

# Simultaneous Supervised Clustering and Feature Selection over a Graph

Xiaotong Shen

University of Minnesota

International Workshop on Perspectives on  
High-Dimensional Data Analysis, June 9-11, 2011

Joint with Hsin-Cheng Huang, Academia Sinica, and Wei  
Pan, U of Minnesota.

# Outline

- 1 Supervised clustering and Feature Selection
- 2 Theory
- 3 Numerical examples

# Introduction

- **Response:**  $\mathbf{Y} = (Y_1, \dots, Y_n)^T$
- **Predictors:**  $p$ -dimensional  $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})^T$ .
- Regression model

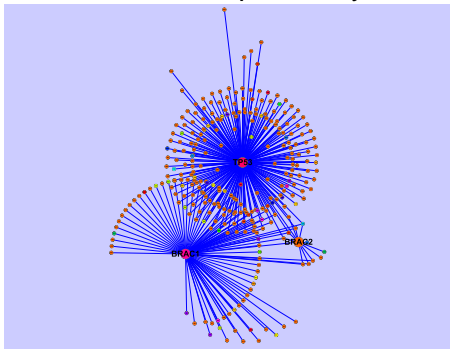
$$Y_i = \mathbf{x}_i^T \boldsymbol{\beta} + \varepsilon_i; \quad \varepsilon_i \sim N(0, \sigma^2); i = 1, \dots, n, \quad (1)$$

where each predictor corresponds to a node in a given undirected graph, and **an edge of the graph indicates possible clustering between two predictors.**

- Identifying two kinds low-dimensional structures simultaneously:
  - **Supervised clustering:** Estimate homogenous and collapsing clusters of predictors.
  - **Feature selection:** Estimate nonzero coefficients of predictors.

# Motivating example

- Identifying subnetworks relevant to breast cancer survival.
- $Y = \log$  survival time,  $X$ : Clinical variables as well as gene expression profiles.
- Gene network: protein-protein interaction network (Chang, et al., 07), and describes dependency structure of genes.



# Simultaneous supervised clustering & feature selection

- **Homogeneity:** Partition  $\{1, \dots, p\}$  into clusters:

$$\beta = (\beta_1, \dots, \beta_p)' \approx (0 \mathbf{1}_{|\mathcal{G}_0|}, \alpha_1 \mathbf{1}_{|\mathcal{G}_1|}, \dots, \alpha_K \mathbf{1}_{|\mathcal{G}_K|})'.$$

- **Goal:** Over the **graph**, estimate true  $\mathcal{G}^0 = (\mathcal{G}_0^0, \mathcal{G}_1^0, \dots, \mathcal{G}_{K_0}^0)'$   
&  $\alpha^0 = (0, \alpha_1^0, \dots, \alpha_{K_0}^0)'$ .

- **Benefits**

- **Structure:** Explore **sparseness & clustering** by leveraging dependency structures given by a graph.
- **Estimation:** Higher accuracy is due to **variance reduction**.
- **Selection:** Overcome **feature selection instability** by grouping and collapsing highly positively correlated predictors, & remove **redundant** clusters by feature selection. Higher accuracy for both.
- **Interpretability:** Simpler model with higher predictive power.

# Objectives and Challenges

- Objectives

- Reconstructing biased OLS based on  $\mathcal{G}^0$ .
- Accurate identification of clusters & parameter estimation/prediction.
- Developing an efficient computational algorithm for large problems.

- Challenges

- More difficult than the problem of feature selection alone & supervised clustering alone.
- Complexity for enumeration over a complete graph is the

Bell number: 
$$B(p) = \frac{1}{e} \sum_{k=0}^{\infty} \frac{k^p}{k!} \approx O\left(e^{e^{p^a}}\right).$$

# Literature

- Graph:
  - Fused Lasso(TSRZK, 05):
$$\lambda_1 \sum_{j=1}^p |\beta_j| + \lambda_2 \sum_{j=1}^p |\beta_j - \beta_{j+1}|. \text{ (QP).}$$
  - LL(08):  $\lambda_1 \sum_{j=1}^p |\beta_j| + \lambda_2 \sum_{j \sim j'}^p \left(\frac{\beta_j}{w_j} - \frac{\beta_{j+1}}{w_{j+1}}\right)^2. \text{ (QP).}$
  - TT(11): Glasso:  $\lambda_1 \sum_{j=1}^p |\beta_j| + \lambda_2 \sum_{j \sim j'}^p |\beta_j - \beta_{j'}|. \text{ (Homotopy).}$
- Non-graph:
  - SH (10)  $\lambda \sum_{j,j'=1}^p J(|\beta_j - \beta_{j'}|). \text{ (Homotopy)}$
  - JKLDY (11)  $\lambda_1 \sum_{j=1}^p |\beta_j| + \lambda_2 \sum_{j,j'=1}^p |\beta_j - \beta_{j'}|. \text{ (QP).}$
- Other types: Clustering in sizes: BR (08), XPS(09). Not a Glasso problem.
- Huge literature for feature selection and encouraging grouping in feature selection.....

# Challenges

- Computation: For large problems, QP is infeasible, and a homotopy method may be inefficient. **Coordinate decent method** breaks down even for F-lasso (special algorithm). Need efficient methods for large  $p$  over an arbitrary undirected graph
- Theory: F-lasso: Rinaldo (2009); clustering: SH (10). Lack of theory to guide practice.

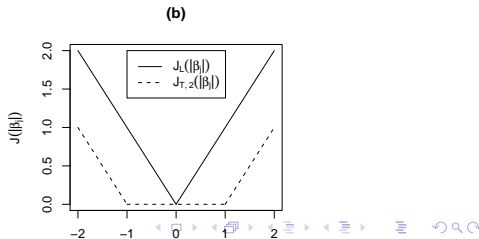
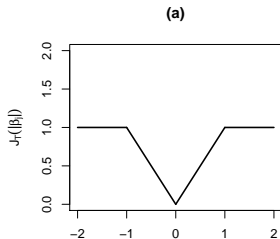


# Constrained Least Squares

- Constrained least squares criterion over a graph:

$$\frac{1}{2n} \sum_{i=1}^n (Y_i - \mathbf{x}_i^T \beta)^2, \quad \text{subj to } \sum_{j=1}^p J(|\beta_j|) \leq s_1,$$
$$\sum_{j < j': (j, j') \in \mathcal{E}} J(|\beta_j - \beta_{j'}|) \leq s_2,$$

- Surrogate of the  $L_0$ :  $J(z) = \min\left(\frac{|z|}{\lambda_3}, 1\right)$ .
- Tuning:  $(s_1, s_2, \lambda_3)$ . Clustering:  $(s_1, s_2)$ ; Threshold:  $\lambda_3 > 0$ .



# Nonconvex minimization

- Theorem: A global minimizer of constrained LS is a local minimizer of  $S(\beta) = (2n)^{-1} \sum_{i=1}^n (Y_i - \mathbf{x}_i^T \beta)^2 + \lambda_1 \sum_{j=1}^p J(|\beta_j|) \lambda_2 \sum_{j < j': (j, j') \in \mathcal{E}} J(|\beta_j - \beta_{j'}|)$ , where  $\lambda_j \rightarrow s_j$ ;  $j = 1, 2$ .
- Local optimality:  $j = 1, \dots, p$ ,

$$-(\mathbf{x}^{(j)})^T (\mathbf{Y} - \mathbf{X}\beta) + \frac{\lambda_1}{\lambda_3} b_j + \frac{\lambda_2}{\lambda_3} \sum_{j': (j, j') \in \mathcal{E}} b_{jj'} = 0, \quad (2)$$

where  $b_j = \text{sign}(\beta_j) I(|\beta_j| < \lambda_3)$  if  $\beta_j \neq 0$ ,  
 $b_{jj'} = \text{sign}(\beta_j - \beta_{j'}) I(|\beta_j - \beta_{j'}| < \lambda_3)$  if  $\beta_j - \beta_{j'} \neq 0$ , &  $b_j = \emptyset$  if  $|\beta_j| = \lambda_3$ , &  $b_{jj'} = \emptyset$  if  $|\beta_j - \beta_{j'}| = \lambda_2$ , are the regular subdifferentials of  $\min(|\beta_j|, \lambda_3)$  &  $\min(|\beta_j - \beta_{j'}|, \lambda_3)$  at  $\beta_j$ .

- **Strategy:** Obtaining a local minimizer of (2) through DC programming, the augmented Lagrangian and coordinate decedent methods.

# Difference Convex Programming

- Decompose  $S = S_1 - S_2$  into a difference of two convex functions.
- Construct a sequence of **upper convex approximations** iteratively by replacing  $S_2$  at iteration  $m$ , by its affine minorization at iteration  $m - 1$ :

$$S^{(m)}(\beta) = (2n)^{-1} \sum_{i=1}^n (Y_i - \mathbf{x}_i^T \beta)^2 + \lambda_1 \sum_{j:j \in \mathcal{F}^{(m-1)}} |\beta_j| + \lambda_2 \sum_{j < j': (j,j') \in \mathcal{E}^{(m-1)}} |\beta_j - \beta_{j'}|, \quad (3)$$

- $\mathcal{E}^{(m-1)} = \{(j, j') \in \mathcal{E}, |\hat{\beta}_j^{(m-1)} - \hat{\beta}_{j'}^{(m-1)}| < \lambda_3\}$  &  
 $\mathcal{F}^{(m-1)} = \{j \in \mathcal{F}, |\hat{\beta}_j^{(m-1)}| < \lambda_3\}$ .
- $\hat{\beta}^{(m-1)}$  is the minimizer of  $S^{(m-1)}(\beta)$ .
- $\hat{\beta}$ : DC estimate  $\hat{\beta}^{(m)}$  after convergence.

## Solution of (3)

- Major challenges:
  - (Stationary points) Coordinate decent method fails for (3).
  - (Graph structure+overcomplete) A graph can be arbitrary.
- Efficient method: Introduce slack variables  $\beta_{jj'} = \beta_j - \beta_{j'}$  for  $j \neq j'$  for an equivalent **augmented** problem of (3) in  $\zeta = (\beta_1, \dots, \beta_p, \beta_{12}, \dots, \beta_{1p}, \dots, \beta_{(p-1)p})^T$ :

$$\begin{aligned} \tilde{S}^{(m)}(\zeta) = & (2n)^{-1} \sum_{i=1}^n (Y_i - \mathbf{x}_i^T \beta)^2 + \lambda_2 \sum_{j < j' : (j, j') \in \mathcal{E}^{(m-1)}} |\beta_{jj'}| \\ & + \lambda_1 \sum_{j: j \in \mathcal{F}^{(m-1)}} |\beta_j|, \end{aligned} \quad (4)$$

subj to linear constraints  $\mathbf{A}\zeta = \mathbf{0}$ .

# Augmented Lagrangian + coordinate decent for (4)

- **Augmented Lagrangian for (4)** by solving its unconstrained version iteratively: At iteration  $t$ , minimize

$$\begin{aligned}\bar{S}^{(m)}(\zeta) = & \tilde{S}^{(m)}(\beta) + \sum_{j < j': (j, j') \in \mathcal{E}} \tau_{jj'}^{(t)} (\beta_j - \beta_{j'} - \beta_{jj'}) \\ & + \frac{1}{2} \sum_{j < j': (j, j') \in \mathcal{E}} \nu_{jj'}^{(t)} (\beta_j - \beta_{j'} - \beta_{jj'})^2,\end{aligned}\quad (5)$$

- $(\tau_{jj'}^{(t)}, \nu_{jj'}^{(t)})$  are Lagrangian multipliers for  $\mathbf{A}\zeta = \mathbf{0}$  and for expediting convergence.
- Solve (5) through analytic updating formula and coordinate decent method.

# Path-following Algorithm

For given  $\lambda_3$ ,

- Step 1** (Initialization) Specify evaluation points for tuning parameters. Supply a good initial estimate  $\hat{\beta}^{(0)}$ . Set tolerance error for convergence,  $\tau_{jj'}^{(0)} = 1$  &  $\nu_{jj'}^{(0)} = 1$ .
- Step 2** (Iteration) Iteration begins with  $m = 1$ . At iteration  $m$ , compute  $\hat{\beta}^{(m)}$  by solving (5) through coordinate descent over active sets.
- Step 3** (Stopping) Terminate when  $S(\hat{\beta}^{(m-1)}) - S(\hat{\beta}^{(m)}) \leq 0$ . The estimate  $\hat{\beta} = \hat{\beta}^{(m_0)}$ , where  $m_0$  is the termination index.

# Computational properties

- DC programming converges **fast** and **finitely**. This is due to the three non-differentiable points of  $J(\cdot)$ .
- The augmented Lagrangian method converges super-linearly.
- The coordinate descent method is efficient when integrated with the augmented Lagrangian method.
- Can handle a problem of size  $p = 3000 - 4000$  easily—complexity for constraint terms is order of  $p^2$ .

# Notation

- $C_{\min} : \inf_{\mathcal{G} \neq \mathcal{G}^0} \frac{1}{n} \|(I - \mathbf{P}_{\mathcal{G} \setminus \mathcal{G}_0^0}) \mathbf{X}_{\mathcal{G} \setminus \mathcal{G}_0^0} \beta_{\mathcal{G} \setminus \mathcal{G}_0^0}\|^2$  for  $\mathcal{G}$  induced by  $\mathcal{E}$ ,  $\mathbf{P}$  is the projection for collapsed predictors over  $\mathcal{G} \setminus \mathcal{G}_0^0$ .  
 $C_{\min}$ : describes the least favorable situation in the KL-loss.
- $\mathbf{X}_{\mathcal{G} \setminus \mathcal{G}_0^0}$  &  $\beta_{\mathcal{G} \setminus \mathcal{G}_0^0}$ : design matrix of predictors & coefficient vector over  $\mathcal{G} \setminus \mathcal{G}_0^0$ ,  $\|\cdot\|$  is the Eucli-norm in  $\mathcal{R}^n$ .
- Oracle estimator  $\hat{\beta}^{ols}$ :  
 $(\hat{\beta}_1^{ols}, \dots, \hat{\beta}_p^{ols})^T = (0_{|\mathcal{G}_0^0|}, \hat{\alpha}_1^{ols} \mathbf{1}_{|\mathcal{G}_1^0|}, \dots, \hat{\alpha}_{K^0}^{ols} \mathbf{1}_{|\mathcal{G}_{K^0}^0|})^T$  given  $\mathcal{G}^0$ ;  $\hat{\alpha}^{ols} \equiv (\hat{\alpha}_1^{ols}, \dots, \hat{\alpha}_{K^0}^{ols})^T = (\mathbf{X}_{\mathcal{G}^0 \setminus \mathcal{G}_0^0}^T \mathbf{X}_{\mathcal{G}^0 \setminus \mathcal{G}_0^0})^{-1} \mathbf{X}_{\mathcal{G}^0 \setminus \mathcal{G}_0^0}^T \mathbf{Y}$ .
- Graph parameters:  $\bar{K}$ —max # clusters allowed;  $S^*$ —number of possible distinct clusters, for the given graph.



# Global minimizer

Let  $\hat{\beta}^{gl}$  be a global minimizer of constrained LS. Let  $p_0$  be # non-zero predictors.

## Theorem

Assume path connectivity for  $j, j' \in \mathcal{G}_k^0$  over  $\mathcal{E}$ . If

$$(s_1, s_2) = (p - p_0, \sum_{(j, j') \in \mathcal{E}} I(|\beta_j^0 - \beta_{j'}^0| \neq 0)),$$

$$\lambda_3 \leq 2\sigma \sqrt{\frac{\log p}{2np^3 \lambda_{\max}(\mathbf{X}^T \mathbf{X})}}, \text{ then}$$

$$P(\hat{\beta}^{gl} \neq \hat{\beta}^{ols}) \leq \exp\left(-\frac{n}{10\sigma^2} \left(C_{\min} - 10\sigma^2 \frac{2\log p + \bar{K} + 2\log S^*}{n}\right)\right).$$

Under condition:  $C_{\min} \geq d_1 \sigma^2 \frac{2\log p + \bar{K} + 2\log S^*}{n}$  for  $d_1 > 10$ , there exist tuning parameter values such that **oracle properties (A)-(D) hold**.

## Global minimizer—continued

Estimate:  $\hat{\beta} = \hat{\beta}_{\hat{\mathcal{G}}}$ ;  $\hat{\mathcal{G}} = (\hat{\mathcal{G}}_0 = \mathbf{0}, \dots, \hat{\mathcal{G}}_K)$ . Truth:  $\beta^0 = \beta_{\mathcal{G}^0}^0$ ,  $\mathcal{G}^0 = (\mathcal{G}_0^0 = \mathbf{0}, \dots, \mathcal{G}_{K^0}^0)$ . As  $n, p \rightarrow \infty$ ,

(A) (Clustering consistency)  $P(\hat{\mathcal{G}} \neq \mathcal{G}^0) \rightarrow 0$ .

(B) (Parameter estimation) For any  $\beta^0$ ,

$$n^{-1} E \|\mathbf{X}(\hat{\beta} - \beta^0)\|^2 = \sigma^2 \frac{K_0}{n}.$$

(C) (Normality) With  $W_{\hat{A}} = \sigma^2 (\mathbf{X}_{\hat{\mathcal{G}} \setminus \hat{\mathcal{G}}^0}^T \mathbf{X}_{\hat{\mathcal{G}} \setminus \hat{\mathcal{G}}^0})^{-1}$ ;  $I_{\mathcal{G}_0 \setminus \mathcal{G}_0^0}$  the identity matrix

$$W_{\hat{\mathcal{G}} \setminus \hat{\mathcal{G}}^0}^{-1/2} (\hat{\beta}_{\hat{\mathcal{G}} \setminus \hat{\mathcal{G}}^0} - \beta_{\hat{\mathcal{G}} \setminus \hat{\mathcal{G}}^0}^0) \sim N(0, I_{\mathcal{G}_0 \setminus \mathcal{G}_0^0}).$$

(D) (Uniformity) (A)-(C) hold uniformly over

$B_0(u, l) \equiv \{\beta \in \mathcal{R}^p : K \leq u, C_{\min} \geq l\}$ : a  $L_0$ -ball of radius  $u > 0$  & resolution level  $l > 0$ . Then  $\hat{\beta}$  is asym minimax.

## Comments and Local minimizer

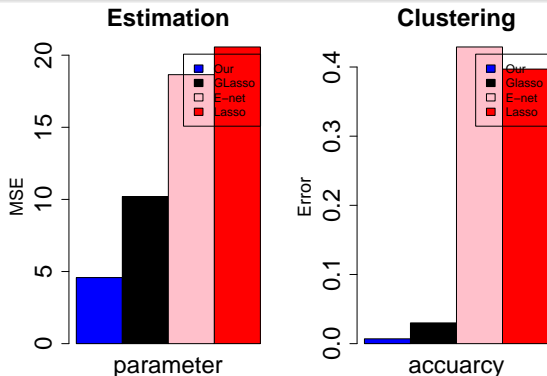
- $\bar{K}$  and  $S^*$  need to be computed for a given graph. Small graph: Fused,  $\bar{K} \leq K_0$ ;  $S^* \leq p_0^{K_0}$ ,  $K_0$  is the number of true clusters.
- Under smaller conditions, any local minimizer, particularly the one computed from the algorithm, has oracle properties (A)-(D).

# Numerical examples

- Ex:** (Graph, Li & Li, 08). Consider a network consisting of 200 subnetworks, each with one transcription factor (TF) and its 10 regulatory target genes. In (1), a predictor of each target gene and the TF follows a bivariate normal distribution with correlation  $\rho = 0.7$ , and target genes are independent  $N(0, 1)$ 's, conditional on the TF;  $n = 100$ ,  $p = 2200$ ,  $\sigma^2 = \sum_{j=1}^p \beta_j^2 / 2$ .
- $\beta =$ 

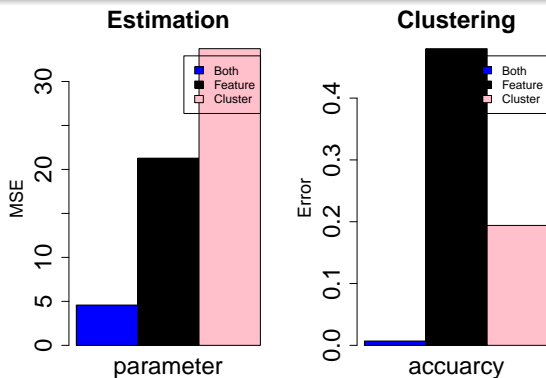
$$\left( 5, \underbrace{\frac{5}{\sqrt{10}}, \dots, \frac{5}{\sqrt{10}}}_{10}, -5, \underbrace{\frac{-5}{\sqrt{10}}, \dots, \frac{-5}{\sqrt{10}}}_{10}, 3, \underbrace{\frac{3}{\sqrt{10}}, \dots, \frac{3}{\sqrt{10}}}_{10}, -3, \underbrace{\frac{-3}{\sqrt{10}}, \dots, \frac{-3}{\sqrt{10}}}_{10}, \underbrace{0, \dots, 0}_{p-44} \right)^T.$$
- Mean squares error:** averaged over 100 replications.
- Selection error & grouping error:** % committing an error.
- Tuning parameters:** are estimated by minimizing MSE over a set of grid points.
- Comparison:** GLasso, Elastic Net, & Lasso.

## Comparison: Example



- Our method outperforms in both parameter estimation and accuracy of selection. Interestingly, Elastic net performs better than Lasso in parameter estimation but worst than selection.
- The average number of DC iterations is about 4.

# Simultaneous supervised clustering & feature selection



- Simultaneous supervised clustering & feature selection performs better than either one alone. Perform relatively better when  $p$  gets large.

# Network-based eQTL analysis

- Goal: Identify genomic loci (Expression quantitative trait loci , called eQTLs) linked to gene expression traits. Improve power in detecting eQTLs for a group of co-regulated genes.
- Mouse dataset in Lan et al. (06): 60 F2 mice from B6 and BTBR founder strains, where the B6 and BTBR strains are diabetes resistant and non-resistant, respectively. About 45000 gene expression traits are measured with genotypes of 194 markers distributed across the mouse genome with an average marker interval of approximately 10cm.

# Network-based eQTL analysis-continued

- Model:  $\mathbf{Y}_g$  : expression of gene  $g$  of 60 mice, linking to  $\mathbf{X}_0$ :

$$\mathbf{Y}_g = \mathbf{X}_0 \beta_g + \epsilon_g; \quad \epsilon_g \sim N(0, \sigma^2 \mathbf{I}); \quad g = 1, \dots, G, \quad (6)$$

where  $\mathbf{X}_0$  is genotypes of 194 markers across 60 mouse genomes, with each genotype taking  $-1, 0, 1$  indicating one of three alleles.

- Estimation of nonzero coefficients  $\beta_g$ .
- GPCR (G protein-coupled receptor) co-expression subnetwork of 17 positively correlated genes based on Ghazalour et al. (06).



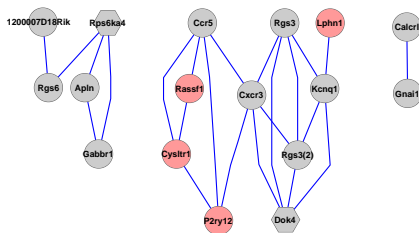
## Network-based eQTL analysis-continued

- It is biologically reasonable to assume that genes connected in a co-expression network are likely to share some common eQTLs; i.e, if two genes are connected, their expressions are likely to be associated with the genotypes at some common genomic loci.
- Combined model:

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}),$$

- Our method identifies 13 out 17 genes, as opposed to 2 by Lasso. The result is in agreement with that in previous studies (Lan et al., 06).

# Network-based eQTL analysis-continued



# Take Away Messages

- Supervised clustering and feature selection can reduce estimation variance while retaining the roughly the same amount of bias, leading to better predictive accuracy.
- The method identifies and collapses highly positive correlated predictors in a process of selection.
- Further develop methods for time varying networks.
- Study other types of clustering, e.g., coefficients of similar size not value, which involves the absolute values.