

A multiattribute Gaussian graphical model for inference of multiscale regulatory network in breast cancer

Julien Chiquet, Guillem Rigai, Martina Sundqvist

September 29, 2017

1 Introduction

Write the damn thing

[...]

A deeper generalization of GGM comes by integrating multiple types of data measured from diverse platforms, what is sometimes referred to as *horizontal* integration: not only does this mean a better treatment of the heterogeneity of the data, but it also makes the network reconstruction more robust.

2 Background

Gaussian Graphical Models (GGMs) Lauritzen [1996], Whittaker [1990] are a very convenient tool for describing the patterns at play in complex data sets. Indeed, through the notion of partial correlation, they provide a well-studied framework for spotting direct relationships between variables, and thus reveal the latent structure in a way that can be easily interpreted. Application areas are very broad and include for instance gene regulatory network inference in biology (using gene expression data) as well as spectroscopy, climate studies, functional magnetic resonance imaging, etc. Estimation of GGMs in a sparse, high-dimensional setting has thus received much attention recently. This section provides an overview of this hot and competitive research field of statistical learning. I mainly focus on the state-of-the-art ℓ_1 -regularization methods and their most recent striking variants, insisting on their computational and statistical properties. This provides the reader with the necessary material to approach the second section of this chapter dedicated to my personal contributions to this field.

2.1 Basics on Gaussian graphical models

Let $\mathcal{P} = \{1, \dots, p\}$ be a set of fixed vertices and $X = (X_1, \dots, X_p)^\top$ a random vector describing a signal over this set. The vector $X \in \mathbb{R}^p$ is assumed to be multivariate Gaussian with unknown mean and unknown covariance matrix $\Sigma = (\Sigma_{ij})_{(i,j) \in \mathcal{P}^2}$. No loss of generality is involved when centering X , so we may assume that $X \sim \mathcal{N}(\mathbf{0}_p, \Sigma)$. The covariance matrix Σ , equal to $\mathbb{E}(XX^\top)$ under

the assumption that X is centered, belongs to the set \mathcal{S}_p^+ of positive definite symmetric matrices of size $p \times p$.

Graph of conditional dependencies. GGMs endow Gaussian random vectors with a graphical representation \mathcal{G} of their *conditional dependency structure*: two variables i and j are linked by an undirected edge (i, j) if, conditional on all other variables indexed by $\mathcal{P} \setminus \{i, j\}$, random variables X_i and X_j remain or become dependent. Thanks to the Gaussian assumption, conditional independence actually boils down to a zero conditional covariance $\text{cov}(X_i, X_j | X_{\mathcal{P} \setminus \{i, j\}})$, or equivalently to a zero partial correlation which we denote by ρ_{ij} , the latter being a normalized expression of the former.

Concretely, the inference of a GGM is based upon a classical result originally emphasized in Dempster [1972] stating that partial correlations ρ_{ij} are actually proportional to the corresponding entries in the *inverse* of the covariance matrix $\Sigma^{-1} = \Theta$, also known as the *concentration matrix*. More precisely, we have

$$\rho_{ij} = -\Theta_{ij} / \sqrt{\Theta_{ii}\Theta_{jj}}, \quad \Theta_{ii} = \text{Var}(X_i | X_{\mathcal{P} \setminus i})^{-1}; \quad (1)$$

thus Θ directly describes the conditional dependency structure of X . Indeed, after a simple rescaling, Θ can be interpreted as the adjacency matrix of an undirected weighted graph representing the partial covariance (or correlation) structure between variables X_1, \dots, X_p . Formally, we denote by $\mathcal{G} = (\mathcal{P}, \mathcal{E})$ this graph, the edges of which are characterized by

$$(i, j) \in \mathcal{E} \Leftrightarrow \Theta_{ij} \neq 0, \quad \forall (i, j) \in \mathcal{P}^2 \text{ such that } i \neq j.$$

In words, \mathcal{G} has no self-loop and contains all edges (i, j) such that Θ_{ij} is nonzero. Therefore recovering nonzero entries of Θ is equivalent to inferring the graph of conditional dependencies \mathcal{G} , and the correct identification of nonzero entries is the main issue in this framework.

Maximum Likelihood inference. GGMs fall into the family of exponential models for which the whole range of classical statistical tools applies. As soon as the sample size n is greater than the number p of variables, the likelihood admits a unique maximum over \mathcal{S}_p^+ , defining a maximum likelihood estimator (MLE): suppose we observe a sample $\{X^1, \dots, X^n\}$ composed of n i.i.d. copies of X , stored row-wise once centered in a matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$ such that $(X^i)^\top$ is the i th row of \mathbf{X} . The empirical covariance matrix is denoted by $\mathbf{S}_n = \mathbf{X}^\top \mathbf{X} / n$. Maximizing the likelihood is equivalent to

$$\hat{\Theta}^{\text{mle}} = \arg \max_{\Theta \in \mathcal{S}_p^+} \log \det(\Theta) - \text{Tr}(\Theta \mathbf{S}_n). \quad (2)$$

When $n > p$, Problem (2) admits a unique solution equal to \mathbf{S}_n . The scaled empirical covariance matrix \mathbf{S}_n follows a Wishart distribution while its inverse \mathbf{S}_n^{-1} follows an inverse Wishart distribution with computable parameters.

There are two major limitations with the MLE regarding the objective of graph reconstruction by recovering the pattern of zeroes in Θ . First, it provides an estimate of the saturated graph: all variables are connected to each other; second, we need n to be larger than p to be able to even define this estimator, which is rarely the case in genomics. In any case, the need for regularization

and feature selection is huge. A natural assumption is that the true set of direct relationships between the variables remains small, that is, the true underlying graph is sparse (say, of the order of p rather than the order of p^2). Sparsity makes estimation feasible in the case where $n < p$ since we can concentrate on sparse or shrinkage estimators with fewer degrees of freedom than in the original problem. Henceforth, the question of selecting the correct set of edges in the graph is treated as a question of variable selection.

High-dimensional inference of GGM. The different methods for the inference of sparse GGMs in high-dimensional settings fall into roughly three categories. The first contains constraint-based methods, performing statistical tests Castelo and Roverato [2006], Drton and Perlman [2007, 2008], Kiiveri [2011], Wille and Bühlmann [2006]. However, they either suffer from the excessive computational burden Castelo and Roverato [2006], Wille and Bühlmann [2006] or strong assumptions Drton and Perlman [2007, 2008] that correspond to regimes never attained in real situations. The second of these categories is composed of Bayesian approaches, see for instance Dobra et al. [2004], Jones et al. [2005], Rau et al. [2012], Schwaller et al. [2015]. However, constructing priors on the set of concentration matrices is not a trivial task and the use of MCMC procedures limits the range of applications to moderate-sized networks. The third category contains regularized estimators, which add a penalty term to the likelihood in order to reduce the complexity or degrees of freedom of the estimator and more generally regularize the problem: throughout this chapter, I focus on methods of this kind. More precisely, I focus on ℓ_1 -regularized procedures, which are freed from any test procedure – and thus multiple testing issues – since they directly perform estimation and selection of the most significant edges by zeroing entries in the estimator of Θ . The remainder of this section is dedicated to a quick review of the state-of-the-art methods of this kind.

2.2 Sparse methods for GGM inference

The idea underlying sparse methods for GGM is the same as for the Lasso in linear regression (see Example ??, Section ??): it basically uses ℓ_1 -regularization as a convex surrogate of the ideal but computationally intensive ℓ_0 -regularized problem:

$$\arg \max_{\Theta \in \mathcal{S}_p^+} \log \det(\Theta) - \text{Tr}(\Theta \mathbf{S}_n) - \lambda \|\Theta\|_{\ell_0}. \quad (3)$$

Problem (3) achieves a trade-off between the maximization of the likelihood and the sparsity of the graph within a single optimization problem. The penalty term can also be interpreted as a log prior on the coefficients in a Bayesian perspective. BIC or AIC criteria are special cases of such ℓ_0 regularized problems, except that the maximization is made upon a restricted subset of candidates $\{\tilde{\Theta}_1, \dots, \tilde{\Theta}_m\}$ and the choice of λ is fixed ($\log(n)$ for BIC and $1/2$ for AIC). Actually solving (3) would require the exploration of all possible 2^p graphs. On the contrary, by preserving the convexity of the optimization problem, ℓ_1 -regularization paves the way to fast algorithms. For the price of a little bias on all the coefficients, we get to shrink some coefficients to exactly 0, operating selection and estimation in one single step as hoped in Problem (3).

Graphical-Lasso. The criterion optimized by the graphical-Lasso was simultaneously proposed in Yuan and Lin [2007] and Banerjee et al. [2008]. It corresponds to the estimator obtained by fitting the ℓ_1 -penalized Gaussian log-likelihood, *i.e.* the tightest convex relaxation of (3):

$$\hat{\Theta}_\lambda^{\text{glasso}} = \arg \max_{\Theta \in S_p^+} \log \det(\Theta) - \text{Tr}(\Theta \mathbf{S}_n) - \lambda \|\Theta\|_{\ell_1}. \quad (4)$$

In this regularized problem, the ℓ_1 -norm drives some coefficients of Θ to zero. The non-negative parameter λ tunes the global amount of sparsity: the larger the λ , the fewer edges in the graph. A large enough penalty level produces an empty graph. As λ decreases towards zero, the estimated graph tends towards the saturated graph and the estimated concentration matrix tends towards the usual MLE (2). By construction, this approach guarantees a well-behaved estimator of the concentration matrix *i.e.* sparse, symmetric and positive-definite, which is a great advantage of this method.

Ever since Criterion (4) was proposed, many efforts have been dedicated to developing efficient algorithms for its optimization. In the original proposal of Banerjee et al. [2008], it is shown that solving for one row of matrix Θ in (4) while keeping other rows fixed boils down to a Lasso problem. The global problem is solved by cycling over the matrix rows until convergence. Thus, if one considers that L passes over the whole matrix are needed to reach convergence, a rough estimation of the overall cost is of the order of $Lp \times$ (cost for solving for one row). With a block-coordinate update each iteration over a row has $\mathcal{O}(p^3)$ complexity and their implementation is $\mathcal{O}(Lp^4)$ for L sweeps over the whole matrix $\hat{\Theta}$. In Banerjee et al. [2008] again, a rigorous analysis is conducted in Nesterov's framework Nesterov [2005] showing that the complexity for a single λ reaches $\mathcal{O}(p^{4.5}/\varepsilon)$ where ε is the desired accuracy of the final estimate.

The *Graphical-Lasso* algorithm of Friedman et al. [2008] follows the same line but builds on a coordinate descent algorithm to solve each underlying Lasso problem. While no precise complexity analysis is possible with these methods, empirical results tend to show that this algorithm is faster than the original proposal of Banerjee et al. [2008]. Additional insights on the convergence of the graphical-Lasso are provided in Mazumder and Hastie, simultaneously with Witten et al. [2011], showing how to take advantage of the problem sparsity by decomposing (4) into block diagonal problems depending on λ : this considerably reduces the computational burden in practice. Implementations of the graphical-Lasso algorithm are available in the R-packages **glasso**, **huge** Zhao et al. [2014], or **simone** ???. The most recent notable efforts related to the optimization of (4) are due to Hsieh et al. [2014, 2013] and the QUIC (then BIG&QUIC) algorithm, a quadratic approximation which allows (4) to be solved up to $p = 1,000,000$ with a super-linear rate of convergence and with bounded memory. The R-package **quic** implements the first version of this algorithm.

On the statistical side, the most striking results are due to Ravikumar et al. [2011]: they show that selection consistency of the estimator defined by (4) – that is, recovery of the true underlying graphical structure –, is met in the sub-Gaussian case when, for an appropriate choice of λ , the sample size n is of the same order as $\mathcal{O}(d^2 \log(p))$, where d is the highest degree in the target graph. Additional conditions on the empirical covariance between relevant and

irrelevant features are required, known as the “irrepresentability conditions” in the Lasso case. Such statistical results are important since they provide insights on the “data” situations where such methods may either be successful or completely hopeless. More on this is discussed in Verzelen [2012]. For instance, this should prevent blindly applying the graphical-Lasso in situations where the sample size n is too small compared to p . Similarly, when the presence of hub nodes with high degree is suspected, the estimated graph should be interpreted with care.

Neighborhood selection. This approach, proposed in Meinshausen and Bühlmann [2006], determines the graph of conditional dependencies by solving a series of p independent Lasso problems, successively estimating the neighborhoods of each variable and then applying a final reconciliation step as post-treatment to recover a symmetric adjacency matrix. Concretely, a given column \mathbf{X}_j of the data matrix is “explained” by the remaining columns $\mathbf{X}_{\setminus j}$ corresponding to the remaining variables: the set $\text{ne}(j)$ of neighbors of variable j in the graph \mathcal{G} is estimated by the support of the vector solving

$$\hat{\beta}_j = \arg \min_{\beta \in \mathbb{R}^{p-1}} \frac{1}{2n} \|\mathbf{X}_j - \mathbf{X}_{\setminus j} \beta\|_{\ell_2}^2 + \lambda \|\beta\|_{\ell_1}. \quad (5)$$

Indeed, if each row of \mathbf{X} is drawn from a multivariate Gaussian $\mathcal{N}(\mathbf{0}, \Theta^{-1})$, then the best linear approximation of \mathbf{X}_j by $\mathbf{X}_{\setminus j}$ is given by

$$\mathbf{X}_j = \sum_{k \in \text{ne}(j)} \beta_{jk} \mathbf{X}_k = - \sum_{k \in \text{ne}(j)} \frac{\Theta_{jk}}{\Theta_{jj}} \mathbf{X}_k, \quad (6)$$

thus coefficients β_j and column Θ_j – once its diagonal elements are removed – share the same support. By support, we mean the set of nonzero coefficients. Adjusting (5) for each $j = 1, \dots, p$ allow us to reconstruct the full graph \mathcal{G} . Because the neighborhoods of the p variable are selected separately, a post symmetrization must be applied to manage inconsistencies between edge selections; Meinshausen and Bühlmann [2006] suggests AND or OR rules.

Let us fill the gap with Criterion (4). First, note that the p regression problem can be rewritten as a unique matrix problem, where \mathbf{B} contains p vectors $\beta_j, j = 1, \dots, p$:

$$\hat{\mathbf{B}}^{\text{ns}} = \arg \min_{\mathbf{B} \in \mathbb{R}^{p \times p}, \text{diag}(\mathbf{B}) = \mathbf{0}_p} \frac{1}{2} \text{Tr}(\mathbf{B}^\top \mathbf{S}_n \mathbf{B}) - \text{Tr}(\mathbf{B}^\top \mathbf{S}_n) + \lambda \|\mathbf{B}\|_{\ell_1}. \quad (7)$$

In fact, it can be shown Rocha et al. [2008], [?], Ravikumar et al. [2010] that the optimization problem (7) corresponds to the minimization of a penalized, negative *pseudo*-likelihood: the joint distribution of X is approximated by the product of the p distributions of the p variables conditional on the other ones, that is

$$\log \mathbb{P}(\mathbf{X}; \Theta) = \sum_{j=1}^p \sum_{i=1}^n \log \mathbb{P}(X_j^i | X_{\setminus j}^i; \Theta_j).$$

This pseudo-likelihood is based upon the (false) assumption that conditional distributions are independent. Moreover, all variables are assumed to share the

same variance in this formulation. Building on these remarks, Rocha et al. [2008] amend criterion (7) by the adjunction of an additional symmetry constraint, and introduce additional parameters to account for different variances between the variables.

Concerning the computational aspect, this approach has very efficient implementation as it basically boils down to solving p Lasso problems. Suppose for instance that the target neighborhood size is k per variable: fitting the whole solution path of a Lasso problem using the Lars algorithm can be done in $\mathcal{O}(npk)$ complexity Bach et al. [2012]. This must be multiplied by p for the whole network, yet we underline that a parallel implementation is straightforward in this case. This makes this approach quite competitive, especially when coupled with additional bootstrap or resampling techniques Meinshausen and Bühlmann [2010].

On the statistical side, neighborhood selection has been reported to be sometimes empirically more accurate in terms of edge detection than is the graphical-Lasso Villers et al. [2008], Rocha et al. [2008] on certain types of data. This is somewhat supported by the statistical analysis of Ravikumar et al. [2011], who show that under the classical irrepresentability conditions for the Lasso Zhao and Yu [2006], Meinshausen and Bühlmann [2006] and for an appropriate choice of λ , neighborhood selection achieves selection consistency with high probability when the sample size n is of the order of $\mathcal{O}(d \log(p))$ with d the maximal degree of the target graph \mathcal{G} . This is to be compared with the $\mathcal{O}(d^2 \log(p))$ required by the graphical-Lasso (even if the corresponding “irrepresentability conditions” are not strictly comparable). A rough explanation for this difference on the asymptotic is that the graphical-Lasso intends to estimate the concentration matrix on top of selecting the nonzero entries, while neighborhood selection focuses on the selection problem.

Constrained ℓ_1 -minimization for inverse matrix estimation (CLIME).

The CLIME estimator has been proposed by Cai et al. [2011] and is designed to avoid the cumbersome “irrepresentability conditions” required for the Graphical-Lasso and the neighborhood selection approaches, while providing statistical guarantees on the support recovery.

The definition of CLIME builds on the remark that the solution to (4) must verify the following first order optimality condition – or subgradient equations:

$$\left(\hat{\Theta}^{\text{glasso}}\right)^{-1} - \mathbf{S}_n = \lambda \mathbf{\Gamma}, \quad \text{with } \mathbf{\Gamma}_{ij} \begin{cases} = \text{sign}(\hat{\Theta}_{ij}^{\text{glasso}}) & \text{if } \hat{\Theta}_{ij}^{\text{glasso}} \neq 0, \\ \in [-1, 1] & \text{otherwise.} \end{cases}$$

This suggests the optimization problem

$$\underset{\Theta \in \mathcal{S}_p^+}{\text{minimize}} \|\Theta\|_{\ell_1}, \quad \text{s.c.} \quad \|\Theta^{-1} - \mathbf{S}_n\|_{\infty} \leq \lambda,$$

which is too hard to solve. Removing the positive-definite requirement and multiplying the constraint by Θ , we encounter the problem solved by CLIME:

$$\hat{\Theta}_{\lambda}^{\text{clime}} = \arg \min_{\Theta \in \mathcal{M}_{pp}} \|\Theta\|_{\ell_1}, \quad \text{s.c.} \quad \|\mathbf{I}_p - \Theta \mathbf{S}_n\|_{\infty} \leq \lambda. \quad (8)$$

This estimator is not necessarily symmetric and a post-treatment is required as for neighborhood selection. But it can also be easily distributed for each column

of Θ_j , which requires the resolution of a linear program of complexity $\mathcal{O}(p^2k)$, with k the targeted number of neighbors per variable. This is slightly more demanding than neighborhood-selection but remains extremely competitive.

On the statistical side, the CLIME estimator achieves selection consistency at a rate comparable to that of the Graphical-Lasso Cai et al. [2011], and is better in its adaptive (weighted) version. Its great advantage is that no particular assumption like a irrepresentability condition – which can never be established in practice – is required for the data matrix \mathbf{X} . This method is distributed via the R-package **fastclime**, and an implementation Wang et al. [2013] is reported to solve for problems with millions of features.

Sparse Partial Correlation Estimation (SPACE). In Peng et al. [2009], the gap is completely filled between linear regression, Gaussian graphical model and neighborhood selection with a method that directly penalizes the partial correlations within the linear model. Indeed, by combining firstly Relationship 1 between the partial correlations and the concentration matrix, and secondly, Relationship 6 between the coefficients in linear regression and concentration matrix, one has

$$\mathbf{X}_j = \sum_{k \in \text{ne}(j)} \beta_{jk} \mathbf{X}_k + \varepsilon = \sum_{k \in \text{ne}(j)} \rho_{jk} \sqrt{\frac{\Theta_{kk}}{\Theta_{jj}}} \mathbf{X}_k + \varepsilon,$$

which suggests the following optimization problem

$$(\hat{\boldsymbol{\rho}}_{\lambda}^{\text{space}}, \text{diag}(\boldsymbol{\Theta})) = \arg \min_{\boldsymbol{\rho} \in \mathbb{R}^{p(p-1)}, \text{diag}(\boldsymbol{\Theta})} \frac{1}{2} \sum_{j=1}^p \omega_j \left\| \mathbf{X}_j - \sum_{k=1}^p \rho_{jk} \sqrt{\frac{\Theta_{kk}}{\Theta_{jj}}} \mathbf{X}_k \right\|_{\ell_2}^2 + \lambda \|\boldsymbol{\rho}\|_{\ell_1}, \quad (9)$$

where $\boldsymbol{\rho}$ is a vector containing all the pairwise partial correlations, $\text{diag}(\boldsymbol{\Theta})$ contains the diagonal elements of $\boldsymbol{\Theta}$, that is to say, the partial covariances of all the variables, and finally ω_j are some positive (given) weights.

Although the optimization of (9) is more demanding than is neighborhood selection, the problem is jointly convex in $(\text{diag}(\boldsymbol{\Theta}), \boldsymbol{\rho})$. When $\text{diag}(\boldsymbol{\Theta})$ is fixed, the problem has the same complexity as does neighborhood selection, and the authors claim that only a couple of iterations alternating over each of the two parameters $(\text{diag}(\boldsymbol{\Theta}), \boldsymbol{\rho})$ are needed for convergence. It thus remains a lot more efficient than the graphical-Lasso. On top of that, the method intrinsically imposes symmetry over the partial correlations $\boldsymbol{\rho}$. In short, it embeds the computational advantage of neighborhood selection while estimating the conditional variance as in the graphical-Lasso. It is available in the R-package **space**. Further refinements and statistical analyses have been recently proposed in Khare et al. [2014].

Model selection issues. Up to this point, we have completely avoided the fundamental model selection issue, that is, the choice of the tuning parameter λ , which is at play in all the sparse methods mentioned thus far. The first possibility is to rely on information criteria of the form

$$\text{IC}_{\lambda} = -2\log\text{lik}(\hat{\boldsymbol{\Theta}}_{\lambda}; \mathbf{X}) + \text{pen}(\text{df}(\hat{\boldsymbol{\Theta}}_{\lambda})),$$

where “pen” is a function penalizing the model complexity, described by df, the degrees of freedom of the current estimator. We meet the AIC by choosing $\text{pen}(x) = 2x$ and the BIC by choosing $\text{pen}(x) = \log(n)x$. However, AIC and BIC are based upon assumptions which are not suited to high-dimensional settings (see Giraud et al. [2012b]). Moreover, the notion of degrees of freedom for sparse methods has to be specified, not to mention that one has to adapt these criteria to the case of GGMs. An example of a criterion meeting these prerequisites is the extended BIC for sparse GGMs Foygel and Drton [2010]:

$$\text{EBIC}_\gamma(\hat{\Theta}_\lambda) = -2\log\text{lik}(\hat{\Theta}_\lambda; \mathbf{X}) + |\mathcal{E}_\lambda|(\log(n) + 4\gamma\log(p)), \quad (10)$$

where the function df is equal to $|\mathcal{E}|$, the total number of edges in the inferred graph. The parameter $\gamma \in [0, 1]$ is used to adjust the tendency of the usual BIC – recovered for $\gamma = 0$ – to choose overly dense graphs in the high-dimensional setting. Further justification can be found in Foygel and Drton [2010]. A competing approach, designed to compare a family of GGM – possibly inferred with different methods –, is GGMSelect Giraud et al. [2012a], Giraud [2008].

Another possibility is to rely on resampling/subsampling procedures to select a set of edges which are robust to small variations of the sample. The most popular approach is the *Stability Selection* procedure proposed in Meinshausen and Bühlmann [2010], also related to the bootstrapped procedure of Bach [2008]. A similar approach, called StaRS (Stability approach to Regularization Selection) is developed specifically in the context of GGM in Liu et al. [2010]. The basic idea is as follows: for a given range of the tuning parameter $\Lambda = [\lambda_{\min}, \lambda_{\max}]$, the same method is fitted on many subsamples (with or without replacement) with size, say $n/2$. The idea is then to construct a score indexed on Λ that measures stability – or instability – of the selected variables. The selected edges are those matching a given score, for which the probability of false discovery is controlled. This requires an additional threshold in place of a choice of λ , but the authors in Meinshausen and Bühlmann [2010], Liu et al. [2010] claim that such a threshold is typically much less sensitive than is the tuning parameter λ . An application of such resampling techniques to the inference of biological networks has been pursued with success in Haury et al. [2012], advocating for the use of stability methods on real problems.

A final possibility — that remains somewhat confidential while writing these lines — is to rely on sparse procedures which are less sensitive to λ : among these, we may cite the “scaled-Lasso” Sun and Zhang [2012] for linear regression, adapted to the context of network inference in a neighborhood-selection-like fashion in Sun and Zhang [2013].

Extensions towards non Gaussian settings. As hopefully illustrated throughout this section, sparse GGM is a mature and well controlled framework, with solid contributions both on the statistical and the computational sides. There is also expanding innovative literature tending to broaden the applicability of GGMs, especially to overcome the Gaussian assumption. Indeed, particularly in genomics, there is a growing interest for the multivariate modeling of discrete random vectors, as sequencing techniques provide us with count data. In this perspective, some attempts were made for a Poisson version of the above techniques: in Allen and Liu [2012] for instance, the neighborhood selection approach is extended to a sparse generalized linear model setup; still, inter-

pretability of the inferred network is questionable, as a null partial correlation does not mean conditional dependency in the non-Gaussian case. In a recent paper Yang et al. [2013], a review of existing Poisson graphical models is provided, where the notion of conditional dependency is more carefully specified.

Finally, there is much interest for pretreatment methods which change the original data into more “Gaussian” data via simple transformations. Hence, we can still take advantage of the well-controlled sparse GGM framework. A successful work based on Gaussian copulas is the nonparanormal distribution developed in Liu et al. [2009]. It is implemented within the R-package **huge**, at a negligible cost compared to that of the inference process itself.

3 Accounting for multiscale data: multi-attribute GGM

We now place ourselves in the situation where, for our collection of features \mathcal{P} , we observe not one but several attributes. The question at hand remains the same, that is to say, unraveling strong interactions between these features according to the observation of their attributes. Such networks are known as “association networks”, which are systems of interacting elements, where a link between two different elements indicates a sufficient level of similarity between element attributes. In this section, we are interested in reconstructing such networks based upon n observations of a set of K attributes of the p elements composing the vertices of the network. To this end, we propose a natural generalization of sparse GGM to sparse *multi-attribute* GGMs.

Remark. *This work was planned for submission as a research paper with Christophe Ambroise and Eric Kolaczyk when we came across an independent work Kolar et al. [2014] on the [arXiv](#) that proposes nearly the same approach. We somewhat gave up this project in its original form, which I choose to include in this manuscript as it brings interesting and renewed questions on GGMs.*

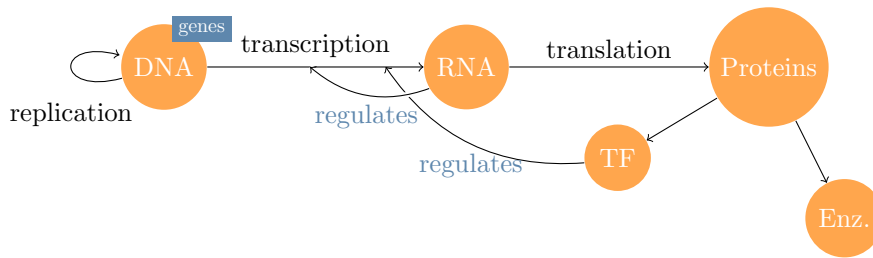


Figure 1: Basic example of a multi-attribute network in genomics: activity of a gene can be measured at the transcriptomic and proteomic levels, and gene regulation affected accordingly

Why multi-attribute networks? The need for multi-attribute networks is relevant in many application fields, but seems particularly applicable in genomics. Indeed, with the plurality of emerging technologies and sequencing

techniques, it is possible to record many signals related to the same set of biological features at various scales or locations of the cell. Consider for instance the simplifying – still hopefully didactic – central dogma of molecular biology, sketched in Figure 1: basically, expression of a gene encoding for a protein can be measured either at the transcriptome level, in terms of its quantity of mRNA, or at the protein level, in terms of the concentration of the associated protein. Still, different technologies are used to measure either the transcriptome or the proteome, typically, microarray or sequencing technology for gene expression levels and cytometric or spectrometric experiments for protein concentrations. Although these signals are very heterogeneous (different levels of noise, count vs. continuous data, etc.), they do share commonality as they undergo common biological processes. We then put an edge in the network if it is supported in both spaces (gene and protein spaces). Our hope is that molecular profiles combined on the same set of biological samples can be *synergistic*, in order to identify a “consensus” and hopefully more robust network.

Multi-attribute GGM. Let $\mathcal{P} = \{1, \dots, p\}$ be a set of variables of interest, each of them having some K attributes. Consider the random vector $X = (X_1, \dots, X_p)^\top$ such as $X_i = (X_{i1}, \dots, X_{iK})^\top \in \mathbb{R}^K$ for $i \in \mathcal{P}$. The vector $X \in \mathbb{R}^{pK}$ describes the K recorded signals for the p features. We assume that X is a multivariate centered Gaussian vector, that is, $X \sim \mathcal{N}(\mathbf{0}, \Sigma)$, with covariance and concentration matrices defined block-wise

$$\Sigma = \begin{bmatrix} \Sigma_{11} & & \Sigma_{1p} \\ & \ddots & \\ \Sigma_{p1} & & \Sigma_{pp} \end{bmatrix}, \quad \Theta = \begin{bmatrix} \Theta_{11} & & \Theta_{1p} \\ & \ddots & \\ \Theta_{p1} & & \Theta_{pp} \end{bmatrix}, \quad \Sigma_{ij}, \Theta_{ij} \in \mathcal{M}_{K,K}, \quad \forall (i, j) \in \mathcal{P}^2,$$

where $\mathcal{M}_{a,b}$ is the set of real-valued matrices with a rows, b columns. Such a multi-attribute framework has been studied in Katenka and Kolaczyk [2012] with a reconstruction method based upon canonical correlations in order to test dependencies between pairs (i, j) at the attribute level using covariance. Here, we propose to rely on partial correlations in a multivariate framework rather than (canonical) correlations to describe relationships between the features, and thus extend GGM to a multi-attribute framework. The objective is to define a “canonical” version of partial correlations. In our setting, the target network $\mathcal{G} = (\mathcal{P}, \mathcal{E})$ is defined as the multivariate analog of the conditional graph for univariate GGM, that is

$$(i, j) \in \mathcal{E} \Leftrightarrow \Theta_{ij} \neq \mathbf{0}_{KK}, \quad \forall i \neq j. \quad (11)$$

In words, there is no edge between two variables i and j when their attributes are all conditionally independent.

A multivariate version of neighborhood selection. Our idea for performing sparse multi-attribute GGM inference is to define a multivariate analog of the neighborhood selection approach Meinshausen and Bühlmann [2006] (see Section 2.2, Equations (5) and (7)). Indeed, it seems to be the most natural and convenient setup toward multivariate generalization. Nevertheless, we think that the graphical-Lasso (4), CLIME (8) or SPACE (9) settings may have a close equivalent multi-attribute version.

To this end, we look at the multivariate analog of equation (6): in a multivariate linear regression setup, it is a matter of straightforward algebra to see that the conditional distribution of $X_j \in \mathbb{R}^K$ on the other variables is

$$X_j | X_{\setminus j} = x \sim \mathcal{N}(-\Theta_{jj}^{-1} \Theta_{j \setminus j} x, \Theta_{jj}^{-1}) .$$

Equivalently, letting $\mathbf{B}_j^T = -\Theta_{jj}^{-1} \Theta_{j \setminus j}$, one has

$$X_j | X_{\setminus j} = \mathbf{B}_j^T X_{\setminus j} + \varepsilon_j \quad \varepsilon_j \sim \mathcal{N}(\mathbf{0}, \Theta_{jj}^{-1}), \quad \varepsilon_j \perp X,$$

where $\mathbf{B}_j \in \mathcal{M}_{(p-1)K, K}$ is defined block-wise

$$\mathbf{B}_j = \begin{bmatrix} \mathbf{B}_j^{(1)} \\ \vdots \\ \mathbf{B}_j^{(p-1)} \end{bmatrix} = \Theta_{jj}^{-1} \times \begin{bmatrix} \Theta_j^{(1)} \\ \vdots \\ \Theta_j^{(p-1)} \end{bmatrix},$$

and where each $\mathbf{B}_j^{(i)}$ is a $K \times K$ matrix which links attributes of variables (i, j) . We see that recovering the support of \mathbf{B}_j block-wise is equivalent to reconstructing the network defined in (11). Estimation of \mathbf{B}_j is thus typically achieved through sparse methods. To this end, we consider an i.i.d. sample $\{X^\ell\}_{\ell=1}^n$ of X such that each attribute is observed n times for the p variable, each X^n being a pK -size row vector staked in a $\mathcal{M}_{n, pK}$ data matrix \mathbf{X} , so that $\mathbf{X}_j \in \mathcal{M}_{n, K}$ is a real-value, $n \times K$ block matrix containing the data related to the j th variable:

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}^1 \\ \vdots \\ \mathbf{x}^n \end{bmatrix} = [\mathbf{X}_1 \quad \dots \quad \mathbf{X}^p] = \begin{bmatrix} X_1^{11} & X_1^{1K} & \dots & X_1^{p1} & \dots & X_1^{pK} \\ \vdots & \vdots & \dots & \vdots & \vdots & \vdots \\ X_n^{11} & X_n^{1K} & \dots & X_n^{p1} & \dots & X_n^{pK} \end{bmatrix}.$$

Using these notations, a direct generalization of the neighborhood selection is to predict for each $j = 1, \dots, p$ the data block \mathbf{X}_j by regressing on $\mathbf{X}_{\setminus j}$. In matrix form, this can be written as the optimization problem

$$\arg \min_{\mathbf{B}_j \in \mathbb{R}} J(\mathbf{B}_j), \quad J(\mathbf{B}_j) = \frac{1}{2n} \|\mathbf{X}_j - \mathbf{X}_{\setminus j} \mathbf{B}_j\|_F^2 + \lambda \Omega(\mathbf{B}_j), \quad (12)$$

where $\|\mathbf{A}\|_F = \sqrt{\sum_{i,j} \mathbf{A}_{ij}^2}$ is the Frobenius norm of matrix \mathbf{A} and Ω is a penalty which constrains \mathbf{B}_j block-wise.

Choosing a penalizer. Various choices for Ω in (12) seem relevant: by simply setting $\Omega_0(A) = \sum_{i,j} |A_{i,j}|$, we just encourage sparsity among the \mathbf{B}_i and thus do not couple the attributes. A clever choice would be to activate a set of attributes all together: hence, the group is defined by all the K attributes between variables i and j , therefore the penalizer turns to a group-Lasso like penalty

$$\Omega_1(\mathbf{B}_j) = \sum_{i \in \mathcal{P} \setminus j} \|\mathbf{B}_j^{(i)}\|_F, \quad (13)$$

in which case convex analysis and subdifferential calculus (see Boyd and Vandenberghe [2006]) can be used to show that a \mathbf{B}_i is optimal for Problem (12) if and only if

$$\begin{cases} \forall i : \mathbf{B}_j^{(i)} \neq 0, & \left(\mathbf{S}_{ij} + \frac{\lambda}{\|\mathbf{B}_j^{(i)}\|_F} I \right)^{-1} \mathbf{S}_{ij} = \mathbf{B}_j^{(i)}, \\ \forall i : \mathbf{B}_j^{(i)} = \mathbf{0}_{KK}, & \|\mathbf{S}_{ij}\|_F \leq \lambda \end{cases}, \quad (14)$$

where $\mathbf{S}_{ij} \in \mathcal{M}_{KK}$ is a $K \times K$ block in the empirical covariance matrix $\mathbf{S}_n = n^{-1} \mathbf{X}^\top \mathbf{X}$, which shows the same block-wise decomposition as $\mathbf{\Sigma}$ or $\mathbf{\Theta}$. This paves the way for an optimization algorithm like block-coordinate descent which we implemented, although we omit details here.

At the time we were working on this model, another idea that we had in mind — although we did not push too far — was to propose a penalty based upon the nuclear norm $\|\mathbf{A}\|_* = \sum_j \nu_j$, where (ν_1, \dots, ν_p) is the vector of singular values of \mathbf{A} . This somewhat penalizes the rank of a matrix, which would be desirable for matrix $\mathbf{\Theta}_{ij}$ when many attributes are shared between (i, j) . We thus might define a penalty on $\mathbf{\Theta}$ in place of \mathbf{B}_j , with something like

$$\Omega_1(\mathbf{B}_j) = \sum_{i \in \mathcal{P} \setminus j} \|\mathbf{\Theta}_{ij}\|_*. \quad (15)$$

However, this idea remains only at the feasible stage for now.

4 Numerical study

We propose a simple simulation to illustrate the interest of using multi-attribute networks and the efficiency of our proposal. The simulations are set up as follows:

1. Draw a random undirected network with p nodes from the Erds-Renyi model;
2. Expand the associated adjacency matrix to multivariate space with

$$\mathbf{A} = (\mathbf{A} + I) \otimes \mathbb{I}_{K \times K};$$

3. Compute $\mathbf{\Theta}$ a positive definite approximation of \mathbf{A} by replacing null and negative eigenvalues by a small constant;
4. Control the difficulty of the problem with $\gamma > 0$ such that $\mathbf{\Theta} = \mathbf{\Theta} + \gamma I$;
5. Draw an i.i.d. sample \mathbf{X} of $X \sim \mathcal{N}(0, \mathbf{\Theta}^{-1})$.

We choose small networks with $p = 20$, with 20 edges on average and vary n from $p/2$ to $2p$. We consider cases where the number of attributes is $K = 2$ or $K = 4$. We either apply the usual neighborhood selection procedure on each dimension separately, or its multi-attribute counterpart with group-like penalty (13) on the multivariate data. We compute the AUC for each method and replicate the experiment 50 times. On Figure 2, it is clear that aggregation improves upon single-attribute methods.

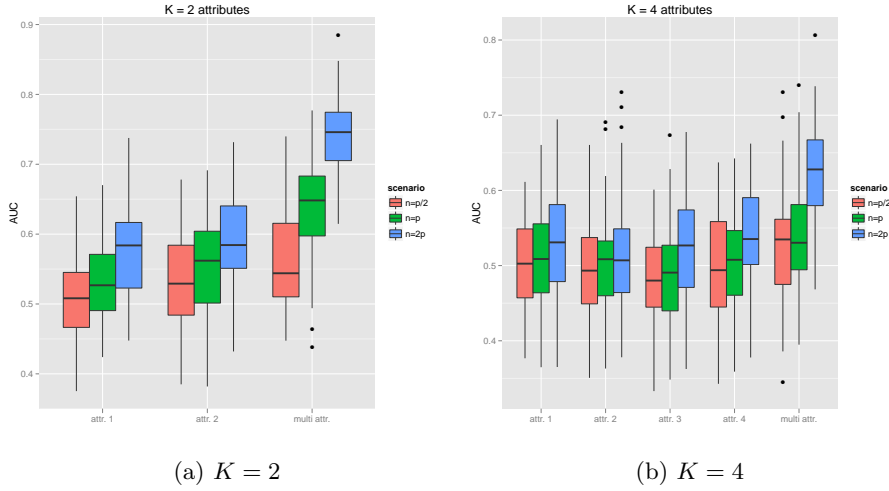


Figure 2: Simple simulation study for the multi-attribute network inference problem: the multivariate procedure improves over the univariate procedures in every situation when networks are close for each attribute.

Illustration: Gene/Protein regulatory network inference. As an illustration, we applied our sparse multi-attribute GGM approach to the NCI-60 cancer line data set. This data set consists in molecular profiles on a panel of 60 diverse human cancer cell lines. We use both protein and gene profiling experiments. For the former, we have samples for 92 antibodies from reverse-phase lysate arrays (RPLA); for the latter, expression is measured for 9,000 RNA with Human Genome U95 affymetrix. A *consensus set* composed of 91 protein and the corresponding gene profiles is retained for the $n = 60$ samples.

We infer a sparse GGM on each attribute (gene and protein), separately to start with, and then on its multi-attribute version. We do this on a large grid of the tuning parameter and thus have three families of networks indexed by their number of edges. Figure 3 demonstrates that our sparse multi-attribute method capture the characteristics of both univariate networks, as the Jaccard similarity index is high between each uni-attribute network and the multi-attribute network, while it remains low when comparing uni-attribute networks together. This tends to prove that this multi-attribute version proposes a consensus version of the interactions at hand in the cell, and one which is hopefully more robust to noise and small misregulations.

References

- G.I. Allen and Z. Liu. A log-linear graphical model for inferring genetic networks from high-throughput sequencing data. In *Bioinformatics and Biomedicine (BIBM), 2012 IEEE International Conference on*, pages 1–6. IEEE, 2012.
- F. Bach. Bolasso: model consistent lasso estimation through the bootstrap. In *Proceedings of the 25th international conference on Machine learning*, pages 33–40. ACM, 2008.

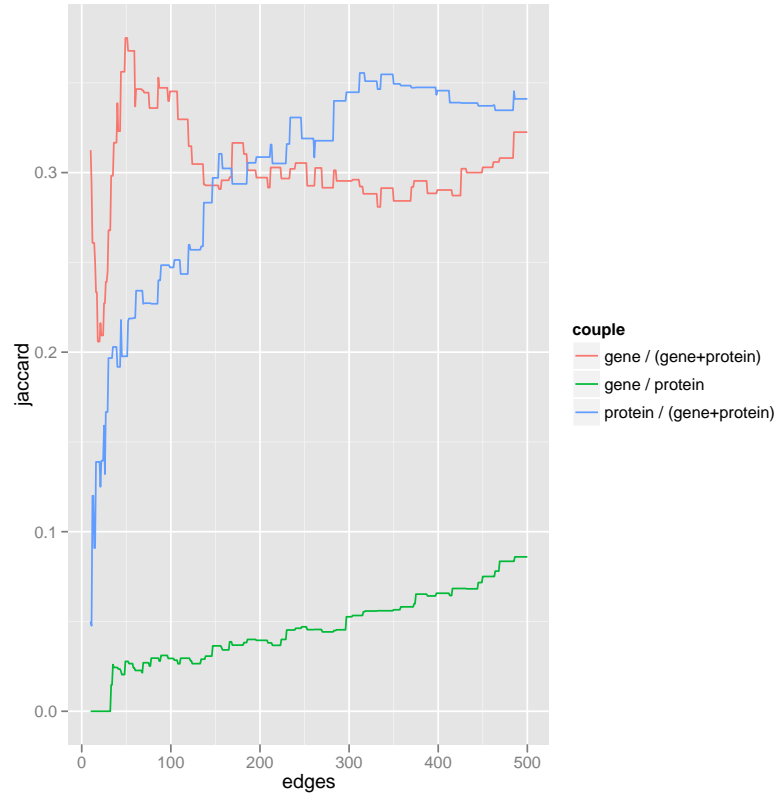


Figure 3: Jaccard's similarity index $J(A, B) = \frac{|A \cap B|}{|A \cup B|}$: multi-attribute network shares a high Jaccard index with both uni-attribute networks.

- F. Bach, R. Jenatton, J. Mairal, and G. Obozinski. Optimization with sparsity-inducing penalties. *Foundations and Trends in Machine Learning*, 4(1):1–106, 2012.
- 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.
- S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, third edition, 2006.
- T. Cai, W. Liu, and X. Luo. A constrained l1 minimization approach to sparse precision matrix estimation. *J. Amer. Statist. Assoc.*, 106:594–607, 2011.
- 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.
- A.P. Dempster. Covariance selection. *Biometrics, Special Multivariate Issue*, 28:157–175, 1972.
- A. Dobra, C. Hans, B. Jones, J. R. Nevins, G. Yao, and M. West. Sparse graphical models for exploring gene expression data. *J. Multivariate Anal.*, 90(1):196–212, 2004.
- M. Drton and M.D. Perlman. Multiple testing and error control in Gaussian graphical model selection. *Statistical Science*, 22:430, 2007.
- M. Drton and M.D. Perlman. A SINful approach to Gaussian graphical model selection. *J. Statist. Plann. Inference*, 138(4):1179–1200, 2008.
- R. Foygel and M. Drton. Extended Bayesian information criteria for Gaussian graphical models. In *Advances in Neural Information Processing Systems (NIPS)*, pages 2020–2028, 2010.
- J. Friedman, T. Hastie, and R. Tibshirani. Sparse inverse covariance estimation with the graphical lasso. *Biostatistics*, 9(3):432–441, 2008.
- C. Giraud. Estimation of Gaussian graphs by model selection. *Electronic Journal of Statistics*, 2:542–563, 2008.
- C. Giraud, S. Huet, and N. Verzelen. Graph selection with GGMselect. *SAGMB*, 11(3):1–50, 2012a.
- C. Giraud, S. Huet, and N. Verzelen. High-dimensional regression with unknown variance. *Statist. Sci.*, 27(4):500–518, 2012b.
- A.-C. Haury, F. Mordelet, P. Vera-Licona, and J.-P. Vert. Tigress: trustful inference of gene regulation using stability selection. *BMC systems biology*, 6(1):145, 2012.
- C.-J. Hsieh, M.A. Sustik, I.S. Dhillon, P. K Ravikumar, and R. Poldrack. Big & quic: Sparse inverse covariance estimation for a million variables. In *Advances in Neural Information Processing Systems (NIPS)*, pages 3165–3173, 2013.

- C.-J. Hsieh, M.A. Sustik, I.S. Dhillon, and P. Ravikumar. Quic: quadratic approximation for sparse inverse covariance estimation. *J. Mach. Learn. Res.*, 15(1):2911–2947, 2014.
- B. Jones, C. Carvalho, A. Dobra, C. Hans, C. Carter, and M. West. Experiments in stochastic computation for high-dimensional graphical models. *Statist. Sci.*, 20(4):388–400, 2005.
- N. Katenka and E.D. Kolaczyk. Inference and characterization of multi-attribute networks with application to computational biology. *Ann. Appl. Stat.*, 6(3):1068–1094, 2012.
- K. Khare, S.-Y. Oh, and B. Rajaratnam. A convex pseudo-likelihood framework for high dimensional partial correlation estimation with convergence guarantees. *J. R. Statist. Soc. B*, 2014.
- H. Kiiveri. Multivariate analysis of microarray data: differential expression and differential connection. *BMC Bioinformatics*, 12(1):42, 2011. ISSN 1471-2105. doi: 10.1186/1471-2105-12-42. URL <http://www.biomedcentral.com/1471-2105/12/42>.
- M. Kolar, H. Liu, and E.P. Xing. Graph estimation from multi-attribute data. *J. Mach. Learn. Res.*, 15(1):1713–1750, 2014.
- S.L. Lauritzen. *Graphical models*, volume 17 of *Oxford Statistical Science Series*. Clarendon Press, New York, 1996. Oxford Science Publications.
- H. Liu, J. Lafferty, and L. Wasserman. The nonparanormal: Semiparametric estimation of high dimensional undirected graphs. *J. Mach. Learn. Res.*, 10:2295–2328, 2009.
- H. Liu, K. Roeder, and L. Wasserman. Stability approach to regularization selection (stars) for high dimensional graphical models. In *Advances in neural information processing systems (NIPS)*, pages 1432–1440, 2010.
- R. Mazumder and T. Hastie. The graphical lasso: New insights and alternatives. *Electronic journal of statistics*, 6:2125.
- N. Meinshausen and P. Bühlmann. High-dimensional graphs and variable selection with the lasso. *Ann. Statist.*, 34(3):1436–1462, 2006.
- N. Meinshausen and P. Bühlmann. Stability selection. *Journal of the Royal Statistical Society, Series B*, 72:417–473, 2010.
- Y. Nesterov. Smooth minimization of non-smooth functions. *Mathematical programming*, 103(1):127–152, 2005.
- J. Peng, P. Wang, N. Zhou, and J. Zhu. Partial correlation estimation by joint sparse regression models. *J. Amer. Statist. Assoc.*, 104(486):735–746, 2009.
- A. Rau, F. Jaffrézic, J.-L. Foulley, and R.W. Doerge. Reverse engineering gene regulatory networks using approximate Bayesian computation. *Statistics and Computing*, 22(6):1257–1271, 2012.

- P. Ravikumar, M. J. Wainwright, and J. Lafferty. High-dimensional Ising model selection using ℓ_1 -regularized logistic regression. *Ann. Stat.*, 38:1287–1319, 2010.
- P. Ravikumar, M. Wainwright, G. Raskutti, and B. Yu. High-dimensional covariance estimation by minimizing ℓ_1 -penalized log-determinant divergence. *Electronic Journal of Statistics*, 5:935–980, 2011.
- G.V. Rocha, P. Zhao, and B. Yu. A path following algorithm for sparse pseudo-likelihood inverse covariance estimation (SPLICE), 2008.
- L. Schwaller, S. Robin, and M. Stumpf. Bayesian inference of graphical model structures using trees. *arXiv preprint arXiv:1504.02723*, 2015.
- T. Sun and C.-H. Zhang. Scaled sparse linear regression. *Biometrika*, page ass043, 2012.
- T. Sun and C.-H. Zhang. Sparse matrix inversion with scaled lasso. *J. Mach. Learn. Res.*, 14(1):3385–3418, 2013.
- N. Verzelen. Minimax risks for sparse regressions: Ultra-high-dimensional phenomena. *Elec. Journal Stat.*, 6:38–90, 2012.
- F. Villers, B. Schaeffer, C. Bertin, and S. Huet. Assessing the validity domains of graphical Gaussian models in order to infer relationships among components of complex biological systems. *Stat. Appl. Genet. Mol. Biol.*, 7(2), 2008.
- H. Wang, A. Banerjee, C.-J. Hsieh, P. K Ravikumar, and I.S. Dhillon. Large scale distributed sparse precision estimation. In C.J.C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems (NIPS)*, pages 584–592. Curran Associates, Inc., 2013. URL <http://papers.nips.cc/paper/5037-large-scale-distributed-sparse-precision-estimation.pdf>
- J. Whittaker. *Graphical models in applied multivariate statistics*. Wiley series in probability and mathematical statistics: Probability and mathematical statistics. Wiley, 1990. ISBN 9780471917502.
- A. Wille and P. Bühlmann. Low-order conditional independence graphs for inferring genetic networks. *Stat. Appl. Genet. Mol. Biol.*, 5(1), 2006.
- D.M. Witten, J.H. Friedman, and N. Simon. New insights and faster computations for the graphical lasso. *J. Comput. Graph. Statist.*, 20(4):892–900, 2011.
- E. Yang, P. Ravikumar, G.I. Allen, and Z. Liu. On poisson graphical models. In *Advances in Neural Information Processing Systems (NIPS)*, pages 1718–1726, 2013.
- M. Yuan and Y. Lin. Model selection and estimation in the Gaussian graphical model. *Biometrika*, 94(1):19–35, 2007.
- P. Zhao and B. Yu. On model selection consistency of Lasso. *J. Mach. Learn. Res.*, 7:2541–2563, 2006.

T. Zhao, H. Liu, K. Roeder, J. Lafferty, and L. Wasserman.
huge: High-dimensional Undirected Graph Estimation, 2014. URL
<http://CRAN.R-project.org/package=huge>. R package version 1.2.6.