

# Stat 577 “Statistical Learning Theory”

## T4: LASSO and SCAD estimator

Xiongzhi Chen

Washington State University

# Subgradient and subdifferential

# Convex set and convex function

- ▶ A set  $U \subseteq \mathbb{R}^n$  is convex if

$$\lambda x + (1 - \lambda) y \in U \text{ for } \forall x, y \in U \text{ and } \forall \lambda \in [0, 1]$$

- ▶ A function  $f : U \rightarrow \mathbb{R}$  is convex if

$$f(\lambda x + (1 - \lambda) y) \leq \lambda f(x) + (1 - \lambda) f(y)$$

for all  $x, y \in U$  and  $\lambda \in [0, 1]$

- ▶ A convex set does not have to be open or closed
- ▶ A convex function is differentiable almost everywhere
- ▶ Always assume  $U \subseteq \mathbb{R}^n$  to be non-empty

# Convex function: equivalent definitions

Let  $U \subseteq \mathbb{R}^n$  be open and convex. Equivalent definitions of  $f : U \rightarrow \mathbb{R}$  being convex:

1.  $f$  has convex *epigraph*

$$\text{epi}(f) = \{(x, z) \in \mathbb{R}^n \times \mathbb{R} : z \geq f(x), x \in U\}$$

2.  $f$  is the *upper envelope* of *affine minorant*, i.e.,  $\forall x_0 \in U$ ,

$$f(x_0) = \sup \{a(x_0) : a \leq f, a : x \mapsto c_a + \langle \vartheta_a, x \rangle\}$$

Each  $a$  induces a “supporting hyperplane” to  $\text{epi}(f)$  at  $(x, f(x))$

3. For  $n = 1$ , the function

$$F(x, y) = \frac{f(x) - f(y)}{(x - y)}, x \neq y$$

is non-decreasing in  $x$  (for fixed  $y$ ) and in  $y$  (for fixed  $x$ )

# Univariate convex function

- ▶ Let  $U = (a, b) \subseteq \mathbb{R}$ . Equivalent definition of a convex  $f : U \rightarrow \mathbb{R}$ :
  - ▶  $f$  has left derivative  $f'_-$  and right derivative  $f'_+$  on  $(a, b)$  (which implies  $f$  is continuous on  $(a, b)$ ), and
  - ▶ For all  $x_1 < x_2$  and  $x_1, x_2 \in (a, b)$ ,

$$f'_-(x_1) \leq f'_+(x_1) \leq \frac{f(x_2) - f(x_1)}{x_2 - x_1} \leq f'_-(x_2) \leq f'_+(x_2)$$

- ▶ Corollary:

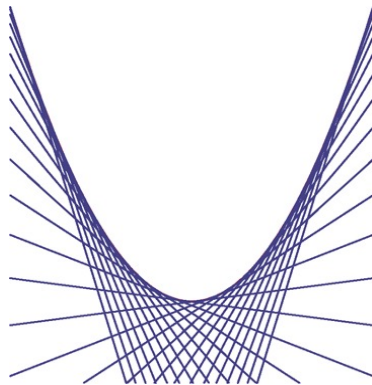
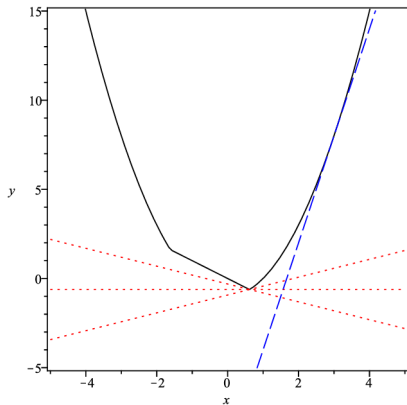
$$f'_-(x) = \sup_{y < x} \frac{f(y) - f(x)}{y - x} \quad \& \quad f'_+(x) = \inf_{y > x} \frac{f(y) - f(x)}{y - x}$$

- ▶ Subdifferential: for  $x_0 \in (a, b)$ ,

$$\begin{aligned} \partial f(x_0) &:= \{\alpha \in \mathbb{R} : f(x) - f(x_0) \geq \alpha(x - x_0), \forall x \in U\} \\ &= [f'_-(x_0), f'_+(x_0)] \end{aligned}$$

# Convex function and envelope: illustration

► *Wiki (Envelope)*



# Subgradient and subdifferential

Let  $U \subseteq \mathbb{R}^n$  be open and convex and  $f : U \rightarrow \mathbb{R}$  be convex.

- ▶ Recall “supporting hyperplanes” induced by  $a$  in

$$f(x) = \sup \{a(x) : a(x) = c_a + \langle \vartheta_a, x \rangle, a(x) \leq f(x)\} \quad (1)$$

- ▶ Subdifferential of  $f$  at  $x_0 \in U$  is defined as

$$\partial f(x_0) := \{v \in \mathbb{R}^n : f(x) - f(x_0) \geq \langle v, x - x_0 \rangle, \forall x \in U\},$$

and  $\partial f(x_0) \neq \emptyset$  whenever  $f$  is convex and continuous at  $x_0$ . An element of  $\partial f(x_0)$  is a “subgradient” of  $f$  at  $x_0$

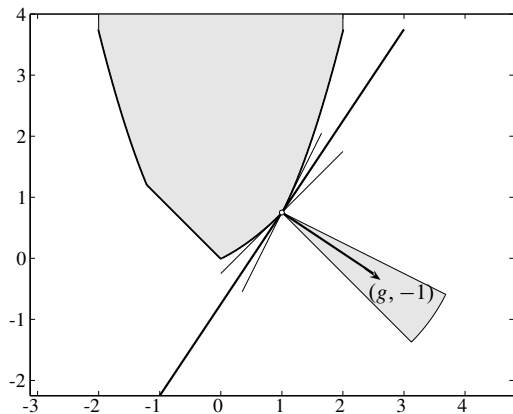
- ▶ Subdifferential  $\partial f(x_0)$  can be derived from the famous *Theorem on Separation of Convex Sets* and is closely related to supporting hyperplanes

# Subgradient: illustration

- ▶ Subgradient  $g$  of  $f$  at  $x_0$ :

$$f(x) - f(x_0) \geq \langle g, x - x_0 \rangle, \forall x \in U$$

- ▶ Image credit: “Nonlinear optimization” by Andrzej Ruszczyński



# Karush-Kuhn-Tucker (KKT) conditions

- ▶ Let  $U \subseteq \mathbb{R}^n$  be open and convex. Given  $f : U \rightarrow \mathbb{R}$ , consider

$$\min_{x \in U} f(x) \quad \text{subject to} \quad \begin{cases} g_i(x) \leq 0, 1 \leq i \leq p' \\ h_j(x) = 0, 1 \leq j \leq q' \end{cases} \quad \text{for } x \in U$$

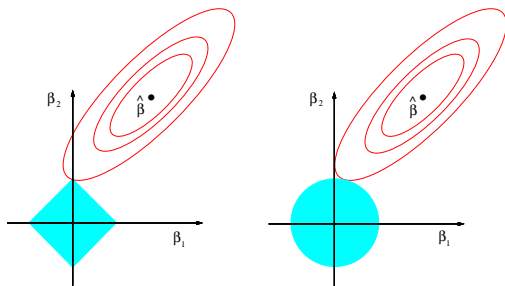
- ▶ If  $f$  and  $\{g_i\}_{i=1}^{p'}$  are convex on  $U$  and  $\{h_j\}_{j=1}^{q'}$  are linear on  $U$ , and
- ▶ If *Slater condition* holds, i.e., there exists  $x_0 \in U$  such that  $g_i(x_0) < 0$  for  $1 \leq i \leq p'$  and  $h_j(x_0) = 0$  for  $1 \leq j \leq q'$
- ▶ Then  $\hat{x}$  is a minimizer iff  $h_j(\hat{x}) = 0$  for  $1 \leq j \leq q'$  and there exist  $\hat{\lambda}_i \geq 0$  for  $1 \leq i \leq p'$  and  $\hat{\mu}_j \in \mathbb{R}$  for  $1 \leq j \leq q'$  such that

$$\begin{cases} 0 \in \partial f(\hat{x}) + \sum_{i=1}^{p'} \hat{\lambda}_i \partial g_i(\hat{x}) + \sum_{j=1}^{q'} \hat{\mu}_j \partial h_j(\hat{x}) \\ \hat{\lambda}_i g_i(\hat{x}) = 0 \text{ for } 1 \leq i \leq p' \end{cases}$$

# Solution for LASSO

# LASSO and Ridge estimates

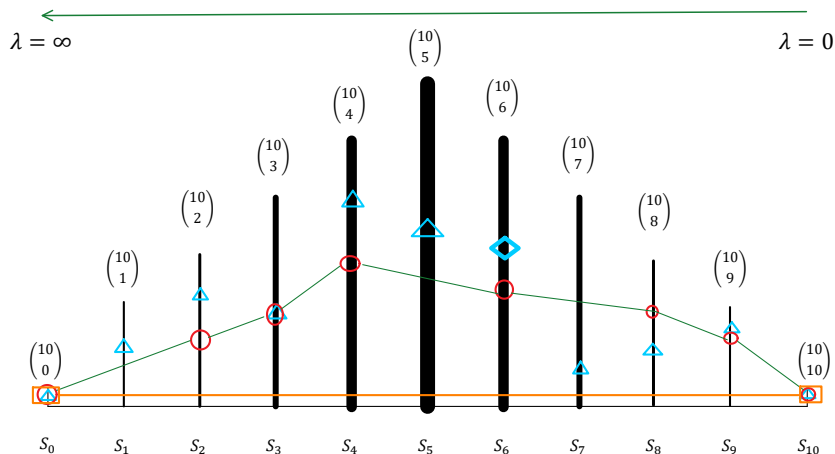
- ▶ Contour of RSS can only meet vertex on  $l_1$ -sphere (due to subdifferential) since  $l_1$ -norm is singular at 0, leading to variable selection via hard thresholding
- ▶ Contour of RSS can meet any point on  $l_2$ -sphere since  $l_2$ -norm is nonsingular at 0, leading to shrinkage only



**FIGURE 3.11.** Estimation picture for the lasso (left) and ridge regression (right). Shown are contours of the error and constraint functions. The solid blue areas are the constraint regions  $|\beta_1| + |\beta_2| \leq t$  and  $\beta_1^2 + \beta_2^2 \leq t^2$ , respectively, while the red ellipses are the contours of the least squares error function.

# Trajectory of selected models

- ▶ “Circle”, “Rectangle” and “Triangle”: models selected by LASSO, Ridge and BSS; “Diamond”: best model by BSS



# LASSO and optimization

- ▶ Model  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  with  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^T \in \mathbb{R}^p$  and LASSO

$$\begin{cases} \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 & \text{subject to } \sum_{i=1}^p |\beta_i| \leq t; \\ h_1(\boldsymbol{\beta}; \lambda) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda \sum_{i=1}^p |\beta_i| \end{cases}$$

- ▶ Subgradient of  $g(x) = |x|$  at  $x = 0$  is  $\partial g(0) = [-1, 1]$
- ▶ By KKT conditions,  $\hat{\boldsymbol{\beta}}_L(\lambda)$  is the optimizer iff  $0 \in \partial h_1(\hat{\boldsymbol{\beta}}_L(\lambda); \lambda)$ , where

$$\partial h_1(\boldsymbol{\beta}; \lambda) = \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} - \mathbf{X}^T \mathbf{y} + \lambda \mathbf{c},$$

and  $\mathbf{c} = (c_1, \dots, c_p)^T$  and

$$\begin{cases} c_i = \text{sgn}(\beta_i) & \text{if } \beta_i \neq 0 \\ c_i \in [-1, 1] & \text{if } \beta_i = 0 \end{cases}$$

# Orthogonal design and solution

- ▶ If  $\mathbf{X}^T \mathbf{X} = \mathbf{I}$ , then LSE  $\hat{\boldsymbol{\beta}} = \mathbf{X}^T \mathbf{y} = (\hat{\beta}_1, \dots, \hat{\beta}_p)$  and  $\hat{\boldsymbol{\beta}}_L(\lambda) = (\hat{\beta}_{L,1}(\lambda), \dots, \hat{\beta}_{L,p}(\lambda))^T$  satisfies

$$\hat{\boldsymbol{\beta}}_L(\lambda) - \hat{\boldsymbol{\beta}} + \lambda \mathbf{c} = 0 \text{ with } \begin{cases} c_i = \text{sgn}(\hat{\beta}_{L,i}(\lambda)) & \text{if } \hat{\beta}_{L,i}(\lambda) \neq 0 \\ c_i \in [-1, 1] & \text{if } \hat{\beta}_{L,i}(\lambda) = 0 \end{cases}$$

- ▶ Thresholding:  $\hat{\beta}_{L,i}(\lambda) = 0 \iff \hat{\beta}_i - \lambda c_i = 0, |c_i| \leq 1 \iff |\hat{\beta}_i| \leq \lambda$
- ▶ Soft thresholding and equivalences:  $\hat{\beta}_{L,i}(\lambda) \neq 0 \iff$ 
  - ▶  $\hat{\beta}_{L,i}(\lambda) - \hat{\beta}_i + \lambda \text{sgn}(\hat{\beta}_{L,i}(\lambda)) = 0$
  - ▶  $\hat{\beta}_{L,i}(\lambda) = \hat{\beta}_i - \lambda$  if  $\hat{\beta}_{L,i}(\lambda) > 0$  but  $\hat{\beta}_{L,i}(\lambda) = \hat{\beta}_i + \lambda$  if  $\hat{\beta}_{L,i}(\lambda) < 0$
  - ▶  $\hat{\beta}_i > \lambda$  and  $\hat{\beta}_{L,i}(\lambda) = \hat{\beta}_i - \lambda$  or  $\hat{\beta}_i < -\lambda$  and  $\hat{\beta}_{L,i}(\lambda) = \hat{\beta}_i + \lambda$
  - ▶  $\hat{\beta}_{L,i}(\lambda) = \text{sgn}(\hat{\beta}_i) (|\hat{\beta}_i| - \lambda)$
- ▶ Summary:  $\hat{\beta}_{L,i}(\lambda) = \text{sgn}(\hat{\beta}_i) (|\hat{\beta}_i| - \lambda)_+$

# Visualization of estimate for orthogonal design

- ▶ If  $\mathbf{X}^T \mathbf{X} = \mathbf{I}$ , then  $\hat{\boldsymbol{\beta}} = \mathbf{X}^T \mathbf{y}$  and  $\hat{\beta}_{L,i}(\lambda) = \text{sgn}(\hat{\beta}_i) (|\hat{\beta}_i| - \lambda)_+$ 
  - ▶  $\hat{\beta}_{L,i}(\lambda)$  is always biased when  $|\hat{\beta}_i| > \lambda > 0$ , i.e., *intrinsic bias*
  - ▶ For Gaussian error, explicit distribution  $F_{i,\lambda}$  for  $\hat{\beta}_{L,i}(\lambda)$  for each *deterministic*  $\lambda$  can be obtained. However, once optimal  $\lambda$  is chosen from data, such is often impossible to obtain

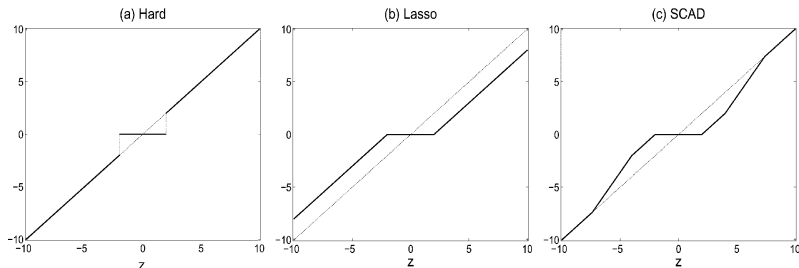


Figure 2. Plot of Thresholding Functions for (a) the Hard, (b) the Soft, and (c) the SCAD Thresholding Functions With  $\lambda = 2$  and  $a = 3.7$  for SCAD.

# Consistency of LASSO

# Intrinsic bias of LASSO estimate

If  $\mathbf{X}^T \mathbf{X} = \mathbf{I}$ , then LSE  $\hat{\beta} = \mathbf{X}^T \mathbf{y}$  and  $\hat{\beta}_{L,i}(\lambda) = \text{sgn}(\hat{\beta}_i) (|\hat{\beta}_i| - \lambda)_+$  and

- ▶  $\hat{\beta}_{L,i}(\lambda)$  is biased when  $|\hat{\beta}_i| > \lambda > 0$
- ▶  $\hat{\beta}_{L,i}(\lambda)$  is biased when  $0 < |\hat{\beta}_i| < \lambda$  for  $\lambda > 0$
- ▶ Bias caused by  $\frac{d}{dx} \varrho_\lambda(|x|) > 0$  for large  $|x|$  and non-differentiability of  $\varrho_\lambda(|x|)$  at  $x = 0$

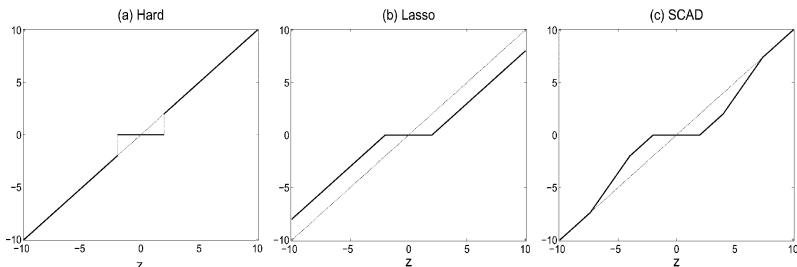


Figure 2. Plot of Thresholding Functions for (a) the Hard, (b) the Soft, and (c) the SCAD Thresholding Functions With  $\lambda = 2$  and  $a = 3.7$  for SCAD.

# Bias of LASSO estimate

If  $\mathbf{X}^T \mathbf{X} = \mathbf{I}$ , then  $\hat{\boldsymbol{\beta}} = \mathbf{X}^T \mathbf{y}$  and  $\hat{\beta}_{L,i}(\lambda) = \text{sgn}(\hat{\beta}_i) (|\hat{\beta}_i| - \lambda)_+$

- ▶ Assume  $|\hat{\beta}_i| > \lambda > 0$ . How can the bias of  $\hat{\beta}_{L,i}(\lambda)$  be 0?
- ▶ Assume  $0 \leq |\hat{\beta}_i| < \lambda$  for  $\lambda > 0$ . How can the bias of  $\hat{\beta}_{L,i}(\lambda)$  be 0?
- ▶ Assume  $\lambda = 0$ . Which  $\hat{\beta}_{L,i}(\lambda)$  among  $1 \leq i \leq p$  has 0 bias?
- ▶ Under what additional conditions will  $\hat{\boldsymbol{\beta}}_L(\lambda)$  be consistent?

# LASSO estimate for general design

- ▶ Let  $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ ,  $\hat{\boldsymbol{\beta}}_L(\lambda) = (\hat{\beta}_{L,1}(\lambda), \dots, \hat{\beta}_{L,p}(\lambda))^T$  and

$$\begin{cases} c_i = \text{sgn}(\hat{\beta}_{L,i}(\lambda)) & \text{if } \hat{\beta}_{L,i}(\lambda) \neq 0 \\ c_i \in [-1, 1] & \text{if } \hat{\beta}_{L,i}(\lambda) = 0 \end{cases}$$

- ▶ Let  $\tilde{\mathbf{x}}_j$  be the  $j$ th column of  $\mathbf{X}$ . The following are equivalent:

- ▶  $\mathbf{X}^T \mathbf{X} \hat{\boldsymbol{\beta}}_L(\lambda) - \mathbf{X}^T \mathbf{y} + \lambda \mathbf{c} = 0 \iff$

- ▶  $\hat{\boldsymbol{\beta}}_L(\lambda) = \hat{\boldsymbol{\beta}} - \lambda (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{c} = (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{y} - \lambda \mathbf{c}) \iff$

- ▶  $\begin{cases} \tilde{\mathbf{x}}_i^T (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}_L(\lambda)) = \lambda \text{sgn}(\hat{\beta}_{L,i}(\lambda)) & \text{if } \hat{\beta}_{L,i}(\lambda) \neq 0 \\ |\tilde{\mathbf{x}}_i^T (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}_L(\lambda))| \leq \lambda & \text{if } \hat{\beta}_{L,i}(\lambda) = 0 \end{cases}$

- ▶ LASSO estimate in general *shrinkage LSE* but is *very complicated*

# Framework for theoretical investigation

- ▶ Sparsity  $\|\beta\|_0 = s$ :  $\beta_1 = (\beta_1, \dots, \beta_s)^T$ ;  $\beta_0 = 0^T$ ;  $\mathcal{A} = \{1, \dots, s\}$
- ▶ Split:  $\mathbf{X} = (\mathbf{X}_1, \mathbf{X}_0)$  with  $\mathbf{X}_1 = (\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_s)$  and  $\mathbf{X}_0 = (\tilde{\mathbf{x}}_{s+1}, \dots, \tilde{\mathbf{x}}_p)$ , i.e.,

$$\beta = \begin{pmatrix} \beta_1 \\ \beta_0 \end{pmatrix} \mapsto \mathbf{X} = (\mathbf{X}_1, \mathbf{X}_0)$$

- ▶ Recall  $\hat{\beta}_L(\lambda) = (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{y} - \lambda \mathbf{c})$  or

$$\begin{cases} \tilde{\mathbf{x}}_i^T (\mathbf{y} - \mathbf{X} \hat{\beta}_L(\lambda)) = \lambda \operatorname{sgn}(\hat{\beta}_{L,i}(\lambda)) & \text{if } \hat{\beta}_{L,i}(\lambda) \neq 0 \\ \left| \tilde{\mathbf{x}}_i^T (\mathbf{y} - \mathbf{X} \hat{\beta}_L(\lambda)) \right| \leq \lambda & \text{if } \hat{\beta}_{L,i}(\lambda) = 0 \end{cases} \quad (2)$$

- ▶ Estimated sparsity  $\|\hat{\beta}_L(\lambda)\|_0 = \hat{s}$  and  $\hat{\mathcal{A}} = \{i : \hat{\beta}_{L,i}(\lambda) \neq 0\}$
- ▶ Strategy on theory: investigate (2)

# Intuition on consistency of LASSO

- ▶ Recall  $\hat{\beta}_L(\lambda) = (\hat{\beta}_{L,1}(\lambda), \dots, \hat{\beta}_{L,p}(\lambda))^T$  satisfies

$$\hat{\beta}_L(\lambda) = (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{y} - \lambda \mathbf{c}), \quad (3)$$

$$\text{where } \begin{cases} c_i = \text{sgn}(\hat{\beta}_{L,i}(\lambda)) & \text{if } \hat{\beta}_{L,i}(\lambda) \neq 0 \\ c_i \in [-1, 1] & \text{if } \hat{\beta}_{L,i}(\lambda) = 0 \end{cases}$$

- ▶ Recall  $\beta_1 = (\beta_1, \dots, \beta_s)^T$ ,  $\beta_0 = 0^T$ ,  $\mathcal{A} = \{1, \dots, s\}$  and splitting

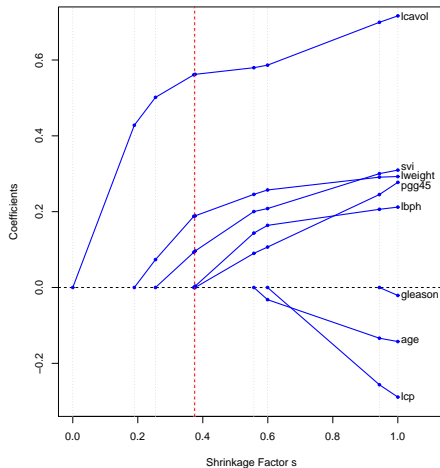
$$\beta = \begin{pmatrix} \beta_1 \\ \beta_0 \end{pmatrix} \mapsto \mathbf{X} = (\mathbf{X}_1, \mathbf{X}_0) \mapsto \mathbf{C} = \mathbf{X}^T \mathbf{X} = \begin{pmatrix} \mathbf{X}_1^T \mathbf{X}_1 & \mathbf{X}_1^T \mathbf{X}_0 \\ \mathbf{X}_0^T \mathbf{X}_1 & \mathbf{X}_0^T \mathbf{X}_0 \end{pmatrix}$$

- ▶ Note  $\mathbf{c}$  is a function of  $\hat{\beta}_L(\lambda)$ . So,  $\mathbf{X}_1$  and  $\mathbf{X}_0$  have to satisfy some conditions for  $\hat{\beta}_L(\lambda)$  to achieve selection consistency and/or estimation consistency

# LASSO consistency: iff condition

- ▶ Recall  $\hat{\beta}_L(\lambda) = (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{y} - \lambda \mathbf{c})$ , where  $c_i = \text{sgn}(\hat{\beta}_{L,i}(\lambda))$  if  $\hat{\beta}_{L,i}(\lambda) \neq 0$  but  $|c_i| \leq 1$  if  $\hat{\beta}_{L,i}(\lambda) = 0$
- ▶ Recall  $\mathbf{C} = \mathbf{X}^T \mathbf{X} = \begin{pmatrix} \mathbf{X}_1^T \mathbf{X}_1 & \mathbf{X}_1^T \mathbf{X}_0 \\ \mathbf{X}_0^T \mathbf{X}_1 & \mathbf{X}_0^T \mathbf{X}_0 \end{pmatrix} = \begin{pmatrix} \mathbf{C}_{11} & \mathbf{C}_{12} \\ \mathbf{C}_{21} & \mathbf{C}_{22} \end{pmatrix}$
- ▶ *Irrepresentability* is about  $\mathbf{C}_{21} \mathbf{C}_{11}^{-1} \mathbf{s}$ , where  $\mathbf{s} = \text{sgn}(\beta_1)$  and about  $\|\beta_1\|_1$  versus  $\|\lambda \mathbf{C}_{11}^{-1} \mathbf{s}\|_1$
- ▶ “Necessary and sufficient conditions for variable selection consistency of the LASSO in high dimensions” by Soumendra N. Lahiri, *Annals of Statistics*, 2021
  - ▶ Namely, lower irrepresentable condition + upper irrepresentable condition == variable selection consistency
  - ▶ *Surprise*: LASSO cannot simultaneously achieve variable selection consistency and  $\sqrt{n}$ -consistency

# LASSO solution path



**FIGURE 3.10.** Profiles of lasso coefficients, as the tuning parameter  $t$  is varied. Coefficients are plotted versus  $s = t / \sum_1^p |\beta_j|$ . A vertical line is drawn at  $s = 0.36$ , the value chosen by cross-validation. Compare Figure 3.8 on page 65; the lasso profiles hit zero, while those for ridge do not. The profiles are piece-wise linear, and so are computed only at the points displayed; see Section 3.4.4 for details.

- ▶ LSE  $\hat{\beta} = (\hat{\beta}_1, \dots, \hat{\beta}_p)^T$
- ▶  $\hat{\beta}_L$  LASSO estimate
- ▶  $t = \left\| \hat{\beta}_L(\lambda) \right\|_1$
- ▶  $\hat{s} = 0.36$  chosen by 10-fold cross-validation
- ▶ Caution: magnitudes of  $\hat{\beta}_{L,i}$ 's not necessarily monotone decreasing in  $s$ , even though they are for this example

# SCAD solution

# SCAD penalty and estimate for orthogonal design

- ▶ SCAD penalty is a *quadratic spline* and is *not convex*
- ▶ Derivative of SCAD penalty for  $a > 2$  is

$$e'_\lambda(\theta) = \lambda \left[ 1_{\{\theta \leq \lambda\}}(\theta) + \frac{(a\lambda - \theta)_+}{(a-1)\lambda} 1_{\{\theta > \lambda\}}(\theta) \right], \theta > 0$$

- ▶ If  $\mathbf{X}^T \mathbf{X} = \mathbf{I}$ , then LSE  $\hat{\boldsymbol{\beta}} = \mathbf{X}^T \mathbf{y}$  and SCAD estimate

$$\hat{\beta}_{S,i}(\lambda) = \begin{cases} \operatorname{sgn}(\hat{\beta}_i) (|\hat{\beta}_i| - \lambda)_+ & \text{if } |\hat{\beta}_i| \leq 2\lambda \\ \frac{(a-1)\hat{\beta}_i - \operatorname{sgn}(\hat{\beta}_i) a\lambda}{a-2} & \text{if } 2\lambda < |\hat{\beta}_i| \leq a\lambda \\ \hat{\beta}_i & \text{if } |\hat{\beta}_i| > a\lambda \end{cases}$$

and LASSO estimate

$$\hat{\beta}_{L,i}(\lambda) = \operatorname{sgn}(\hat{\beta}_i) (|\hat{\beta}_i| - \lambda)_+$$

- ▶ Note the differences between “soft shrinkage” parts

# Visualizing estimates for orthogonal design

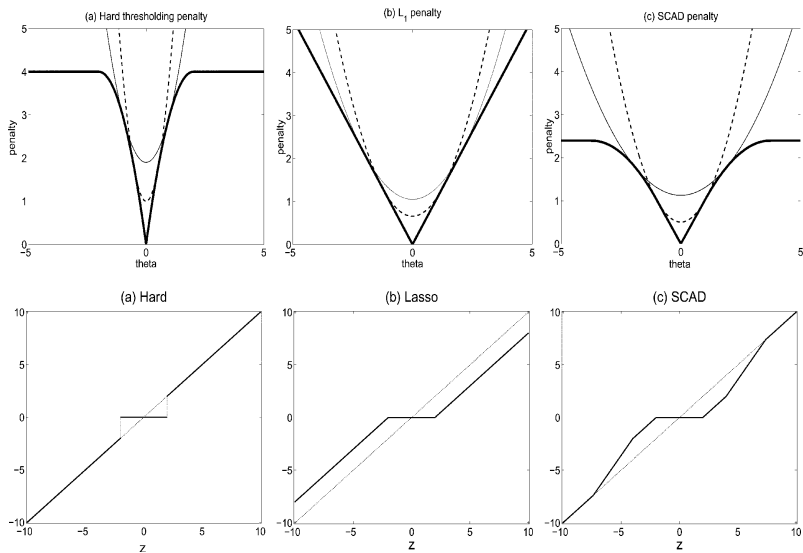


Figure 2. Plot of Thresholding Functions for (a) the Hard, (b) the Soft, and (c) the SCAD Thresholding Functions With  $\lambda = 2$  and  $a = 3.7$  for SCAD.

# SCAD estimator in general

- ▶ Since SCAD penalty  $\varrho_\lambda(|x|)$  is *not convex*, KKT conditions are *usually only necessary*
- ▶ Penalized RSS as objective function

$$h_S(\boldsymbol{\beta}; \lambda) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \sum_{i=1}^p \varrho_\lambda(|\beta_i|)$$

and optimizer  $\hat{\boldsymbol{\beta}}_S(\lambda) = (\hat{\beta}_{S,1}(\lambda), \dots, \hat{\beta}_{S,p}(\lambda))^T$

- ▶ Let  $\tilde{\mathbf{x}}_j$  be the  $j$ th column of  $\mathbf{X}$ . Necessary condition on  $\hat{\boldsymbol{\beta}}_S(\lambda)$ :

$$\begin{cases} \tilde{\mathbf{x}}_i^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}_S(\lambda)) = \text{sgn}(\hat{\beta}_{S,i}(\lambda)) \varrho'_\lambda(|\hat{\beta}_{S,i}(\lambda)|) \text{ if } \hat{\beta}_{S,i}(\lambda) \neq 0 \\ \left| \tilde{\mathbf{x}}_i^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}_S(\lambda)) \right| \leq \lambda \text{ if } \hat{\beta}_{S,i}(\lambda) = 0 \end{cases}$$

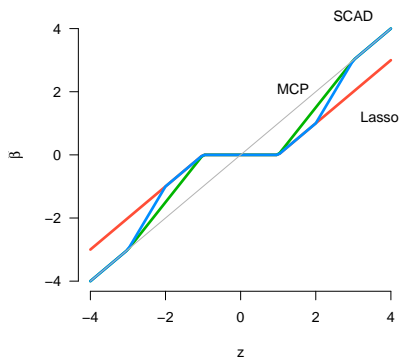
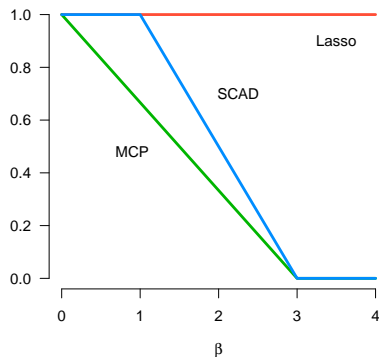
since  $\varrho_\lambda(|x|) = \lambda|x|$  for  $x$  close to 0

# Non-uniqueness of SCAD optimizer

- ▶ Since SCAD penalty  $\rho_\lambda(|x|)$  is not convex, optimizer  $\hat{\beta}_S(\lambda)$  is not unique in general
- ▶ Despite nearly unbiasedness of  $\hat{\beta}_S(\lambda)$ , SCAD estimate is not as popular as LASSO estimate since the former is harder to compute than the latter
- ▶ Through years of research, Prof. Jianqing Fan and his collaborators found that, in some situations, any local optimizer  $\hat{\beta}_S(\lambda)$  is good enough, thus avoiding global search of optimizers of  $h_S(\beta; \lambda)$
- ▶ SCAD estimate, if it is good, often outperforms LASSO estimate in terms of accuracy
- ▶ *We know little about performances of high-dimensional models under non-convex penalties*

# Minimax concave penalty (MCP)

- ▶ MCP further improves SCAD; see Zhang (2010)
- ▶ MCP is not convex, and it reduces bias in estimating moderate sized effects, compared to SCAD and LASSO



# Consistency of SCAD estimator

# Framework on theoretical investigation

- ▶ Recall SCAD estimate  $\hat{\beta}_S(\lambda)$  as

$$\begin{cases} \tilde{\mathbf{x}}_i^T (\mathbf{y} - \mathbf{X}\hat{\beta}_S(\lambda)) = \text{sgn}(\hat{\beta}_{S,i}(\lambda)) \varrho'_\lambda(|\hat{\beta}_{S,i}(\lambda)|) \text{ if } \hat{\beta}_{S,i}(\lambda) \neq 0 \\ \left| \tilde{\mathbf{x}}_i^T (\mathbf{y} - \mathbf{X}\hat{\beta}_S(\lambda)) \right| \leq \lambda \text{ if } \hat{\beta}_{S,i}(\lambda) = 0 \end{cases}$$

- ▶ Selection consistency is similar to that of LASSO
- ▶ Estimation consistency should be better than that of LASSO
- ▶ The above representation of  $\hat{\beta}_S(\lambda)$  can be the starting point of theoretical investigation on  $\hat{\beta}_S(\lambda)$
- ▶ An alternative framework has been provided by, e.g., Fan and Li (2001) and Fan and Peng (2004) and Fan and Lv (2011)

# Penalized log-likelihood

- ▶ Recall penalized RSS as objective function

$$h_S(\boldsymbol{\beta}; \lambda) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \sum_{i=1}^p \varrho_\lambda(|\beta_i|)$$

- ▶ To accommodate model

$$\begin{cases} \Pr(Y \in A|X) = \int_A f(g^{-1}(X^T \boldsymbol{\beta}), y) dy \\ g(E(Y|X)) = X^T \boldsymbol{\beta} \end{cases}$$

and asymptotic study, consider *penalized log-likelihood*

$$h(\boldsymbol{\beta}; \lambda) = -n^{-1} \sum_{i=1}^n \ln f(g^{-1}(\mathbf{x}_i^T \boldsymbol{\beta}), y_i) + \sum_{i=1}^p \varrho_\lambda(|\beta_i|)$$

for i.i.d. observations  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$  from  $(X, Y)$

# Oracle property of optimizer

- ▶ Sparsity  $\|\beta\|_0 = s$ :  $\beta_1 = (\beta_1, \dots, \beta_s)^T$ ;  $\beta_0 = 0^T$ ;  $\mathcal{A} = \{1, \dots, s\}$
- ▶ Optimizer  $\hat{\beta}(\lambda) = (\hat{\beta}_1(\lambda), \dots, \hat{\beta}_p(\lambda))^T$  such that

$$\hat{\beta}(\lambda) \in \operatorname{argmax}_{\beta \in \mathbb{R}^p} h(\beta; \lambda)$$

- ▶ Oracle property of  $\hat{\beta}(\lambda) = (\hat{\beta}_1(\lambda)^T, \hat{\beta}_2(\lambda)^T)^T$  as

$$\left\{ \begin{array}{l} \text{Matching sparsity: } \hat{\beta}_2(\lambda) = \beta_0 = 0 \\ \text{Asymptotic Normality: } \sqrt{n}A_n \left( \hat{\beta}_1(\lambda) - \beta_1 + \mathbf{b}_n \right) \\ \rightsquigarrow \text{Normal}(0, \mathbf{I}) \end{array} \right.$$

- ▶ Oracle property is stronger than

$$\sqrt{n}A_n \left( \hat{\beta}(\lambda) - \beta + \mathbf{b}_n \right) \rightsquigarrow \text{Normal}(0, \mathbf{I})$$

# Alternative framework on consistency

- ▶ Existence of  $\sqrt{n}$ -consistent estimate: Pick suitable  $\alpha_n = O(n^{-1/2})$ . For any  $\varepsilon > 0$ , there exists constant  $C > 0$  such that

$$\Pr \left( \sup_{\|\mathbf{u}\|=C} h(\boldsymbol{\beta} + \alpha_n \mathbf{u}; \lambda) < h(\boldsymbol{\beta}; \lambda) \right) > 1 - \varepsilon$$

iff with probability tending to one,  $h(\cdot; \lambda)$  has a local maximum  $\hat{\boldsymbol{\beta}}(\lambda)$  in the ball

$$B_C(\boldsymbol{\beta}) = \{\boldsymbol{\beta} + \alpha_n \mathbf{u} : \|\mathbf{u}\| \leq C\} \quad \text{and} \quad \|\hat{\boldsymbol{\beta}}(\lambda) - \boldsymbol{\beta}\| = O_{\Pr}(\alpha_n)$$

- ▶ “Matching sparsity lemma”: for any given  $\tilde{\boldsymbol{\beta}} = (\tilde{\boldsymbol{\beta}}_1^T, \tilde{\boldsymbol{\beta}}_2^T)^T$  such that  $\|\tilde{\boldsymbol{\beta}}_1 - \boldsymbol{\beta}_1\| = O_{\Pr}(n^{-1/2})$  and any constant  $C > 0$ ,

$$h \left( \left( \tilde{\boldsymbol{\beta}}_1^T, 0^T \right)^T ; \lambda \right) = \max_{\|\tilde{\boldsymbol{\beta}}_2\| \leq Cn^{-1/2}} h \left( \left( \tilde{\boldsymbol{\beta}}_1^T, \tilde{\boldsymbol{\beta}}_2^T \right)^T ; \lambda \right)$$

# Summary on alternative framework on consistency

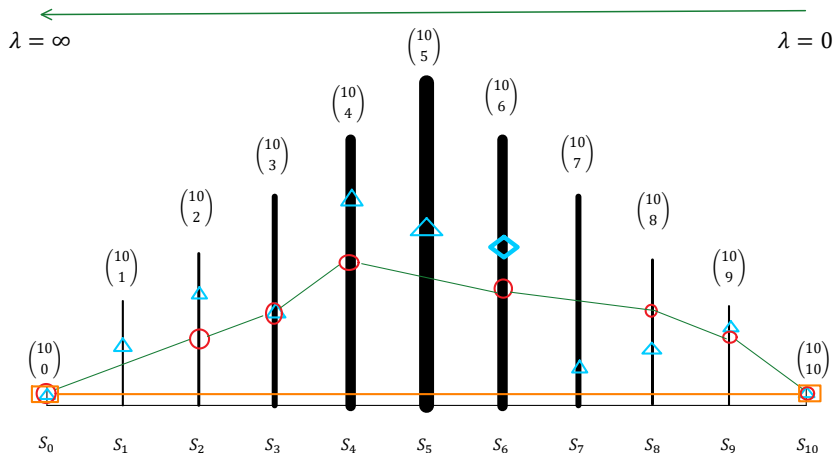
- ▶ Existence of  $\sqrt{n}$ -consistent estimate and matching sparsity lemma
  - ▶ form the backbone of this framework
  - ▶ will eventually lead to oracle property
  - ▶ rely heavily on properties of *sample size, number of coefficients, sparsity, magnitudes of non-zero coefficients, likelihood function, penalty, and tuning parameter* (which is true for other methods)
  - ▶ rely heavily on how the above *seven elements interact with each other* (which is true for other methods)
- ▶ The general principles of Fan and Li (2001) to achieve oracle property work for the regime of a diverging number of parameters (see, e.g., Fan and Peng (2004) and Fan and Lv (2011))
- ▶ Note that properties of design matrix has been incorporated into conditions that govern how likelihood function, penalty and tuning parameter should interact with each other

# Traditional versus modern

- ▶ Traditional variable/model selection *penalizes/controls model size*
  - ▶ Search algorithm: Best subset selection or forward/backward stepwise selection, all driven by, e.g., RSS
  - ▶ Model selection criteria: AIC; BIC; Mallows'  $C_p$ ; adjusted R-squared; *cross-validation*
- ▶ Modern variable/model selection *penalizes/controls effect size*
  - ▶ Search algorithm: "homotopy" via tuning parameter, i.e., the family of models searched is  $\{\mathcal{M}_\lambda : \lambda \geq 0\}$
  - ▶ Model selection criterion: maximizing penalized/regularized (log)-likelihood function
  - ▶ Tuning parameter selection via *cross-validation and estimated prediction risk*
  - ▶ Implementation: pathwise coordinate optimization

# Search algorithm for model space

- ▶ “Circle”, “Rectangle” and “Triangle”: models selected by LASSO, Ridge and BSS; “Diamond”: best model by BSS



# Visualizing estimates for orthogonal design

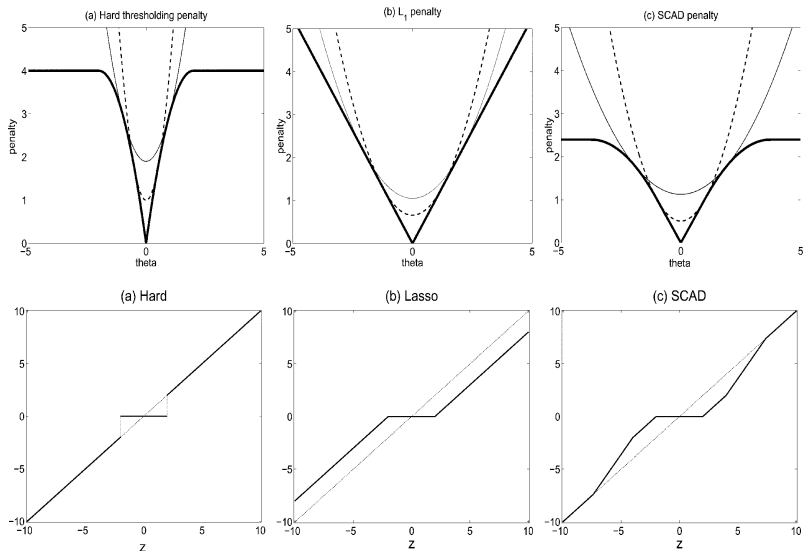


Figure 2. Plot of Thresholding Functions for (a) the Hard, (b) the Soft, and (c) the SCAD Thresholding Functions With  $\lambda = 2$  and  $a = 3.7$  for SCAD.

- Fan, J. and Li, R. (2001). Variable selection via nonconcave penalized likelihood and its oracle properties. *J. Amer. Statist. Assoc.*, 96:1348–1360.
- Fan, J. and Lv, J. (2011). Non-concave penalized likelihood with np-dimensionality. *IEEE Trans Inf Theory*, 57(8):5467–5484.
- Fan, J. and Peng, H. (2004). Nonconcave penalized likelihood with a diverging number of paramters. *Ann. Statist.*, 32(2):928–961.
- Zhang, C.-H. (2010). Nearly unbiased variable selection under minimax concave penalty. *Ann. Statist.*, 38(2):894–942.