

STATISTICS 593C: Spring, 2007

Model Selection and Regularization

Jon A. Wellner

Lecture 16 (May 17): In this (half) lecture I will briefly discuss the paper by Meinshausen and Bühlmann (2006) on covariance selection in Gaussian graphical models via the lasso

Suppose that $\underline{X} = (X_1, \dots, X_p) \sim N_p(\mu, \Sigma)$. Consider a graph $\mathcal{G} = (\Gamma, E)$ associated with the random vector X as follows:

$\Gamma = \Gamma(n) = \{1, \dots, p = p_n\}$ is the set of *nodes* in the graph,

$E \subset \Gamma \times \Gamma$ is the set of *edges* in the graph, and

$(a, b) \in E$ if and only if $X_a \not\perp X_b \mid \{X_k : k \in \Gamma \setminus \{a, b\}\}$.

Thus if $(a, b) \in E^c$ it follows that

$$X_a \perp X_b \mid \{X_k : k \in \Gamma \setminus \{a, b\}\},$$

and this corresponds to a 0 in Σ^{-1} . (Here $X_a \perp X_b \mid W$ means “ X_a and X_b are conditionally independent given the collection of random variables W .”)

Now suppose that $\underline{X}_1, \dots, \underline{X}_n$ are i.i.d. as \underline{X} . The goal of “covariance selection” as introduced by Dempster (1963) is to estimate (or “discover” or “learn”) the graph \mathcal{G} from the data. See e.g. Buhl (1993) for methods based on maximum likelihood and see Heckerman, Chickering, Meek, Rounthwaite, and Kadie (2000). The goal of Meinshausen and Bühlmann is to use the lasso to study covariance selection when $p = n^\alpha$ with (possibly) $\alpha > 1$.

Define

$ne_a =$ smallest subset of $\Gamma \setminus \{a\}$ such that

X_a is conditionally independent of all other variables ;

that is,

$$X_a \perp \{X_k : k \in \Gamma \setminus cl_a\} \mid X_{ne_a}$$

where $cl_a \equiv ne_a \cup \{a\}$.

Assumptions:

A1 $p = O(n^\alpha)$ for some $0 < \alpha < \infty$.

A2 $Var(X_a) = 1$ for all $a \in \Gamma$ and there exists a constant $v^2 > 0$ such that $Var(X_a \mid X_{\Gamma \setminus \{a\}}) \geq v^2$ for all $a \in \Gamma$.

A3 $\max\{|\text{ne}_a| : a \in \Gamma\} = O(n^\kappa)$ for some $0 \leq \kappa < 1$. (This is a *sparsity assumption*.)

Before we introduce the rest of the assumptions needed, we define the lasso type procedures introduced by Meinshausen and Bühlmann (2006). For an arbitrary subset $\mathcal{A} \subset \Gamma$, let

$$\begin{aligned}\theta^{a,\mathcal{A}} &\equiv \operatorname{argmin}_{\theta: \theta_k=0, k \notin \mathcal{A}} E \left(X_a - \sum_{k \in \Gamma} \theta_k X_k \right)^2, \\ \theta^{a,\Gamma \setminus \{a\}} &\equiv \theta^a = (\theta_b^a, b \in \Gamma).\end{aligned}$$

Then it turns out that

$$\theta_b^a = -\frac{\Sigma_{a,b}^{-1}}{\Sigma_{aa}^{-1}}.$$

Therefore $\text{ne}_a = \{b \in \Gamma : \theta_b^a \neq 0\}$.

This suggests defining the lasso estimator of θ^a and hence the lasso estimator of ne_a as follows:

$$\begin{aligned}\hat{\theta}^{a,\lambda} &\equiv \operatorname{argmin}_{\theta: \theta_a=0} \{n^{-1} \|X_{n,a} - \underline{X}_n \theta\|_2^2 + \lambda \|\theta\|_1\}, \\ \hat{\text{ne}}_a^\lambda &\equiv \{b \in \Gamma : \hat{\theta}_b^{a,\lambda} \neq 0\}.\end{aligned}$$

where $\underline{X}_n \equiv \sum_1^n X_i$.

Before stating the main result of Meinshausen and Bühlmann, we need three more assumptions:

Assumptions, continued:

A4 $\|\theta^{a,\text{ne}_b \setminus \{a\}}\|_1 \leq \eta$ for all $a, b \in \Gamma$ such that $(a, b) \in E$.

A5 If $\pi_{a,b}$ denotes the partial correlation between X_a and X_b (i.e. the correlation after correcting for linear effects from all the other variables), then

$$|\pi_{a,b}| \geq \delta n^{-(1-\xi)/2} \quad \text{for all } (a, b) \in E$$

for some $\delta > 0$ and $\xi > \kappa$.

A6 (Neighborhood stability) There is a $\delta < 1$ such that for all $a, b \in \Gamma$ with $b \notin \text{ne}_a$,

$$S_a(b) \equiv \sum_{k \in \text{ne}_a} \text{sign}(\theta_k^{a,\text{ne}_a}) \theta_b^{k,\text{ne}_a}$$

satisfies $|S_a(b)| < \delta$.

With this preparation we can state the main result of M & B (2006):

Theorem. Suppose that A1-A6 hold. Let $\lambda_n \sim dn^{-(1-\epsilon)/2}$ with some $\kappa < \epsilon < \xi$ and $d > 0$. Then there exists a $c > 0$ such that for all $a \in \Gamma$,

$$\begin{aligned}P(\hat{\text{ne}}_a^\lambda \subset \text{ne}_a) &= 1 - O(e^{-cn^\epsilon}) \quad \text{as } n \rightarrow \infty, \quad \text{and} \\ P(\text{ne}_a \subset \hat{\text{ne}}_a^\lambda) &= 1 - O(e^{-cn^\epsilon}) \quad \text{as } n \rightarrow \infty.\end{aligned}$$