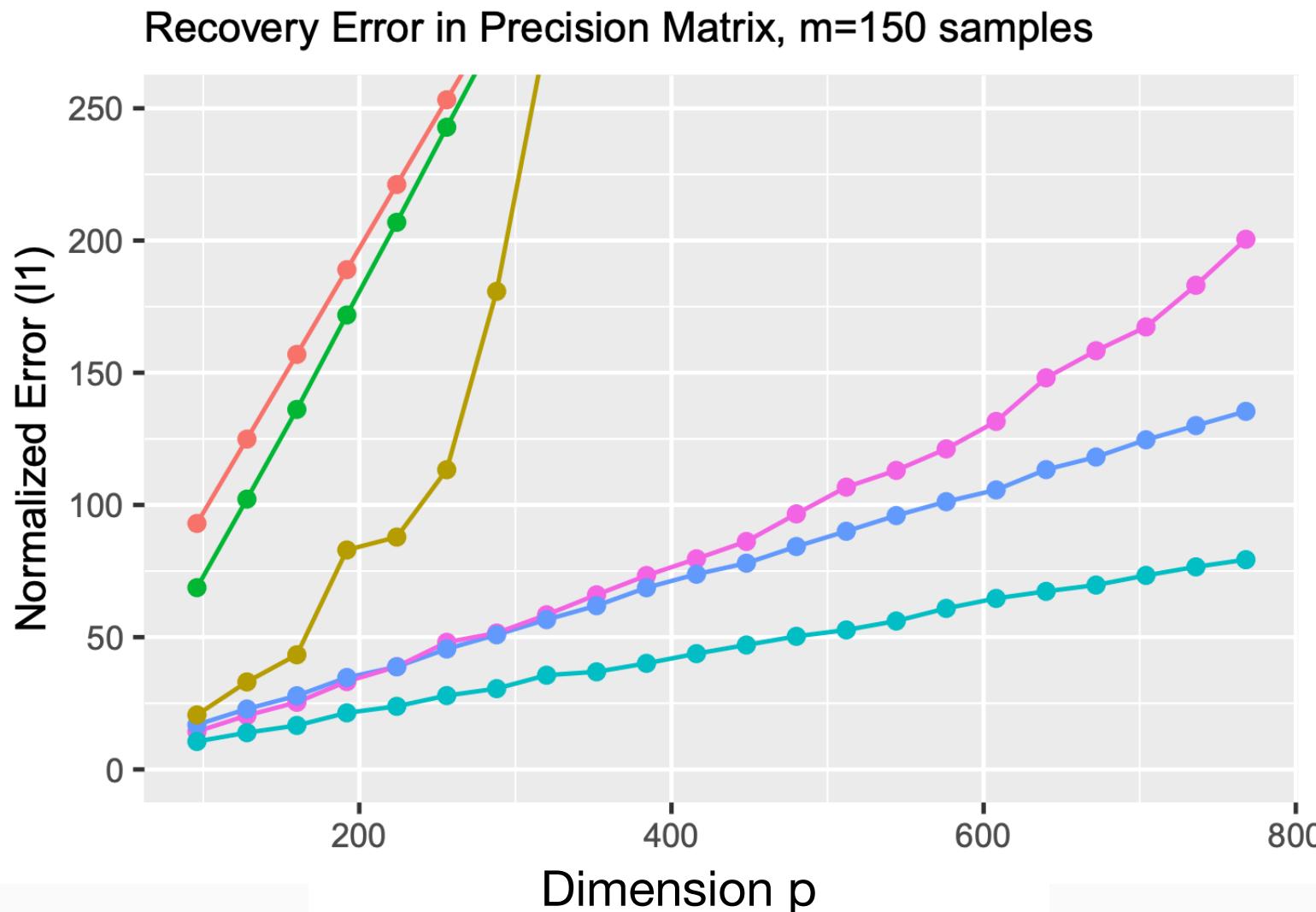


greedy: 600



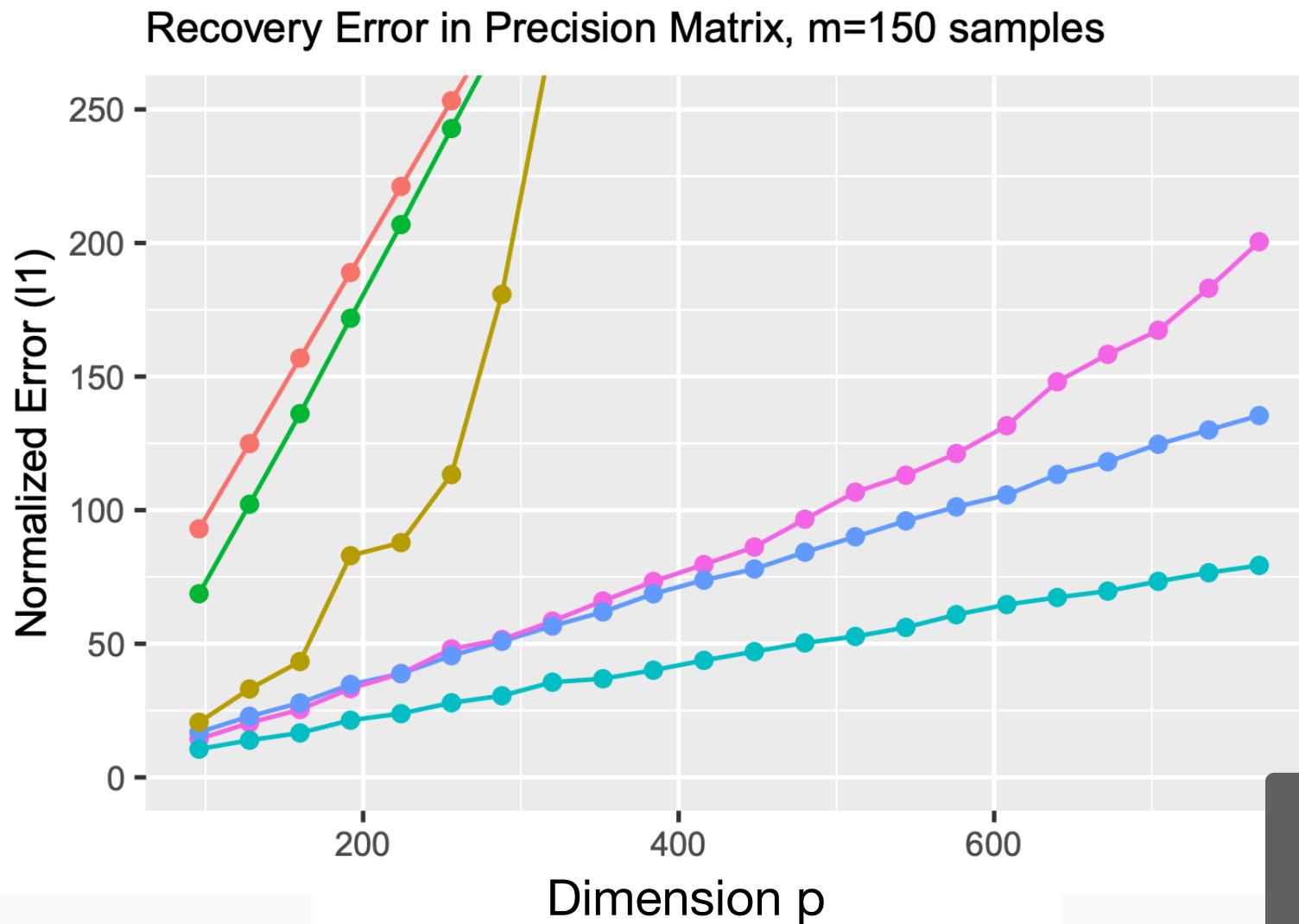
Final Outputs

A Simple Challenge: Path + Clique



GreedyPrune has best error

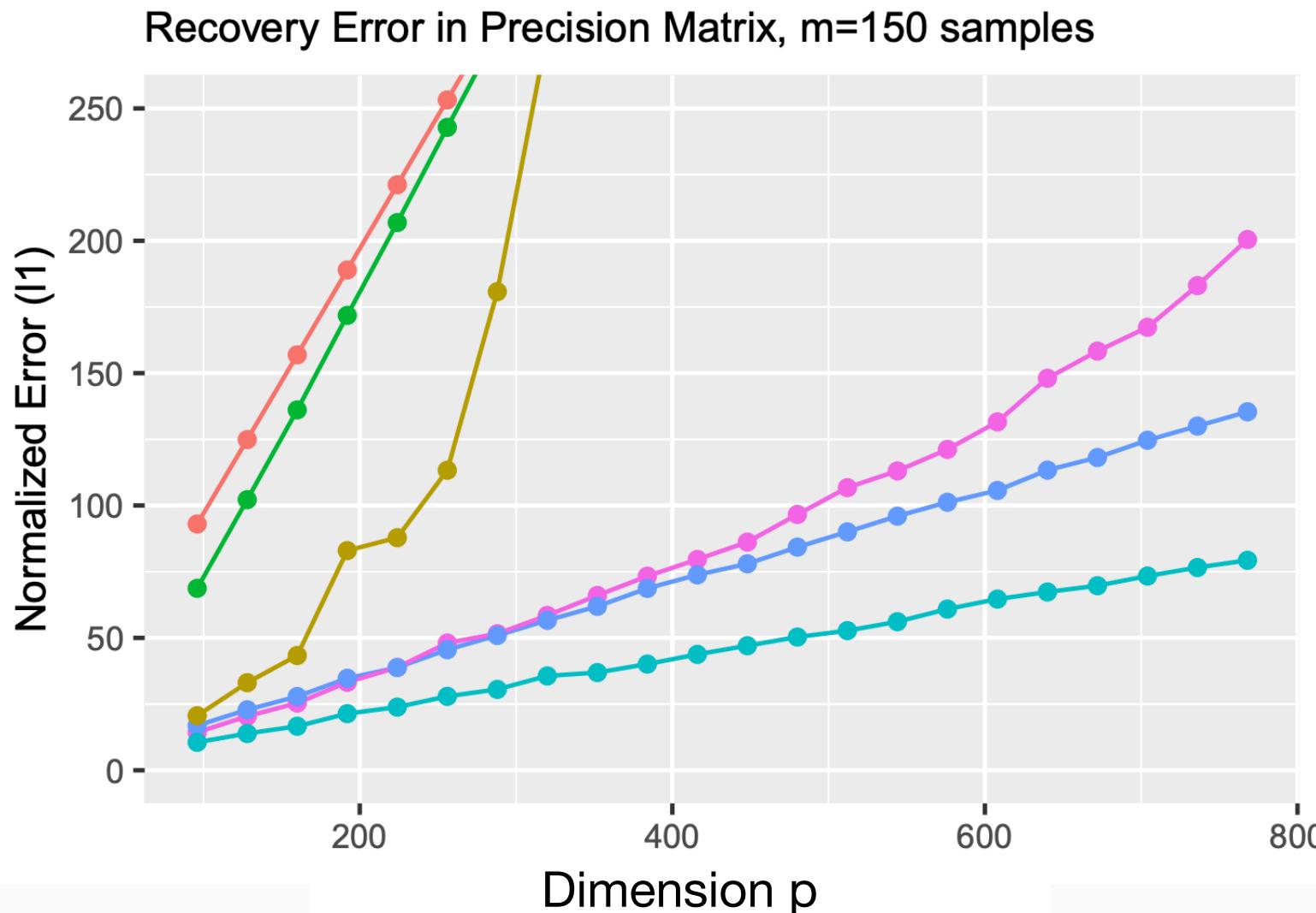
A Simple Challenge: Path + Clique



Meinhausen-Buhlmann
estimator (Lasso based)

GreedyPrune has best error

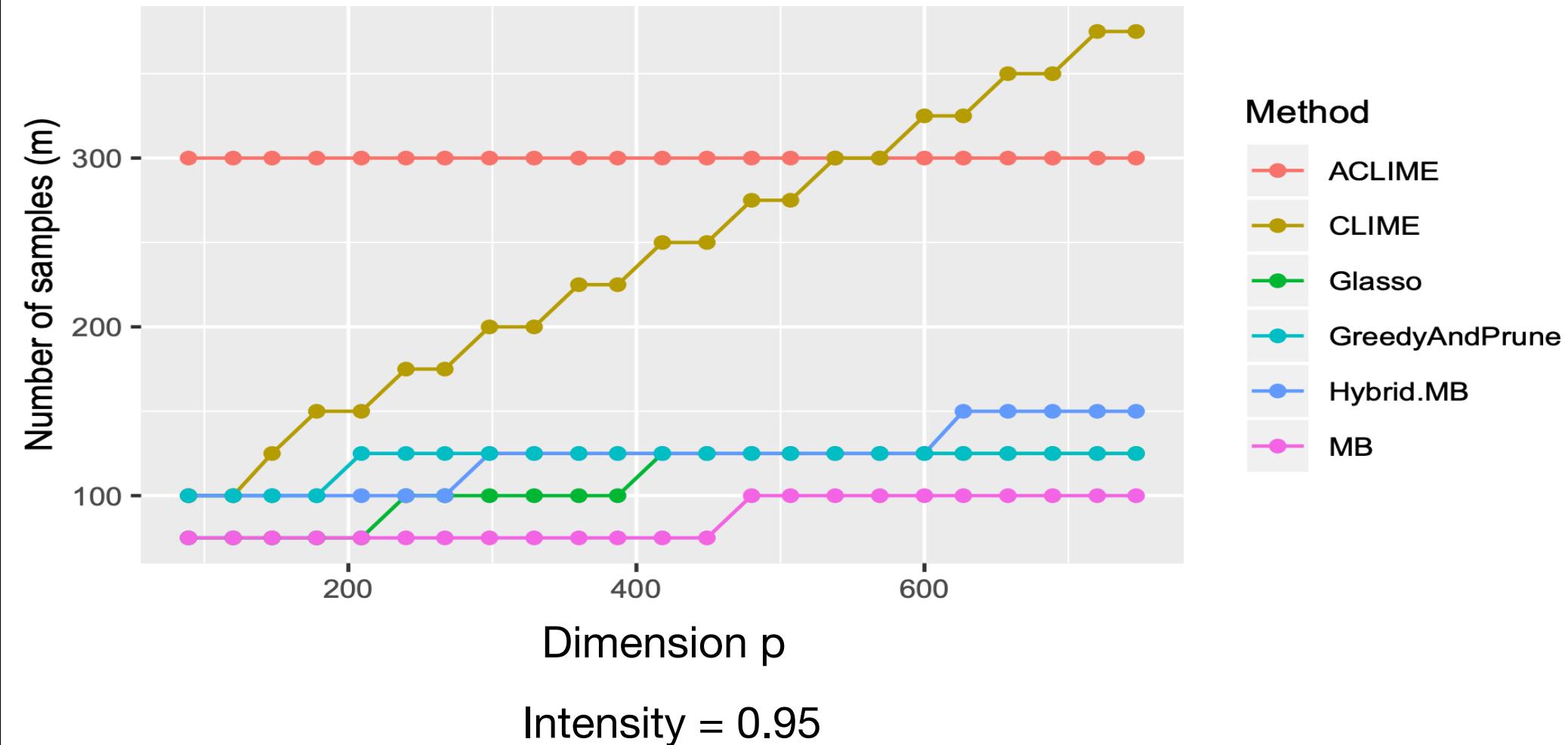
A Simple Challenge: Path + Clique



GreedyPrune has best error

A Simple Challenge: Path + Clique

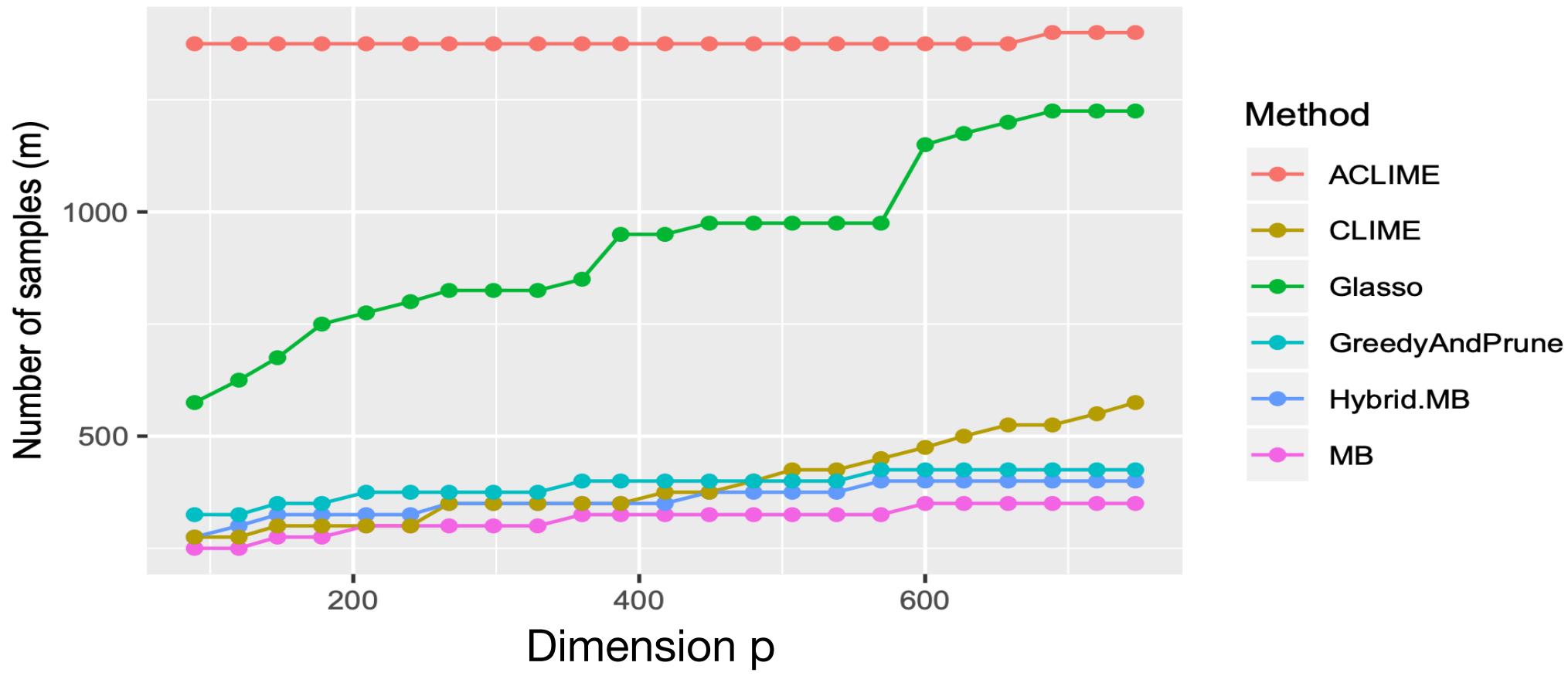
Sample Complexity for Edge Recovery



GreedyPrune needs very few samples.
CLIME grows nearly linearly ...

A Simple Challenge: Path + Clique

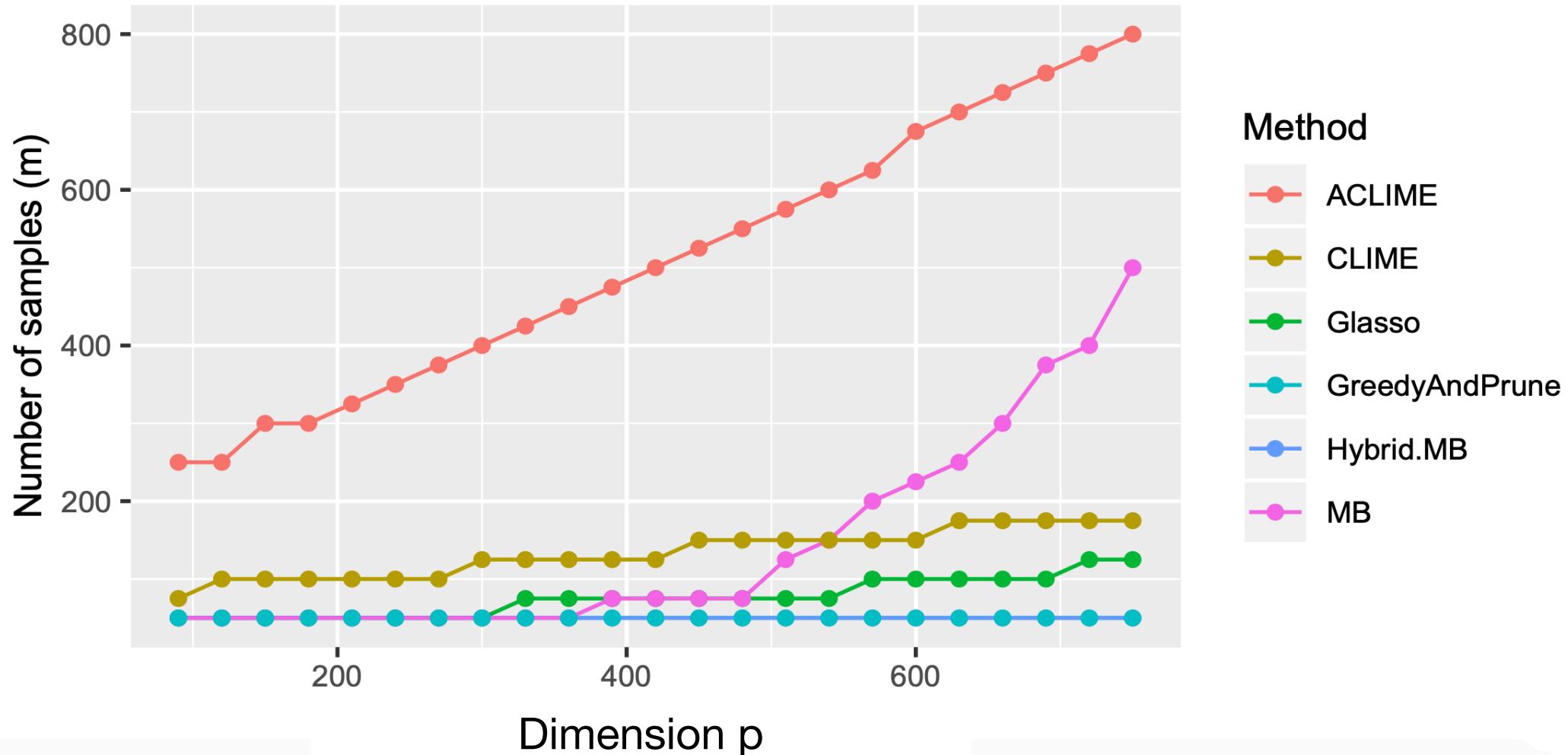
Sample Complexity for Edge Recovery



GreedyPrune needs very few samples.
GLASSO grows nearly linearly ...

A Simple Challenge: Random walk

Sample Complexity for Edge Recovery



GreedyPrune needs very few samples.
MB grows nearly linearly ...

GreedyPrune Summary

KKMM: **GreedyPrune** learns attractive models with $\tilde{O}(d \log p/\kappa^2)$ samples and quadratic run-time.

KKMM: **GreedyPrune** learns walk-summable models with $O(d^2 \log p/\kappa^6)$ samples and quadratic run-time.

GreedyPrune Summary

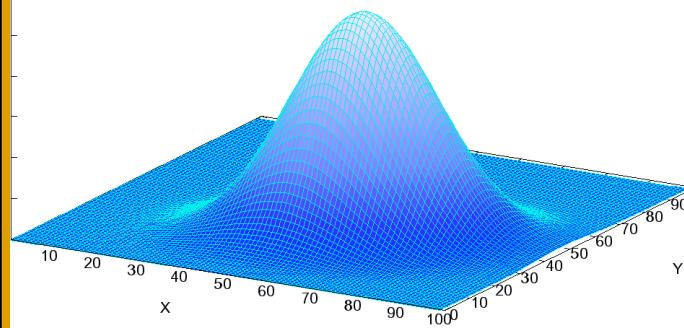
KKMM: **GreedyPrune** learns attractive models with $\tilde{O}(d \log p/\kappa^2)$ samples and quadratic run-time.

KKMM: **GreedyPrune** learns walk-summable models with $O(d^2 \log p/\kappa^6)$ samples and quadratic run-time.

Also ...

- Recovers guarantees of GLASSO, CLIME
- Empirically better
- Non-Gaussian distributions: Can learn precision matrices if good tail behaviour

GGMs



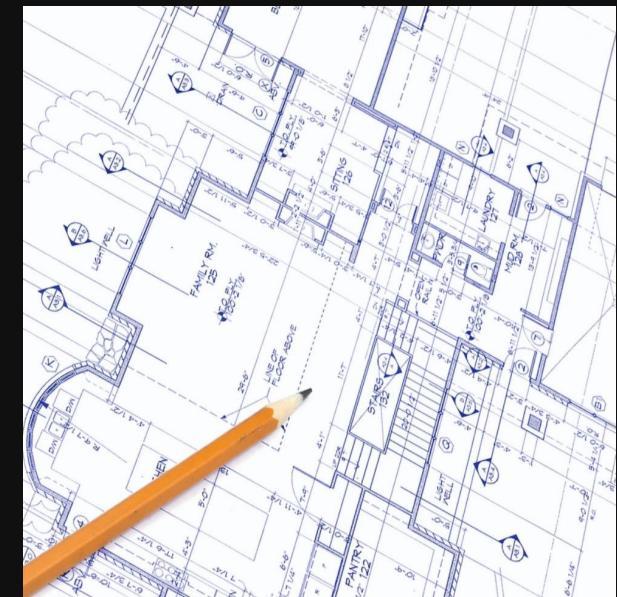
Motivation

Algorithm

```
While |S| < t :  
Add  $\arg \min_j Var(X_i | X_{S \cup j})$ .
```

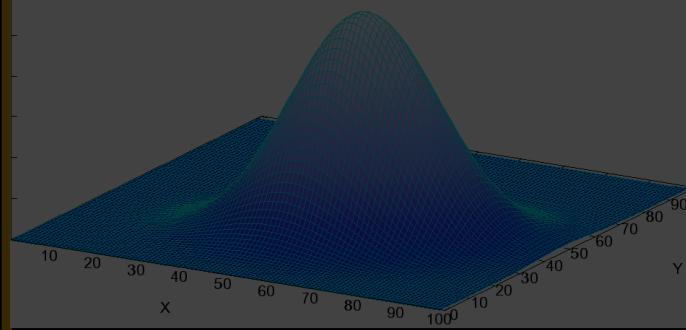
Special Models

Analysis



Attractive. SDD

GGMs



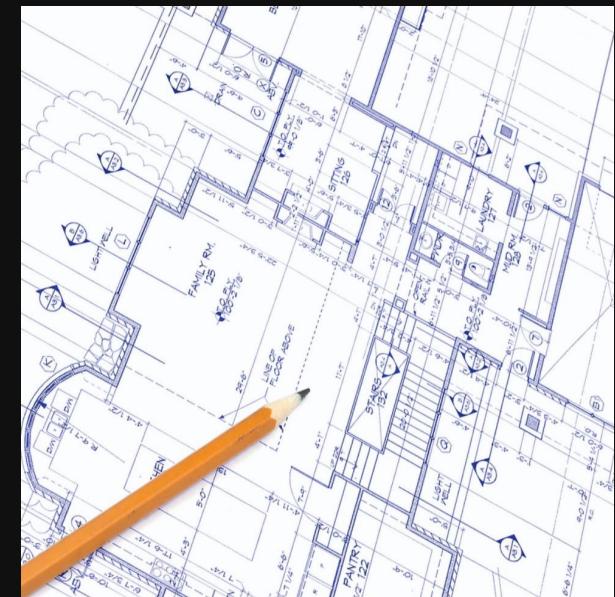
Motivation

Algorithm

```
While  $|S| < t$  :  
Add  $\arg \min_j \text{Var}(X_i | X_{S \cup j})$ .
```

Special Models

Analysis



Attractive. SDD

Summary

Main Challenge: Can we recover structure of GGM of degree d , pairwise-correlation κ with run-time $p^{o(d)}$ and sample complexity $n \approx O_{d,\kappa}(poly(\log p))$?

Very well-studied, practically important!

Summary

Main Challenge: Can we recover structure of GGM of degree d , pairwise-correlation κ with run-time $p^{o(d)}$ and sample complexity $n \approx O_{d,\kappa}(poly(\log p))$?

Very well-studied, practically important!

Today: A simple greedy algorithm solves interesting special classes

Many Questions for GGMs ...

- What other classes can we solve efficiently without condition number assumptions?
- Testing if a model is correct?
- Hidden variables?
- Computational hardness?

Many Questions for GGMs ...

- What other classes can we solve efficiently without condition number assumptions?
- Testing if a model is correct?
- Hidden variables?
- Computational hardness?

Conjecture: No $p^{o(d)}$ algorithm to recover structure of general GGMs with $\text{poly}(d, 1/\kappa, \log p)$ samples.

Bigger Picture

Can we learn sparse dependency graphs from few samples?

(aka learning Markov random fields, undirected graphical models)

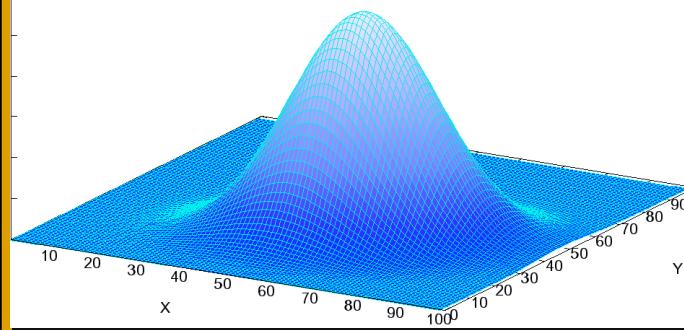
Bigger Picture

Can we learn sparse dependency graphs from few samples?

(aka learning Markov random fields, undirected graphical models)

Lots of work ... [Bresler10], [KM17], [HKM17], ...

GGMs



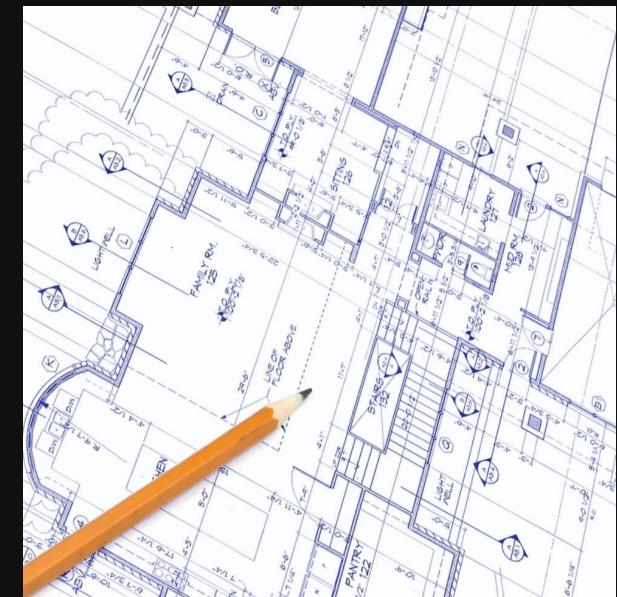
Motivation

Algorithm

```
While  $|S| < t$  :  
Add  $\arg \min_j \text{Var}(X_i | X_{S \cup j})$ .
```

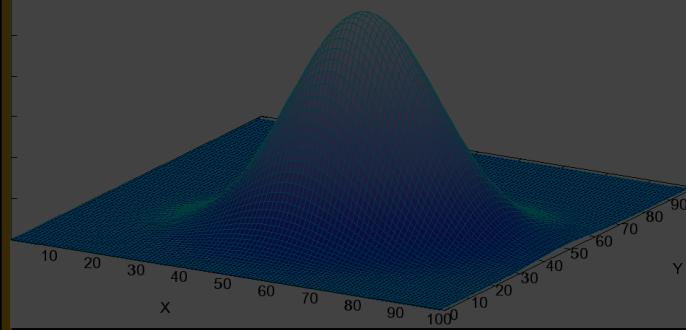
Special Models

Analysis



Attractive. SDD

GGMs



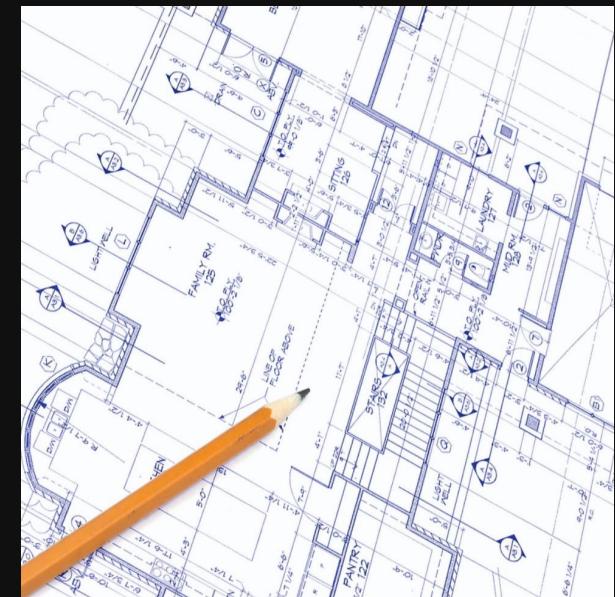
Motivation

Algorithm

```
While  $|S| < t$  :  
Add  $\arg \min_j \text{Var}(X_i | X_{S \cup j})$ .
```

Special Models

Analysis



Attractive. SDD

Learning GGMs Greedily

Input: Samples from a sparse GGM $\sim X$.

Output: Dependency graph of X .

(d: max-degree; p: dimension;
n: num. samples; κ : pairwise correlation)

Learning GGMs Greedily

Input: Samples from a sparse GGM $\sim \mathbf{X}$.

Output: Dependency graph of \mathbf{X} .

(d : max-degree; p : dimension;

n : num. samples; κ : pairwise correlation)

Thm 1: GreedyPrune learns attractive models with $\tilde{O}(d \log p / \kappa^2)$ samples and quadratic run-time.

Learning GGMs Greedily

Input: Samples from a sparse GGM $\sim \mathbf{X}$.

Output: Dependency graph of \mathbf{X} .

(d : max-degree; p : dimension;
 n : num. samples; κ : pairwise correlation)

Thm 1: **GreedyPrune** learns attractive models with
 $\tilde{O}(d \log p / \kappa^2)$ samples and quadratic run-time.

Thm 2: **GreedyPrune** learns walk-summable models
with $O(d^2 \log p / \kappa^6)$ samples and quadratic run-time.

Analysis for Attractive: Supermodularity

Fix vertex 1. For $S \subseteq \{2, \dots, n\}$, define

$$f(S) = \text{Var}(X_1 \mid X_S).$$

Analysis for Attractive: Supermodularity

Fix vertex 1. For $S \subseteq \{2, \dots, n\}$, define

$$f(S) = \text{Var}(X_1 | X_S).$$

Claim [MJW06, MS12]: X is attractive GGM. Then, $f(\cdot)$ is monotonically decreasing and a supermodular function.

Analysis for Attractive: Supermodularity

Fix vertex 1. For $S \subseteq \{2, \dots, n\}$, define

$$f(S) = \text{Var}(X_1 | X_S).$$

Claim [MJW06, MS12]: X is attractive GGM. Then, $f(\cdot)$ is monotonically decreasing and a **supermodular** function.

Supermodular: $S \subset T$,

$$f(S) - f(S \cup \{j\}) \geq f(T) - f(T \cup \{j\}).$$

Analysis for Attractive: Supermodularity

Fix vertex 1. For $S \subseteq \{2, \dots, n\}$, define

$$f(S) = \text{Var}(X_1 | X_S).$$

Claim [MS, MJW]: X is attractive GGM. Then, $f(\cdot)$ is monotonically decreasing and a **supermodular** function.

Why useful?

- Optimizing supermodularity very well understood
- Greedy algorithm finds minimizer

Analysis for Attractive: Supermodularity

Fix vertex 1. For $S \subseteq \{2, \dots, n\}$, define

$$f(S) = \text{Var}(X_1 | X_S).$$

Claim [MS, MJW]: X is attractive GGM. Then, $f(\cdot)$ is monotonically decreasing and a **supermodular** function.

Are we done?

Analysis for Attractive: Supermodularity

Fix vertex 1. For $S \subseteq \{2, \dots, n\}$, define

$$f(S) = \text{Var}(X_1 | X_S).$$

Claim [MS, MJW]: X is attractive GGM. Then, $f(\cdot)$ is monotonically decreasing and a **supermodular** function.

Are we done?

Almost ...

- Only have estimates for f
- Crucial: missing a vertex means noticeably away from optimum

Analysis for Attractive: Supermodularity

Fix vertex 1. For $S \subseteq \{2, \dots, n\}$, define

$$f(S) = \text{Var}(X_1 | X_S).$$

Claim [MS, MJW]: X is attractive GGM. Then, $f(\)$ is monotonically decreasing and a **supermodular** function.

Analysis for Attractive: Supermodularity

Fix vertex 1. For $S \subseteq \{2, \dots, n\}$, define
 $f(S) = \text{Var}(X_1 | X_S)$.

Claim [MS, MJW]: X is attractive GGM. Then, $f(\)$ is monotonically decreasing and a **supermodular** function.

Proof: ...

Analysis for Attractive: Supermodularity

Fix vertex 1. For $S \subseteq \{2, \dots, n\}$, define
 $f(S) = \text{Var}(X_1 | X_S)$.

Claim [MS, MJW]: X is attractive GGM. Then, $f(\cdot)$ is monotonically decreasing and a **supermodular** function.

Proof: ...

- Assume $\Theta_{ii} = 1$.

Analysis for Attractive: Supermodularity

Fix vertex 1. For $S \subseteq \{2, \dots, n\}$, define
 $f(S) = \text{Var}(X_1 | X_S)$.

Claim [MS, MJW]: X is attractive GGM. Then, $f(\)$ is monotonically decreasing and a **supermodular** function.

Proof: ...

- Assume $\Theta_{ii} = 1$.
- Linear algebra:

$$\text{Var}(X_1 | X_S) = \sum_{i_1, i_2, \dots, i_k \notin S} (-\Theta_{ii_1})(-\Theta_{i_1 i_2}) \cdots (-\Theta_{i_k i})$$

Analysis for Attractive: Supermodularity

Fix vertex 1. For $S \subseteq \{2, \dots, n\}$, define
 $f(S) = \text{Var}(X_1 | X_S)$.

Claim [MS, MJW]: X is attractive GGM. Then, $f(\)$ is monotonically decreasing and a **supermodular** function.

Proof: ...

-

Product over cycles not touching S ...

$$\text{Var}(X_1 | X_S) = \sum_{i_1, i_2, \dots, i_k \notin S} (-\Theta_{ii_1})(-\Theta_{i_1 i_2}) \cdots (-\Theta_{i_k i})$$

Analysis for Attractive: Supermodularity

Fix vertex 1. For $S \subseteq \{2, \dots, n\}$, define
 $f(S) = \text{Var}(X_1 | X_S)$.

Claim [MS, MJW]: X is attractive GGM. Then, $f(\)$ is monotonically decreasing and a **supermodular** function.

Proof: ...

-

Product over cycles not touching S ...

$$\text{Var}(X_1 | X_S) = \sum_{i_1, i_2, \dots, i_k \notin S} (-\Theta_{ii_1})(-\Theta_{i_1 i_2}) \cdots (-\Theta_{i_k i})$$

- Thus: $f(S) - f(S \cup \{j\})$ is sum on cycles that touch j but don't touch S .

Analysis for Attractive: Supermodularity

Fix vertex 1. For $S \subseteq \{2, \dots, n\}$, define
 $f(S) = \text{Var}(X_1 | X_S)$.

Claim [MS, MJW]: X is attractive GGM. Then, $f(\cdot)$ is monotonically decreasing and a **supermodular** function.

Supermodular: $S \subset T$,

$$f(S) - f(S \cup \{j\}) \geq f(T) - f(T \cup \{j\}).$$

$$\text{Var}(X_1 | X_S) = \sum_{i_1, i_2, \dots, i_k \notin S} (-\Theta_{ii_1})(-\Theta_{i_1 i_2}) \cdots (-\Theta_{i_k i})$$

- **Thus:** $f(S) - f(S \cup \{j\})$ is sum on cycles that touch j but don't touch S .

Analysis for Attractive: Supermodularity

Fix vertex 1. For $S \subseteq \{2, \dots, n\}$, define
 $f(S) = \text{Var}(X_1 | X_S)$.

Claim [MS, MJW]: X is attractive GGM. Then, $f(\)$ is monotonically decreasing and a **supermodular** function.

Supermodular: $S \subset T$,

$$f(S) - f(S \cup \{j\}) \geq f(T) - f(T \cup \{j\}).$$

$$\text{Var}(X_1 | X_S) = \sum_{i_1, i_2, \dots, i_k \notin S} (-\Theta_{ii_1})(-\Theta_{i_1 i_2}) \cdots (-\Theta_{i_k i})$$

- **Thus:** $f(S) - f(S \cup \{j\})$ is sum on cycles that touch j but don't touch S .
- Is decreasing with S , so f is supermodular.

Analysis for Walk-Summable

Fix vertex 1. For $S \subseteq \{2, \dots, n\}$, define

$$f(S) = \text{Var}(X_1 \mid X_S).$$

No longer supermodular!
(Not even weakly-supermodular ...)

Analysis for Walk-Summable

Fix vertex 1. For $S \subseteq \{2, \dots, n\}$, define

$$f(S) = \text{Var}(X_1 | X_S).$$

No longer supermodular!
(Not even weakly-supermodular ...)

- Walk-summable = Symmetric Diagonally Dominant (SDD)
- SDD implies neighbors have noticeable effect

Analysis for Walk-Summable

Def: Θ is SDD if for all i

$$|\Theta_{ii}| \geq \sum_{j \neq i} |\Theta_{ij}|$$

- **Walk-summable = Symmetric Diagonally Dominant (SDD)**
- **SDD implies neighbors have noticeable effect**

Analysis for Walk-Summable

Def: Θ is SDD if for all i

$$|\Theta_{ii}| \geq \sum_{j \neq i} |\Theta_{ij}|$$

Claim [RTT09]: Θ walk-summable iff a diagonal $D \geq 0$,
 $D\Theta D$ is SDD.

- Walk-summable = Symmetric Diagonally Dominant (SDD)
- SDD implies neighbors have noticeable effect

Analysis for Walk-Summable

Def: Θ is SDD if for all i

$$|\Theta_{ii}| \geq \sum_{j \neq i} |\Theta_{ij}|$$

Claim [RTT09]: Θ walk-summable iff a diagonal $D \geq 0$,
 $D\Theta D$ is SDD.

Suffices to show greedy works for SDD!

- Walk-summable = Symmetric Diagonally Dominant (SDD)
- SDD implies neighbors have noticeable effect

Analysis: Bounded Conditional Variances

Lemma: If Θ is SDD, then for any vertex i , there exists a neighbor j such that $Var(X_i | X_j) \leq 4Var(X_i | X_{-i})/\kappa^2$.

Analysis: Bounded Conditional Variances

There exists a neighbor, that is comparable to the entire neighborhood in conditioning ...

Lemma: If Θ is SDD, then for any vertex i , there exists a neighbor j such that $Var(X_i | X_j) \leq 4Var(X_i | X_{-i})/\kappa^2$.

Analysis: Bounded Conditional Variances

Why useful ...

- $\text{Var}(X_i)$ can be very large!
- But conditioning on one neighbor brings the variance down.
- We can detect it with few samples.

Lemma: If Θ is SDD, then for any vertex i , there exists a neighbor j such that $\text{Var}(X_i | X_j) \leq 4\text{Var}(X_i | X_{-i})/\kappa^2$.

Analysis: Bounded Conditional Variances

Proof ...

- Triangle inequality of **effective resistance** metric on graphs
- Properties of conditional variance ...

Lemma: If Θ is SDD, then for any vertex i , there exists a neighbor j such that $Var(X_i | X_j) \leq 4Var(X_i | X_{-i})/\kappa^2$.

Analysis: Bounded Conditional Variances

For intuition: Assume $\Theta_{ii} = 1, \forall i$.

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

So $i \sim j \Rightarrow |\Theta_{ij}| \geq \kappa$.

Analysis: Bounded Conditional Variances

For intuition: Assume $\Theta_{ii} = 1, \forall i$.

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

So $i \sim j \Rightarrow |\Theta_{ij}| \geq \kappa$.

Lemma 1: If Θ is SDD, then conditional variances are bounded — $i \sim j \Rightarrow \text{Var}(X_i | X_j) \leq 1/\kappa$.

Analysis: Bounded Conditional Variances

For intuition: Assume $\Theta_{ii} = 1, \forall i$.

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

So $i \sim j \Rightarrow |\Theta_{ij}| \geq \kappa$.

Lemma 1: If Θ is SDD, then conditional variances are bounded — $i \sim j \Rightarrow \text{Var}(X_i | X_j) \leq 1/\kappa$.

Analysis: Bounded Conditional Variances

For intuition: Assume $\Theta_{ii} = 1, \forall i$.

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

So $i \sim j \Rightarrow |\Theta_{ij}| \geq \kappa$.

Lemma 1: If Θ is SDD, then conditional variances are bounded — $i \sim j \Rightarrow \text{Var}(X_i | X_j) \leq 1/\kappa$.

Example: In random walk model

$$\text{Var}(X_{i+1} | X_i) = 1; \quad \text{Var}(X_{i+1}) = i + 1!$$

Analysis: Bounded Conditional Variances

For intuition: Assume $\Theta_{ii} = 1, \forall i$.

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

So $i \sim j \Rightarrow |\Theta_{ij}| \geq \kappa$.

Lemma 1: If Θ is SDD, then conditional variances are bounded — $i \sim j \Rightarrow \text{Var}(X_i | X_j) \leq 1/\kappa$.

Not true for general sparse precision matrices ...

Analysis: Bounded Conditional Variances

For intuition: Assume $\Theta_{ii} = 1, \forall i$.

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

So $i \sim j \Rightarrow |\Theta_{ij}| \geq \kappa$.

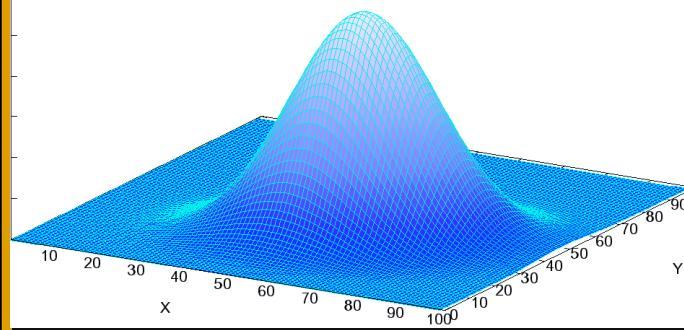
Lemma 1: If Θ is SDD, then conditional variances are bounded — $i \sim j \Rightarrow \text{Var}(X_i | X_j) \leq 1/\kappa$.

“Proof ...”:

- If Laplacian ...

$$\text{Var}(X_i | X_j) = \frac{1}{2} R_{eff}(i, j) \leq \frac{1}{|\Theta_{ij}|} \leq 1/\kappa.$$

GGMs



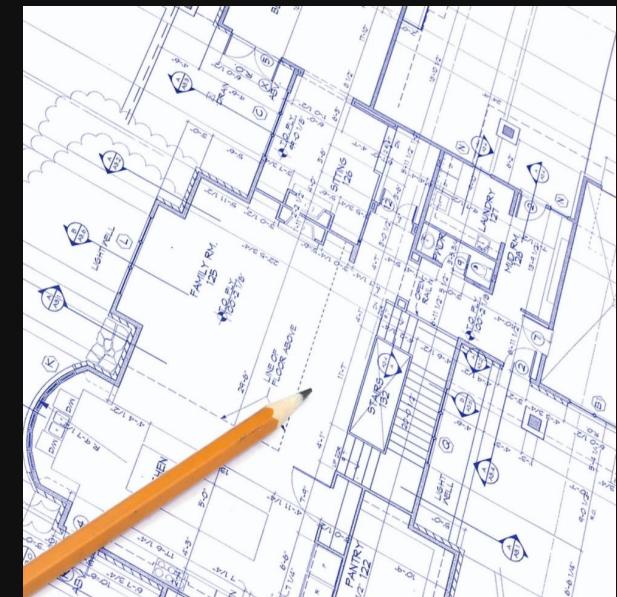
Motivation

Algorithm

```
While  $|S| < t$  :  
Add  $\arg \min_j \text{Var}(X_i | X_{S \cup j})$ .
```

Special Models

Analysis



Attractive. SDD