

# Functional Sparse Estimation of Time Varying Graphical Model

Meilei Jiang, Yufeng Liu  
Department of Statistics and Operations Research  
University of North Carolina at Chapel Hill

April 12, 2016

## 1 Introduction

Graphical models are quite useful in many domains to uncover the dependence structure among observed variables. Typically, we consider a  $p$ -dimensional multivariate normal distributed random variable

$$\mathbf{X} = (X_1, \dots, X_p) \sim \mathbb{N}(\mathbf{0}, \mathbf{\Sigma}),$$

where  $p$  is the number of features. Then a useful graph of these  $p$  features can be constructed based on there conditional dependence structure. More precisely, we can construct a Gaussian graphical model

$$\mathcal{G} = (V, E), \text{ where } V = \{1, \dots, p\} \text{ is the set of nodes,} \\ \text{and } E = \{(j, l) | X_j \text{ is conditionally dependent with } X_l, \text{ given } X_{V/\{j, l\}}\}.$$

Let  $\mathbf{\Omega} = \mathbf{\Sigma}^{-1} = (\omega_{j,l})_{1 \leq j, l \leq p}$  be the precision matrix. Then  $X_j$  and  $X_l$  are conditionally dependent given other features if and only if  $\omega_{jl} = 0$ . Therefore, estimating the covariance matrix and precision matrix of  $X$  is equivalent to estimate the structure of Gaussian graphical model  $\mathcal{G}$ . More discussion can be found in (Lauritzen, 1996).

### 1.1 Estimation Sparse Precision Matrix $\mathbf{\Omega}$

Given a random sample  $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(n)}$  of  $\mathbf{X}$ , we aim to estimate  $\mathbf{\Omega}$  and recover its support, i.e. the corresponding undirected Gaussian graph. When  $n > p$ , a nature estimator of  $\mathbf{\Omega}$  can be  $\hat{\mathbf{\Omega}}_n = \hat{\mathbf{\Sigma}}_n^{-1}$ , where  $\hat{\mathbf{\Sigma}}_n = \sum_{k=1}^n (\mathbf{X}^{(k)} - \bar{\mathbf{X}}_n)(\mathbf{X}^{(k)} - \bar{\mathbf{X}}_n)'$  is the sample covariance matrix. In the case  $n < p$ , which is quite common in the many applications, the estimation of  $\mathbf{\Omega}$  is much more challenging since  $\hat{\mathbf{\Sigma}}_n$  is no longer invertible.

There are lots of literatures discussing about the estimation of sparse precision matrix  $\mathbf{\Omega}$  in high dimension low sample size settings, i.e.  $p > n$ . Generally speaking, there are three main approaches.

1. **Covariance selection approach.** There is a connection between linear regression and prediction matrix  $\mathbf{\Omega}$ :

$$\mathbf{X}_j = \mathbf{X}_{-j}\boldsymbol{\beta}_j + \boldsymbol{\varepsilon}_j = \sum_{l \neq j} \mathbf{X}_l \beta_{jl} + \boldsymbol{\varepsilon}_j \quad (1)$$

It could be shown that  $\beta_{jl} = \omega_{jl}/\omega_{jj}$ . Thus estimating  $\boldsymbol{\beta}_j$  can identify the support of  $j$ th row of  $\mathbf{\Omega}$ . Meinshausen and Bühlmann (2006) applied LASSO penalty (Tibshirani, 1996) on multivariate regression 1 to estimate the support of  $\mathbf{\Omega}$  row by row. Peng et al. (2012) considered a joint sparse regression to estimate the support of  $\mathbf{\Omega}$  together through an active-shooting algorithm. Yuan (2010) applied Danzig selector (Candes and Tao, 2007) to the problem 1 to estimate each column of  $\mathbf{\Omega}$ .

2. **Penalized likelihood approach.** Another nature way is to estimate  $\mathbf{\Omega}$  is the penalized likelihood approach. The log-likelihood of the parameters in  $\mathbf{\Omega}$  is as following:

$$l(\mathbf{X}^{(i)}, 1 \leq i \leq n | \mathbf{\Omega}) = -\text{tr}(\mathbf{\Omega} \hat{\boldsymbol{\Sigma}}_n) + \log |\mathbf{\Omega}| \quad (2)$$

In order to have a sparse estimation of  $\mathbf{\Omega}$ , different penalized likelihood estimator are considered. Yuan and Lin (2007) proposed the *max-det problem* to solve the LASSO-type estimator.

$$\hat{\mathbf{\Omega}}_L = \arg \min \text{tr}(\mathbf{\Omega} \hat{\boldsymbol{\Sigma}}_n) - \log |\mathbf{\Omega}| + \lambda \|\mathbf{\Omega}\|_1 \quad (3)$$

and the non-negative garrote-type estimator.

$$\begin{aligned} \hat{\mathbf{\Omega}}_G &= \arg \min \text{tr}(\mathbf{\Omega} \hat{\boldsymbol{\Sigma}}_n) - \log |\mathbf{\Omega}| + \lambda \sum_{j \neq i} \frac{\omega_{jl}}{\tilde{\omega}_{jl}} \\ \text{subject to } &\frac{\omega_{jl}}{\tilde{\omega}_{jl}} \geq 0, \mathbf{\Omega} \text{ p.d.} \end{aligned} \quad (4)$$

Banerjee et al. (2008) used a block coordinate descent algorithm to solve (3). And then Friedman et al. (2008) proposed the *graphical lasso* algorithm for (3) based on least square lasso type estimator, which is simple and fast. Rothman et al. (2008) proposed the sparse permutation invariant covariance estimator (SPICE) which could extend to the  $\ell_q$ -type penalized likelihood estimator. Fan et al. (2009), Lam et al. (2009) studied the penalized likelihood estimator with the smoothly clipped absolute deviation (SCAD) penalty and the adaptive LASSO penalty.

3.  **$\ell_1$  constrained minimization approach.** Cai et al. (2011) performed a constrained  $\ell_1$  minimization approach to estimate sparse precision matrix (CLIME).

$$\begin{aligned} \hat{\boldsymbol{\Omega}}_C &= \arg \min \|\boldsymbol{\Omega}\|_1 \\ \text{subject to } &\|\boldsymbol{\Omega}\hat{\boldsymbol{\Sigma}} - \mathbf{I}\|_\infty \leq \lambda_n \end{aligned} \quad (5)$$

Cai et al. (2016) proposed adaptive constrained  $\ell_1$  minimization estimator (ACLIME), which achieved the optimal minimax rate of convergence.

## 1.2 Heterogeneous Data And Time Varying Graphical Model

The methods aforementioned focus on estimating a single Gaussian graph by assuming samples are identically distributed. However, in many applications it is more realistic to assume that data are heterogeneous due to batch effects or latent factors. Guo et al. (2011) reparameterized the off-diagonal entry  $\omega_{jl} = \theta_{jl}\gamma_{jl}^k$  and estimated them through a penalized likelihood with the hierarchical penalty on common structures  $\theta_{jl}$  and individual structure  $\gamma_{jl}^k$ . Danaher et al. (2014) propose the *joint graphical lasso*, to estimate multiple graphical models corresponding to distinct but related conditions. The *joint graphical lasso* utilized fused lasso and group lasso on log-likelihood to force similarity among graphs. Lee and Liu (2015) proposed a method to estimate the common structure and unique structure through the constrained  $\ell_1$  minimization.

In many cases the sample indexes have orders, e.g. time, and the corresponding graphs evolve through the order. In such cases it could be quite interesting to estimate a time varying graphical model.

$$\mathbf{X}(t) \sim \mathbb{N}(\mathbf{0}, \boldsymbol{\Sigma}(t)) \quad (6)$$

Zhou et al. (2010) developed a nonparametric framework for estimating time varying graphical model by kernel smoothing and  $\ell_1$  penalty. Zhou's model assumed that the observations  $\mathbf{X}^t$  are independent and changed smoothly.

$$\begin{aligned} \hat{\boldsymbol{\Omega}}(\tau) &= \arg \min_{\boldsymbol{\Omega}} \left\{ \text{tr}(\boldsymbol{\Omega}\hat{\boldsymbol{\Omega}}(\tau)) - \log |\boldsymbol{\Omega}| + \lambda \|\boldsymbol{\Sigma}^-\|_1 \right\} \\ \text{where } \hat{\boldsymbol{\Omega}}(\tau) &= \sum_i \omega_i^\tau \mathbf{X}^i (\mathbf{X}^i)', \text{ and } \omega_i^\tau = \frac{K_h(t_i - \tau)}{\sum_{i'} K_h(t_{i'} - \tau)} \end{aligned} \quad (7)$$

Lu et al. (2015) proposed a dynamic nonparanormal graphical model, which is more robust, by estimating a weighted Kendall's tau correlation matrix. These approaches generated estimation of similar graphs by weighting the samples among different times.

### 1.3 Varying Coefficient Model And Sparse Derivatives

If we are interested in estimating a dynamic graph which is smoothing over a time region, a functional model which consider each feature as functional data over time is very attractive. Since the structure of Gaussian graph can be recovered through regression approach, we can consider dynamic graphical model under the context of varying coefficient model (Hastie and Tibshirani, 1993; Fan and Zhang, 2008). Typically, we are interested in the basis expansion approach to estimate the varying coefficient models (Huang et al., 2002, 2004).

$$Y(t) = X(t)^T \beta(t) + \varepsilon(t)$$

Then we are looking at the following model

$$X_j(t) = \mathbf{X}'_{-j}(t) \beta_j(t) + \varepsilon_j(t) = \sum_{l \neq j} X_l(t) \beta_{jl}(t) + \varepsilon_j(t). \quad (8)$$

Assume that we collect sample at  $t_1, \dots, t_n$ , denote  $X_j^i = X_j(t_i)$ . Kolar et al. (2009, 2010); Kolar and Xing (2011, 2012) proposed a local linear regression approach with “kernel  $\ell_1$ ” penalty to estimate the smoothly varying graph,

$$\hat{\beta}_j(\tau) = \arg \min_{\beta \in \mathbb{R}^{p-1}} \sum_i (X_j^i - \sum_{l \neq j} X_l^i \beta_l)^2 \omega_i^T + \lambda |\beta|_1 \quad (9)$$

and total variation penalty to estimate graph with jumps.

$$\{\hat{\beta}_j(t_1), \dots, \hat{\beta}_j(t_n)\} = \arg \min_{\beta(t_i), i \leq n} \sum_i (X_j^i - \sum_{l \neq j} X_l^i \beta_l(t_i))^2 + \lambda_1 \sum_i |\beta(t_i)|_1 + \lambda_2 \sum_{i=2}^n |\beta(t_i) - \beta(t_{i-1})|_1 \quad (10)$$

Moreover, in order to estimate a sparse graph, we need to gain sparse functional coefficient  $\beta(t)$ . Based on the idea of FLiRTI in James, Wang and Zhu’s paper James et al. (2009), we put penalty on the derivative matrices of  $\beta(t)$ . Kim et al. (2009) also put the  $l_1$  penalty on derivative of coefficient function in the trend filtering problem. Tibshirani et al. (2014) applied the generalized lasso (Tibshirani and Taylor, 2011) to solve the optimization problem in the trendfilter. We will apply similar algorithm to solve our optimization problem.

## 2 Methodology

### 2.1 Functional Undirected Graph

Consider  $p$  smooth functional continuous variables  $\{X_1(t), X_2(t) \dots, X_p(t)\}$  on the ‘time’ domain  $\mathcal{T}$ . On each  $t$ , we assume

$$(X_1(t), X_2(t) \dots, X_p(t)) \sim \mathcal{N}(\mathbf{0}, \Sigma(t)).$$

Define the undirected graph

$$\begin{aligned} G(t) &= \{V, E(t)\}, \text{ where } V = \{1, \dots, p\}, \\ \text{and } E(t) &= \{(j, l) \in V^2 : \text{Cov}[X_j(t), X_l(t) | X_k(t), k \neq j, l] \neq 0, j \neq l\}. \end{aligned} \quad (11)$$

Namely,  $G(t)$  is the Gaussian graphical model at each  $t$ . This model is quite flexible, which allows  $G(t)$  evolve over time and includes the time dependence. Assume that data are observed at  $t_1, \dots, t_n$  and at each time point  $t$  we have  $n_t$  samples.

## 2.2 Functional Nearest Neighborhood Selection

Consider the following functional linear model

$$X_j^r(t) = \mathbf{X}_{-j}^r(t)^T \boldsymbol{\beta}_j(t) + \varepsilon_j^r(t) = \sum_{l \neq j} X_l^r(t) \beta_{jl}(t) + \varepsilon_j^r(t),$$

$$\text{where } \mathbf{X}_{-j}^r(t) = (X_l^r(t))_{l \neq j} \in \mathbb{R}^{(p-1) \times 1}, r = 1, \dots, n_t, t = t_1, \dots, t_n, j = 1, \dots, p. \quad (12)$$

For each functional coefficient  $\beta_{jl}(t)$ , we consider the basis  $\mathbf{B}_{jl}(t) = (B_{jl1}(t), \dots, B_{jlk_{jl}}(t))$  and then

$$\beta_{jl}(t) = \sum_{s=1}^{k_{jl}} B_{jls}(t) \gamma_{jls} + e_{jl}(t) = \mathbf{B}_{jl}(t) \boldsymbol{\gamma}_{jl}$$

Thus Equation (12) can be represented as

$$X_j^r(t) = \sum_{l \neq j} \sum_{s=1}^{k_{jl}} X_l^r(t) B_{jls}(t) \gamma_{jls} + \tilde{\varepsilon}^r(t),$$

$$\text{where } \tilde{\varepsilon}^r(t) = \sum_{l \neq j} X_l^r(t) e_{jl}^r(t) + \varepsilon_j^r(t), r = 1, \dots, n_t, t = t_1, \dots, t_n, j = 1, \dots, p. \quad (13)$$

As seen in Equation (13), our model is quite flexible since the basis of each functional coefficient can be different.

To get the matrix form of Equation (13), denote

$$\begin{aligned}
\mathbf{B}(t) &= \text{diag}\{\mathbf{B}_{jl}(t)\} \in \mathbb{R}^{(p-1) \times \sum_{l \neq j} k_{jl}}, \\
\mathbf{U}_j^r(t) &= \mathbf{B}(t)^T X_{-j}^r(t) \in \mathbb{R}^{\sum_{l \neq j} k_{jl} \times 1}, \\
\mathbf{U}_j(t) &= (\mathbf{U}_j^1(t), \dots, \mathbf{U}_j^{n_t}(t))^T \in \mathbb{R}^{n_t \times \sum_{l \neq j} k_{jl}}, \\
\mathbf{U}_j &= (\mathbf{U}_j(t_1), \dots, \mathbf{U}_j(t_n))^T \in \mathbb{R}^{\sum_{t=1}^n n_t \times \sum_{l \neq j} k_{jl}}, \\
\mathbf{X}_j(t) &= (X_j^1(t), \dots, X_j^{n_t}(t))^T \in \mathbb{R}^{n_t \times 1}, \\
\mathbf{X}_j &= (X_j(t_1), \dots, X_j(t_n))^T \in \mathbb{R}^{\sum_{t=1}^n n_t \times 1}, \\
\boldsymbol{\varepsilon}_j(t) &= (\varepsilon_j^1(t), \dots, \varepsilon_j^{n_t}(t))^T \in \mathbb{R}^{n_t \times 1}, \\
\boldsymbol{\varepsilon}_j &= (\varepsilon_j(t_1), \dots, \varepsilon_j(t_n))^T \in \mathbb{R}^{\sum_{t=1}^n n_t \times 1}, \\
\boldsymbol{\gamma}_j &= (\gamma_{jl})_{l \neq j} \in \mathbb{R}^{\sum_{l \neq j} k_{jl} \times 1}.
\end{aligned}$$

Then the Equation (13) can be expressed as

$$\mathbf{X}_j = \mathbf{U}_j \boldsymbol{\gamma}_j + \boldsymbol{\varepsilon}_j, j = 1, \dots, p. \quad (14)$$

### 2.3 Control The Sparsity Of Derivatives

In Model (14), we want to estimate a sparse graph and interpretable coefficient functions. For each  $i$  and  $l \neq j$ , we want to control the sparsity of  $\beta_{jl}^{(m)} = \frac{d^m}{dt^m} \beta_{jl}(t) \approx \frac{d^m}{dt^m} \mathbf{B}_{jl}(t)^T \boldsymbol{\gamma}_{jl}$  for some  $m$ . Say assume  $\beta_{jl}^{(0)}(t) = 0$  and  $\beta_{jl}^{(2)}(t) = 0$  in large area, then  $\beta_{jl}(t)$  is zero in many region and linear in the left regions.

Let

$$\mathbf{A}_{jl} = [D^m \mathbf{B}_{jl}(t_1) \quad \dots \quad D^m \mathbf{B}_{jl}(t_n)]^T \in \mathbb{R}^{n \times k_{jl}}, \quad (15)$$

where  $D^m$  is the  $m$ th finite difference operator, i.e.,  $D\mathbf{B}_{jl}(t_k) = [\mathbf{B}_{jl}(t_k) - \mathbf{B}_{jl}(t_{k-1})]/[t_k - t_{k-1}]$ ,  $D^2\mathbf{B}_{jl}(t_k) = [D\mathbf{B}_{jl}(t_k) - D\mathbf{B}_{jl}(t_{k-1})]/[t_k - t_{k-1}]$ , etc.

Next, set

$$\boldsymbol{\eta}_{jl} = \mathbf{A}_{jl} \boldsymbol{\gamma}_{jl} \in \mathbb{R}^{n \times 1} \quad (16)$$

Then  $\boldsymbol{\eta}_{jl} \approx (\beta_{jl}^{(m)}(t_k))_{1 \leq k \leq n}$ . Moreover, we denote

$$\boldsymbol{\eta}_j = (\boldsymbol{\eta}_{jl})_{l \neq j} = \mathbf{A}_j \boldsymbol{\gamma}_j \in \mathbb{R}^{n(p-1)}, \quad (17)$$

where  $\mathbf{A}_j = \text{diag}(\mathbf{A}_{jl})_{l \neq j} \in \mathbb{R}^{n(p-1) \times \sum_{l \neq j} k_{jl}}$ .

We want to put the sparsity penalty on the  $\boldsymbol{\eta}_j$ . Then Model 14 can be expressed as the following optimization problem, which is a generalized lasso problem.

$$\begin{aligned}
\hat{\boldsymbol{\gamma}}_{j,L} &= \arg \min_{\boldsymbol{\gamma}_j} \frac{1}{2} \|\mathbf{X}_j - \mathbf{U}_j \boldsymbol{\gamma}_j\|_2^2 \\
\text{subject to } &\|\boldsymbol{\eta}_j\|_1 = \|\mathbf{A}_j \boldsymbol{\gamma}_j\|_1 \leq t \\
\text{i.e. } \hat{\boldsymbol{\gamma}}_{j,L} &= \arg \min_{\boldsymbol{\gamma}_j} \frac{1}{2} \|\mathbf{X}_j - \mathbf{U}_j \boldsymbol{\gamma}_j\|_2^2 + \lambda \|\mathbf{A}_j \boldsymbol{\gamma}_j\|_1
\end{aligned} \quad (18)$$

Moreover, if we want to control the sparsity of multiple derivatives of  $\beta_{jl}(t)$ , say we want both  $\beta_{jl}^{(0)}(t) = 0$  and  $\beta_{jl}^{(2)}(t) = 0$  in large area.

There is a connection between Model 18 and fused lasso.

## References

- Banerjee, O., El Ghaoui, L., and d’Aspremont, A. (2008), “Model Selection Through Sparse Maximum Likelihood Estimation for Multivariate Gaussian or Binary Data,” *Journal of Machine Learning Research*, 9, 485–516.
- Cai, T., Liu, W., and Luo, X. (2011), “A constrained  $\ell_1$  minimization approach to sparse precision matrix estimation,” *Journal of the American Statistical Association*, 106, 594–607.
- Cai, T. T., Liu, W., and Zhou, H. H. (2016), “Estimating sparse precision matrix: Optimal rates of convergence and adaptive estimation,” *The Annals of Statistics*, 44, 455–488.
- Candes, E. and Tao, T. (2007), “The Dantzig selector: statistical estimation when  $p$  is much larger than  $n$ ,” *The Annals of Statistics*, 2313–2351.
- Danaher, P., Wang, P., and Witten, D. M. (2014), “The joint graphical lasso for inverse covariance estimation across multiple classes,” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 76, 373–397.
- Fan, J., Feng, Y., and Wu, Y. (2009), “Network exploration via the adaptive LASSO and SCAD penalties,” *The annals of applied statistics*, 3, 521.
- Fan, J. and Zhang, W. (2008), “Statistical methods with varying coefficient models,” *Statistics and its Interface*, 1, 179.
- Friedman, J., Hastie, T., and Tibshirani, R. (2008), “Sparse inverse covariance estimation with the graphical lasso,” *Biostatistics*, 9, 432–441.
- Guo, J., Levina, E., Michailidis, G., and Zhu, J. (2011), “Joint estimation of multiple graphical models,” *Biometrika*, asq060.
- Hastie, T. and Tibshirani, R. (1993), “Varying-coefficient models,” *Journal of the Royal Statistical Society. Series B (Methodological)*, 757–796.
- Huang, J. Z., Wu, C. O., and Zhou, L. (2002), “Varying-coefficient models and basis function approximations for the analysis of repeated measurements,” *Biometrika*, 89, 111–128.
- (2004), “Polynomial spline estimation and inference for varying coefficient models with longitudinal data,” *Statistica Sinica*, 763–788.

- James, G. M., Wang, J., and Zhu, J. (2009), “Functional linear regression that’s interpretable,” *The Annals of Statistics*, 2083–2108.
- Kim, S.-J., Koh, K., Boyd, S., and Gorinevsky, D. (2009), “ $\ell_1$  Trend Filtering,” *SIAM review*, 51, 339–360.
- Kolar, M., Song, L., Ahmed, A., and Xing, E. P. (2010), “Estimating time-varying networks,” *The Annals of Applied Statistics*, 94–123.
- Kolar, M., Song, L., and Xing, E. P. (2009), “Sparsistent learning of varying-coefficient models with structural changes,” in *Advances in Neural Information Processing Systems*, pp. 1006–1014.
- Kolar, M. and Xing, E. P. (2011), “On time varying undirected graphs,” *Journal of Machine Learning Research*, 15, 407–415.
- (2012), “Estimating networks with jumps,” *Electronic journal of statistics*, 6, 2069.
- Lam, C., Fan, J., et al. (2009), “Sparsistency and rates of convergence in large covariance matrix estimation,” *The Annals of Statistics*, 37, 4254–4278.
- Lauritzen, S. L. (1996), *Graphical models*, Clarendon Press.
- Lee, W. and Liu, Y. (2015), “Joint Estimation of Multiple Precision Matrices with Common Structures,” *Journal of Machine Learning Research*, 16, 1035–1062.
- Lu, J., Kolar, M., and Liu, H. (2015), “Post-regularization Inference for Dynamic Nonparanormal Graphical Models,” *arXiv preprint arXiv:1512.08298*.
- Meinshausen, N. and Bühlmann, P. (2006), “High-dimensional graphs and variable selection with the lasso,” *The annals of statistics*, 1436–1462.
- Peng, J., Wang, P., Zhou, N., and Zhu, J. (2012), “Partial correlation estimation by joint sparse regression models,” *Journal of the American Statistical Association*.
- Rothman, A. J., Bickel, P. J., Levina, E., Zhu, J., et al. (2008), “Sparse permutation invariant covariance estimation,” *Electronic Journal of Statistics*, 2, 494–515.
- Tibshirani, R. (1996), “Regression shrinkage and selection via the lasso,” *Journal of the Royal Statistical Society. Series B (Methodological)*, 267–288.
- Tibshirani, R. J. and Taylor, J. (2011), “THE SOLUTION PATH OF THE GENERALIZED LASSO,” *The Annals of Statistics*, 1335–1371.
- Tibshirani, R. J. et al. (2014), “Adaptive piecewise polynomial estimation via trend filtering,” *The Annals of Statistics*, 42, 285–323.



- Yuan, M. (2010), “High dimensional inverse covariance matrix estimation via linear programming,” *The Journal of Machine Learning Research*, 11, 2261–2286.
- Yuan, M. and Lin, Y. (2007), “Model selection and estimation in the Gaussian graphical model,” *Biometrika*, 94, 19–35.
- Zhou, S., Lafferty, J., and Wasserman, L. (2010), “Time varying undirected graphs,” *Machine Learning*, 80, 295–319.