

Graphical Modeling

Christophe Giraud

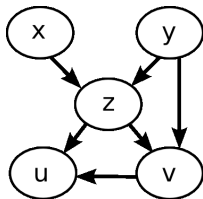
Université Paris-Sud and Ecole Polytechnique
Maths Department

Lecture Notes on High-Dimensional Statistics :

<http://www.cmap.polytechnique.fr/~giraud/MSV/LectureNotes.pdf>

Please ask questions!

Graphical models :



- Graphical modeling is a convenient representation of conditional dependences among random variables
- It is a powerful tool for
 - exploring “direct effect” between variables
 - fast computations in complex models
- It is popular in many different fields, including bioinformatics, computer vision, speech recognition, environmental statistics, economics, social sciences, etc



Two important topics :

- learning graphical models
- learning with graphical models



Two important topics :

- learning graphical models
- ~~learning with graphical models~~

Example 1

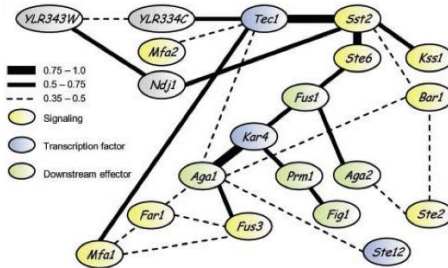


Figure: Learning Biological Regulatory Networks

Seminal reference : Friedman [6]

Example 1

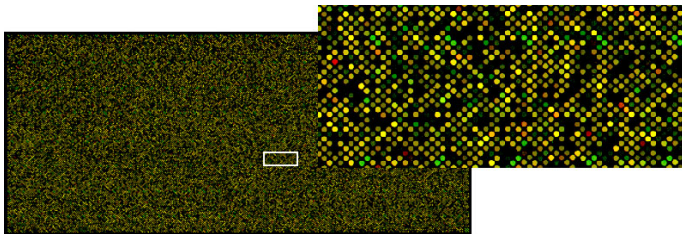
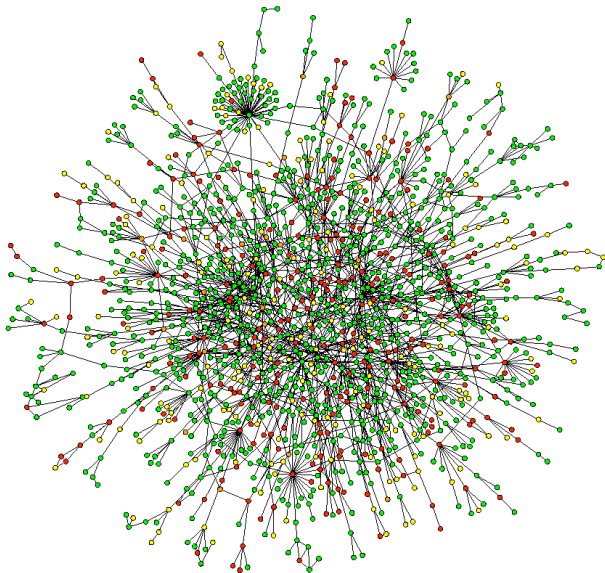


Figure: Learning Gene-Gene Regulatory Networks from microarrays

Seminal reference : Friedman [6]

Example 1



Example 2

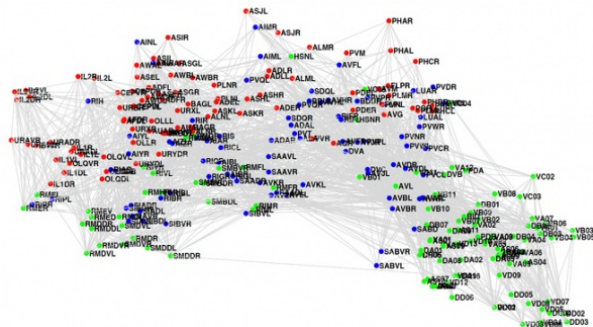


Figure: Learning Brain Connectivity Networks

Example 3

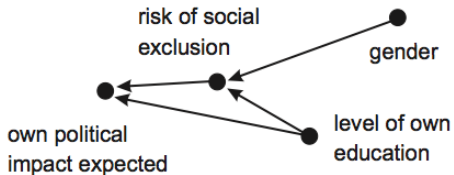


Figure: Learning “direct effects” in Social Sciences

Conditional independence

The concept of conditional dependence is more suited than the concept of dependence in order to catch "direct" dependences between variables

Traffic jams and snowmen are correlated.

But conditionally on snow falls, the size of the traffic jams and the number of snowmen are independent.



Figure: Difference between dependence and conditional dependence

Conditional independence

random variables X and Y are independent conditionally on a variable Z (we write $X \perp\!\!\!\perp Y \mid Z$) if

$$\text{law}((X, Y) \mid Z) = \text{law}(X \mid Z) \otimes \text{law}(Y \mid Z).$$

Characterisation

When the distribution of (X, Y, Z) has a positive density f , then

$$\begin{aligned} X \perp\!\!\!\perp Y \mid Z &\iff f(x, y \mid z) = f(x \mid z)f(y \mid z) \\ &\iff f(x, y, z) = f(x, z)f(y, z)/f(z) \\ &\iff f(x \mid y, z) = f(x \mid z), \end{aligned}$$

Directed acyclic model

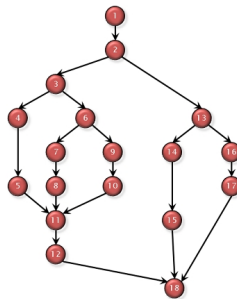
Terminology

Directed graph \vec{g}

set of nodes and arrows

Acyclic

no sequence of arrows forms a loop
in the graph



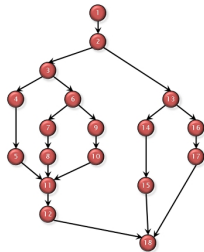
Parents

the parents of a is the set $pa(a)$ of nodes b such that $b \rightarrow a$

Descendent

the descendent of a is the set $de(a)$ of nodes that can be reached from a by following some sequence of arrows.

$X_a \perp\!\!\!\perp \{X_b : a \nrightarrow \dots \nrightarrow b\}$
 conditionally on $\{X_c : c \rightarrow a\}$



Directed acyclic graphical model

The law of the random variable $X = (X_1, \dots, X_p)$ is a graphical model according to the directed acyclic graph \vec{g} if

for all a , $X_a \perp\!\!\!\perp \{X_b, b \notin \text{de}(a)\} \mid \{X_c, c \in \text{pa}(a)\}$

We write $\mathcal{L}(X) \sim \vec{g}$.

Remark: if $\vec{g} \subset \vec{g}'$ and $\mathcal{L}(X) \sim \vec{g}$ then $\mathcal{L}(X) \sim \vec{g}'$.

Warning

There is no unique minimal graph in general!



Be careful with the interpretation of directed graphical models!

Example:

$$X_{i+1} = \alpha X_i + \varepsilon_i \quad \text{with} \quad \varepsilon_i \text{ independent of } X_1, \dots, X_{i-1}.$$

Then, the two graphs

$$1 \rightarrow 2 \rightarrow \dots \rightarrow p \quad \text{and} \quad 1 \leftarrow 2 \leftarrow \dots \leftarrow p$$

are minimal graphs for this model.

The issue of estimating the minimal \vec{g} is ill-posed in this context.

Yet,

- 1 it is very useful for defining / computing laws (next slides)
- 2 it can be used for exploring “causal effect” (last part of the talk)

Here, we assume that \vec{g} is known (from expert knowledge).

Factorization formula

If $\mathcal{L}(X) \sim \vec{g}$, we have

$$f(x_1, \dots, x_p) = \prod_{b=1}^p f(x_b | x_{\text{pa}(b)})$$

Proof: for a leaf p

$$\begin{aligned} f(x_1, \dots, x_p) &= f(x_p | x_1, \dots, x_{p-1}) f(x_1, \dots, x_{p-1}) \\ &= f(x_p | x_{\text{pa}(p)}) f(x_1, \dots, x_{p-1}) \end{aligned}$$

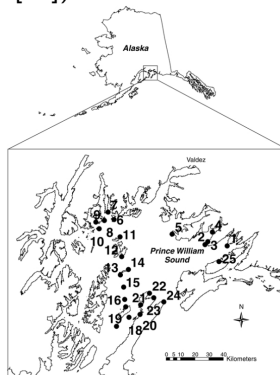
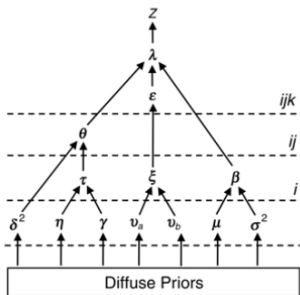
□

- Very useful for defining / computing f !
- Sampling with Gibbs sampler

Examples of applications

- speech recognition
- computer vision
- ecological monitoring
- decision making
- diagnosis
- environmental statistics
- etc

Seals monitoring (Ver Hoef and Frost [17])



Ice streams (Berliner *et al.*)

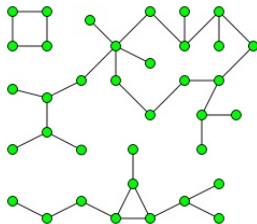
http://www.stat.osu.edu/~sses/collab_ice.html

Non-directed model

Non-directed graph

set of nodes and edges

The nodes are labelled by $1, \dots, p$.



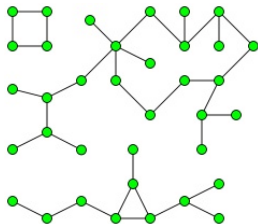
Neighbors

The neighbors of a are the nodes in $ne(a) = \{b : b \stackrel{g}{\sim} a\}$

Class of a

We set $cl(a) = ne(a) \cup \{a\}$

X_a independent from $\{X_b : b \approx a\}$
conditionally on $\{X_c : c \sim a\}$



Non-directed graphical model

The law of the random variable $X = (X_1, \dots, X_p)$ is a graphical model according to the non-directed graph g if

$$\text{for all } a: X_a \perp\!\!\!\perp \{X_b, b \notin \text{cl}(a)\} \mid \{X_c, c \in \text{ne}(a)\}.$$

We write $\mathcal{L}(X) \sim g$.

Remark: if $g \subset g'$ and $\mathcal{L}(X) \sim g$ then $\mathcal{L}(X) \sim g'$.

Minimal graph

When X has a positive density there exists a unique minimal graph g_* such that $\mathcal{L}(X) \sim g_*$.

In the following, we call simply “graph of X ” the minimal graph g_* such that $\mathcal{L}(X) \sim g_*$.

Our main goal will be to learn g_* from data.

Connection with directed acyclic graphs

Moral graph

The moral graph g^m associated to a directed acyclic graph \vec{g} is obtained by

- setting an edge between each parents of each nodes
- replacing arrows by edges

Proposition

$$\mathcal{L}(X) \sim_{\vec{g}} \implies \mathcal{L}(X) \sim g^m$$

Proof :

- 1 From the factorization formula, $\exists g_1, g_2$ such that

$$f(x) = g_1(x_a, x_{ne^m(a)}) g_2(x_{nn^m(a)}, x_{ne^m(a)})$$

where $nn^m(a) = \{1, \dots, p\} \setminus cl^m(a)$.

- 2 This ensures $X_a \perp\!\!\!\perp X_{nn^m(a)} \text{ given } X_{ne^m(a)}$.



Hammersley-Clifford formula

For a random variable X with positive density f

$$\mathcal{L}(X) \sim g \iff f(x) \propto \exp \left(\sum_{c: c \in \text{cliques}(g)} g_c(x_c) \right).$$

Proof : based on Möbius inversion formula.

Questions so far?

Gaussian graphical models

In the remaining $X \sim \mathcal{N}(0, \Sigma)$, with Σ non singular.

Reminder on Gaussian distribution (1/2)

Proposition 1 : Gaussian conditioning

We consider two sets $A = \{1, \dots, k\}$ and $B = \{1, \dots, p\} \setminus A$, and a Gaussian random vector $X = \begin{bmatrix} X_A \\ X_B \end{bmatrix} \in \mathbb{R}^p$ with $\mathcal{N}(0, \Sigma)$ distribution.

We write $K = \begin{bmatrix} K_{AA} & K_{AB} \\ K_{BA} & K_{BB} \end{bmatrix}$ for Σ^{-1} and K_{AA}^{-1} for the inverse $(K_{AA})^{-1}$ of K_{AA} .

Then, we have

$$\text{Law}(X_A | X_B) = \mathcal{N}(-K_{AA}^{-1}K_{AB}X_B, K_{AA}^{-1})$$

which means

$$X_A = -K_{AA}^{-1}K_{AB}X_B + \varepsilon_A,$$

with $\varepsilon_A \sim \mathcal{N}(0, K_{AA}^{-1})$ independent of X_B .

Proof: We have for some function f that we do not need to make explicit

$$g(x_A|x_B) = g(x_A, x_B)/g(x_B) = \exp\left(-\frac{1}{2}x_A^T K_{AA}x_A - x_A^T K_{AB}x_B\right) f(x_B).$$

As a consequence,

$$g(x_A|x_B) \propto \exp\left(-\frac{1}{2}(x_A + K_{AA}^{-1}K_{AB}x_B)^T K_{AA}(x_A + K_{AA}^{-1}K_{AB}x_B)\right)$$

where the factor of proportionality does not depend on x_A .

We recognize the density of the Gaussian $\mathcal{N}(-K_{AA}^{-1}K_{AB}x_B, K_{AA}^{-1})$ law. \square

Ref: Lauritzen [11]

Partial correlation

For any $a, b \in \{1, \dots, p\}$, we have

$$\text{cor}(X_a, X_b | X_c : c \neq a, b) = \frac{-K_{a,b}}{\sqrt{K_{aa} K_{bb}}}.$$

Proof : The previous proposition with $A = \{a, b\}$ and $B = A^c$ gives

$$\text{cov}(X_A | X_B) = \begin{pmatrix} K_{aa} & K_{ab} \\ K_{ab} & K_{bb} \end{pmatrix}^{-1} = \frac{1}{K_{aa}K_{bb} - K_{ab}^2} \begin{pmatrix} K_{bb} & -K_{ab} \\ -K_{ab} & K_{aa} \end{pmatrix}.$$

Plugging this formula in the definition of the partial correlation gives the result. \square

Reading the graph g on K

From K to g

We set $K = \Sigma^{-1}$ and define the graph g by

$$a \overset{g}{\sim} b \iff K_{a,b} \neq 0. \quad (1)$$

GGM and precision matrix

For the graph g defined by (1), we have

- 1 $\mathcal{L}(X) \sim g$ and g is minimal.
- 2 There exists $\varepsilon_a \sim \mathcal{N}(0, K_{aa}^{-1})$ independent of $\{X_b : b \neq a\}$ such that

$$X_a = - \sum_{b \in \text{ne}(a)} \frac{K_{ab}}{K_{aa}} X_b + \varepsilon_a.$$

Proof. We apply Proposition 1 :

- 1 We set $A = \{a\} \cup \text{nn}(a)$ and $B = \text{ne}(a)$, where $\text{nn}(a) = \{1, \dots, p\} \setminus \text{cl}(a)$. The precision matrix restricted to A is

$$K_{AA} = \begin{pmatrix} K_{aa} & 0 \\ 0 & K_{\text{nn}(a)\text{nn}(a)} \end{pmatrix} \text{ so its inverse is}$$

$$(K_{AA})^{-1} = \begin{pmatrix} K_{aa}^{-1} & 0 \\ 0 & (K_{\text{nn}(a)\text{nn}(a)})^{-1} \end{pmatrix}.$$

The above Lemma ensures that the law of $X_{\{a\} \cup \text{nn}(a)}$ given $X_{\text{ne}(a)}$ is Gaussian with covariance matrix $(K_{AA})^{-1}$ so X_a and $X_{\text{nn}(a)}$ are independent conditionally on $X_{\text{ne}(a)}$.

- 2 The second point is obtained with $A = \{a\}$ and $B = A^c$.

Goal

From a n -sample X_1, \dots, X_n i.i.d. with distribution $\mathcal{N}(0, \Sigma)$, we want to estimate the (minimal) graph g such that $\mathcal{L}(X) \sim g$.

The above results suggest 3 estimations strategies:

- 1 by estimating the partial correlations + multiple testing
- 2 by a sparse estimation of K
- 3 by a regression approach

Estimation with partial correlation (1/3)

Reminder 1

$$a \stackrel{g}{\sim} b \iff \rho_{a,b} := \text{cor}(X_a, X_b | X_c : c \neq a, b) \neq 0$$

Reminder 2

$$\rho_{a,b} = \frac{-K_{a,b}}{\sqrt{K_{aa} K_{bb}}}$$

Partial covariance estimation

For $n > p$, we estimate $\rho_{a,b}$ by

$$\hat{\rho}_{ab} = \frac{-[\hat{\Sigma}^{-1}]_{ab}}{\sqrt{[\hat{\Sigma}^{-1}]_{aa} [\hat{\Sigma}^{-1}]_{bb}}},$$

where $\hat{\Sigma}$ is the empirical covariance.

Estimation with partial correlation (2/3)

Under the null hypothesis

when $\rho_{a,b} = 0$ and $n > p - 2$, we have

$$\hat{t}_{a,b} := \sqrt{n - 2 - p} \times \frac{\hat{\rho}_{ab}}{\sqrt{1 - \hat{\rho}_{ab}^2}} \sim \text{Student}(n - p - 2).$$

Estimation procedure

- 1 Compute the $\hat{t}_{a,b}$
- 2 Apply a multiple testing thresholding

Weakness

- when $p > n - 2$ the procedure cannot be applied
- when $n > p$ but $n - p$ small, $\hat{t}_{a,b}$ has a large variance and the procedure is powerless

Solution 1: Shrinking the conditioning



work with $\widehat{\text{cor}}(X_a, X_b | X_c : c \in S)$ with S small

Ref: Wille and Bühlmann [19], Castelo and Roverato [2], Spirtes et al. [16] or Kalisch and Bühlmann [8].



$\widehat{\text{cor}}(X_a, X_b | X_c : c \in S)$ is stable when S is small



it is unclear what we estimate at the end (in general)

Solution 2 : Sparse estimation of K

The instability for large p comes from the instability of $\widehat{\Sigma}^{-1}$ for estimating K .



Build a more stable estimator of K capitalizing on its sparsity.

Sparse estimation of K (1/2)

The likelihood of a $p \times p$ positive symmetric matrix $K \in \mathcal{S}_p^+$ is

$$\text{Likelihood}(K) = \prod_{i=1}^n \sqrt{\frac{\det(K)}{(2\pi)^p}} \exp\left(-\frac{1}{2} X_i^T K X_i\right).$$

Negative log-likelihood

The negative-log-likelihood

$$K \rightarrow -\frac{n}{2} \log(\det(K)) + \frac{n}{2} \langle K, \hat{\Sigma} \rangle_F$$

is convex.

Graphical Lasso : sparse estimation of K

$$\hat{K}_\lambda = \operatorname{argmin}_{K \in \mathcal{S}_p^+} \left\{ -\frac{n}{2} \log(\det(K)) + \frac{n}{2} \langle K, \hat{\Sigma} \rangle_F + \lambda \sum_{a \neq b} |K_{ab}| \right\}$$

Sparse estimation of K (2/2)

- Efficient optimization algorithms.
Ref: Friedman *et al.* [5], Banerjee *et al.* [1]
- Poor empirical results reported by Villers *et al.* [18]
- Theoretical guaranties under some “compatibility conditions” hard to check/interpret (by Ravikumar *et al.* [14])



keep the (good) idea of exploiting the sparsity, but move to the more classical regression framework.

Definitions

- Θ = the set of $p \times p$ matrices with zero on the diagonal
- θ : matrix in Θ defined by $\theta_{ab} = -K_{ab}/K_{bb}$ for $a \neq b$

Characterization

$$\theta = \operatorname{argmin}_{\theta \in \Theta} \|\Sigma^{1/2}(I - \theta)\|_F^2$$

Proof:

$\mathbb{E}[X_a | X_b : b \neq a] = \sum_b \theta_{ba} X_b$ since $X_a = \sum_b \theta_{ba} X_b + \varepsilon_a$. So:

$$\begin{aligned} \theta &= \operatorname{argmin}_{\theta \in \Theta} \mathbb{E} \left[\sum_{a=1}^p \left(X_a - \sum_{b:b \neq a} \theta_{ba} X_b \right)^2 \right] \\ &= \operatorname{argmin}_{\theta \in \Theta} \mathbb{E} [\|X - \theta^T X\|^2] = \operatorname{argmin}_{\theta \in \Theta} \|\Sigma^{1/2}(I - \theta)\|_F^2 \end{aligned}$$

Replacing Σ by $\widehat{\Sigma}$, we obtain

$$\langle (I - \theta), \widehat{\Sigma}(I - \theta) \rangle_F = \frac{1}{n} \|\mathbf{X}(I - \theta)\|_F^2.$$

Estimation procedure



$$\widehat{\theta}_\lambda = \operatorname{argmin}_{\theta \in \Theta} \left\{ \frac{1}{n} \|\mathbf{X} - \mathbf{X}\theta\|_F^2 + \lambda \Omega(\theta) \right\}$$

with $\Omega(\theta)$ enforcing coordinate sparsity.

Examples :

- 1 ℓ^1 penalty : $\Omega(\theta) = \sum_{a \neq b} |\theta_{ab}|$
- 2 ℓ^1/ℓ^2 penalty : $\Omega(\theta) = \sum_{a < b} \sqrt{\theta_{ab}^2 + \theta_{ba}^2}$

With the ℓ^1 penalty : (Meinshausen and Bühlmann [13])

☺ We can split the minimization into p problems in \mathbb{R}^{p-1}

$$[\hat{\theta}_{ba}^{\ell^1}]_{b:b \neq a} = \operatorname{argmin}_{\beta \in \mathbb{R}^{p-1}} \left\{ \frac{1}{n} \|\mathbf{X}_a - \sum_b \beta_b \mathbf{X}_b\|^2 + \lambda |\beta|_{\ell^1} \right\}$$

Very efficient algorithms by coordinate descent.

☹ No constraint enforces that $\hat{\theta}_{ab}^{\ell^1} \neq 0$ when $\hat{\theta}_{ba}^{\ell^1} \neq 0$.

\implies choose an arbitrary decision rule to build \hat{g} from $\hat{\theta}^{\ell^1}$.

Examples:

- 1 set an edge between $a \sim b$ in \hat{g} when either $\hat{\theta}_{ab}^{\ell^1} \neq 0$ or $\hat{\theta}_{ba}^{\ell^1} \neq 0$.
- 2 set an edge $a \sim b$ in \hat{g} when both $\hat{\theta}_{ab}^{\ell^1} \neq 0$ and $\hat{\theta}_{ba}^{\ell^1} \neq 0$.

Regression approach (4/4)

With the ℓ^1/ℓ^2 penalty :

☺ Symmetric zeros

⇒ no ambiguity to define \hat{g} from $\hat{\theta}_\lambda^{\ell^1/\ell^2}$

☹ Computational cost

The minimization problem cannot be split into p subproblems and it is less easy to minimize it in large dimensions.

Algorithm : iterate on couple (a, b) until convergence

1 set $\Delta = \begin{pmatrix} \Delta_{ab} \\ \Delta_{ba} \end{pmatrix}$ with $\Delta_{ab} = \frac{1}{n} \mathbf{X}_a^T (\mathbf{X}_b - \sum_{k \neq a, b} \hat{\theta}_{kb} \mathbf{X}_k)$.

2 set

$$\begin{pmatrix} \hat{\theta}_{ab} \\ \hat{\theta}_{ba} \end{pmatrix} \leftarrow \left(1 - \frac{\lambda}{2\|\Delta\|} \right)_+ \begin{pmatrix} \Delta_{ab} \\ \Delta_{ba} \end{pmatrix}.$$

A series of papers [20, 4, 15] investigate the Bayesian approach.

Issues

- 1 design of sensible priors
- 2 efficient posterior sampling

To the best of my knowledge, cannot handle large dimensional problems

Conclusion

- we have the choice between multiple procedures
- for each procedure, there is at least one (non-scale free) tuning parameter to choose

⇒ we need a selection criterion

We have a collection \mathcal{G} of graphs.

Unbiased risk estimation

$$AIC = -2 \log(L(g)) + 2|g|$$

Bayesian criterion

$$-2 \log(\mathbb{P}(g|\mathbf{X})) \stackrel{n \rightarrow \infty}{\approx} BIC = -2 \log(L(g)) + |g| \log(n) - 2 \log(\mathbb{P}(g))$$

Only mathematically grounded in asymptotic setting :
 p fixed and $n \rightarrow \infty$.

Cross-Validation schemes

train	train	train	train	test
train	train	train	test	train
train	train	test	train	train
train	test	train	train	train
test	train	train	train	train

Figure recursive data splitting for 5-fold Cross-Validation

No guaranty in high-dimensional settings : $p \gg n$ or $p \approx n$.

GGMselect

R package (available on <http://cran.r-project.org/>) which

- 1 generates a collection $\hat{\mathcal{G}}$ of candidates graphs according to the above procedures (+ some variants)
- 2 selects “the best” graph among $\hat{\mathcal{G}}$

Quality criterion

For a graph g

$$\begin{aligned}\text{MSEP}(g) &= \text{Mean Square Error of Prediction related to } g \\ &= \text{bias}(g) + \text{variance}(g)\end{aligned}$$

where

- $\text{bias}(g)$ quantifies how important are the missing edges
- $\text{variance}(g)$ is roughly proportional to the number of edges in g divided by n .

Why MSEP?

It is a way to quantify the importance of each edge.

Ideal

Select $g^* = \operatorname{argmin}\{\operatorname{MSEP}(g) : g \in \hat{\mathcal{G}}\}$
→ g^* unknown!

Selection criterion

”select \hat{g} which minimizes some penalized empirical MSEP”

where the penalty term:

- roughly penalizes each node of \hat{g} according to its degree (number of edges),
- is based on quantiles of Fisher random variables.

Theorem : oracle-like inequality for GGMselect

$$\text{If } \max_{g \in \hat{\mathcal{G}}} \{\text{deg}(g)\} \leq \rho \frac{n}{2(1.1 + \sqrt{\log p})^2}, \quad \text{for some } \rho < 1,$$

then the estimated graph \hat{g} fulfills

$$\text{MSEP}(\hat{g}) \leq c_\rho \mathbb{E} \left[\inf_{g \in \hat{\mathcal{G}}} \{\text{bias}(g) + \log(p) \text{var}(g)\} \right] + R_n,$$

where $R_n = O(\text{Tr}(\Sigma)e^{-c'_\rho n} + \text{CVar}(\Sigma) \log(p)/n)$

with $\text{CVar}(\Sigma) = \sum_a (\Sigma_{aa}^{-1})^{-1}$.

Ref: Giraud *et al.* [7]

Theorem : oracle-like inequality for GGMselect

$$\text{If } \max_{g \in \hat{\mathcal{G}}} \{\text{deg}(g)\} \leq \rho \frac{n}{2(1.1 + \sqrt{\log p})^2}, \quad \text{for some } \rho < 1,$$

then the estimated graph \hat{g} fulfills

$$\text{MSEP}(\hat{g}) \leq c_\rho \mathbb{E} \left[\inf_{g \in \hat{\mathcal{G}}} \{\text{bias}(g) + \log(p) \text{var}(g)\} \right] + R_n,$$

where $R_n = O(\text{Tr}(\Sigma)e^{-c'_\rho n} + \text{CVar}(\Sigma) \log(p)/n)$

with $\text{CVar}(\Sigma) = \sum_a (\Sigma_{aa}^{-1})^{-1}$.

Ref: Giraud *et al.* [7]

Optimality?

- Optimal selection criterion?
 - "minimal" size of the penalty to avoid overfitting
 - minimax estimation rates when $\hat{\mathcal{G}}$ contains good graphs
- What about the condition on the degree? ($n/2 \log p$)
unavoidable, otherwise estimation rate gets worse.

Gaussianity?

Hammersley-Clifford formula

For a random variable X with positive density f

$$\mathcal{L}(X) \sim g \iff f(x) \propto \exp \left(\sum_{c: c \in \text{cliques}(g)} g_c(x_c) \right).$$

Gaussianity?



Data transformation: $f_j(X_j) := \Phi^{-1}(F_j(X_j)) \sim \mathcal{N}(0, 1)$

Assumption: $(f_1(X_1), \dots, f_p(X_p)) \sim \mathcal{N}(0, \Sigma)$

Key point: $\text{graph}(X_1, \dots, X_p) = \text{graph}(f_1(X_1), \dots, f_p(X_p))$

Estimation: work with $\hat{f}_j(X_j) = \Phi^{-1}(\hat{F}_j(X_j))$ for some estimator \hat{F}_j .

Ref: Data transformations proposed by Lafferty *et al.* [10]

Hidden variables?

we may only observe part of the relevant variables:

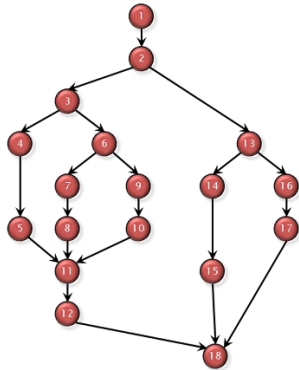
$X = \begin{pmatrix} X_O \\ X_H \end{pmatrix} \sim \mathcal{N} \left(0, \begin{pmatrix} \Sigma_{OO} & \Sigma_{HO} \\ \Sigma_{OH} & \Sigma_{HH} \end{pmatrix} \right)$ with X_O observed and X_H unobserved.

We only have access to $(\Sigma_{OO})^{-1} = K_{OO} - K_{OH}(K_{HH})^{-1}K_{HO}$

Ref: Chandrasekaran *et al.* [3] proposes a sparse + low rank estimation to recover K_{OO} when H is small

Back to directed models

$X_a \perp\!\!\!\perp \{X_b : a \nrightarrow \dots \nrightarrow b\}$
conditionally on $\{X_c : c \rightarrow a\}$



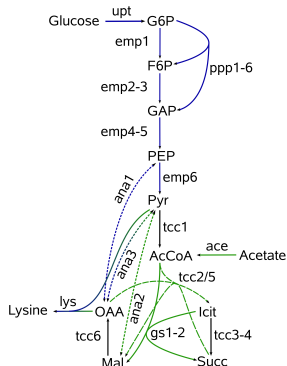
“Causal” inference

Setting

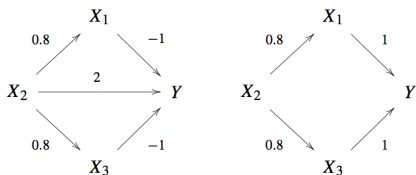
- We have p covariables X_1, \dots, X_p and Y a variable of interest.
- $[X, Y]$ is a Gaussian graphical model according to \vec{g} .
No arrows from Y to X_1, \dots, X_p

Example

- Y is the end product of a metabolic network
- X_1, \dots, X_p are protein abundances



“Causal” inference



Direct versus “Causal” effect

- **Direct effect** given by θ_a in

$$\mathbb{E}[Y|X_1, \dots, X_p] = \sum_b \theta_b X_b$$

- **Causal effect** (relative to \vec{g}) given by β_a in

$$\mathbb{E}[Y|X_a, X_b : b \in \text{pa}(a)] = \beta_a X_a + \sum_{b \in \text{pa}(a)} \beta_b X_b$$

Main issue

Causal effects are defined relative to \vec{g} and there is no unique minimal directed graph...



- find all the DAG $\vec{g}_{(1)}, \dots, \vec{g}_{(m)}$ such that $\mathcal{L}(X) \sim \vec{g}_{(k)}$
- compute a lower bound of the causal effect:

$$\beta_* = \min \{ \beta_{(1)}, \dots, \beta_{(m)} \}$$

PCalg

R package (available on <http://cran.r-project.org/>) which

- estimates the DAG $\vec{g}_{(1)}, \dots, \vec{g}_{(m)}$ from the data
- estimates $\beta_{(1)}, \dots, \beta_{(m)}$ and β_*

Ref:

- Kalish *et al* [9]
- Maathuis *et al.* [12]

Principle of PC algorithm

Init : $g =$ complete graph

Iterate :

- for $a = 1, \dots, p$, for $b \in \text{ne}(a)$: remove $a - b$ if $\widehat{\text{cor}}(X_a, X_b) < t_0$
- for $a = 1, \dots, p$, for $b \in \text{ne}(a)$: remove $a - b$ if $\exists c_1 \in \text{ne}(a)$ such that $\widehat{\text{cor}}(X_a, X_b | X_{c_1}) < t_1$
- for $a = 1, \dots, p$, for $b \in \text{ne}(a)$: remove $a - b$ if $\exists c_1, c_2 \in \text{ne}(a)$ such that $\widehat{\text{cor}}(X_a, X_b | X_{c_1}, X_{c_2}) < t_2$
- ...

Output : skeleton of the DAGs

Last step: compute $\vec{g}_{(1)}, \dots, \vec{g}_{(m)}$ from the skeleton

Riboflavin prediction

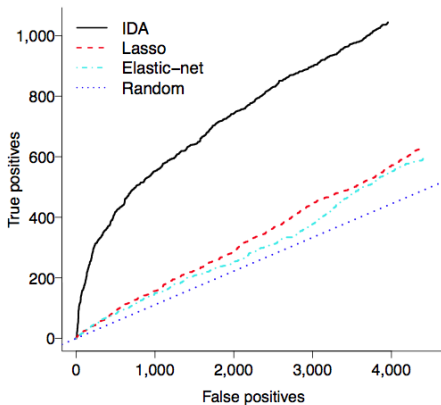


Figure: Important genes for riboflavin production

Ref : Kalish *et al* [9]

- Be aware of over-interpretation : we cannot reliably infer causal networks on i.i.d. data
- Relies on uncheckable assumptions
- But, seems promising for ranking covariables

Statistics in high-dimensional setting

- Despite theorem, do not trust too much statistical inferences in high-dimensional setting $n \ll p$
ex: gene pre-selection, metagenes, etc
- It is not a validation tool, but rather a good tool for providing good hints
- Requires experimental validations.

References on high-dimensional statistics:

- Lecture Notes on High-Dimensional Statistics
<http://www.cmap.polytechnique.fr/~giraud/MSV/LectureNotes.pdf>
- The Element of Statistical Learning
by Hastie, Tibshirani, Friedman
www-stat.stanford.edu/~tibs/ElemStatLearn/

Thank you!

- [1] O. Banerjee, L. El Ghaoui, and A. d'Aspremont.
Model selection through sparse maximum likelihood estimation for multivariate Gaussian or binary data.
[J. Mach. Learn. Res.](#), 9:485–516, 2008.
- [2] R. Castelo and A. Roverato.
A robust procedure for Gaussian graphical model search from microarray data with p larger than n .
[J. Mach. Learn. Res.](#), 7:2621–2650, 2006.
- [3] V. Chandrasekaran, P. Parrilo, and A. Willsky.
Latent variable graphical model selection via convex optimization.
[Annals of Statistics](#), 2012.

- [4] P. Dellaportas, P. Giudici, and G. Roberts.
Bayesian inference for nondecomposable graphical Gaussian models.
[Sankhyā](#), 65(1):43–55, 2003.

- [5] J. Friedman, T. Hastie, and R. Tibshirani.
Sparse inverse covariance estimation with the graphical lasso.
[Biostatistics](#), 9(3):432–441, 2008.

- [6] Nir Friedman.
Inferring Cellular Networks Using Probabilistic Graphical Models.
[Science](#), 303(5659):799–805, 2004.

- [7] C. Giraud, S. Huet, and N. Verzelen.
Graph selection with GGMselect.
[Stat. Appl. Genet. Mol. Biol.](#), 11(3):1–50, 2012.

- [8] M. Kalisch and P. Bühlmann.
Robustification of the pc-algorithm for directed acyclic graphs.
[J. Comput. Graph. Statist.](#), 17(4):773–789, 2008.

- [9] Markus Kalisch, Martin Mächler, Diego Colombo, Marloes H. Maathuis, and Peter Bühlmann.
Causal inference using graphical models with the r package pcalg.
[Journal of Statistical Software](#), 47(11):1–26, 5 2012.

- [10] John Lafferty, Han Liu, and Larry Wasserman.
Sparse nonparametric graphical models.
[Statistical Science](#), 27:519–537, 2012.

- [11] Steffen L. Lauritzen.
[Graphical Models](#).
Oxford University Press, 1996.

- [12] Marloes H. Maathuis, Markus Kalisch, and Peter Bühlmann.
Estimating high-dimensional intervention effects from observational data.
[Ann. Statist.](#), 37(6A):3133–3164, 2009.
- [13] N. Meinshausen and P. Bühlmann.
High-dimensional graphs and variable selection with the lasso.
[Ann. Statist.](#), 34(3):1436–1462, 2006.
- [14] Pradeep Ravikumar, Martin J. Wainwright, Garvesh Raskutti, and Bin Yu.
High-dimensional covariance estimation by minimizing, ℓ^1 -penalized log-determinant divergence.
[Electronic Journal of Statistics](#), 5:935–980, 2011.

- [15] J.G. Scott and C. M. Carvalho.
Feature-inclusion stochastic search for gaussian graphical models.
[J. Comp. Graph. Statist.](#), 17:790–808, 2009.
- [16] P. Spirtes, C. Glymour, and R. Scheines.
[Causation, prediction, and search.](#)
Adaptive Computation and Machine Learning. MIT Press,
Cambridge, MA, second edition, 2000.
- [17] J. Verhoef and K. Frost.
A bayesian hierarchical model for monitoring harbor seal changes in
prince william sound, alaska.
[Environmental and Ecological Statistics](#), 10:201–219, 2003.

- [18] Fanny Villers, Brigitte Schaeffer, Caroline Bertin, and Sylvie Huet.
Assessing the Validity Domains of Graphical Gaussian Models in Order to Infer Relationships among Components of Complex Biological Systems.
[Statistical Applications in Genetics and Molecular Biology](#), 7, 2008.
- [19] A. Wille and P. Bühlmann.
Low-order conditional independence graphs for inferring genetic networks.
[Stat. Appl. Genet. Mol. Biol.](#), 5:Art. 1, 34 pp. (electronic), 2006.
- [20] F. Wong, C. K. Carter, and R. Kohn.
Efficient estimation of covariance selection models.
[Biometrika](#), 90(4):809–830, 2003.