

# Model Selection And Sharp Asymptotic Minimaxity

Zheyang Wu    Huibin (Harrison) Zhou

Biostatistics Division, Yale School of Medicine  
Department of Statistics  
Yale University

May 16, 2009

# Problem description

- **Multivariate normal mean problem:**

$$y_i = \theta_i + \sigma_n z_i, \quad z_i \stackrel{i.i.d.}{\sim} N(0, 1), \quad i = 1, \dots, n.$$

- **Sparse parameter space:**  $\theta \in \Theta_{n,p}(\eta_n)$ ,  $\Theta_{n,p}(\eta_n)$  is one of

$$l_0 \text{ ball} : \quad \{\theta \in \mathbb{R}^n : \|\theta\|_0 \leq \eta_n n\}$$

$$l_p \text{ ball} : \quad \left\{ \theta \in \mathbb{R}^n : \sum_{i=1}^n |\theta_i|^p \leq \eta_n^p n, 0 < p < 2 \right\}$$

$$m_p \text{ ball} : \quad \left\{ \theta \in \mathbb{R}^n : |\theta|_{[k]} \leq \eta_n \left( \frac{n}{k} \right)^{1/p}, 0 < p < 2, k = 1, \dots, n \right\}$$

where  $\eta_n \in [n^{-1} \log^\gamma n, b_2 n^{-b_3}]$ ,  $\gamma > 4.5$ .

- **Goal:** to estimate  $\theta = (\theta_1, \dots, \theta_n)$ .

# Main theorems

- Penalized estimation:  $\hat{\theta} = \arg \min_{\|\mu\|_0 \leq n / \log n} \left[ \|y - \mu\|_2^2 + \text{Pen}(\|\mu\|_0) \right]$ .
- Main theorems: Let  $\|\mu\|_0 = k$ ,  $\text{Pen}(k) = \sum_{i=1}^k u_i^2$ ,  $R_n(\Theta_{n,p}(\eta_n))$  be the minimax risk,

## Theorem

For some  $0 < \varepsilon < 1$ , and any constant  $c'$ , let

$2 \log \frac{n}{i} - (1 - \varepsilon) \log \log \frac{n}{i} \leq u_i^2 \leq 2 \log \frac{n}{i} + c' \log \log n$ , then

$\sup_{\theta \in \Theta_{n,p}(\eta_n)} \mathbb{E} \|\hat{\theta} - \theta\|_2^2 \sim R_n(\Theta_{n,p}(\eta_n))$ .

## Theorem

Let  $u_i^2 = c_n \log \frac{n}{i}$ ,  $c_n \rightarrow c > 2$ , then

$\sup_{\theta \in \Theta_{n,p}(\eta_n)} \mathbb{E} \|\hat{\theta} - \theta\|_2^2 \sim c^* R_n(\Theta_{n,p}(\eta_n))$ , where  $c^* = \frac{c}{2}$  for  $l_0$  ball,

$c^* = \left(\frac{c}{2}\right)^{1-p/2}$  for  $l_p$  ball and  $m_p$  ball.

## Interpretation

- L0 norm leads to minimaxity if  $\text{Pen}(k)$  is in a sharp range of  $2k \log \frac{n}{k}$ .
- L0 norm leads to non-minimaxity if  $\text{Pen}(k)$  is not dominated by  $2k \log \frac{n}{k}$ .

## Significance

- Minimax estimation procedure for  $\text{Pen}(k)$ 
  - Foster and Stine (1999):  $2 \sum_{i=1}^k \log \frac{n}{i}$ .
  - George and Foster (2000):  $2 \sum_{i=1}^k \log \left( \frac{n+1}{i} - 1 \right)$ .
  - Berge and Massart (2001):  $2k \log \frac{n}{k}$ .
  - ABDJ (2006):  $\sum_{i=1}^k \left( \Phi^{-1} \left( \frac{q_n^i}{2n} \right) \right)^2$ .
- Non-minimax estimation procedure for  $\text{Pen}(k)$ 
  - Tibshirani and Knight (1999):  $\sum_{i=1}^k \log \frac{n}{i}$ .
  - Abramovich, Grinshtein, Pensky (2007):  $ck \log \frac{n}{k}$ , some  $c > 2$ .