

Two-Layer Heterogeneity Model for Massive Data

Ching-Wei Cheng
(Ongoing joint work with Prof. Guang Cheng)

August 9, 2016
(BaT Group Meeting)

Heterogeneity is a promising yet challenging feature in **Big Data**

Example 1 (US Big-Data Health Network Launches Aspirin Study)

- ▶ Initiation of a 10-million pilot study that aim to investigate the use of aspirin to prevent heart disease
- ▶ **Various healthcare data** (e.g., insurance claims, blood tests, medical histories, etc.) will be collected from as many as 30 million people in the United States through a large healthcare network

<http://www.nature.com/news/us-big-data-health-network-launches-aspirin-study-1.15675>

- ▶ To fix terminology...
 - ▶ Observation \subseteq Unit \subseteq Cluster/subgroup \subseteq Full data (massive)
- ▶ Different data units may be endowed with different features (**Heterogeneity**)
 - ▶ However, some data units may be similar enough to treated homogeneous
 - ▶ Data similarity implies potential clustering effects among the data units, and thus encourages the development of **multi-layer heterogeneity models**

An Illustration of Nestedly Heterogeneous Data

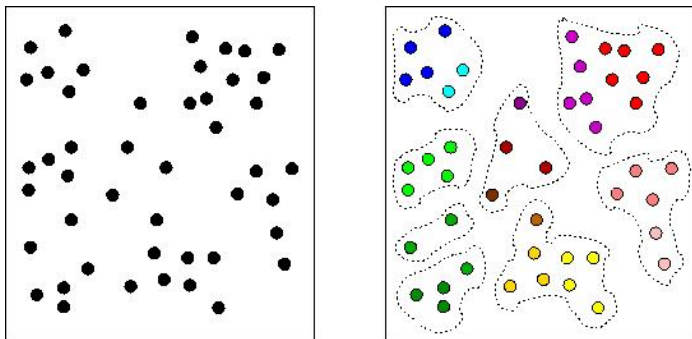


Figure 1 : Every dot is a data unit; colors and dashed lines indicate possible nested clustering structure

In this talk, we focus on linear model family to illustrate our idea

- ▶ Linear models $\mathcal{Y} = \mathcal{X}^\top \beta + \epsilon$
 - ▶ Simple and easy to interpret, but homogeneous
- ▶ Linear fixed-effect models $\mathcal{Y} = \mathcal{X}^\top \beta + \mathcal{Z}^\top \theta_i + \epsilon$
 - ▶ Fixed effects θ_i indicates (finite) heterogeneity among units
- ▶ Linear mixed-effects models (LMMs) $\mathcal{Y} = \mathcal{X}^\top \beta + \mathcal{Z}^\top \vartheta_i + \epsilon$ with $E[\vartheta_i] = \theta_i$ and $\text{Cov}(\vartheta_i) = \sigma_u^2 \mathbf{I}$
 - ▶ A well-known powerful tool for grouped (e.g., longitudinal, panel, cross-sectional) data
 - ▶ Random effects ϑ_i accounts for unit heterogeneity
 - ▶ Relax independence assumption
- ▶ Write $\vartheta_i = \theta_i + \mathbf{u}_i$
 - ▶ Fixed θ_i (between-unit heterogeneity)
 - ▶ Random \mathbf{u}_i with $E[\mathbf{u}_i] = \mathbf{0}$ and $\text{Cov}(\mathbf{u}_i) = \sigma_u^2 \mathbf{I}$
- ▶ We then consider **some of θ_i 's are identical**
 - ▶ Resulting in, if unit i belongs to subgroup s ,

$$\begin{aligned}\mathcal{Y} &= \mathcal{X}^\top \beta + \mathcal{Z}^\top \theta_i + \mathcal{Z}^\top \mathbf{u}_i + \epsilon \\ &= \mathcal{X}^\top \beta + \mathcal{Z}^\top \alpha_s + \mathcal{Z}^\top \mathbf{u}_i + \epsilon\end{aligned}$$

Outline

The Model

- Practical Model Formulation

- Oracular Model Representations

- A CD Fusion Approach for Massive Data

Theoretical Properties

- An Intermediate Estimator

- Oracle Properties

Simulation

Discussion

The Model

- ▶ Two-Layer heterogeneity model (THM):

$$\mathcal{Y} = \mathcal{X}^\top \boldsymbol{\beta} + \mathcal{Z}^\top \boldsymbol{\theta}_i + \mathcal{Z}^\top \mathbf{u}_i + \epsilon \quad (1)$$

- ▶ $\boldsymbol{\beta}$ is the p -vector of common fixed effects across all units
- ▶ $\boldsymbol{\theta}_i$ is the q -vector of unit-specific fixed effects of unit i
- ▶ \mathbf{u}_i is the q -vector of unit-specific random effects with $E[\mathbf{u}_i] = \mathbf{0}$ and $\text{Cov}(\mathbf{u}_i) = \sigma_u^2 \mathbf{I}$
- ▶ ϵ_i is the error n_i -vector with $E[\epsilon_i] = \mathbf{0}$ and $\text{Cov}(\epsilon_i) = \sigma_\epsilon^2$
- ▶ Massive data $\{(\mathbf{y}_i, \mathbf{x}_i, \mathbf{z}_i, L_i)\}_{i=1}^M$, aggregated from M data units
 - ▶ $\mathbf{y}_i = (y_{i1}, \dots, y_{in_i})^\top$ is a n_i -vector of responses
 - ▶ $\mathbf{x}_i(\mathbf{z}_i)$ is an $n_i \times p$ ($n_i \times q$) design matrix
 - ▶ L_i is the (latent) subgroup label for unit i
 - ▶ Total sample size $N = \sum_{i=1}^M n_i$

The Model

Practical Model Formulation

► Unit level:

$$\mathbf{y}_i = \mathbf{x}_i\boldsymbol{\beta} + \mathbf{z}_i\boldsymbol{\theta}_i + \mathbf{z}_i\mathbf{u}_i + \boldsymbol{\epsilon}_i, \quad i = 1, \dots, M \quad (2)$$

- Define $\mathbf{W}_i = \text{Cov}(\mathbf{y}_i)^{-1} = (\sigma_\epsilon^2 \mathbf{I}_{n_i} + \sigma_u^2 \mathbf{z}_i \mathbf{z}_i^\top)^{-1}$

► Full-data level:

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\Theta} + \mathbf{Z}\mathbf{U} + \boldsymbol{\mathcal{E}} \quad (3)$$

- $\mathbf{Y}_{(N \times 1)} = (\mathbf{y}_1^\top, \dots, \mathbf{y}_M^\top)^\top$
- $\mathbf{X}_{(N \times p)} = (\mathbf{x}_1^\top, \dots, \mathbf{x}_M^\top)^\top$
- $\mathbf{Z}_{(N \times Mq)} = \text{bdiag}(\mathbf{z}_1, \dots, \mathbf{z}_M)$
- $\boldsymbol{\Theta}_{(Mq \times 1)} = (\boldsymbol{\theta}_1^\top, \dots, \boldsymbol{\theta}_M^\top)^\top$
- $\mathbf{U}_{(Mq \times 1)} = (\mathbf{u}_1^\top, \dots, \mathbf{u}_M^\top)^\top$
- $\boldsymbol{\mathcal{E}}_{(N \times 1)} = (\boldsymbol{\epsilon}_1^\top, \dots, \boldsymbol{\epsilon}_M^\top)^\top$
- Define $\mathbf{W} = \text{Cov}(\mathbf{Y})^{-1} = \text{bdiag}(\mathbf{W}_i)$

The Model

Oracular Model Representations

When subgroup labels L_i 's are known, model (1) can be expressed as

$$\begin{aligned}\mathcal{Y} &= \mathcal{X}^\top \boldsymbol{\beta} + \mathcal{Z}^\top \boldsymbol{\theta}_i + \mathcal{Z}^\top \mathbf{u}_i + \epsilon \\ &= \mathcal{X}^\top \boldsymbol{\beta} + \mathcal{Z}^\top \boldsymbol{\alpha}_s + \mathcal{Z}^\top \mathbf{u}_i + \epsilon,\end{aligned}\tag{4}$$

where $\boldsymbol{\alpha}_s$ is the q -vector of group-specific fixed effects of subgroup s

- ▶ $\boldsymbol{\alpha}_{(Sq \times 1)} = (\boldsymbol{\alpha}_1^\top, \dots, \boldsymbol{\alpha}_S^\top)^\top$
- ▶ Define

$$\mathcal{M}_{\mathcal{G}} = \{ \boldsymbol{\Theta} \in \mathbb{R}^{Mq} : \boldsymbol{\theta}_i = \boldsymbol{\theta}_j \text{ for any } L_i = L_j \}$$

- ▶ For each $\boldsymbol{\Theta} \in \mathcal{M}_{\mathcal{G}}$, we have $\boldsymbol{\Theta} = \mathbf{A}\boldsymbol{\alpha}$, where \mathbf{A} is a $(Mq) \times (Sq)$ matrix with the (i, s) block being \mathbf{I}_q if $L_i = s$, and \mathbf{O}_q otherwise.
 - ▶ An example of matrix \mathbf{A} will be delivered later...

The Model

Oracular Model Representations

- **Unit level:**

$$\mathbf{y}_i = \mathbf{x}_i\boldsymbol{\beta} + \mathbf{z}_i\boldsymbol{\alpha}_s + \mathbf{z}_i\mathbf{u}_i + \epsilon_i, \quad i = 1, \dots, M. \quad (5)$$

- **Full-data level:** (Recall that $\boldsymbol{\Theta} = \mathbf{A}\boldsymbol{\alpha}$)

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{A}\boldsymbol{\alpha} + \mathbf{Z}\mathbf{U} + \boldsymbol{\mathcal{E}}, \quad (6)$$

- Oracle estimator, defined via GLS:

$$\begin{pmatrix} \hat{\boldsymbol{\beta}}_{\text{OR}} \\ \hat{\boldsymbol{\alpha}}_{\text{OR}} \end{pmatrix} = [(\mathbf{X}, \mathbf{Z}\mathbf{A})^\top \mathbf{W}(\mathbf{X}, \mathbf{Z}\mathbf{A})]^{-1} (\mathbf{X}, \mathbf{Z}\mathbf{A})^\top \mathbf{W}\mathbf{Y} \quad (7)$$

The Model

Oracular Model Representations

Example 2

Suppose $M = 5, S = 2$, $\theta_1 = \theta_2 = \alpha_1$ and $\theta_3 = \theta_4 = \theta_5 = \alpha_2$. Then we have

$$\Theta_{(5q) \times 1} = \begin{pmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \\ \theta_4 \\ \theta_5 \end{pmatrix}_{(5q) \times 1} = \begin{pmatrix} \alpha_1 \\ \alpha_1 \\ \alpha_2 \\ \alpha_2 \\ \alpha_2 \end{pmatrix}_{(5q) \times 1} = \begin{bmatrix} I_q & & & & \\ I_q & & & & \\ & I_q & & & \\ & I_q & & & \\ & I_q & & & \end{bmatrix}_{(5q) \times (2q)} \begin{pmatrix} \alpha_1 \\ \alpha_2 \end{pmatrix}_{(2q) \times 1}$$

and

$$ZA = \begin{bmatrix} z_1 & & & & \\ & z_2 & & & \\ & & z_3 & & \\ & & & z_4 & \\ & & & & z_5 \end{bmatrix} \begin{bmatrix} I_q & & & & \\ I_q & & & & \\ & I_q & & & \\ & I_q & & & \\ & I_q & & & \end{bmatrix} = \begin{bmatrix} z_1 & & & & \\ & z_2 & & & \\ & & z_3 & & \\ & & & z_4 & \\ & & & & z_5 \end{bmatrix},$$

The Model

A CD Fusion Approach for Massive Data

When dealing with massive data, GLS may be computationally unavailable...

- ▶ Due to massive sample size, direct estimation towards the full-data model (3) is usually unavailable
 - ▶ GLS approach can be applied for LMMs, but may be computationally infeasible due to massive sample size
- ▶ Starting with unit GLS estimates

$$\begin{pmatrix} \hat{\beta}_i \\ \hat{\theta}_i \end{pmatrix} = [(\mathbf{x}_i, \mathbf{z}_i)^\top \mathbf{W}_i (\mathbf{x}_i, \mathbf{z}_i)]^{-1} (\mathbf{x}_i, \mathbf{z}_i)^\top \mathbf{W}_i \mathbf{y}_i, \quad i = 1, \dots, M \quad (8)$$

- ▶ If σ_ϵ^2 and σ_u^2 are unknown, they can be consistently estimated by restricted maximum likelihood (REML) method
- ▶ Throughout this talk, we assume σ_ϵ^2 and σ_u^2 are known for simplicity
- ▶ We adopt the **confidence distribution (CD)** concept to merge the unit estimates
- ▶ To discover potential subgroups, we consider **pairwise concave fusion penalty** as in Ma and Huang (2016)

The Model

A CD Fusion Approach for Massive Data

Confidence Distribution (CD)

- ▶ A CD can be viewed as
 - ▶ “A sample-dependent distribution function that can represent confidence intervals of all levels for a parameter of interest”
 - ▶ “The frequentist distribution estimator of a parameter”
- ▶ See Xie and Singh (2013) for a comprehensive review for the CD development
- ▶ The CD concept has been shown effective for combining information from independent sources (unit estimates)
- ▶ We are going to adopt the CD approach proposed by Liu *et al.* (2015) for multi-parameter estimates

The Model

A CD Fusion Approach for Massive Data

A **CD fusion estimator** can be formed as follows:

- ▶ Given $(\hat{\beta}_i^\top, \hat{\theta}_i^\top)^\top \xrightarrow{D} \mathcal{N}((\beta_0^\top, \theta_{i,0}^\top)^\top, [(\mathbf{x}_i, \mathbf{z}_i)^\top \mathbf{W}_i(\mathbf{x}_i, \mathbf{z}_i)]^{-1})$
- ▶ Following the CD concept, the **CD densities** $h_i(\beta, \theta_i)$ can be defined as the density function of $\mathcal{N}((\hat{\beta}_i^\top, \hat{\theta}_i^\top)^\top, [(\mathbf{x}_i, \mathbf{z}_i)^\top \mathbf{W}_i(\mathbf{x}_i, \mathbf{z}_i)]^{-1})$
- ▶ The combined CD density can be defined by

$$h(\beta, \Theta) = \prod_{i=1}^M h_i(\beta, \theta_i)$$

- ▶ A CD estimator of (β, Θ) is defined as $\arg \max_{\beta, \Theta} \log h(\beta, \Theta)$ (yet not what we want)
- ▶ Recall that some underlying values of θ_i 's are assumed identical, and so we incorporate a **pairwise concave fusion penalty** into the objective function for a fusion estimation

The Model

A CD Fusion Approach for Massive Data

We propose a **CD fusion estimator** defined by

$$\begin{pmatrix} \check{\beta}(\lambda) \\ \check{\Theta}(\lambda) \end{pmatrix} = \arg \min_{\beta \in \mathbb{R}^p, \Theta \in \mathbb{R}^{Mq}} Q_N^{\text{CD}}(\beta, \Theta), \quad (9)$$

where

$$\begin{aligned} Q_N^{\text{CD}}(\beta, \Theta) = & \frac{1}{2} \sum_{i=1}^M \begin{pmatrix} \hat{\beta}_i - \beta \\ \hat{\theta}_i - \theta_i \end{pmatrix}^\top (x_i, z_i)^\top W_i(x_i, z_i) \begin{pmatrix} \hat{\beta}_i - \beta \\ \hat{\theta}_i - \theta_i \end{pmatrix} \\ & + \sum_{1 \leq i < j \leq M} p_\gamma(\|\theta_i - \theta_j\|; \lambda), \end{aligned} \quad (10)$$

- ▶ First term of the R.H.S. comes from the simplified $-\log h(\beta, \Theta)$ by omitting additive constant terms, due to asymptotic normality of unit GLS estimates
- ▶ $p_\gamma(t; \lambda)$ is a concave penalty function with a tuning parameter $\lambda > 0$ and a parameter $\gamma > 0$ which is associated with the concavity of the penalty

The Model

Remarks on the CD fusion estimator:

- The objective function (10) can be rewritten as

$$\begin{aligned} Q_N^{\text{CD}}(\boldsymbol{\beta}, \boldsymbol{\Theta}) = & \frac{1}{2} \sum_{i=1}^M (\hat{\boldsymbol{\beta}}_i - \boldsymbol{\beta})^\top \mathbf{x}_i^\top \mathbf{W}_i \mathbf{x}_i (\hat{\boldsymbol{\beta}}_i - \boldsymbol{\beta}) \\ & + \frac{1}{2} \sum_{i=1}^M (\hat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}_i)^\top \mathbf{z}_i^\top \mathbf{W}_i \mathbf{z}_i (\hat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}_i) \\ & + \sum_{1 \leq i < j \leq M} p_\gamma(\|\boldsymbol{\theta}_i - \boldsymbol{\theta}_j\|; \lambda). \end{aligned} \quad (11)$$

- Minimizer of $\boldsymbol{\beta}$ is actually free of $\boldsymbol{\Theta}$, resulting in

$$\check{\boldsymbol{\beta}}(\lambda) \equiv \check{\boldsymbol{\beta}} = \left(\sum_{i=1}^M \mathbf{x}_i^\top \mathbf{W}_i \mathbf{x}_i \right)^{-1} \left(\sum_{i=1}^M \mathbf{x}_i^\top \mathbf{W}_i \mathbf{x}_i \hat{\boldsymbol{\beta}}_i \right). \quad (12)$$

- To find $\check{\boldsymbol{\Theta}}(\lambda)$, we reformulate the objective function as an augmented Lagrangian problem, and establish an [alternating direction method of multipliers \(ADMM, Boyd et al. \(2010\)\)](#) algorithm

The Model

Remarks on $p_\gamma(t; \lambda)$:

- ▶ L_1 or Lasso penalty produces biased estimates and thus may not correctly form the clusters
 - ▶ Tends to result in either a large number of subgroups or no subgroup on the solution path
- ▶ MCP and SCAD penalty are more appropriate
 - ▶ Enjoy sparsity (on pairwise distances) as the L_1 penalty so that they automatically fuse some of θ_i 's together
 - ▶ Due to their unbiasedness property, they do not shrink large estimated parameters
 - ▶ Remain unbiased in ADMM iterations for concave optimization solvers

Theoretical Properties

An Intermediate Estimator

- ▶ We introduce an **intermediate estimator**:

$$\begin{pmatrix} \hat{\beta}(\lambda) \\ \hat{\Theta}(\lambda) \end{pmatrix} = \arg \min_{\beta \in \mathbb{R}^p, \Theta \in \mathbb{R}^{Mq}} Q_N(\beta, \Theta), \quad (13)$$

where

$$\begin{aligned} Q_N(\beta, \Theta) = & \frac{1}{2}(\mathbf{Y} - \mathbf{X}\beta - \mathbf{Z}\Theta)^\top \mathbf{W}(\mathbf{Y} - \mathbf{X}\beta - \mathbf{Z}\Theta) \\ & + \sum_{1 \leq i < j \leq M} p_\gamma(\|\theta_i - \theta_j\|; \lambda), \end{aligned} \quad (14)$$

- ▶ A GLS version of the estimator proposed by Ma and Huang (2016)
- ▶ Our analysis strategy:
 - ▶ **Prove the CD fusion estimator is equivalent to the intermediate estimator**
 - ▶ Show oracle properties of the proposed CD fusion estimator through the intermediate estimator

Theoretical Properties

An Intermediate Estimator

Theorem 2.1 (Equivalence between the proposed CD fusion estimator and the intermediate estimator)

Let

$$\begin{pmatrix} \check{\beta}(\lambda) \\ \check{\Theta}(\lambda) \end{pmatrix} = \arg \min_{\beta \in \mathbb{R}^p, \Theta \in \mathbb{R}^{Mq}} Q_N^{\text{CD}}(\beta, \Theta)$$

and

$$\begin{pmatrix} \hat{\beta}(\lambda) \\ \hat{\Theta}(\lambda) \end{pmatrix} = \arg \min_{\beta \in \mathbb{R}^p, \Theta \in \mathbb{R}^{Mq}} Q_N(\beta, \Theta),$$

where the objective functions $Q_N^{\text{CD}}(\beta, \Theta)$ and $Q_N(\beta, \Theta)$ are defined in (10) and (14), respectively. Then we have

$$P \left(\begin{pmatrix} \check{\beta} \\ \check{\Theta}(\lambda) \end{pmatrix} = \begin{pmatrix} \hat{\beta}(\lambda) \\ \hat{\Theta}(\lambda) \end{pmatrix} \right) = 1.$$

Theoretical Properties

An Intermediate Estimator

Remarks on Theorem 2.1:

- ▶ The intermediate estimator uses the full data at once, and thus can be treated as the **IPD (individual participant data)** estimator, which is taken as the **gold standard** in meta-analysis
- ▶ Liu *et al.* (2015) showed that the CD estimator is asymptotically equivalent to the IPD estimator
- ▶ Our equivalence result is stronger (non-asymptotic)
 - ▶ May result from the LMM structure and GLS estimation approach

Theoretical Properties

Oracle Properties

When the true subgroup membership L_i 's are known, define the following:

- ▶ $\mathcal{G}_s = \{(\mathbf{y}_i, \mathbf{x}_i, \mathbf{z}_i, L_i) : L_i = s\}, s = 1, \dots, S$
- ▶ $M_s = |\mathcal{G}_s|$, the number of units in subgroup s
 - ▶ $M_{\min} = \min_{1 \leq s \leq S} M_s$
 - ▶ $M_{\max} = \max_{1 \leq s \leq S} M_s$
- ▶ $g_s = \sum_{i: L_i = s} n_i$, the number of observations in subgroup s
 - ▶ $g_{\min} = \min_{1 \leq s \leq S} g_s$
 - ▶ $g_{\max} = \max_{1 \leq s \leq S} g_s$
- ▶ The **oracle estimator** is defined from model (6) by

$$\begin{pmatrix} \hat{\beta}_{\text{OR}} \\ \hat{\alpha}_{\text{OR}} \end{pmatrix} = [(\mathbf{X}, \mathbf{Z}\mathbf{A})^\top \mathbf{W}(\mathbf{X}, \mathbf{Z}\mathbf{A})]^{-1} (\mathbf{X}, \mathbf{Z}\mathbf{A})^\top \mathbf{W}\mathbf{Y} \quad (15)$$

Theoretical Properties

Oracle Properties

Theorem 2.2 (Asymptotics for the oracle estimator)

Assume regularity conditions. If $g_{\min} \gg \sqrt{(p + Sq)N \log N}$ and $p + Sq = o(N^{1 \wedge \tau})$, we have

$$\left\| \begin{pmatrix} \hat{\beta}_{\text{OR}} - \beta_0 \\ \hat{\Theta}_{\text{OR}} - \Theta_0 \end{pmatrix} \right\| \leq \phi_N = O_P \left(g_{\min}^{-1} \sqrt{(p + Sq)N \log N} \right), \quad (16)$$

$$\|\hat{\beta}_{\text{OR}} - \beta_0\| = O_P \left(\sqrt{\frac{p \log N}{N}} \right), \quad (17)$$

$$\|\hat{\alpha}_{\text{OR}} - \alpha_0\| = O_P \left(\sqrt{\frac{Sq \log g_{\min}}{g_{\min}}} \right). \quad (18)$$

Moreover, for any sequence of $(p + Sq)$ -vectors $\{\mathbf{a}_N\}$, we have

$$\sigma_N^{-1}(\mathbf{a}_N) \mathbf{a}_N^\top \begin{pmatrix} \hat{\beta}_{\text{OR}} - \beta_0 \\ \hat{\alpha}_{\text{OR}} - \alpha_0 \end{pmatrix} \xrightarrow{D} \mathcal{N}(0, 1), \quad (19)$$

where $\sigma_N(\mathbf{a}_N) = \left(\mathbf{a}_N^\top [(X, ZA)^\top W(X, ZA)]^{-1} \mathbf{a}_N \right)^{1/2}$.

Theoretical Properties

Oracle Properties

- For $S \geq 2$, let

$$\Delta_N = \min_{L_i=s, L_j=s', s \neq s'} \|\boldsymbol{\theta}_i - \boldsymbol{\theta}_j\| = \min_{s \neq s'} \|\boldsymbol{\alpha}_s - \boldsymbol{\alpha}_{s'}\|$$

denote the minimal signal measure for the unit-specific fixed effects

Theorem 2.3 (Oracle property)

Suppose regularity conditions and $S \geq 2$. If $\Delta_N \gg a\phi_N$, $\lambda \gg \phi_N$, where a is defined in the assumption for penalty functions and ϕ_N is given in

Theorem 2.2, then there exists a local minimizer $(\check{\boldsymbol{\beta}}^\top, \check{\boldsymbol{\Theta}}(\lambda)^\top)^\top$ of the objective function $Q_N^{\text{CD}}(\boldsymbol{\beta}, \boldsymbol{\Theta})$ given in (10) satisfying

$$P \left(\begin{pmatrix} \check{\boldsymbol{\beta}} \\ \check{\boldsymbol{\Theta}}(\lambda) \end{pmatrix} = \begin{pmatrix} \hat{\boldsymbol{\beta}}_{\text{OR}} \\ \hat{\boldsymbol{\Theta}}_{\text{OR}} \end{pmatrix} \right) \rightarrow 1.$$

Theoretical Properties

Oracle Properties

Corollary 2.1

Under the conditions in Theorem 2.3, we have for any $(p + Sq)$ -vector \mathbf{a}_N ,

$$\sigma_N^{-1}(\mathbf{a}_N) \mathbf{a}_N^\top \begin{pmatrix} \check{\boldsymbol{\beta}} - \boldsymbol{\beta}_0 \\ \check{\boldsymbol{\alpha}}(\lambda) - \boldsymbol{\alpha}_0 \end{pmatrix} \xrightarrow{D} \mathcal{N}(0, 1),$$

where $\sigma_N(\mathbf{a}_N)$ is given in Theorem 2.2. Moreover, we have for any $\mathbf{a}_{N1} \in \mathbb{R}^p$ and $\mathbf{a}_{N2} \in \mathbb{R}^{Sq}$,

$$\begin{aligned} \sigma_{N1}^{-1}(\mathbf{a}_{N1}) \mathbf{a}_{N1}^\top (\check{\boldsymbol{\beta}} - \boldsymbol{\beta}_0) &\xrightarrow{D} \mathcal{N}(0, 1), \\ \sigma_{N2}^{-1}(\mathbf{a}_{N2}) \mathbf{a}_{N2}^\top (\check{\boldsymbol{\alpha}}(\lambda) - \boldsymbol{\alpha}_0) &\xrightarrow{D} \mathcal{N}(0, 1), \end{aligned}$$

where

$$\begin{aligned} \sigma_{N1}(\mathbf{a}_{N1}) &= \left\{ \mathbf{a}_{N1}^\top (\mathbf{X}^\top \mathbf{Q}_{\mathbf{Z}\mathbf{A}} \mathbf{X})^{-1} \mathbf{a}_{N1} \right\}^{1/2}, \\ \sigma_{N2}(\mathbf{a}_{N2}) &= \left\{ \mathbf{a}_{N2}^\top [(\mathbf{Z}\mathbf{A})^\top \mathbf{Q}_{\mathbf{X}} (\mathbf{Z}\mathbf{A})]^{-1} \mathbf{a}_{N2} \right\}^{1/2}. \end{aligned}$$

Theoretical Properties

Oracle Properties

Remarks on Theorem 2.3 and Corollary 2.1:

- ▶ As θ_i 's can be consistently estimated by $\hat{\theta}_i(\lambda)$'s, the subgroup membership L_i 's and the cluster-specific fixed effects α_s 's can be obtained consistently as well
- ▶ The asymptotic variances $\sigma_{N1}(\mathbf{a}_{N1})$ and $\sigma_{N2}(\mathbf{a}_{N2})$ are derived from $\text{Cov}[(\hat{\beta}_{\text{OR}}^\top, \hat{\alpha}_{\text{OR}}^\top)^\top]$
 - ▶ $(\hat{\beta}_{\text{OR}}^\top, \hat{\alpha}_{\text{OR}}^\top)^\top$ is formed via GLS, which produces **best linear unbiased estimator (BLUE)**
 - ▶ $\sigma_{N1}(\mathbf{a}_{N1})$ and $\sigma_{N2}(\mathbf{a}_{N2})$ are the smallest asymptotic variances of $\mathbf{a}_{N1}^\top \check{\beta}(\lambda)$ and $\mathbf{a}_{N2}^\top \check{\alpha}(\lambda)$ for any $\mathbf{a}_{N1} \in \mathbb{R}^p$ and $\mathbf{a}_{N2} \in \mathbb{R}^{Sq}$
 - ▶ In summary, Corollary 2.1 suggests that the optimal inference be made by the proposed CD fusion estimator

Simulation

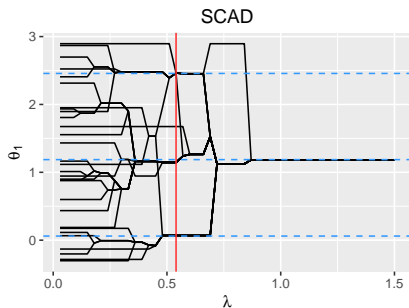
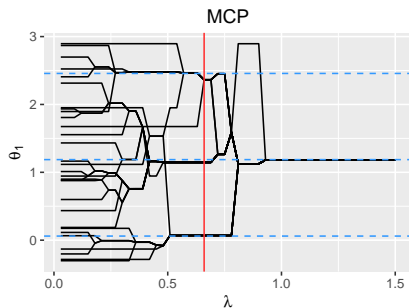
Tuning Parameter Selection

- ▶ Tuning parameter selection:
 - ▶ To select λ , we apply the modified BIC (Wang *et al.*, 2009) for diverging parameter dimension
 - ▶ $\check{\lambda} = \arg \min_{\lambda} \text{BIC}(\lambda)$
 - ▶ $\gamma = 3$
- ▶ **Solution path analysis** for $\theta_{s,1}$:
 - ▶ Blue dashed lines are $\hat{\alpha}_{\text{OR},s,1}$'s
 - ▶ Red solid line indicates $\check{\lambda}$ (via grid search)

Simulation

Solution Path Analysis

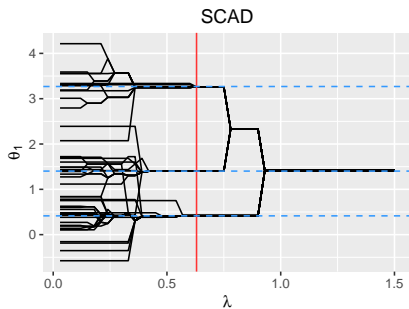
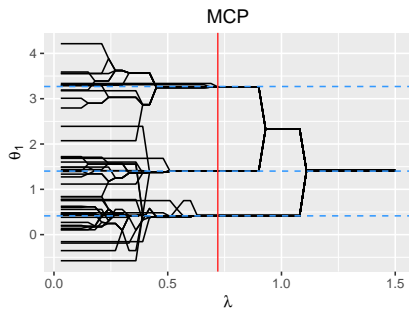
$$n_i \equiv n = 256, M = 30, S = 3; p = 5, q = 4$$



Simulation

Solution Path Analysis

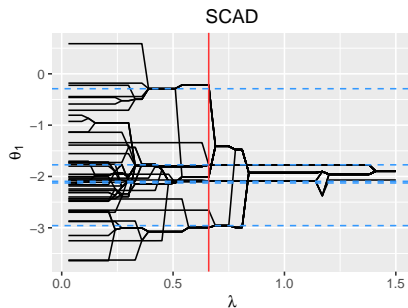
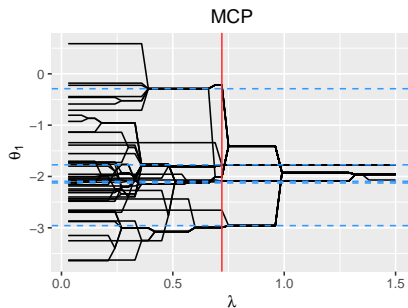
$$n_i \equiv n = 256, M = 50, S = 3; p = 5, q = 4$$



Simulation

Solution Path Analysis

$$n_i \equiv n = 256, M = 50, S = 5; p = 5, q = 4$$



Discussion

- ▶ The theoretical results are based on **fixed design**
 - ▶ Currently trying to move to **random design** regime, which is more appropriate for Big Data
 - ▶ Need to deal with lots of probability statements (e.g., bounded columns, eigenvalue bounds, projection matrices, etc.)
- ▶ The exhaustive pairwise fusion penalty is too costly when M is large (i.e., $M(M-1)/2$ pairs)
 - ▶ Have tried vectorization, i.e., there exists a sparse matrix B s.t. $B^\top \Theta$ includes all pairs $\theta_i - \theta_j$, but seems not good enough
 - ▶ May consider the aCARDS proposed by Ke *et al.* (2015)
 - ▶ Make a rough segmentation (e.g., $M^{0.7}$ segments)
 - ▶ Use a **hybrid pairwise fusion penalty** to consider **between-segment** and **within-segment** penalties
- ▶ More evaluation measures for clustering and estimation in simulation

Selected References

- Boyd, S., Parikh, N., Chu, E., Peleato, B. and Eckstein, J. (2010) Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends in Machine Learning*, **3**, 1–122. URL <http://dx.doi.org/10.1561/22000000016>.
- Ke, Z. T., Fan, J. and Wu, Y. (2015) Homogeneity pursuit. *Journal of the American Statistical Association*, **110**, 175–194. URL <http://dx.doi.org/10.1080/01621459.2014.892882>.
- Liu, D., Liu, R. Y. and Xie, M. (2015) Multivariate meta-analysis of heterogeneous studies using only summary statistics: Efficiency and robustness. *Journal of the American Statistical Association*, **110**, 326–340. URL <http://dx.doi.org/10.1080/01621459.2014.899235>.
- Ma, S. and Huang, J. (2016) A concave pairwise fusion approach to subgroup analysis. *Journal of the American Statistical Association*. URL <http://dx.doi.org/10.1080/01621459.2016.1148039>. To appear.
- Wang, H., Li, B. and Leng, C. (2009) Shrinkage tuning parameter selection with a diverging number of parameters. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **71**, 671–683. URL <http://dx.doi.org/10.1111/j.1467-9868.2008.00693.x>.
- Xie, M. and Singh, K. (2013) Confidence distribution, the frequentist distribution estimator of a parameter: A review. *International Statistical Review*, **81**, 3–39. URL <http://dx.doi.org/10.1111/insr.12000>.