UPS Delivers Optimal Phase Diagram in Variable Selection Non-optimal Regions for L¹ and L⁰-Penalization Methods

Jiashun Jin (STAT. Carnegie Mellon)

With Pengsheng Ji (Math. Cornell)

June 9, 2011

$$Y = X\beta + z,$$
 $X = X_{n,p},$ $z \sim N(0, I_n)$

• $p \gg n \gg 1$

- β is sparse: many coordinates are 0
- ▶ Gram matrix (of the column) X'X
 - has unit diagonals
 - is sparse (few large coordinates in each row)
- ► Ex. Compressive Sensing, Genomics

$$\frac{1}{2} \|Y - X\beta\|_2^2 + \frac{\lambda^2}{2} \|\beta\|_0$$

- L⁰-penalization method
- Variants: Cp, AIC, BIC, RIC
- Computationally challenging

Mallows (1973), Akaike (1974), Schwartz (1978), Foster & George (1994)



$$\frac{1}{2} \|\boldsymbol{Y} - \boldsymbol{X}\boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_1$$

- L¹-penalization method; Basis Pursuit
- Widely used
 - computationally efficient even when p is large
 - if sufficiently sparse and noiseless, equivalent to L⁰-penalization

Chen et al. (1998); Tibshirani (1996); Donoho (2006)

$$ig|(X'X)(i,j)ig| \geq rac{1}{\log(p)}, \qquad {\sf say}$$

• G is sparse in many applications

Signal sparsity and graph sparsity

- Despite its sparsity, G is usually complicate
- \blacktriangleright Denote the support of β by

$$S = S(\beta) = \{1 \leq i \leq p, \ \beta_i \neq 0\}$$

Restricting nodes to S forms a subgraph G_S

Key insight: G_S splits into many small-size component that are disconnected to each other

Component: a maximal connected subgraph



Univariate Penalization Screening (UPS)



U-step: screen by univariate thresholding P-step: clean by penalized MLE



$$Y = X\beta + z,$$
 $X = X_{n,p},$ $z \sim N(0, I_n).$

• Let
$$\tilde{Y} = X'Y$$
. For a threshold $t_p^* > 0$ TBD,

keep the
$$j$$
-th variable $\iff \widetilde{Y}_j > t_p^*$

Denote the set of survived indices by

$$\mathcal{U}_{p}(t_{p}^{*})=\{1\leq j\leq p, \hspace{0.1cm} \widetilde{Y}_{j}>t_{p}^{*}\}.$$

Donoho (2006), Fan and Lv (2008), Genovese et al. (2010) Use one-sided thresholding for simplicity

Jiashun Jin (STAT. Carnegie Mellon) UPS Delivers Optimal Phase Diagram in Variable Selection

Two important properties of U-step

 $\mathcal{U}_{p}(t_{p}^{*}) = \{1 \leq j \leq p, \ \tilde{Y}_{j} > t_{p}^{*}\}:$ survived indices

If both signals and Graph G are sparse:

 Sure Screening (SS): U_p(t^{*}_p) retains all but a small proportion of signals

Separable After Screening (SAS): U_p(t^{*}_p) splits into many components, each is small in size, and different ones are disconnected
 Original problem reduces to many small-size regression problems

Reduce to many small-size regression

$$Y = X'Y,$$
 $\mathcal{I}_0 \subset \mathcal{U}_p(t_p^*):$ a component
 $(X'X)^{\mathcal{I}_0}:$ row restriction; $(X'X)^{\mathcal{I}_0,\mathcal{I}_0}:$ row & column restriction

• Restrict regression to \mathcal{I}_0

$$\begin{split} Y &= X\beta + z \implies \tilde{Y} = X'X\beta + X'z \\ &\implies \tilde{Y}^{\mathcal{I}_0} = (X'X\beta)^{\mathcal{I}_0} + (X'z)^{\mathcal{I}_0} \end{split}$$

•
$$(X'z)^{\mathcal{I}_0} \sim N(0, (X'X)^{\mathcal{I}_0, \mathcal{I}_0})$$
 since $z \sim N(0, I_n)$

- Key: $(X'X\beta)^{\mathcal{I}_0} \approx (X'X)^{\mathcal{I}_0,\mathcal{I}_0}\beta^{\mathcal{I}_0}$
- Result: many small-size regression:

$$ilde{Y}^{\mathcal{I}_0} pprox \mathsf{N}ig((X'X)^{\mathcal{I}_0,\mathcal{I}_0}eta^{\mathcal{I}_0},(X'X)^{\mathcal{I}_0,\mathcal{I}_0}ig)$$

Reduce to small-size regression, II

Why
$$(X'X\beta)^{\mathcal{I}_0} = (X'X)^{\mathcal{I}_0}\beta \approx (X'X)^{\mathcal{I}_0,\mathcal{I}_0}\beta^{\mathcal{I}_0}$$
?
 $\begin{bmatrix} (X'X)^{\mathcal{I}_0,\mathcal{I}_0} \blacksquare (X'X)^{\mathcal{I}_0,\mathcal{J}_0} \blacksquare \dots \end{bmatrix} \begin{bmatrix} \beta^{\mathcal{I}_0} \\ 0 \\ \beta^{\mathcal{J}_0} \\ 0 \\ \dots \end{bmatrix}$

- $\mathcal{I}_0, \mathcal{J}_0 \subset \mathcal{U}_p(t_p^*)$: components
- By SS property, $\beta^{\blacksquare} = 0$
- By SAS property, $(X'X)^{\mathcal{I}_0,\mathcal{J}_0}pprox 0$

P-step

$$Y = X\beta + z, \qquad z \sim N(0, I_n)$$

• $ilde{Y} = X'Y; \quad \mathcal{I}_0 \subset \mathcal{U}_p(t_p^*): \text{ a component}$

• $\beta^{\mathcal{I}_0}$: restricting rows of β to \mathcal{I}_0

• M_0 : restricting rows/columns of X'X to \mathcal{I}_0

Fixing (λ^{ups}, u^{ups}) ,

• estimate $\beta^{\mathcal{I}_0}$ via minimizing

$$[\tilde{Y}^{\mathcal{I}_0} - M_0\beta^{\mathcal{I}_0}]'M_0^{-1}[\tilde{Y}^{\mathcal{I}_0} - M_0\beta^{\mathcal{I}_0}] + (\lambda^{ups})^2 \|\beta^{\mathcal{I}_0}\|_0,$$

where $\beta^{\mathcal{I}_0}$ takes values from $\{0, u^{ups}\}$

•
$$j \notin \mathcal{U}_p(t_p^*)$$
: set $\hat{\beta}_j = 0$

Sparse model

$$Y = X\beta + z, \qquad z \sim N(0, I_n)$$

Sparse signals:

$$\beta_j \stackrel{iid}{\sim} (1 - \epsilon_p) \nu_0 + \epsilon_p \pi_p$$

where ν_0 is point mass at 0 and

$$\epsilon_{p} = p^{-\vartheta}, \qquad 0 < \vartheta < 1$$

Measuring errors with Hamming distance:

$$\operatorname{Hamm}_{p}(\hat{\beta}; \epsilon_{p}, \pi_{p} | X) = E\left[\sum_{j=1}^{p} \mathbb{1}\{\operatorname{sgn}(\hat{\beta}_{j}) \neq \operatorname{sgn}(\beta_{j})\}\right]$$

$$eta_j \stackrel{\textit{iid}}{\sim} (1 - \epsilon_p)
u_0 + \epsilon_p \pi_p, \qquad \epsilon_p = p^{-\vartheta}$$

Theorem 1. Fix r > 0 and let $\tau_p = \sqrt{2r \log p}$. If the support of π_p is contained in $[-\tau_p, 0) \cup (0, \tau_p]$, then as $p \to \infty$,

$$\frac{\operatorname{Hamm}_{p}(\hat{\beta};\epsilon_{p},\pi_{p}|X)}{p\epsilon_{p}} \gtrsim \begin{cases} L_{p} \cdot p^{-(r-\vartheta)^{2}/(4r)}, & r > \vartheta, \\ 1, & r < \vartheta, \end{cases}$$

where L_p is a generic multi-log(p) term.

Random design model

$$Y = X\beta + z, \qquad X = \begin{pmatrix} X'_1 \\ \dots \\ X'_n \end{pmatrix}, \qquad X_i \stackrel{iid}{\sim} N(0, \frac{1}{n}\Omega)$$

• Ex: Compressive Sensing, Computer Security Dinur and Nissim (2004), Nowak et al. (2007)

Additional Assumptions

$$\beta_j \stackrel{iid}{\sim} (1 - \epsilon_p)\nu_0 + \epsilon_p \pi_p, \quad \epsilon_p = p^{-\vartheta}, \quad \tau_p = \sqrt{2r \log p}$$
$$\eta = \min\{\frac{2\vartheta}{r}, \ 1 - \frac{\vartheta}{r}, \ \sqrt{2 - 2\omega_0} - 1 + \frac{\vartheta}{r}\}$$

- Support of $\pi_p \subset [\tau_p, (1+\eta)\tau_p], r > \vartheta$ (*)
- $\sum_{j=1}^p |\Omega(i,j)|^\gamma \leq C$, $\gamma \in (0,1)$
- $\max\{\|U(\Omega)\|_1, \|U(\Omega)\|_\infty\} \le \omega_0 < 1/2$ (*)
- $\Omega(i,j) \geq 0$ or $r/\vartheta \leq 3 + 2\sqrt{2}$ (*)

(*): relaxable; Jin and Zhang (2011), in progress

Theorem 2. Set
$$t_p^* = \sqrt{2q \log p}$$
, $q \in (0, \frac{(\vartheta + r)^2}{4r}]$, $\lambda^{ups} = \sqrt{2\vartheta \log p}$, and $u^{ups} = \tau_p$. As $p \to \infty$,

- Both SS and SAS property hold
- UPS achieves optimal rate of convergence

$$\frac{E_{\Omega}[\operatorname{Hamm}_{p}(\hat{\beta}^{ups};\epsilon_{p},\pi_{p}|X)]}{p\epsilon_{p}} \leq L_{p}p^{-\frac{(r-\vartheta)^{2}}{4r}}$$

UPS delivers optimal phase diagram



 $\epsilon_p = p^{-\vartheta}; \quad \tau_p = \sqrt{2r\log p}; \quad \text{line: } r = \vartheta; \quad \text{curve: } r = (1 + \sqrt{1 - \vartheta})^2$

Jiashun Jin (STAT. Carnegie Mellon) UPS Delivers Optimal Phase Diagram in Variable Selection

Non-optimality of lasso & subset selection

$$\epsilon_{p} = p^{-\vartheta}, \qquad \tau_{p} = \sqrt{2r\log p}, \qquad X_{i} \stackrel{iid}{\sim} N(0, \frac{1}{n}\Omega)$$

$$\beta_j = \begin{cases} \tau_p, \text{ prob. } \epsilon_p, \\ 0, \text{ prob. } 1 - \epsilon_p \end{cases}$$
$$\Omega(i,j) = 1\{i = j\} + a \cdot 1\{|i - j| = 1\}, \quad a \in (0, 1/2)$$

- UPS achieves the optimal phase diagram
- The lasso and the subset selection do not, even with tuning parameters set ideally

Non-optimal region (lasso)



$$\begin{split} \epsilon_{\rho} &= \rho^{-\vartheta}, \quad \tau_{\rho} = \sqrt{2r\log \rho}, \quad a = 0.4; \quad \lambda_{\rho}^{lasso} = \max\{\frac{r+\vartheta}{2r}, \ \frac{1}{1+\sqrt{\frac{1-a}{1+a}}}\}\tau_{\rho}; \\ \text{Two lines:} \ r &= \left(\frac{1+\sqrt{1-a^2}}{a}\right)\vartheta; \ r = (1+\sqrt{\frac{1+a}{1-a}})^2(1-\vartheta) \end{split}$$

Non-optimal region (subset selection)



$$\begin{split} \epsilon_{p} &= p^{-\vartheta}, \quad \tau_{p} = \sqrt{2r\log p}, \quad a = 0.4; \quad \lambda_{p}^{SS} = \min\{\frac{r+\vartheta}{2r}, \frac{r(1-a^{2})+2\vartheta}{2r\sqrt{1-a^{2}}}\}\tau_{p};\\ \text{Line:} \quad r &= \left((2-\sqrt{1-a^{2}})/[\sqrt{1-a^{2}}-(1-a^{2})]\right)\vartheta;\\ \text{Curve:} \quad \sqrt{r} &= \frac{1}{v}\left(\sqrt{1-2\vartheta}+\sqrt{1-2\vartheta+\vartheta v}\right), \quad v \equiv 2\sqrt{1-a^{2}}-1 \end{split}$$

Understanding lasso & subset selection

$$Y = X\beta + z \implies \tilde{Y} \approx N(\Omega\beta, \Omega)$$

Major components:

- Signal singleton: $\beta_j = \tau_p$
- Signal pairs: $\beta_j = \beta_{j+1} = \tau_p$
- Pure noise: $\beta_j = (\Omega\beta)_j = 0$
- Fake signals: $\beta_j = 0$, $(\Omega\beta)_j \neq 0$

Key Insight:

- The lasso is too **loose** on fake signals
- Subset selection is too harsh on signal pairs

How lasso/subset selection works

$$Y = X\beta + z \implies \tilde{Y} \approx N(\Omega\beta, \Omega)$$



For properly set λ_p^{lasso} , many entries of β are estimated 0; remaining entries break into small clusters, where lasso applies independently

Simulation

Hamming error/ $(p\epsilon_p)$



Left to right: $\vartheta = 0.2, 0.45, 0.6.$ x-axis: τ_p . $p = 10^4$, n = 1600, $\epsilon_p = p^{-\vartheta}$, $\pi_p = \nu_{\tau_p}$, Ω is tridiagonal (a = 0.45)

Hamming errors:

ϑ	τ_p	5	6	7	8	9	10	11
.25	UPS	49	11.1	1.79	0.26	0.02	0	0
	lasso	186.7	99.35	58.26	38.53	25.97	18.18	12.94
.50	UPS	10.06	2.11	0.37	0.09	0	0	0
	lasso	16.36	5.11	1.47	0.51	0.28	0.33	0.26
.65	UPS	5.49	1.29	0.33	0.06	0	0	0
	lasso	7.97	2.43	0.69	0.18	0.07	0.03	.02

p = 2000, Ω tridiagonal (a = 0.45), $\epsilon_p = p^{-\vartheta}$, $\pi_p = \nu_{\tau_p}$.

Theoretic boundaries for exact recovery: $\tau_p \ge 8.01$, 7.32, 7.01

- Suggested UPS as a two-stages procedure for variable selection
- Proved that UPS achieves the optimal phase diagram
- Identified non-optimal regions for
 - the lasso
 - the subset selection
- Exposed intuition for the non-optimality of penalization methods

Alphabetically:

Peter Bickel, Tony Cai, David Donoho, Jianqing Fan, Stephen Fienberg, Runze Li, Larry Wasserman, Cun-Hui Zhang for inspiration and discussion