

Stat 577 “Statistical Learning Theory”

T3: Designing a penalty and Ridge estimator

Xiongzhi Chen

Washington State University

Regularization and shrinkage

Secrets behind shrinkage

▶ Stein's phenomenon:

- ▶ Given a single observation \mathbf{z}_0 from $\mathbf{z} \sim \text{Normal}(\boldsymbol{\theta}, \mathbf{I})$ with $\mathbf{z} \in \mathbb{R}^p$. If $p \geq 3$, then the LSE $\hat{\boldsymbol{\theta}} = \mathbf{z}_0$ of $\boldsymbol{\theta}$ is inadmissible under the loss $E \left[\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|^2 \right]$ compared to the estimator James-Stein estimator

$$\hat{\boldsymbol{\theta}}_{\text{JS}} = \mathbf{z}_0 \left[1 - (p - 2) \|\mathbf{z}_0\|^{-2} \right]$$

▶ Hodges's super-efficiency phenomenon:

- ▶ Given n i.i.d. observations $\{z_i\}_{i=1}^n$ from $z \sim \text{Normal}(\theta, 1)$ with $z \in \mathbb{R}$. Hodges estimator

$$\hat{\theta}_{n,H} = \begin{cases} \bar{z}_n & \text{if } |\bar{z}_n| \geq n^{-1/4} \\ 0 & \text{if } |\bar{z}_n| < n^{-1/4} \end{cases} \Rightarrow \begin{cases} n^\alpha (\hat{\theta}_{n,H} - \theta) \rightsquigarrow 0, \forall \alpha \in \mathbb{R} \\ \sqrt{n} (\hat{\theta}_{n,H} - \theta) \sim \sqrt{n} (\bar{z}_n - \theta) \end{cases}$$

Norms and sparsity

Some norms on Euclidean space

For a vector $\mathbf{x} = (x_1, \dots, x_n)^T \in \mathbb{R}^n$, we have

l_0 -“norm”:	$\ \mathbf{x}\ _0 = \sum_{i=1}^n \mathbf{1}_{\{x_i \neq 0\}}(x_i)$
l_p -norm:	$\ \mathbf{x}\ _p = \left(\sum_{i=1}^n x_i ^p\right)^{1/p}, 0 < p < \infty$
l_∞ -norm:	$\ \mathbf{x}\ _\infty = \max_{1 \leq i \leq n} x_i $

- ▶ $\|\mathbf{x}\|_0$ is the number of nonzero entries of \mathbf{x} , Euclidean norm of \mathbf{x} is $\|\mathbf{x}\|_2$, and l_0 -“norm” is not a norm
- ▶ Define $0^0 = 0$. Then

$$\lim_{p \rightarrow 0^+} \|\mathbf{x}\|_p = \|\mathbf{x}\|_0 \quad \text{and} \quad \lim_{p \rightarrow \infty} \|\mathbf{x}\|_p = \|\mathbf{x}\|_\infty$$

- ▶ Different l_p -norms usually induce different l_p -balls

Norms: axioms and properties

- ▶ Definition of *norm* $\|\cdot\|_* : \mathcal{V} \rightarrow \mathbb{R}$:
 - ▶ Positive-definiteness: $\|\mathbf{x}\|_* \geq 0, \forall \mathbf{x}$ and $\|\mathbf{x}\|_* = 0$ iff $\mathbf{x} = 0$
 - ▶ Triangular inequality: $\|\mathbf{x} + \mathbf{y}\|_* \leq \|\mathbf{x}\|_* + \|\mathbf{y}\|_*$
 - ▶ Absolute homogeneity: $\|\lambda \mathbf{x}\|_* = |\lambda| \|\mathbf{x}\|_*, \forall \lambda \in \mathbb{R}$ and $\forall \mathbf{x}$
- ▶ $\|\cdot\|_*$ must be convex and Lipschitz, i.e.,

$$\begin{cases} \|\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}\|_* \leq |\lambda| \|\mathbf{x}\|_* + |1 - \lambda| \|\mathbf{y}\|_* \\ \left| \|\mathbf{x}\|_* - \|\mathbf{y}\|_* \right| \leq \|\mathbf{x} - \mathbf{y}\|_* \end{cases}$$

- ▶ l_p -norm, $0 < p \leq \infty$, are norms due to *Minkowski inequality*

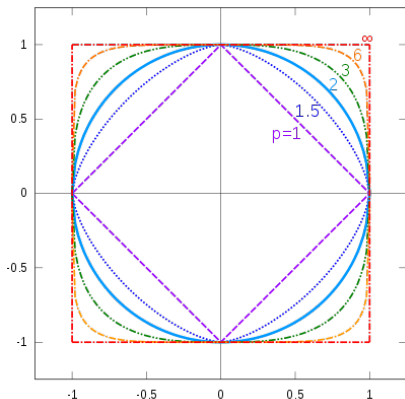
$$\|\mathbf{x} + \mathbf{y}\|_p \leq \|\mathbf{x}\|_p + \|\mathbf{y}\|_p$$

(which can be proved by *Hölder inequality*)

- ▶ l_0 -“norm” is *not a norm* since it is NOT homogeneous

Geometry: visualization

- Unit l_p -ball $B_p = \{ \mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\|_p \leq 1 \}$ in \mathbb{R}^2 (Credit: Wiki)



Linear model with scalar response

- ▶ Predictor $X = (X_1, \dots, X_p)^T \in \mathbb{R}^p$ and response $Y \in \mathbb{R}$
- ▶ Coefficient vector $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^T \in \mathbb{R}^p$ and population model

$$Y = \beta_0 + \sum_{j=1}^p X_j \beta_j + \varepsilon, \quad (1)$$

where $\varepsilon \in \mathbb{R}$ is the random error term $E(\varepsilon) = 0$ and $\text{var}(\varepsilon) = \sigma^2$

- ▶ Data version: set $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_N)^T \in \mathbb{R}^N$ and each $\varepsilon_i \sim \varepsilon$, then

$$\mathbf{y} = \beta_0 + \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

- ▶ Variable/model selection for (1) is the same as retaining X_j whose $\beta_j \neq 0$

Searching through model space

Always keep β_0 in $Y = \beta_0 + \sum_{j=1}^p X_j \beta_j + \varepsilon$

- ▶ Continuous in model size (and *controlling model size*):
 - ▶ Exhaustive: Best subset selection (BSS)
 - ▶ Non-exhaustive: Backward-stepwise selection; forward-stepwise selection
- ▶ Discontinuous in model size (and *controlling coefficient size*):
 - ▶ Non-exhaustive: regularization such as LASSO and “least angle regression (LAR)”
 - ▶ Non-exhaustive: implicit regularization via use of priors
- ▶ Hybrid:
 - ▶ Warm start without regularization and then search model space via regularization (such as “two-stage methods” including adaptive LASSO of Hui ZOU)

Concepts on sparsity

For an unknown coefficient vector $\beta = (\beta_1, \dots, \beta_p)^T \in \mathbb{R}^p$, a general concept of “sparsity” ϑ is

“the proportion of entries of β whose magnitude exceed a specific threshold”

- ▶ ϑ for variable selection (in linear model): $\|\beta\|_0 \leq s$
- ▶ ϑ originated from quantum physics: $\|\beta\|_q \leq s$ for $0 < q \leq \infty$

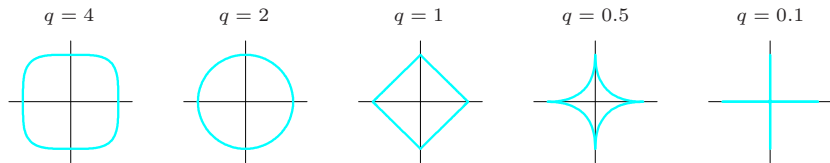


FIGURE 3.12. Contours of constant value of $\sum_j |\beta_j|^q$ for given values of q .

False discovery rate

Variable selection and multiple testing

- ▶ Model: $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^T \in \mathbb{R}^p$ and

$$Y = \beta_0 + \sum_{j=1}^p X_j \beta_j + \varepsilon \Leftrightarrow \mathbf{y} = \beta_0 + \mathbf{X}\boldsymbol{\beta} + \varepsilon$$

- ▶ Sparsity: $\|\boldsymbol{\beta}\|_0 = s$; e.g., $\min_{1 \leq i \leq s} |\beta_j| > 0$ and $\max_{s < j \leq p} |\beta_j| = 0$
- ▶ Hypotheses: $H_{j0} : \beta_j = 0$ versus $H_{j1} : \beta_j \neq 0$ for $1 \leq j \leq p$
- ▶ Classification table for H_{j0} :

	Truth	
	$\beta_j = 0$	$\beta_j \neq 0$
Selected (reject H_{j0})	False positive	True positive
Excluded (retain H_{j0})	True negative	False negative

Variable selection and false discovery rate

- ▶ Classification table for all $\{H_{j0}\}_{j=1}^p$ upon variable selection:

	Null is true	Null is false	Total
Reject null	V	S	R
Retain null	U	T	$p - R$
Total	$p - s$	s	p

- ▶ The false discovery rate (FDR) α of the selection procedure \mathcal{R} :

$$\alpha(\mathcal{R}) = \mathbb{E} \left[\frac{V}{\max\{R, 1\}} \right]$$

- ▶ The proportion of true nulls π_0 and the proportion of false nulls π_1 are defined respectively as

$$\pi_0 = \frac{p - s}{p} \quad \text{and} \quad \pi_1 = \frac{s}{p}$$

Power of a variable selection procedure

- ▶ Classification table of selection procedure \mathcal{R} :

	Null is true	Null is false	Total
Reject null	V	S	R
Retain null	U	T	$p - R$
Total	$p - s$	s	p

- ▶ FDR of the selection procedure \mathcal{R} :

$$\alpha = \mathbb{E} \left[\frac{V}{\max\{R, 1\}} \right]$$

- ▶ Power of \mathcal{R} is measured by “true discovery rate (TDR)”

$$\omega = \mathbb{E} [S/s]$$

- ▶ Other measures available on selection performance

The Benjamini-Hochberg procedure

- ▶ Let p_i be the p-value associated with H_{i0}
- ▶ Let $\{p_{(i)}\}_{i=1}^m$ be the order statistics of $\{p_i\}_{i=1}^m$ such that $p_{(j)} \leq p_{(j+1)}$
- ▶ Set

$$r = \max \left\{ 1 \leq k \leq m : p_{(k)} \leq \alpha k m^{-1} \right\}$$

if r is well-defined, then reject $H_{(i)}$ such that $p_{(i)} \leq p_{(r)}$; otherwise, no rejections are made

Theorem (Benjamini and Yekutieli (2001))

If the joint distribution of p-values satisfy PRDS on the subset I_0 of true null hypotheses and null p-values are super-uniform, then $\hat{\alpha} \leq \pi_0 \alpha$, with equality for independent p-values with continuous, uniform null distribution.

Target of variable selection

- ▶ Model $Y = \beta_0 + \sum_{j=1}^p X_j \beta_j + \varepsilon$
- ▶ Truth $\mathcal{A} = \{j : \beta_j \neq 0\}$ with $|\mathcal{A}| = s$
- ▶ Hypotheses $H_{j0} : \beta_j = 0$ versus $H_{j1} : \beta_j \neq 0$ for $1 \leq j \leq p$
- ▶ Selection $\mathcal{I} = \{j : H_{j0} \text{ is rejected}\}$
- ▶ Performance criteria:
 - ▶ “Strong selection consistency”: $\Pr(\mathcal{I} = \mathcal{A}) \rightarrow 1$
 - ▶ “Weak selection consistency”: $\Pr(\mathcal{I} \supseteq \mathcal{A}) \rightarrow 1$
 - ▶ “Estimation consistency”: $\mathbb{E} \left[\left\| \hat{\beta} - \beta \right\| \right] \rightarrow 0$
 - ▶ “Oracle property”: $\mathcal{I} = \mathcal{A}$ and

$$\mathbf{A}_n \left[\sqrt{n} \left(\hat{\beta}_{\mathcal{I}} - \beta_{\mathcal{A}} + \mathbf{b}_n \right) \right] \rightsquigarrow \text{Normal}(0, \mathbf{I})$$

Some relationships

- ▶ FDR and TDR of the selection procedure \mathcal{R} :

$$\alpha = \mathbb{E} \left[\frac{V}{\max\{R, 1\}} \right] \quad \text{and} \quad \omega = \mathbb{E} [S/s]$$

- ▶ Strong or weak selection consistency:

$$\Pr(\mathcal{I} = \mathcal{A}) \rightarrow 1 \quad \text{or} \quad \Pr(\mathcal{I} \supseteq \mathcal{A}) \rightarrow 1$$

- ▶ Oracle property and estimation consistency:

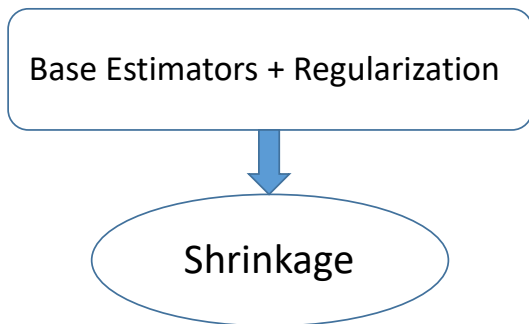
$$\begin{cases} \mathcal{I} = \mathcal{A} \text{ and } \mathbf{A}_n \left[\sqrt{n} \left(\hat{\boldsymbol{\beta}}_{\mathcal{I}} - \boldsymbol{\beta}_{\mathcal{A}} + \mathbf{b}_n \right) \right] \rightsquigarrow \text{Normal}(0, \mathbf{I}) \\ \mathbb{E} \left[\left\| \hat{\boldsymbol{\beta}} - \boldsymbol{\beta} \right\| \right] \rightarrow 0 \end{cases}$$

- ▶ Do you know their relationships?

Shrinkage-inducing penalty

General principle to induce shrinkage

- ▶ Base estimators: least squares estimate, maximum likelihood estimator, moment estimator, or kernel estimator
- ▶ Regularization: implicit/explicit and/or deterministic/random regularization



- ▶ RSS: e.g., $\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2$ or

$$\sum_{i=1}^N (y_i - f(\mathbf{x}_i))^2 = \|\mathbf{y} - f(\mathbf{X})\|^2$$

- ▶ Penalized RSS:

$$\text{Penalized RSS} = \text{RSS} + \text{Penalty}$$

- ▶ Example models: additive model

$$E(Y|X) = f(X) + \varepsilon$$

Penalized likelihood

- ▶ Likelihood: $h(y; \text{other parameters})$ and $l(\mathbf{y}; \text{other parameters})$

$$l(\mathbf{y}; *) = \prod_{i=1}^n h(y_i; *)$$

- ▶ Penalized likelihood:

Penalized likelihood = Likelihood + Penalty

- ▶ Example models:

$g[E(Y|X)] = f(X)$ and g is a link function

Designing a penalty

A *shrinkage-inducing* penalty should produce an estimator that possesses

- ▶ **Continuity**: estimator is continuous in data input to avoid instability in **model prediction**
- ▶ for purpose of *variable selection*
 1. **Unbiasedness**: estimator is nearly unbiased when true unknown parameter is large to avoid unnecessary **modeling bias**
 2. **Sparsity**: estimator incorporates a thresholding rule, which automatically sets small estimated coefficients to zero to reduce **model complexity**

Various penalties

Recall $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^T \in \mathbb{R}^p$ and penalty $\varrho : \mathbb{R}^p \rightarrow \mathbb{R}_{\geq 0}$.

- ▶ Separable penalty: $\varrho(\boldsymbol{\beta}) = \sum_{i=1}^p \varrho_i(\beta_i)$ with $\varrho_i : \mathbb{R} \rightarrow \mathbb{R}_{\geq 0}$, such as
 - ▶ the l_q -penalty

$$\varrho(\boldsymbol{\beta}) = \|\boldsymbol{\beta}\|_q^q = \sum_{i=1}^p |\beta_i|^q, 0 \leq q < \infty,$$

including BSS, LASSO and Ridge

- ▶ Non-separable penalty: $\varrho(\boldsymbol{\beta}) \neq \sum_{i=1}^p \varrho_i(\beta_i)$, such as
 - ▶ Benjamini-Hochberg (BH) procedure; see Abramovich et al. (2006)
- ▶ Note: BH procedure induces a random penalty

Orthogonal design and penalized RSS

- ▶ Model: $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ with $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^T \in \mathbb{R}^p$
- ▶ Assume $\mathbf{X}^T\mathbf{X} = \mathbf{I}$. Then LSE $\hat{\boldsymbol{\beta}} = \mathbf{X}^T\mathbf{y}$ and

$$\text{RSS}(\boldsymbol{\beta}) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 = \|\mathbf{y} - \hat{\mathbf{y}}\|^2 + \|\mathbf{z} - \boldsymbol{\beta}\|^2,$$

where $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$ and $\mathbf{z} = \mathbf{X}^T\mathbf{y}$ since $\mathbf{y} - \hat{\mathbf{y}} \perp \langle \mathbf{X} \rangle$ and $\mathbf{X}^T\mathbf{X} = \mathbf{I}$

- ▶ Further, RSS with *separable, piecewise smooth penalty* is

$$\begin{aligned} \text{pRSS}(\boldsymbol{\beta}) &:= \frac{1}{2}\text{RSS}(\boldsymbol{\beta}) + \sum_{i=1}^p \varrho_\lambda(|\beta_i|) \\ &= \frac{1}{2}\|\mathbf{y} - \hat{\mathbf{y}}\|^2 + \frac{1}{2}\|\mathbf{z} - \boldsymbol{\beta}\|^2 + \sum_{i=1}^p \varrho_\lambda(|\beta_i|) \end{aligned} \quad (2)$$

Componentwise optimization

- ▶ Let $\boldsymbol{\theta} = (\theta_1, \dots, \theta_p)^T$ and $\mathbf{z} = (z_1, \dots, z_p)^T$. We have equivalent objective function to pRSS ($\boldsymbol{\beta}$) as

$$u(\boldsymbol{\theta}) = \frac{1}{2} \|\mathbf{z} - \boldsymbol{\theta}\|^2 + \sum_{i=1}^p \varrho_\lambda(|\theta_i|) = \frac{1}{2} \sum_{i=1}^p (z_i - \theta_i)^2 + \sum_{i=1}^p \varrho_\lambda(|\theta_i|)$$

- ▶ First order condition for minimizer $\hat{\boldsymbol{\theta}} = (\hat{\theta}_1, \dots, \hat{\theta}_p)^T$ from

$$a(\theta) = \frac{1}{2} (z - \theta)^2 + \varrho_\lambda(|\theta|) \quad \text{and} \quad u(\boldsymbol{\theta}) = \sum_{i=1}^p a_i(\theta)$$

gives

$$\frac{\partial a(\theta)}{\partial \theta} = \begin{cases} \text{sgn}(\theta) [|\theta| + \varrho'_\lambda(|\theta|)] - z & \text{if } \theta \neq 0 \\ \text{subdifferential } \partial a(\theta) & \text{if } \theta = 0 \end{cases}$$

Shrinkage without hard thresholding

- ▶ Recall $a(\theta) = \frac{1}{2}(z - \theta)^2 + \varrho_\lambda(|\theta|)$ and

$$\frac{\partial a(\theta)}{\partial \theta} = \begin{cases} \operatorname{sgn}(\theta) [|\theta| + \varrho'_\lambda(|\theta|)] - z & \text{if } \theta \neq 0 \\ \text{subdifferential } \partial a(\theta) & \text{if } \theta = 0 \end{cases}$$

- ▶ Let $\varrho_\lambda(|\theta|) = 2^{-1}\lambda\theta^2$ for $\lambda > 0$. Then $\frac{\partial a(\theta)}{\partial \theta}$ exits at 0 and the optimal solution

$$\hat{\theta} = \frac{z}{1 + \lambda} \quad (3)$$

- ▶ Note LSE $\hat{\theta} = z$. Namely, (3) is a “shrinkage” estimator of z (since $|\hat{\theta}| < |z|$ whenever $\lambda > 0$)
- ▶ (3) is the *ridge estimator* of z (when $\mathbf{X}^T \mathbf{X} = \mathbf{I}$)

Shrinkage via hard thresholding

- ▶ Recall $a(\theta) = \frac{1}{2}(z - \theta)^2 + \varrho_\lambda(|\theta|)$ and

$$\frac{\partial a(\theta)}{\partial \theta} = \begin{cases} \operatorname{sgn}(\theta) [|\theta| + \varrho'_\lambda(|\theta|)] - z & \text{if } \theta \neq 0 \\ \text{subdifferential } \partial a(\theta) & \text{if } \theta = 0 \end{cases}$$

- ▶ Since $\partial a(\theta) = \theta - z + \lambda v$ for $|v| \leq 1$ at $\theta = 0$ for penalty

$$\varrho_\lambda(|\theta|) = 2^{-1}\lambda^2 - 2^{-1}(|\theta| - \lambda)^2 1_{\{|\theta| < \lambda\}}(\theta),$$

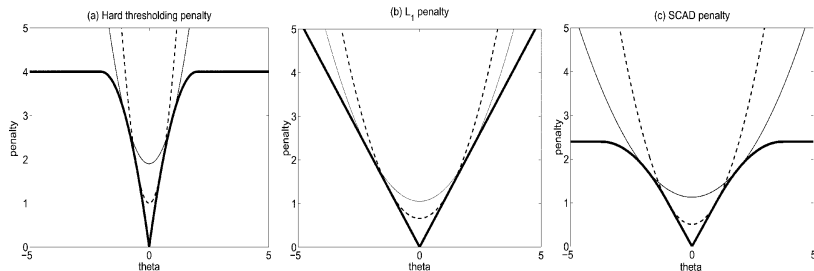
then $\hat{\theta} = 0$ iff $|z| \leq \lambda$, yielding the “hard thresholding” estimator

$$\hat{\theta} = z 1_{\{|z| > \lambda\}}(z) \quad (4)$$

- ▶ (4) is equivalent to best subset selection when $\mathbf{X}^T \mathbf{X} = \mathbf{I}$. Why?

Three penalties for variable selection

- ▶ Hard-thresholding penalty; l_1 -penalty for LASSO, and “smoothly clipped absolute deviation” as SCAD
- ▶ We want unbiasedness, sparsity and continuity of estimate. Which penalty is preferred?



Nearly unbiasedness and thresholding

- ▶ Recall $a(\theta) = \frac{1}{2}(z - \theta)^2 + \varrho_\lambda(|\theta|)$ and

$$\frac{\partial a(\theta)}{\partial \theta} = \begin{cases} \operatorname{sgn}(\theta) [|\theta| + \varrho'_\lambda(|\theta|)] - z & \text{if } \theta \neq 0 \\ \text{subdifferential } \partial a(\theta) & \text{if } \theta = 0 \end{cases}$$

- ▶ *Nearly unbiasedness*: $\hat{\theta} = z$ for $|z|$ large if $\varrho'_\lambda(|\theta|) = 0$ for large $|\theta|$
- ▶ *Thresholding rule*: $\hat{\theta} = 0$ for small $|z|$, i.e.,

$$\hat{\theta} = 0 \text{ when } |z| < \tau_\lambda = \inf_{\theta \neq 0} \{|\theta| + \rho'_\lambda(|\theta|)\}$$

$$\text{iff } \tau_\lambda > 0 \ \& \ \operatorname{arginf}_{\theta \in \mathbb{R}} \{|\theta| + \rho'_\lambda(|\theta|)\} = 0.$$

- ▶ Namely, $\rho_\lambda(|\theta|)$ has to be non-differentiable at $\theta = 0$ to induce thresholding rule

Thresholding rule: visualization

- ▶ $\operatorname{argmin}_{\theta \in \mathbb{R}} \left\{ \frac{1}{2} (z - \theta)^2 + \rho_{\lambda} (|\theta|) \right\}$ when $\tau_{\lambda} = \inf_{\theta \neq 0} \{ |\theta| + \rho'_{\lambda} (|\theta|) \} > 0$
- ▶ $\hat{\theta} (z_1) = 0$ but $\hat{\theta} (z_2) = \sup \left\{ \theta : \frac{\partial a(\theta)}{\partial \theta} = 0 \right\}$ to retain continuity

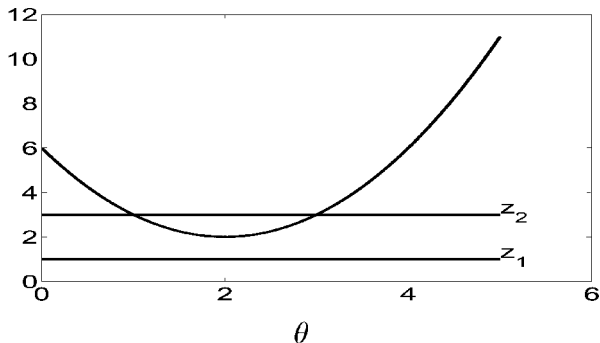


Figure 3. A Plot of $\theta + \rho'_{\lambda}(\theta)$ Against $\theta(\theta > 0)$.

Derivatives of three penalties

- ▶ Hard-thresholding penalty: $\varrho_\lambda (|\theta|) = \lambda^2 - (|\theta| - \lambda)^2 \mathbf{1}_{\{|\theta| < \lambda\}} (\theta)$
- ▶ SCAD penalty: $\varrho'_\lambda (\theta) = \lambda \left[\mathbf{1}_{\{\theta \leq \lambda\}} (\theta) + \frac{(a\lambda - \theta)_+}{(a-1)\lambda} \mathbf{1}_{\{\theta > \lambda\}} (\theta) \right], \theta > 0$
- ▶ l_q -penalty: $\varrho_\lambda (|\theta|) = \lambda |\theta|^q$

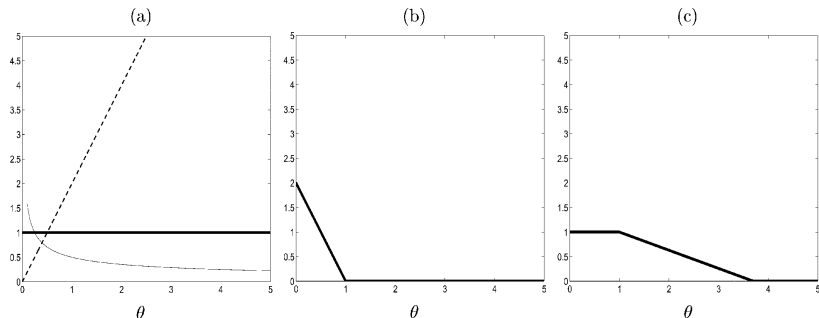
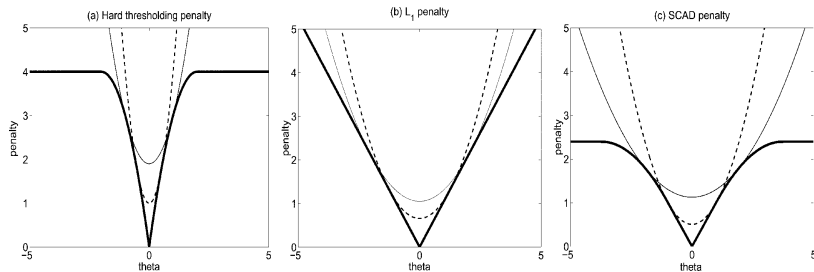


Figure 4. Plot of $p'_\lambda(\theta)$ Functions Over $\theta > 0$ (a) for L_q Penalties, (b) the Hard Thresholding Penalty, and (c) the SCAD Penalty. In (a), the heavy line corresponds to L_1 , the dash-dot line corresponds to L_s , and the thin line corresponds to L_2 penalties.

Recap on three penalties for variable selection

- ▶ Pay attention to the “rounding” parts of a hard thresholding penalty and SCAD penalty
- ▶ None of these penalties are differentiable at 0
- ▶ All are even functions on \mathbb{R} and non-decreasing on $\mathbb{R}_{>0}$



Estimates given by three penalties

When tuning parameter $\lambda > 0$:

- ▶ Hard thresholding estimate is discontinuous in data
- ▶ LASSO estimate is intrinsically biased
- ▶ Both LASSO and SCAD incorporate hard thresholding

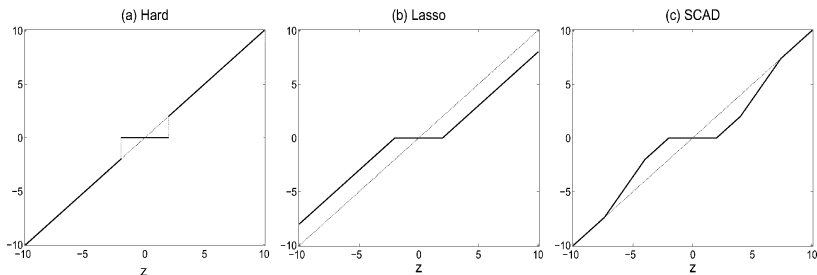


Figure 2. Plot of Thresholding Functions for (a) the Hard, (b) the Soft, and (c) the SCAD Thresholding Functions With $\lambda = 2$ and $a = 3.7$ for SCAD.

Recap on penalties for variable selection: I

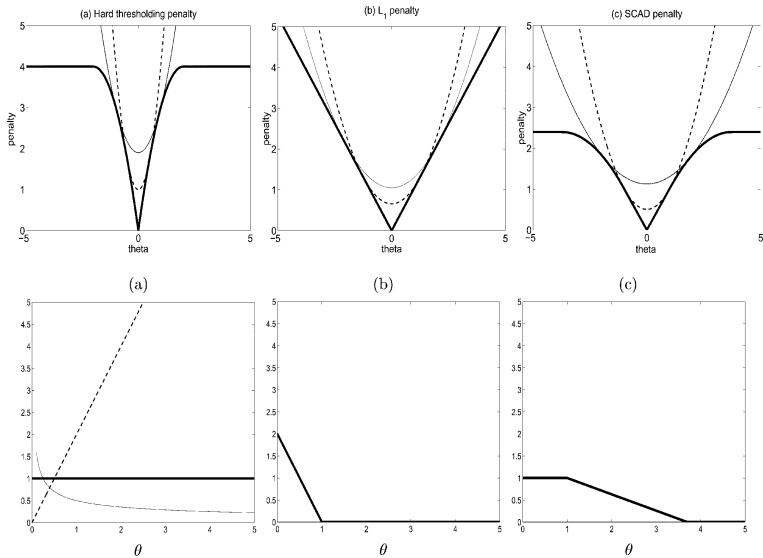


Figure 4. Plot of $p'_j(\theta)$ Functions Over $\theta > 0$ (a) for L_2 Penalties, (b) the Hard Thresholding Penalty, and (c) the SCAD Penalty. In (a), the heavy line corresponds to L_1 , the dash-dot line corresponds to L_2 , and the thin line corresponds to L_2 penalties.

Recap on penalties for variable selection: II

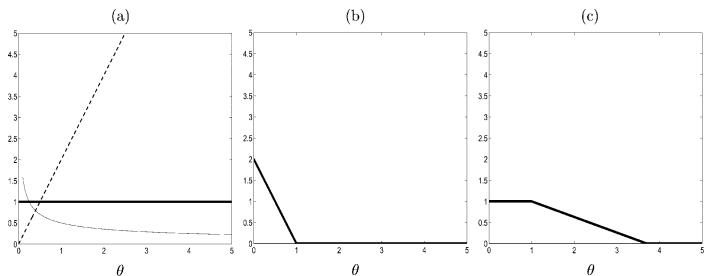


Figure 4. Plot of $p_\lambda(\theta)$ Functions Over $\theta > 0$ (a) for L_∞ Penalties, (b) the Hard Thresholding Penalty, and (c) the SCAD Penalty. In (a), the heavy line corresponds to L_1 , the dash-dot line corresponds to L_∞ , and the thin line corresponds to L_2 penalties.

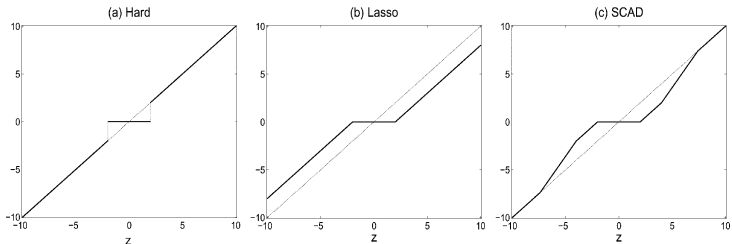


Figure 2. Plot of Thresholding Functions for (a) the Hard, (b) the Soft, and (c) the SCAD Thresholding Functions With $\lambda = 2$ and $a = 3.7$ for SCAD.

Geometry of linear transform

Singular value decomposition

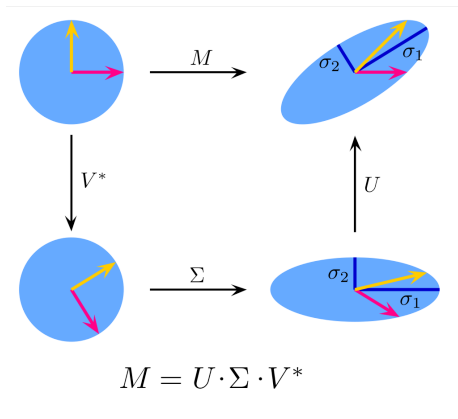
Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ with $r = \text{rank}(\mathbf{A})$. The “singular value decomposition (SVD)” of \mathbf{A} is

$$\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{V}^T = \sum_{i=1}^r d_i \mathbf{u}_i \mathbf{v}_i^T$$

- ▶ $\mathbf{D} = \text{diag}\{d_1, \dots, d_r\}$ and $d_1 \geq \dots \geq d_r > 0$. Each d_i is called a “singular value” of \mathbf{A}
- ▶ $\mathbf{U} \in \mathbb{R}^{m \times r}$ and $\mathbf{V} \in \mathbb{R}^{n \times r}$ such that $\mathbf{U}^T \mathbf{U} = \mathbf{V}^T \mathbf{V} = \mathbf{I}_r$. Further, $\langle \mathbf{U} \rangle = \langle \mathbf{A} \rangle$ and $\langle \mathbf{V} \rangle = \langle \mathbf{A}^T \rangle$
- ▶ \mathbf{u}_i is the i th column of \mathbf{U} and is called a “left singular vector”, and \mathbf{v}_i the i th column of \mathbf{V} and is called a “right singular vector”
- ▶ $(\mathbf{U}, \mathbf{D}, \mathbf{V})$ above can be extended into $\mathbf{U} \in \mathbb{R}^{m \times n}$, $\mathbf{V} \in \mathbb{R}^{n \times n}$ and $\mathbf{D} = \text{diag}\{d_1, \dots, d_n\}$ such that $\mathbf{U}^T \mathbf{U} = \mathbf{V}^T \mathbf{V} = \mathbf{I}_n$ and $d_{r+1} = \dots = d_n = 0$

Geometric meaning of SVD: illustration

- ▶ SVD $M = U\Sigma V^*$ where Σ is diagonal; note the rotations induced by U and V^* and scaling induced by Σ (Credit: Wiki)



Spectral decomposition

Let $\mathbf{A} \in \mathbb{R}^{m \times m}$ with $\mathbf{A} = \mathbf{A}^T$. The “spectral decomposition” of \mathbf{A} is

$$\mathbf{A} = \tilde{\mathbf{V}}\tilde{\mathbf{D}}\tilde{\mathbf{V}}^T = \sum_{i=1}^m \tilde{d}_i \tilde{\mathbf{v}}_i \tilde{\mathbf{v}}_i^T$$

- ▶ $\tilde{\mathbf{D}} = \text{diag} \{ \tilde{d}_1, \dots, \tilde{d}_m \}$ with $|\tilde{d}_1| \geq \dots \geq |\tilde{d}_m|$, and each \tilde{d}_i is called an “eigenvalue” of \mathbf{A}
- ▶ $\tilde{\mathbf{V}} \in \mathbb{R}^{m \times m}$ such that $\tilde{\mathbf{V}}^T \tilde{\mathbf{V}} = \mathbf{I}_m$, and the i th column $\tilde{\mathbf{v}}_i$ of $\tilde{\mathbf{V}}$ is an eigenvector of \mathbf{A} associated with \tilde{d}_i
- ▶ $(\tilde{d}_i, \tilde{\mathbf{v}}_i)$ is called the “ i th eigenpair”

Remark: Spectral decomposition \iff SVD; decompositions hold for Hermitian matrix $\mathbf{A} = \mathbf{A}^*$ with operation “conjugate transpose $*$ ”

Maximal scaling

- ▶ $\mathbf{A} \in \mathbb{R}^{m \times n}$ induces a linear map $\mathbf{A} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ via *right multiplication* $\mathbf{w} = \mathbf{A}\mathbf{v}$
- ▶ Set $\mathbf{S} = \mathbf{A}^T \mathbf{A}$. Then $\mathbf{S} \succeq 0$, i.e., \mathbf{S} is positive semi-definite, and

$$\|\mathbf{w}\|^2 = \|\mathbf{A}\mathbf{v}\|^2 = \mathbf{v}^T (\mathbf{A}^T \mathbf{A}) \mathbf{v} = \mathbf{v}^T \mathbf{S} \mathbf{v}$$

- ▶ SVD $\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{V}^T$ gives $\mathbf{S} = \mathbf{V}\mathbf{D}^2\mathbf{V}^T$ with $\mathbf{D} = \text{diag}\{d_1, \dots, d_n\}$ where

$$d_1 \geq \dots \geq d_n \geq 0$$

- ▶ Set $\|\mathbf{A}\|_{\text{op}} = \max_{1 \leq i \leq n} d_i = d_1$. **Rayleigh's inequality:**

$$\|\mathbf{A}\mathbf{v}\| \leq \|\mathbf{A}\|_{\text{op}} \|\mathbf{v}\|,$$

i.e., $\mathbf{v} \in \mathbb{R}^n$ can be scaled at most $\|\mathbf{A}\|_{\text{op}}$ by \mathbf{A} when mapped to $\mathbf{A}\mathbf{v}$

Rayleigh's inequality

- ▶ Recall $\mathbf{S} = \mathbf{A}^T \mathbf{A} = \mathbf{V} \mathbf{D}^2 \mathbf{V}^T$, $\mathbf{V}^T \mathbf{V} = \mathbf{I}$ and $\mathbf{D} = \text{diag} \{d_1, \dots, d_n\}$
- ▶ Recall $\mathbf{w} = \mathbf{A} \mathbf{v}$. So, $\|\mathbf{w}\|^2 = \mathbf{v}^T \mathbf{V} \mathbf{D}^2 \mathbf{V}^T \mathbf{v}$. Let $\tilde{\mathbf{v}} = \mathbf{V}^T \mathbf{v}$. Then

$\|\tilde{\mathbf{v}}\| = \|\mathbf{v}\|$, i.e., orthogonal transform preserves length

Recall $\|\mathbf{A}\|_{\text{op}} = \max_{1 \leq i \leq n} d_i$. Then

$$\|\mathbf{A} \mathbf{v}\|^2 = \mathbf{v}^T \mathbf{V} \mathbf{D}^2 \mathbf{V}^T \mathbf{v} = \tilde{\mathbf{v}}^T \mathbf{D}^2 \tilde{\mathbf{v}} \leq \|\mathbf{A}\|_{\text{op}}^2 \|\tilde{\mathbf{v}}\|^2 = \|\mathbf{A}\|_{\text{op}}^2 \|\mathbf{v}\|^2$$

with equality iff $\mathbf{v} = c \mathbf{v}_1$ with (d_1, \mathbf{v}_1) and $c \neq 0$

- ▶ So, $\|\mathbf{A} \mathbf{v}\| \leq \|\mathbf{A}\|_{\text{op}} \|\mathbf{v}\|$ and

$$\left(\max_{\|\mathbf{v}\|=1} \|\mathbf{A} \mathbf{v}\|, \operatorname{argmax}_{\|\mathbf{v}\|=1} \|\mathbf{A} \mathbf{v}\| \right) = \left(\|\mathbf{A}\|_{\text{op}}, \mathbf{v}_1 \right)$$

Successive maximal scaling

- ▶ Recall $\mathbf{S} = \mathbf{A}^T \mathbf{A} = \mathbf{V} \mathbf{D}^2 \mathbf{V}^T$ with $\mathbf{V}^T \mathbf{V} = \mathbf{I}$ and

$$\mathbf{D} = \text{diag} \{d_1, \dots, d_n\} \quad \text{and} \quad d_1 \geq \dots \geq d_n \geq 0$$

- ▶ Rayleigh's inequality: $\max_{\|\mathbf{v}\|=1} \|\mathbf{A}\mathbf{v}\| \leq \|\mathbf{A}\|_{\text{op}} = d_1$
- ▶ Recall \mathbf{v}_i as i th column of \mathbf{V} and eigenpair (d_i, \mathbf{v}_i) . Same arguments (plus induction) can show, for $1 \leq k \leq n-1$,

$$\left(\max_{\|\mathbf{v}\|=1, \mathbf{v} \perp \langle \mathbf{v}_1, \dots, \mathbf{v}_k \rangle} \|\mathbf{A}\mathbf{v}\|, \operatorname{argmax}_{\|\mathbf{v}\|=1, \mathbf{v} \perp \langle \mathbf{v}_1, \dots, \mathbf{v}_k \rangle} \|\mathbf{A}\mathbf{v}\| \right) = (d_{k+1}, \mathbf{v}_{k+1})$$

by noting that $\{\mathbf{v}_i\}_{i=1}^n$ are orthonormal and constraint $\|\mathbf{v}\| = 1$

- ▶ The above is called “successive maximal scaling (SMS)” lemma

Successive maximal scaling: sketch of proof

- ▶ Let $(\max_{\|\mathbf{v}\|=1} \|\mathbf{A}\mathbf{v}\|, \operatorname{argmax}_{\|\mathbf{v}\|=1} \|\mathbf{A}\mathbf{v}\|) = (d_1, \mathbf{v}_1)$
- ▶ Recall $\mathcal{S}^{n-1} = \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\| = 1\}$. Since $\{\mathbf{v}_i\}_{i=1}^n$ are orthonormal, then $\|\mathbf{v}\| = 1$ and $\mathbf{v} \perp \mathbf{v}_1$ iff $\mathbf{v} \in \langle \{\mathbf{v}_i\}_{i=2}^n \rangle$ and

$$\exists \mathbf{a} = (a_1, \dots, a_{n-1})^T \in \mathcal{S}^{n-1} \text{ such that } \mathbf{v} = \sum_{i=2}^n a_{i-1} \mathbf{v}_i$$

- ▶ So, $\mathbf{A}\mathbf{v} = \sum_{i=2}^n a_{i-1} \mathbf{A}\mathbf{v}_i$. But $\{\mathbf{v}_i\}_{i=1}^n$ are orthonormal eigenvectors of $\mathbf{S} = \mathbf{A}^T \mathbf{A}$, So,

$$\|\mathbf{A}\mathbf{v}\|^2 = \sum_{i=2}^n a_{i-1}^2 \|\mathbf{A}\mathbf{v}_i\|^2 = \sum_{i=2}^n a_{i-1}^2 d_i^2 \leq d_2^2 \|\mathbf{a}\|^2 = d_2^2,$$

i.e., $(\max_{\|\mathbf{v}\|=1, \mathbf{v} \perp \mathbf{v}_1} \|\mathbf{A}\mathbf{v}\|, \operatorname{argmax}_{\|\mathbf{v}\|=1, \mathbf{v} \perp \mathbf{v}_1} \|\mathbf{A}\mathbf{v}\|) = (d_2, \mathbf{v}_2)$

- ▶ Remaining part done via induction

Orthogonal projections and PCA

Directions of successive maximal variability

- ▶ Fact: if $\mathbf{v} \in \mathcal{S}^{n-1}$ and $\mathbf{x} \in \mathbb{R}^n$, then their inner product $\langle \mathbf{x}, \mathbf{v} \rangle$ is the orthogonal projection of \mathbf{x} onto \mathbf{v}
- ▶ Let $X \in \mathbb{R}^n$ have covariance matrix $\Sigma \in \mathbb{R}^{n \times n}$ with eigenpairs (σ_i, \mathbf{v}_i) . Given any orthonormal $\{\mathbf{u}_i\}_{i=1}^n$, define $y_i = \langle X, \mathbf{u}_i \rangle$. Then $\text{var}(\langle X, \mathbf{u}_i \rangle) = \mathbf{u}_i^T \Sigma \mathbf{u}_i$
- ▶ SMS Lemma (applied to Σ) implies that $\{\mathbf{v}_i\}_{i=1}^n$ are the *directions* for which $\{y_i\}_{i=1}^n$ successively achieve maximal variances, i.e.,

$$\begin{cases} \max_{\|\mathbf{v}\|=1, \mathbf{v} \perp \langle \mathbf{v}_1, \dots, \mathbf{v}_k \rangle} \text{var}(\langle X, \mathbf{v} \rangle) = d_{k+1} \\ \text{argmax}_{\|\mathbf{v}\|=1, \mathbf{v} \perp \langle \mathbf{v}_1, \dots, \mathbf{v}_k \rangle} \text{var}(\langle X, \mathbf{v} \rangle) = \mathbf{v}_{k+1} \end{cases}$$

- ▶ \mathbf{v}_i is called the “*i*th principal component direction”. Further, $\text{cov}(\langle X, \mathbf{v}_i \rangle, \langle X, \mathbf{v}_j \rangle) = 0$ for $i \neq j$ when $E(X) = 0$

Principal components

- ▶ Let $X \in \mathbb{R}^n$ and $\mathbf{X} = (\mathbf{x}_1^T, \dots, \mathbf{x}_m^T)^T$ be the *column-centered* data matrix (i.e., sample mean for entries of each column of \mathbf{X} is 0)
- ▶ SVD $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T$ gives $\mathbf{S} = \mathbf{X}^T\mathbf{X} = \mathbf{V}\mathbf{D}^2\mathbf{V}^T$, and $\mathbf{S} = m^{-1}\widetilde{\text{cov}}(X)$
- ▶ Let $\{(d_i^2, \mathbf{v}_i)\}_{i=1}^n$ be the eigenpairs of \mathbf{S} . Given any orthonormal $\{\tilde{\mathbf{u}}_i\}_{i=1}^n$, define $\mathbf{z}_i = \mathbf{X}\tilde{\mathbf{u}}_i$. Then

$$\mathbf{z}_i = (z_{1i}, \dots, z_{mi})^T = (\langle \mathbf{x}_1, \tilde{\mathbf{u}}_i \rangle, \dots, \langle \mathbf{x}_m, \tilde{\mathbf{u}}_i \rangle)^T$$

and $\widetilde{\text{cov}}(\mathbf{z}_i) = \tilde{\mathbf{u}}_i^T \mathbf{S} \tilde{\mathbf{u}}_i$. (What is the interpretation of \mathbf{z}_i ?)

- ▶ SMS Lemma (applied to \mathbf{S}) implies that $\{\mathbf{v}_i\}_{i=1}^n$ are the *directions* for which $\{\mathbf{z}_i\}_{i=1}^n$ *successively* achieve *maximal sample variances*, i.e.,

$$\begin{cases} \max_{\|\mathbf{v}\|=1, \mathbf{v} \perp \langle \mathbf{v}_1, \dots, \mathbf{v}_k \rangle} \widetilde{\text{var}}(\langle \mathbf{X}, \mathbf{v} \rangle) = m^{-1}d_{k+1}^2 \\ \text{argmax}_{\|\mathbf{v}\|=1, \mathbf{v} \perp \langle \mathbf{v}_1, \dots, \mathbf{v}_k \rangle} \widetilde{\text{var}}(\langle \mathbf{X}, \mathbf{v} \rangle) = \mathbf{v}_{k+1} \end{cases}$$

- ▶ Thus, optimal $\mathbf{z}_i = d_i \mathbf{v}_i$ (why?)

Recap on classic PCA

- ▶ Population version of optimal projections:

- ▶ Let $X \in \mathbb{R}^n$ have covariance matrix $\Sigma \in \mathbb{R}^{n \times n}$ with eigenpairs (σ_i, \mathbf{v}_i)
- ▶ Optimal linear combinations of entries of X are $\{y_i = \langle X, \mathbf{v}_i \rangle\}_{i=1}^n$ when $\{\mathbf{v}_i\}_{i=1}^n$ have to be orthonormal
- ▶ $\{y_i\}_{i=1}^n$ *successively achieve maximal variances and are mutually uncorrelated*

- ▶ Sample version of optimal projections:

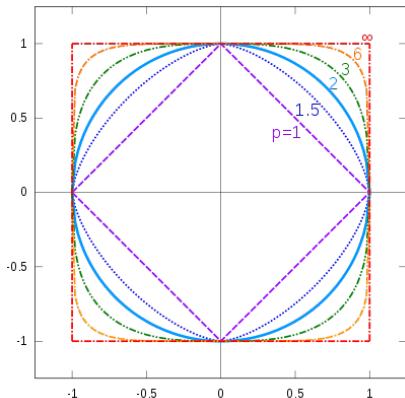
- ▶ Let $X \in \mathbb{R}^n$ and $\mathbf{X} = (\mathbf{x}_1^T, \dots, \mathbf{x}_m^T)^T$ be the *column-centered* data matrix
- ▶ Let $\{(d_i^2, \mathbf{v}_i)\}_{i=1}^n$ be the eigenpairs of $\mathbf{S} = \mathbf{X}^T \mathbf{X}$
- ▶ Optimal linear combinations of columns of \mathbf{X} are $\{\mathbf{z}_i = \mathbf{X} \mathbf{v}_i\}_{i=1}^n$ when $\{\mathbf{v}_i\}_{i=1}^n$ have to be orthonormal
- ▶ $\{\mathbf{z}_i\}_{i=1}^n$ *successively achieve maximal sample variances*

- ▶ Both versions involve orthogonal projections

Solution for ridge regression

Norms and geometry

- ▶ Unit l_p -ball: $B_p = \{x \in \mathbb{R}^n : \|x\|_p \leq 1\}$ (Credit: Wikipedia l_p -space)



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

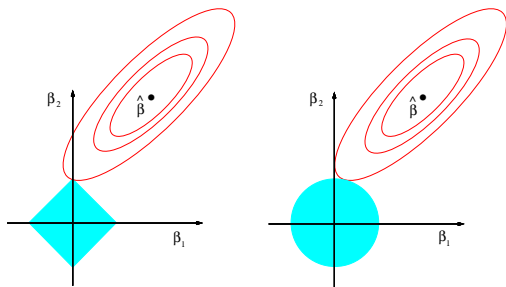


FIGURE 3.11. Estimation picture for the lasso (left) and ridge regression (right). Shown are contours of the error and constraint functions. The solid blue areas are the constraint regions $|\beta_1| + |\beta_2| \leq t$ and $\beta_1^2 + \beta_2^2 \leq t^2$, respectively, while the red ellipses are the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

Linear model with scalar response

- ▶ Predictor $X = (X_1, \dots, X_p)^T \in \mathbb{R}^p$ and response $Y \in \mathbb{R}$
- ▶ Coefficient vector $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^T \in \mathbb{R}^p$ and population model

$$Y = \sum_{j=1}^p X_j \beta_j + \varepsilon,$$

where $\varepsilon \in \mathbb{R}$ is the random error term $E(\varepsilon) = 0$ and $\text{var}(\varepsilon) = \sigma^2$

- ▶ Data version: set $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_N)^T \in \mathbb{R}^N$ and each $\varepsilon_i \sim \varepsilon$, then

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

Ridge regression

- ▶ l_2 -penalized RSS:

$$\begin{aligned}h_2(\boldsymbol{\beta}; \lambda) &= \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda \|\boldsymbol{\beta}\|^2 \\ &= \|\mathbf{y}\|^2 - (\mathbf{X}\boldsymbol{\beta})^T \mathbf{y} + \boldsymbol{\beta}^T (\lambda \mathbf{I} + \mathbf{X}^T \mathbf{X}) \boldsymbol{\beta}\end{aligned}$$

- ▶ If $|\mathbf{X}^T \mathbf{X}| \neq 0$, then unique optimizer

$$\hat{\boldsymbol{\beta}}_R(\lambda) = (\lambda \mathbf{I} + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \text{ for } \lambda \geq 0$$

- ▶ If $\mathbf{X}^T \mathbf{X} = \mathbf{I}$, then $\hat{\boldsymbol{\beta}}_R(\lambda) = (1 + \lambda)^{-1} \mathbf{X}^T \mathbf{y} = (1 + \lambda)^{-1} \hat{\boldsymbol{\beta}}$
- ▶ $\hat{\boldsymbol{\beta}}_R(0) = \hat{\boldsymbol{\beta}}$, i.e., LSE and unbiased when $E(\boldsymbol{\varepsilon}) = 0$
- ▶ $\hat{\boldsymbol{\beta}}_R(\infty) = 0$, i.e., largest bias but smallest “variance”
- ▶ For some $\lambda \in (0, \infty)$, there is a balance between bias and variance when $|\mathbf{X}^T \mathbf{X}| = 0$

Ridge regression: solution path

►
$$\text{df}(\lambda) = \text{trace} \left(\mathbf{X} (\lambda \mathbf{I} + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \right) = \text{trace} (\mathbf{H}_\lambda) \downarrow \lambda$$

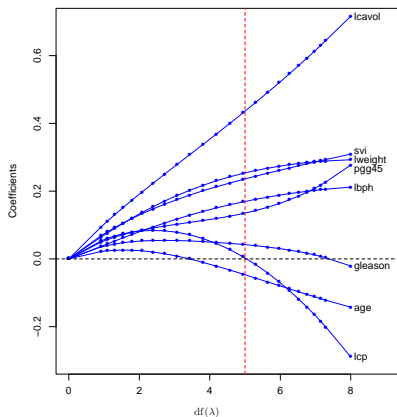
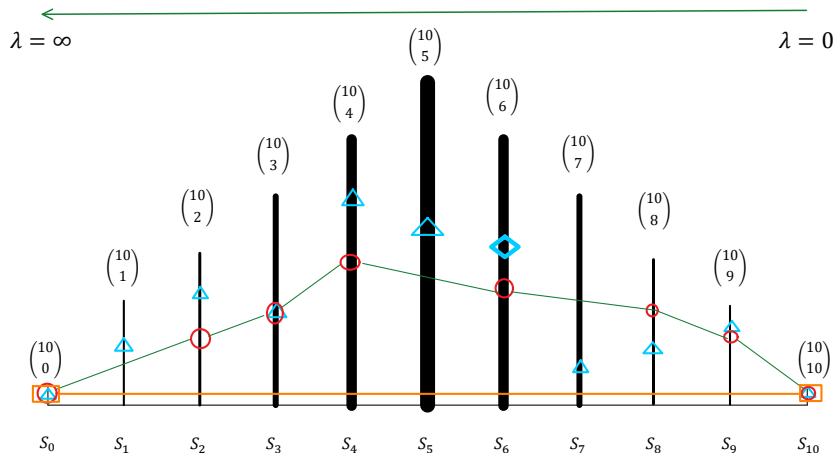


FIGURE 3.8. Profiles of ridge coefficients for the prostate cancer example, as the tuning parameter λ is varied. Coefficients are plotted versus $\text{df}(\lambda)$, the effective degrees of freedom. A vertical line is drawn at $\text{df} = 5.0$, the value chosen by cross-validation.

Trajectory of selected models

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



Adaptive shrinkage for non-orthogonal design

- ▶ Recall ridge estimate $\hat{\beta}_R(\lambda) = (\lambda\mathbf{I} + \mathbf{X}^T\mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ for $\lambda \geq 0$
- ▶ Recall $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T \in \mathbb{R}^{n \times p}$ and $\mathbf{S} = \mathbf{X}^T\mathbf{X}$, where \mathbf{u}_i is the i th column of \mathbf{U} , $\mathbf{U}^T\mathbf{U} = \mathbf{V}^T\mathbf{V} = \mathbf{I}_p$ and

$$\mathbf{D} = \text{diag} \{d_1, \dots, d_n\} \ \& \ d_1 \geq \dots \geq d_n \geq 0$$

- ▶ Adaptive shrinkage on LSE $\hat{\beta} = \hat{\beta}_R(0)$, i.e., $\hat{\beta}_R(\lambda) \neq c_\lambda \hat{\beta}$ for some $c_\lambda \neq 0$ (except when $\mathbf{X}^T\mathbf{X} = \mathbf{I}$)
- ▶ Adaptive shrinkage in prediction:

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\beta}_R(\lambda) = \mathbf{U}\mathbf{D}(\lambda\mathbf{I} + \mathbf{D}^2)^{-1} \mathbf{D}\mathbf{U}^T \mathbf{y} = \sum_{j=1}^p \mathbf{u}_j \frac{d_j^2}{d_j^2 + \lambda} \mathbf{u}_j^T \mathbf{y},$$

i.e., $\hat{\mathbf{y}}$ is a sum of projections of \mathbf{y} onto each \mathbf{u}_j with shrinkage proportional to $d_j^2 (d_j^2 + \lambda)^{-1}$

Adaptive shrinkage for non-orthogonal design

- ▶ Recall $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T$ and $\mathbf{S} = \mathbf{X}^T\mathbf{X} = \mathbf{V}\mathbf{D}^2\mathbf{V}^T$, where \mathbf{u}_i is the i th column of \mathbf{U} and

$$\mathbf{D} = \text{diag} \{d_1, \dots, d_n\} \ \& \ d_1 \geq \dots \geq d_n \geq 0$$

- ▶ Recall $\mathbf{z}_i = \mathbf{X}\mathbf{v}_i = d_i\mathbf{u}_i$, $1 \leq i \leq n$ successively achieve maximal sample variances, i.e., $\widetilde{\text{var}}(\mathbf{z}_i) = m^{-1}d_i^2$. Recall

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}_R(\lambda) = \sum_{j=1}^p \mathbf{u}_j d_j^2 (d_j^2 + \lambda)^{-1} \mathbf{u}_j^T \mathbf{y}$$

- ▶ So, $\hat{\mathbf{y}}$ is a sum of shrunken \mathbf{y} 's in the directions $\{\mathbf{u}_i\}_{i=1}^n$ of successive maximal sample variances where a larger such variance leads to more shrinkage

Adaptive shrinkage in prediction: illustration

► Recall $\hat{\mathbf{y}} = \sum_{j=1}^p \mathbf{u}_j d_j^2 (d_j^2 + \lambda)^{-1} \mathbf{u}_j^T \mathbf{y}$

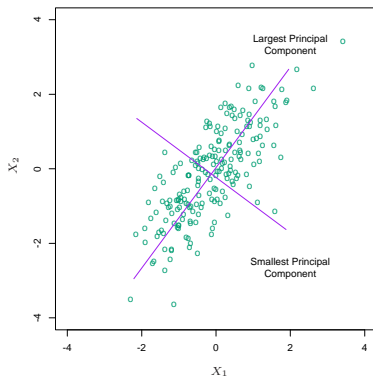


FIGURE 3.9. Principal components of some input data points. The largest principal component is the direction that maximizes the variance of the projected data, and the smallest principal component minimizes that variance. Ridge regression projects \mathbf{y} onto these components, and then shrinks the coefficients of the low-variance components more than the high-variance components.

- Abramovich, F., Benjamini, Y., Donoho, D. L., and Johnstone, I. M. (2006). Adapting to unknown sparsity by controlling the false discovery rate. *Ann. Statist.*, 34(2):584–653.
- Benjamini, Y. and Yekutieli, D. (2001). The control of the false discovery rate in multiple testing under dependency. *Ann. Statist.*, 29(4):1165–1188.