

# Undirected Graphical Model

Guoyu Zhang, Jiaqiang Li

University of Science and Technology of China

December 3, 2021

# Content

## 1 Introduction and preliminaries

- Connecting an undirected graph to random variables
- Examples of Graph Model

## 2 Inferring of GGM

- Estimation of Parameters when Graph Structure is Known
- Estimation of the Graph Structure
- Generalization of Graphical Lasso

# Table of Contents

## 1 Introduction and preliminaries

- Connecting an undirected graph to random variables
- Examples of Graph Model

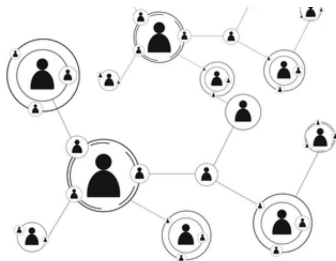
## 2 Inferring of GGM

- Estimation of Parameters when Graph Structure is Known
- Estimation of the Graph Structure
- Generalization of Graphical Lasso

# Motivation

- Many researches concerns to obtain the empirical evidence on variables capturing characteristics, behaviour, abilities, attitudes of people. These variables are properties of observational units.
- The objective of such studies is to improve knowledge about the association structure of the observed units via their properties, or to infer the system of variables via their association structure.

# Social Network



**Figure:** An example of social network. We want to study the relationship of people via their properties, or predict the behaviour of people via their contact network.

# Motivation

- One hopes represent the structure as sparse as possible, i.e. with a few strong associations, which account for the indirect relations in the system.
- Direct and only indirect relations are key notions in studying association structures. To say that a relation between a pair of variables is indirect means that they have a substantial association given the information on other variables.
- This implies the system of variables are generally related, but only a few of these relationships are directly connected without conditioning on extra information.

# Why Graphs?

There are sometimes that what we know or are interested in only concerns about the (conditional) independence among variables, rather than their accurate relationships.

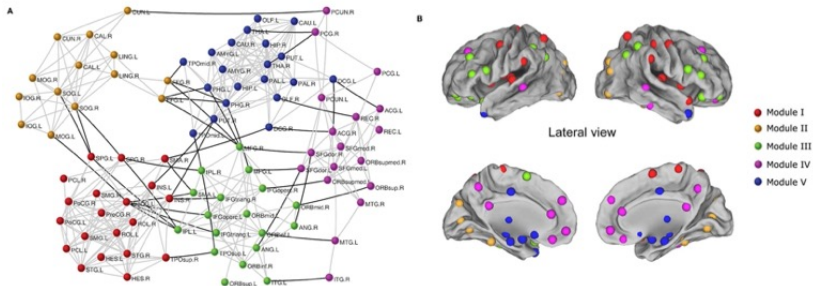


Figure: An example of graphs.

# Introduction

- Graph models show our knowledge of the association structure of a given set of variables, meaning they contain the information of (conditional) independence, while still leaving enough freedom for the exact form of the distribution.
- Undirected graphical models (UGM), to be more specific, are often used in problems in which there is little causal structure to guide the construction of a directed graph but with structure of conditional independence involving.



# Basic Knowledge of Graphs

- Typically, we denote an undirected graph as a tuple  $(V, E)$ .
- $V$  is the set of vertices of the graph, which we will use to represent the variables.
- $E$  is the set of edges of the graph, which we will use to represent the 'dependence among variables' (to be defined later).

Notation: for  $i, j \in V$  we use the symbol  $(i, j)$  to denote the edge from  $i$  to  $j$  whether it is in  $E$  or not.

# Basic Settings

- For a  $d$ -dimensional random vector  $X = (X_1, X_2, \dots, X_d)$ , we need to find a graph  $G = (V, E)$  to describe their probabilistic relationship.
- For each variable  $X_i$ , we use one vertex  $i \in V$  to represent it.
- We hope to use whether there are edges between vertices to reflect the conditional dependence between variables given the remainders.

To achieve this, we need to define firstly what is 'clique'.

# Cliques and Maximal Cliques

## Definition 1

*A clique  $C$  is a subset of vertices  $V$  that are all joined by edges, i.e. for all distinct  $j, k \in C$ , we have  $(i, j) \in E$*

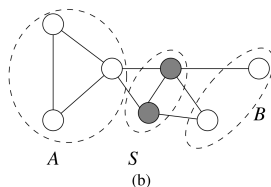
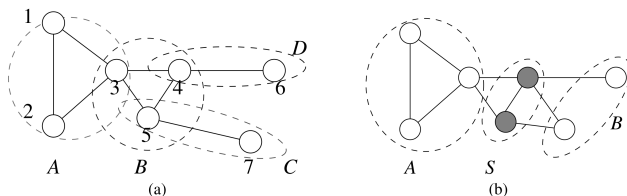
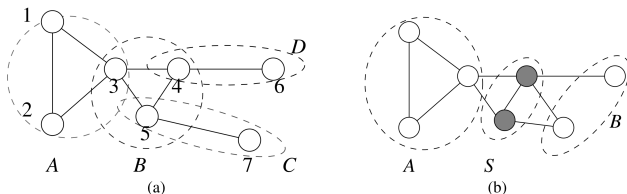


Figure: Illustration of Cliques (a).

# Cliques and Maximal Cliques

## Definition 2

*A maximal clique  $C$  is a clique that is not a subset of any other clique.*



**Figure:** Illustration of maximal Cliques (a).

Note that any function on a clique can be seen as a function on some maximal clique.

# Functions on Cliques

- We use  $\mathcal{C}$  to denote the set of all cliques  $C \subset V$ .
- By using notation  $\psi_C$ , we mean a function on clique  $C$ , i.e.

$$\psi_C : x_C \rightarrow \mathbb{R}.$$

where  $x_C := \{x_i\}_{i \in C}$ . With these preparations, we are ready to give a representation for probabilistic mass via an undirected graph.

# Factorization

## Definition 3

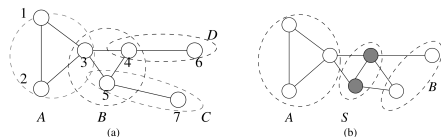
*We call random vector  $(X_1, \dots, X_d)$  factorizes according to graph  $G$  if we have*

$$p(x_1, x_2, \dots, x_d) \propto \prod_{C \in \mathcal{C}} \psi_C(x_C)$$

*where  $\psi_C$  is a non-negative function on clique  $C$ .*

- Note that we can equivalently define the factorization only according to maximal cliques instead of cliques, but sometimes using the 'clique' version will bring us some convenience, as we will see in the example of Gaussian model later.

# An Example



**Figure:** Illustration of factorization (b).

The graph to factorize the random vector does always exist (complete graph), while the density factorized by given graph (a) must have this form

$$p(x_1, \dots, x_7) \propto \psi_{123}(x_1, x_2, x_3) \psi_{345}(x_3, x_4, x_5) \psi_{46}(x_4, x_6) \psi_{57}(x_5, x_7).$$

# A Question

- Till now, we have denoted the distribution via an undirected graph, but why is it that way?
- Recall our purpose: we want to define a way to easily show the conditional dependence among variables, does this definition satisfy our request?

To illustrate this, let us first define the conditional independence 'required' by an undirected graph structures.



# Vertex Cut

For a vertices subset  $S$  of  $V$ , we define the graph of removing  $S$  as  $G(V \setminus S) = (V \setminus S, E(V \setminus S))$ , here the set of edges are defined as

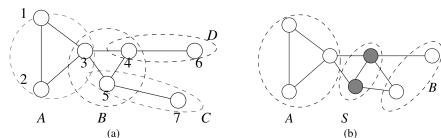
$$E(V \setminus S) = \{(j, k) \in E \mid j, k \in V \setminus S\}.$$

It is natural to expect that if the graph  $G(V \setminus S)$  consists of two distinct sub-graph, they are independent (conditioned on  $S$ ).

# Conditional Independence induced by Undirected Graph

## Definition 4

We say a random variable  $X = (X_1, X_2, \dots, X_d)$  is Markov with respect to the undirected graph  $G$ , if  $G(V \setminus S)$  is two distinct sub-graph  $A$  and  $B$ , we have  $X_A \perp X_B \mid X_S$ , i.e.  $X_A$  and  $X_B$  is independent conditioned on  $X_S$ .



**Figure:** Illustration of independence induced by graph.

# Markov Random Field

- A Markov random field (MRF), is the generalization of Markov chain for the random object defined on a field, which satisfies Markov properties with respect to a given graph of that field.
- The neighbourhood of node is defined via a graph, which encodes the conditional independence of the field of variables.

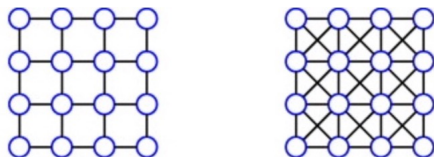


Figure: Illustration of MRF.

# Markov Chains as an example

- We now consider an example: Markov chain. The states are  $1, 2, \dots, k$  and we can easily draw a graph that contains these  $k$  nodes, with edges only existing between  $(i, i + 1)$
- By the graph, we might write the probability as

$$p(x_1, \dots, x_k) \propto p(x_1)p(x_2|x_1) \cdots p(x_k|x_{k-1})$$

which means it factorize according to graph

- Now if we cut the nodes from some  $i$  to  $j$ , the remaining two parts are independent conditioned on the variables  $i$  to  $j$ , which means this it is also Markov with the graph.

Is this just a coincidence that these two properties are satisfied simultaneously? The famous Hammersley–Clifford theorem[6] tells you 'NO'!

# Hammersley–Clifford equivalence

## Theorem 1

*For a given undirected graph  $G$  and random variables  $X = (X_1, X_2, \dots, X_d)$ , with strictly positive density  $p$ , the following two properties are equivalent:*

- (1)  $X$  factorize as the graph  $G$ .*
- (2)  $X$  is Markov with the graph  $G$ .*

This interesting theorem tells us two things:

- Graph models can represent the conditional dependence relationships among variables.
- The only way to define the probability induced by a graph, if request to keep independence, is the way we used above.

# Simple proof for necessity: A closer look at Markov chains

To give the idea of the proof for the necessity, i.e. (2) contains (1), we look closer at the simplest situation of 'Markov chain':

- Consider a 3-state Markov chain, with graph

$$G = (\{1, 2, 3\}, \{(1, 2), (2, 3)\}).$$

- According to the factorization, we can write

$$p(x_1, x_2, x_3) \propto \psi_{12}(x_1, x_2)\psi_{23}(x_2, x_3)$$

■

$$\begin{aligned} p(x_2) &\propto \sum_{x_1, x_3} \psi_{12}(x_1, x_2)\psi_{23}(x_2, x_3) \\ &= \sum_{x_1} \psi_{12}(x_1, x_2) \sum_{x_3} \psi_{23}(x_2, x_3) \end{aligned}$$

# Simple proof for necessity: A closer look at Markov chains



$$\begin{aligned} p(x_1, x_3 | x_2) &\propto \frac{\psi_{12}(x_1, x_2)}{\sum_{x_1} \psi_{12}(x_1, x_2)} \cdot \frac{\psi_{23}(x_2, x_3)}{\sum_{x_3} \psi_{23}(x_2, x_3)} \\ &= p(x_1 | x_2) p(x_3 | x_2) \end{aligned}$$

and the conditional independence is obtained.

# Gaussian Graph Model

We say a  $d$ -dimension random vector  $X$  obeys multivariate normal distribution  $\mathcal{N}(\mu, \Sigma)$ , if  $X$  has the density

$$f(x) = \frac{1}{\sqrt{(2\pi)^d \det(\Sigma)}} e^{-\frac{1}{2}(x-\mu)' \Sigma^{-1}(x-\mu)}$$

From now on, let  $\mu = 0$ .

We use the notation  $\sigma_{ij}$  to denote the  $(i, j)$  element of  $\Sigma$  and  $\sigma^{ij}$  for the  $(i, j)$  element of  $\Sigma^{-1}$ , then we can rewrite the density function as:

$$f(x) = \frac{1}{\sqrt{(2\pi)^d \det(\Sigma)}} e^{-\frac{1}{2} \sum_{i,j} \sigma^{ij} x_i x_j}.$$



# Conditional Independence for multivariate normal distribution

- By using the notations we define above, let us see what will happen if  $\sigma^{12} = 0$ .
- Under this condition, we have

$$f(x) \propto e^{-\frac{1}{2}x_1(\sum_{j \neq 1,2} 2\sigma^{1j}x_j + \sigma^{11}x_1)} e^{-\frac{1}{2}x_2(\sum_{i \neq 1,2} 2\sigma^{2i}x_i + \sigma^{22}x_2)} e^{-\frac{1}{2}\sum_{i,j \neq 1,2} \sigma^{ij}x_i x_j},$$

which exactly tells us that  $X_1 \perp X_2 | X_3, \dots, X_d$ .

- Similarly, we can prove that  $\sigma^{ij} = 0$  if and only if that  $X_i \perp X_j | X_{-i,-j}$ , here  $X_{-i,-j}$  means all the other components except  $i, j$ . Even better, similar conclusion holds for sub-matrix of  $\Sigma^{-1}$ .

# Gaussian Graph Model

With above preparations, we can now define the structure of Gaussian graph models:

- For  $X \sim \mathcal{N}(0, \Sigma)$ , we connect an undirected graph model  $G = (V, E)$  to it
- Vertex  $i$  in  $V$  represent the  $i$ -th component of  $X$ .
- There is an edge between  $i$  and  $j$  if and only if  $\sigma^{ij} \neq 0$ .

Then we can easily show that with these settings,  $X$  factorizes with the graph  $G$ , thus is also Markov with  $G$ .

Note that if we require the factor is defined as functions of maximal cliques, it will be much more inconvenient to represent this.

# Covariance Estimation with Gaussian Graph Model

- A natural question arises: when knowing the structure of Gaussian graph model and samples  $X_1, \dots, X_n$ , how can we estimate the covariance matrix now?
- According to our construction, knowing the structure of Gaussian graph model is equivalent to knowing which  $\sigma^{ij} = 0$  if  $(i, j) \notin E$ .
- At this time, the usual unbiased estimation

$$S = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})' (X_i - \bar{X})$$

where  $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ , might not be our ideal option since the  $(i, j)$ th element of  $\hat{\Theta} := S^{-1}$  may not be 0 if  $(i, j) \notin E$ .

# Covariance Estimation with Gaussian Graph Model

Let  $J = \{(i, j) \mid \sigma^{ij} = 0\}$ , the set of all the conditional independent pairs, and  $I = J^c$ , the remaining set of  $J$ , then one seemingly crazy way to construct our new estimator  $\tilde{S}$  is as follows:

- $W|_I = S|_I$ , i.e.  $W$  is the same as  $S$  at the position of  $I$ .
- The  $(i, j)$  element of the inverse of  $W$  is 0 for any  $(i, j) \in J$ .

Clearly, we might ask these two questions:

- Is there really existing such an estimator that satisfying the above conditions? If yes, will it be unique?
- Why we choose our estimator like this? What properties it might own? [1]

# Covariance Estimation with Gaussian Graph Model

## Theorem 2

*Given positive definite matrices  $L$  and  $M$  defined on the vertices  $V$  of a graph  $G = (V, E)$  there exists a unique positive definite matrix  $K$  such that*

- (1)  $K_{ij} = L_{ij}$  if  $(i, j) \in E$  or  $i=j$ ,*
- (2)  $(K^{-1})_{ij} = M_{ij}$  if  $(i, j) \notin E$  and  $i \neq j$ .*

If there exists one positive definite matrix that agrees with  $S$  at all the positions  $(i, j) \in I$ , there exists one and only one positive definite matrix satisfies additionally that the inverse of which is 0 in position  $J$ .

# Discussion on the above theorem

- Clearly, the above theorem can actually be generalized by not assuming zero on position  $J$ .
- An interesting way to view this theorem is given the marginal distribution of all maximal cliques, is there a joint distribution of all components owning these marginals and graph structure?
- The existence can be proven in a general exponential families setting, but instead, we will give the algorithm of how to find it to prove the existence.

# The Maximum of Likelihood Property

With almost the same observation, we can show the following maximum likelihood property:

## Theorem 3

*Among all the multivariate normal distributions whose inverse covariance matrix are 0 at  $J$ , our choice of  $W$  is the MLE of covariance matrix.*

By writing down the likelihood, we can easily find that the elements of  $S$  of position  $I$  is sufficient statistics, and this theorem just tells us how to use this sufficient statistics to construct MLE.

# The Maximum of Likelihood Property

## Proof.

By Lagrange multiplier, the log-likelihood of the data can be written as

$$l_C(\Theta) = \log \det \Theta - \text{trace}(S\Theta) - \sum_{(i,j) \in J} \gamma_{ij} \theta_{jk}.$$

The gradient equation is

$$\Theta^{-1} - S - \Gamma = 0 \quad (1)$$

$\Gamma$  is a matrix of Lagrange parameters with nonzero values for all pairs not in  $E$ . Therefore,  $W = S + \Gamma$ , satisfying that  $W|_I = S|_I$ .  $\square$



# The general cyclic algorithm

With all these good properties, a natural question might come:  
how can we really construct  $W$  from  $S$ ?

Fortunately, there is actually an algorithm that dealing with such kind of problem:

- Given two positive definite matrices  $P, Q$  of order  $d$ , and two disjoint set  $I, J$  such that  $I \cup J = \{(i, j) \mid 1 \leq i \leq j \leq d\}$
- How can we find a positive definite matrix  $F$  so that

$$F|_I = P|_I, \quad F^{-1}|_J = Q^{-1}|_J$$

The algorithm is called general cyclic algorithm [5].

# General idea for general cyclic algorithm

We now give the main idea for such a powerful algorithm:

- First we need to find some sets  $I_1, I_2, \dots, I_m$  (not necessarily distinct) so that their union is just  $I$ .
- Let  $F_0 = Q^{-1}$ , we try to construct a sequence of  $F_n$  so that

$$F_n|_{I_k} = P|_{I_k}, \quad F_n^{-1}|_{I_k^c} = F_{n-1}^{-1}|_{I_k^c}$$

where  $n = k(\text{mod } m)$ .

- While maintaining the second condition, we cyclic through  $I_m$  so that the first condition is also satisfied eventually.

When applying this to our problem, we either fix the inverse of matrix to be 0 at  $J$  and forcing the  $I$  position to become  $S$ , or firstly fixing the position of  $I$  to be the same as  $S$  and making the inverse matrix to be 0 at  $J$ .

# Ising Model

Consider a vector  $X = (X_1, \dots, X_d)$  of binary random variables, with each  $X_j \in \{0, 1\}$ . Given an undirected graph  $G = (V, E)$ , it posits a factorization of the form

$$p(x_1, \dots, x_d; \theta^*) = \frac{1}{Z(\theta^*)} \exp \left\{ \sum_{j \in V} \theta_j^* x_j + \sum_{(j,k) \in E} \theta_{jk}^* x_j x_k \right\}$$

where the parameter  $\theta_j^*$  is associated with vertex  $j \in V$ , and the parameter  $\theta_{jk}^*$  is associated with edge  $(j, k) \in E$  and  $x_j \in \{-1, 1\}$ .

# Ising Model

The quantity  $Z(\theta^*)$  is a constant that serves to enforce that the probability mass function  $p$  normalizes properly to one; more precisely, we have

$$Z(\theta^*) = \sum_{\mathbf{x} \in \{0,1\}^d} \exp \left\{ \sum_{j \in V} \theta_j^* x_j + \sum_{(j,k) \in E} \theta_{jk}^* x_j x_k \right\}.$$

In general, the parameters  $\theta^*$  would not possess nice properties as Gaussian cases, which requires gradient-based optimization, Monte Carlo sampling or variational algorithm for statistical inference.

# Table of Contents

## 1 Introduction and preliminaries

- Connecting an undirected graph to random variables
- Examples of Graph Model

## 2 Inferring of GGM

- Estimation of Parameters when Graph Structure is Known
- Estimation of the Graph Structure
- Generalization of Graphical Lasso

# Gaussian Graphical Model

Undirected Graphical Model (UGM) characterizes the relation between variables. For a graph  $(V, E)$ ,  $\{i, j\} \notin E$  means then variables **conditionally independent** (CI) given the other variables.

Let  $\vec{X} \sim N_p(\mu, \Sigma)$ , the **precision matrix** is defined as  $\Theta = \Sigma^{-1}$ . If the  $ij$ th component of  $\Theta$  is zero, then variables  $i$  and  $j$  are CI.

- If we have known the structure of UGM related to a Gaussian distribution, can we estimate  $\Sigma$  better?
- Now we focus on extracting the MLE of covariance ( $\mathbf{W}$ ), given the graph constraint  $(V, E)$ .
- Specifically, we can make  $\mathbf{W}_{ij} = \mathbf{S}_{ij}$  if  $\{i, j\} \in E$  or  $i = j$ ,  $\mathbf{W}_{ij}^{-1} = 0$  if  $\{i, j\} \notin E$ , where  $\mathbf{S}$  represent the empirical covariance matrix.

# Multiple Liner Regression

Partition  $\vec{X} = (\vec{Z}, Y)$  where  $\vec{Z} = (X_1, \dots, X_{p-1})$ ,  $Y = X_p$ . Partition  $\Sigma$  as

$$\Sigma = \begin{pmatrix} \Sigma_{ZZ} & \sigma_{ZY} \\ \sigma_{ZY}^T & \sigma_{YY} \end{pmatrix}$$

For simplicity, we denote  $\Sigma_{ZZ}$  as  $\Sigma_{11}$ ,  $\sigma_{ZY}$  as  $\sigma_{12}$ , similar for other elements in  $(X_1, \dots, X_p)$ .

Partition  $\Theta$  in the same way, since  $\Sigma\Theta = I$ , we obtain

$$\theta_{12} = -\theta_{22}\Sigma_{11}^{-1}\sigma_{12} = -\theta_{22}\beta,$$

where  $\beta := \Sigma_{11}^{-1}\sigma_{12}$  is exactly the regression coefficient of  $Y$  onto  $\vec{Z}$ . Hence  $\beta = -\theta_{12}/\theta_{22}$ .

# Multiple Liner Regression

Now we focus on iteratively updating the columns of  $\mathbf{W}$ . Let

$$\mathbf{W} = \begin{pmatrix} \mathbf{W}_{11} & \mathbf{w}_{12} \\ \mathbf{w}_{12}^T & w_{22} \end{pmatrix}.$$

Let  $\hat{\Theta}$  be the estimation of  $\Theta$ . Since  $\mathbf{W}\hat{\Theta} = \mathbf{I}$ ,

$$\mathbf{w}_{12} = -\mathbf{W}_{11}\hat{\theta}_{12}/\hat{\theta}_{22} := \mathbf{W}_{11}\hat{\beta}.$$

Focus on the upper right block of equation (1), we have

$$\mathbf{w}_{12} - \mathbf{s}_{12} - \gamma_{12} = \mathbf{0}.$$

By given  $\mathbf{W}_{11}$ , we can solve  $\hat{\beta}$  from  $\mathbf{W}_{11}\hat{\beta} - \mathbf{s}_{12} - \gamma_{12} = \mathbf{0}$ .



# Modified Regression

Remove the rows and columns that the corresponding predictors don't connect to the respond, we have

$$\mathbf{W}_{11}^* \boldsymbol{\beta}^* - \mathbf{s}_{12}^* = \mathbf{0}, \quad (2)$$

then  $\boldsymbol{\beta}^* = \mathbf{W}_{11}^{*-1} \mathbf{s}_{12}^*$ . Pad the zeros in corresponding locations to give  $\hat{\boldsymbol{\beta}}$ .

# A Modified Regression Algorithm

---

**Algorithm 17.1** *A Modified Regression Algorithm for Estimation of an Undirected Gaussian Graphical Model with Known Structure.*

---

1. Initialize  $\mathbf{W} = \mathbf{S}$ .
  2. Repeat for  $j = 1, 2, \dots, p, 1, \dots$  until convergence:
    - (a) Partition the matrix  $\mathbf{W}$  into part 1: all but the  $j$ th row and column, and part 2: the  $j$ th row and column.
    - (b) Solve  $\mathbf{W}_{11}^* \beta^* - s_{12}^* = 0$  for the unconstrained edge parameters  $\beta^*$ , using the reduced system of equations as in (17.19). Obtain  $\hat{\beta}$  by padding  $\hat{\beta}^*$  with zeros in the appropriate positions.
    - (c) Update  $w_{12} = \mathbf{W}_{11} \hat{\beta}$
  3. In the final cycle (for each  $j$ ) solve for  $\hat{\theta}_{12} = -\hat{\beta} \cdot \hat{\theta}_{22}$ , with  $1/\hat{\theta}_{22} = s_{22} - w_{12}^T \hat{\beta}$ .
-

# Why this algorithm converges to our target?

Algorithm 17.1 is a special case of the general cyclic algorithm [5].

- $\mathbf{W}^n$  is the matrix after  $n$  iterations and  $\mathbf{W}^0 = \mathbf{S}$ .
- Suppose  $\mathbf{W}^n$  is invertible and positive definite, by induction we know  $\mathbf{W}^{n+1}$  is invertible and positive definite.
- Let  $\Theta^n = (\mathbf{W}^n)^{-1}$ , then  $\{\Theta^n\}$  satisfies that:

$$\mathbf{W}_{ij}^n = \mathbf{W}_{ij}^{n-1} = \mathbf{S}_{ij}, \quad (i, j) \notin \tilde{E}_{n'}, \quad (3)$$

$$\Theta_{ij}^n = \Theta_{ij} = 0, \quad (i, j) \in \tilde{E}_{n'}. \quad (4)$$

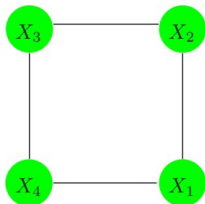
where  $\tilde{E}_k = \{(k, j) | j \in V, (k, j) \notin E\}$ ,  $k = 1, \dots, p$ , is the edge set without connecting to node  $k$ .  $n' = n \pmod{p}$ .

# Iteration Process

- Suppose (3),(4) hold after  $(n - 1)$ th iterations, we focus on the  $n$ -th iteration. Noting that  $\mathbf{W}_{11}$  would not change in  $n$ th iterations. We just need to focus on  $\mathbf{w}_{12}$ .
- Let  $A$  be the node set connecting target, and  $B$  be the node set without linking target, and

$$(\mathbf{w}_{12})_A = (\mathbf{W}_{AA}, \mathbf{W}_{AB}) \times \begin{pmatrix} \mathbf{W}_{AA}^{-1} \mathbf{s}_{12}^* \\ \mathbf{0} \end{pmatrix} = \mathbf{s}_{12}^*.$$

# A Simple Example



$$\mathbf{S} = \begin{pmatrix} 10 & 1 & 5 & 4 \\ 1 & 10 & 2 & 6 \\ 5 & 2 & 10 & 3 \\ 4 & 6 & 3 & 10 \end{pmatrix}$$

$$\hat{\Sigma} = \begin{pmatrix} 10.00 & 1.00 & \mathbf{1.31} & 4.00 \\ 1.00 & 10.00 & 2.00 & \mathbf{0.87} \\ \mathbf{1.31} & 2.00 & 10.00 & 3.00 \\ 4.00 & \mathbf{0.87} & 3.00 & 10.00 \end{pmatrix}, \quad \hat{\Sigma}^{-1} = \begin{pmatrix} 0.12 & -0.01 & \mathbf{0.00} & -0.05 \\ -0.01 & 0.11 & -0.02 & \mathbf{0.00} \\ \mathbf{0.00} & -0.02 & 0.11 & -0.03 \\ -0.05 & \mathbf{0.00} & -0.03 & 0.13 \end{pmatrix}$$

# Estimating Graphs from Data

- Meinshausen et al. (2006) apply lasso to estimate which components of  $\Theta$  are nonzero.
- Fit a lasso regression using each variable as the response and the others as predictors.
- Include an edge  $\{i, j\}$  in the graph if the estimated precision matrix of variable  $i$  on  $j$  is nonzero.

# Graphical Lasso

- Consider maximizing the penalized log-likelihood

$$\max_{\Theta > 0} \{ \log \det \Theta - \text{trace}(\mathbf{S}\Theta) - \lambda \|\Theta\|_1 \}, \quad (5)$$

- The gradient equation (stationary condition) is

$$\Theta^{-1} - \mathbf{S} - \lambda \mathbf{\Gamma} = \mathbf{0},$$

where  $\mathbf{\Gamma} = \text{Sign}(\Theta)$  with  $\text{Sign}(\theta_{jk}) = \text{sign}(\theta_{jk})$  if  $\theta_{jk} \neq 0$ , else  $\text{Sign}(\theta_{jk}) \in [-1, 1]$ , which is the sub-gradient of  $\|\Theta\|_1$ .

- It's easy to show that  $\mathbf{w}_{22} = \mathbf{s}_{22} + \lambda$ .

# Graphical Lasso

- Similar with the previous section

$$\mathbf{W}_{11}\boldsymbol{\beta} - \mathbf{s}_{12} + \lambda \text{Sign}(\boldsymbol{\beta}) = \mathbf{0} \quad (6)$$

- Note that

$$\frac{\partial}{\partial \boldsymbol{\beta}} \left( \frac{1}{2} \boldsymbol{\beta}' \mathbf{W}_{11} \boldsymbol{\beta} - \boldsymbol{\beta}' \mathbf{s}_{12} + \lambda \|\boldsymbol{\beta}\|_1 \right) = \mathbf{W}_{11} \boldsymbol{\beta} - \mathbf{s}_{12} + \lambda \text{Sign}(\boldsymbol{\beta}).$$

The system above is exactly equivalent to a lasso regression:

$$\min_{\boldsymbol{\beta}} \left\{ \frac{1}{2} \|\mathbf{V}^{1/2} \boldsymbol{\beta} - \mathbf{V}^{-1/2} \mathbf{s}_{12}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_1 \right\}, \quad (7)$$

where  $\mathbf{V} = \mathbf{W}_{11}$ .



# Graphical Lasso

Friedman et al. [2] use the pathwise coordinate descent method to solve the modified lasso problem at each stage.

- Focus on the  $j$ th row of (6) and fix  $\hat{\beta}_k, (k \neq j)$ .

$$\mathbf{v}_{jj}\hat{\beta}_j + \sum_{k \neq j} \mathbf{v}_{jk}\hat{\beta}_k - (\mathbf{s}_{12})_j + \lambda \text{Sign}(\hat{\beta}_j) = 0$$

$$\hat{\beta}_j = \lambda [((\mathbf{s}_{12})_j - \sum_{k \neq j} \mathbf{v}_{jk}\hat{\beta}_k) / \lambda - \text{Sign}(\hat{\beta}_j)] / \mathbf{v}_{jj}$$

- Since  $\text{Sign}(t) \in [-1, 1]$  if  $t=0$ , we obtain

$$\hat{\beta}_j = S\left((\mathbf{s}_{12})_j - \sum_{k \neq j} \mathbf{v}_{jk}\hat{\beta}_k, \lambda\right) / \mathbf{v}_{jj},$$

where  $S$  is the soft-threshold operator:

$$S(x, \lambda) = \text{sign}(x)(|x| - \lambda)_+.$$

# Graphical Lasso

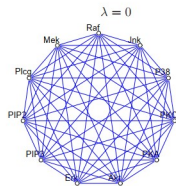
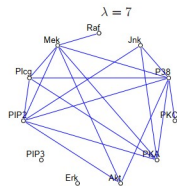
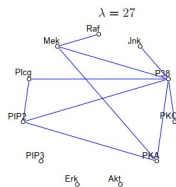
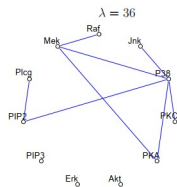
---

**Algorithm 17.2** *Graphical Lasso.*

---

1. Initialize  $\mathbf{W} = \mathbf{S} + \lambda \mathbf{I}$ . The diagonal of  $\mathbf{W}$  remains unchanged in what follows.
  2. Repeat for  $j = 1, 2, \dots, p, 1, 2, \dots, p, \dots$  until convergence:
    - (a) Partition the matrix  $\mathbf{W}$  into part 1: all but the  $j$ th row and column, and part 2: the  $j$ th row and column.
    - (b) Solve the estimating equations  $\mathbf{W}_{11}\beta - s_{12} + \lambda \cdot \text{Sign}(\beta) = 0$  using the cyclical coordinate-descent algorithm (17.26) for the modified lasso.
    - (c) Update  $w_{12} = \mathbf{W}_{11}\hat{\beta}$
  3. In the final cycle (for each  $j$ ) solve for  $\hat{\theta}_{12} = -\hat{\beta} \cdot \hat{\theta}_{22}$ , with  $1/\hat{\theta}_{22} = w_{22} - w_{12}^T \hat{\beta}$ .
-

# Model Selection



# Model Selection

The graph becomes more sparse as the penalty parameter is increased.

We can use BIC for model selection. For a given penalty parameter  $\lambda$ ,

$$BIC(\lambda) = -\log \det(\hat{\Theta}(\lambda)) + \text{trace}(\hat{\Theta} \mathbf{S}) + \frac{\log n}{n} \sum_{i \leq j} \hat{e}_{ij}(\lambda),$$

where  $\hat{e}_{ij} = 0$  if  $\hat{\theta}_{ij} = 0$ , and  $\hat{e}_{ij} = 1$  otherwise [7].

# The dual problem of Graphical Lasso

## Lemma 1

*The primal problem (5) and its stationarity conditions are equivalent to the stationary conditions for the box-constrained SDP (semi-definite programming) [4].*

$$\max_{\|\tilde{\mathbf{r}}\|_{\infty} \leq \lambda} \log \det(\mathbf{S} + \tilde{\mathbf{r}}) + p, \quad (8)$$

*under the transformation  $\mathbf{S} + \tilde{\mathbf{r}} = \mathbf{\Theta}^{-1}$ .*

- Rewrite gradient conditions of (5):  $-(\mathbf{S} + \lambda \mathbf{\Gamma})^{-1} + \mathbf{\Theta} = \mathbf{0}$ .  
Let  $\tilde{\mathbf{r}} = \lambda \mathbf{\Gamma}$  then  $\|\tilde{\mathbf{r}}\|_{\infty} \leq \lambda$ . Let  $\mathbf{P} = \text{abs}(\mathbf{\Theta})$ .

- Such  $\tilde{\Gamma}$ ,  $\mathbf{P}$  satisfy the following equations:

$$\begin{aligned} -(\mathbf{S} + \tilde{\Gamma})^{-1} + \mathbf{P} * \text{sign}(\tilde{\Gamma}) &= \mathbf{0} \\ \mathbf{P} * (\text{abs}(\tilde{\Gamma}) - \lambda \mathbf{1}_p \mathbf{1}_{p'}) &= \mathbf{0} \\ \|\tilde{\Gamma}\|_{\infty} &\leq \lambda, \end{aligned} \tag{9}$$

where  $\mathbf{P}$  is a symmetric  $p \times p$  matrix with non-negative entries and the operator '\*' denotes element-wise product.

In fact, (9) are the KKT conditions for the box-constrained SDP (8).

- Take  $\Theta = \mathbf{P} * \text{sign}(\tilde{\Gamma})$ ,  $\Gamma = \tilde{\Gamma}/\lambda$ , conditions (9) imply  $-(\mathbf{S} + \lambda \Gamma)^{-1} + \Theta = \mathbf{0}$ .

## Lemma 2

Assume  $\mathbf{W}_{11} > 0$ . The stationarity equations

$$-\mathbf{W}_{11}\hat{\beta} + \mathbf{s}_{12} + \lambda\hat{\gamma}_{12} = \mathbf{0}, \quad (10)$$

correspond to the solution of

$$\min_{\beta \in \mathbb{R}^{p-1}} \frac{1}{2} \beta' \mathbf{W}_{11} \beta - \beta' \mathbf{s}_{12} + \lambda \|\beta\|_1. \quad (11)$$

Solving it is equivalent to solving the following box-constrain:

$$\min_{\tilde{\gamma} \in \mathbb{R}^{p-1}, \|\tilde{\gamma}\|_\infty \leq \lambda} \frac{1}{2} (\mathbf{s}_{12} + \tilde{\gamma})' \mathbf{W}_{11}^{-1} (\mathbf{s}_{12} + \tilde{\gamma}), \quad (12)$$

where  $\hat{\beta}$  and  $\tilde{\gamma}_{12}$  are related by

$$\hat{\beta} = \mathbf{W}_{11}^{-1} (\mathbf{s}_{12} + \tilde{\gamma}_{12}). \quad (13)$$

# Proof of Lemma 2

- It's easy to find (10) is the KKT condition for problem (11) since  $\beta = -\theta_{12}/\theta_{22}$ ,  $Sign(\beta) = -Sign(\theta_{12})$ .
- Write (10) as

$$\hat{\beta} - \mathbf{W}_{11}^{-1}(\mathbf{s}_{12} + \lambda \hat{\gamma}_{12}) = \mathbf{0} \quad (14)$$

Suppose  $\hat{\beta}, \hat{\gamma}_{12}$  satisfy (14), then  $\tilde{\gamma}_{12} := \lambda \hat{\gamma}_{12}$ ,  $\tilde{\mathbf{p}}_{12} := \text{abs}(\hat{\beta})$  satisfy

$$\begin{aligned} \mathbf{W}_{11}^{-1}(\mathbf{s}_{12} + \tilde{\gamma}_{12}) + \tilde{\mathbf{p}}_{12} * \text{sign}(\tilde{\gamma}_{12}) &= \mathbf{0} \\ \tilde{\mathbf{p}}_{12} * (\text{abs}(\tilde{\gamma}_{12}) - \lambda \mathbf{1}_{p-1}) &= \mathbf{0} \\ \|\tilde{\gamma}_{12}\|_{\infty} &\leq \lambda, \end{aligned} \quad (15)$$

which are KKT conditions of (12).



- If  $\tilde{\gamma}_{12}, \tilde{\mathbf{p}}_{12}$  satisfy (15), then

$$\hat{\gamma}_{12} = \tilde{\gamma}_{12}/\lambda, \quad \hat{\beta} = -\tilde{\mathbf{p}}_{12} * \text{sign}(\hat{\gamma}_{12})$$

satisfy (14). □

Moreover, consider solving (8) for the block  $\tilde{\gamma}_{12}$  holding the rest of  $\tilde{\Gamma}$  fixed. From (9), we have

$$\begin{aligned} -\theta_{12} + \mathbf{p}_{12} * \text{sign}(\tilde{\gamma}_{12}) &= \mathbf{0} \\ \mathbf{p}_{12} * (\text{abs}(\tilde{\gamma}_{12}) - \lambda \mathbf{1}_{p-1}) &= \mathbf{0} \\ \|\tilde{\gamma}_{12}\|_{\infty} &\leq \lambda. \end{aligned} \tag{16}$$

Since  $\hat{\beta} = -\theta_{12}/\theta_{22}, \theta_{22} > 0$ , let  $\tilde{\mathbf{p}}_{12} = \mathbf{p}_{12}/\theta_{22}$ , (16) are equivalent with conditions (15). By lemma 2, GLASSO solves (15), hence (16). So GLASSO performs dual block updates for the dual of the graphical problem.

## DP-GLASSO

- Instead of updating  $\mathbf{W}$ , now we update  $\Theta$  directly. Similar as before, we obtain

$$\Theta_{11}^{-1} \theta_{12} w_{22} + \mathbf{s}_{12} + \lambda \gamma_{12} = \mathbf{0}, \quad (17)$$

where  $w_{22} = s_{22} + \lambda$  is fixed.

- With  $\alpha = \theta_{12} w_{22}$ ,  $\Theta_{11} > 0$ , By Lemma 2, it's equivalent to solve:

$$\min_{\tilde{\gamma} \in \mathbb{R}^{p-1}, \|\tilde{\gamma}\|_{\infty} \leq \lambda} \frac{1}{2} (\mathbf{s}_{12} + \tilde{\gamma})' \Theta_{11} (\mathbf{s}_{12} + \tilde{\gamma}) \quad (18)$$

The optimal solutions of (17) and (18) are related by

$$\hat{\alpha} = -\Theta_{11} (\mathbf{s}_{12} + \tilde{\gamma}).$$

## DP-GLASSO

- We update  $\hat{\theta}_{12}$  as  $\hat{\alpha}/w_{22}$ .
- $\theta_{22}$  is updated by

$$\hat{\theta}_{22} = \frac{1 - (\mathbf{s}_{12} + \tilde{\gamma})' \hat{\theta}_{12}}{w_{22}}. \quad (19)$$

# DP-GLASSO

---

**Algorithm 3** DP-GLASSO algorithm

---

1. Initialize  $\Theta = \text{diag}(\mathbf{S} + \lambda \mathbf{I})^{-1}$ .
2. Cycle around the columns repeatedly, performing the following steps till convergence:
  - (a) Rearrange the rows/columns so that the target column is last (implicitly).
  - (b) Solve (18) for  $\tilde{\gamma}$  and update

$$\hat{\theta}_{12} = -\Theta_{11}(\mathbf{s}_{12} + \tilde{\gamma})/w_{22}$$

- (c) Solve for  $\theta_{22}$  using (19).
    - (d) Update the working covariance  $\mathbf{w}_{12} = \mathbf{s}_{12} + \tilde{\gamma}$ .
-

# Nonparanormal Case

- Han Liu et al.(2012) propose *the nonparanormal SKEPTIC* to estimate high-dimensional UGM efficiently and robustly [3].
- Their simulation suggests that the nonparanormal SKEPTIC can be used as a safe replaced for Gaussian based estimator, even when the data are truly Gaussian.
- Idea: Use rank-based correlation coefficient estimators (Spearman's rho, Kendall's tau) which are invariant under monotone transformations.

# Nonparanormal

Let  $\vec{f} = \{f_1, \dots, f_d\}$  be a set of monotone univariate functions and let  $\Sigma^0 \in \mathbb{R}^{p \times p}$  be a positive-definite correlation matrix with  $\text{diag}(\Sigma^0) = \mathbf{1}$ .

## Definition 5 (nonparanormal)

We say a  $p$ -dimensional random variable  $\vec{X} = (X_1, \dots, X_p)^T$  has a **nonparanormal** distribution  $\vec{X} \sim \text{NPN}_p(\vec{f}, \Sigma^0)$  if

$$\vec{f}(\vec{X}) := (f_1(X_1), \dots, f_p(X_p))^T \sim N_p(\mathbf{0}, \Sigma^0).$$

- Monotone transformation won't change conditional independence between variables, the graph structure is invariant under such transformations.

# Estimating marginal transformations

Assume  $x^1, \dots, x^n$  are realizations of  $\vec{X} \sim NPN_p(\vec{f}, \Sigma^0)$ , here  $x_j^i = (x_1^i, \dots, x_p^i)$ . Estimating marginal transformations is useful for calculating the likelihood of a nonparnormal fit.

Let  $\tilde{F}_j(t) = \frac{1}{n} \sum_{i=1}^n I(x_j^i \leq t)$  be the emperical distribution of  $X_j$ . Estimate the marginal transformation  $f_j$  by

$$\tilde{f}_j(x) := \Phi^{-1}(\mathcal{T}_{\delta_n}[\tilde{F}_j(x)]),$$

where  $\Phi$  is the distribution function of  $N(0, 1)$  and the truncation operator

$$\mathcal{T}_{\delta_n} := \delta_n \cdot I(x < \delta_n) + x \cdot I(\delta_n \leq x \leq 1 - \delta_n) + (1 - \delta_n) \cdot I(x > 1 - \delta_n),$$

is used to alleviate the instability of  $\Phi^{-1}(t)$  for too large or small  $t$ .

# Normal-score rank correlation

- The normal-score rank correlation coefficient of  $\tilde{f}_j(x)$  could be used directly in glasso. However, it appears that their efficiency cannot be generalized to the high-dimensional setting.
- The reason is that the standard Gaussian quantile function diverges very quickly when it is evaluated at a point close to 1.



# Rank-based Correlation Coefficient Estimators

Let  $r_j^i$  be the rank of  $x_j^i$  among  $x_j^1, \dots, x_j^n$ ,  $\bar{r}_j = \frac{1}{n} \sum_{i=1}^n r_j^i = \frac{n+1}{2}$ .

## ■ Spearman's rho

$$\hat{\rho}_{jk} = \frac{\sum_{i=1}^n (r_j^i - \bar{r}_j)(r_k^i - \bar{r}_k)}{\sqrt{\sum_{i=1}^n (r_j^i - \bar{r}_j)^2 \cdot \sum_{i=1}^n (r_k^i - \bar{r}_k)^2}}$$

## ■ Kendall's tau

$$\hat{\tau}_{jk} = \frac{2}{n(n-1)} \sum_{1 \leq i < i' \leq n} \text{sign}((x_j^i - x_j^{i'})(x_k^i - x_k^{i'})).$$

## ■ The two statistics are the same for $\vec{X}$ and $\vec{f}(\vec{X})$ .

# SKEPTIC

## Lemma 3 (Kendall (1948), Kruskal (1958))

Assume  $\vec{X} \sim \text{NPN}_p(\vec{f}, \Sigma^0)$ , we have  $\Sigma_{jk}^0 = 2 \sin(\frac{\pi}{6} \rho_{jk}) = \sin(\frac{\pi}{2} \tau_{jk})$ .

Motivated by this lemma, define  $\hat{\mathbf{S}}^\rho = [\hat{S}_{jk}^\rho]$  for the unknown correlation matrix  $\Sigma^0$  as follows:

$$\hat{S}_{jk}^\rho = \begin{cases} 2 \sin(\frac{\pi}{6} \hat{\rho}_{jk}), & j \neq k, \\ 1, & j = k. \end{cases}$$

Similarly we can define  $\hat{\mathbf{S}}^\tau = [\hat{S}_{jk}^\tau]$ .

- Estimate  $\hat{\mathbf{S}}$  by SKEPTIC instead of calculating Pearson correlation.

# SKEPTIC with Graphical Lasso

Plug the estimated correlation matrix  $\hat{\mathbf{S}}$  into the graphical lasso:

$$\hat{\Theta}^{g\text{lasso}} = \arg \min_{\Theta > 0} \{ \text{trace}(\hat{\mathbf{S}}\Theta) - \log \det \Theta + \lambda \|\Theta\|_1 \}$$

$\hat{\mathbf{S}}$  may not be positive semi-definite. Even though the problem is still convex, certain algorithms (like the blockwise coordinate descent algorithm) will fail. However, other algorithms that don't have positive semi-definite assumption on  $\hat{\mathbf{S}}$  (like the two-metric projected Newton method) can still work.

# References I

- [1] A. P. Dempster. “Covariance Selection”. In: *Biometrics* 28.1 (1972), pp. 157–175. ISSN: 0006341X, 15410420. URL: <http://www.jstor.org/stable/2528966>.
- [2] Jerome H Friedman. *The elements of statistical learning: Data mining, inference, and prediction*. [springer open](#), 2017.
- [3] Han Liu et al. “High-dimensional semiparametric Gaussian copula graphical models”. In: *The Annals of Statistics* 40.4 (2012), pp. 2293–2326.
- [4] Rahul Mazumder and Trevor Hastie. “The graphical lasso: New insights and alternatives”. In: *Electronic journal of statistics* 6 (2012), p. 2125.

# References II

- [5] T. P. Speed and H. T. Kiiveri. “Gaussian Markov Distributions over Finite Graphs”. In: *The Annals of Statistics* 14.1 (1986), pp. 138–150. ISSN: 00905364. URL: <http://www.jstor.org/stable/2241271>.
- [6] M.J. Wainwright. *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 2019. ISBN: 9781108498029. URL: <https://books.google.nl/books?id=IluHDwAAQBAJ>.
- [7] Ming Yuan and Yi Lin. “Model selection and estimation in the Gaussian graphical model”. In: *Biometrika* 94.1 (2007), pp. 19–35.