

3.4 Shrinkage Methods

exclude β_0 from the penalty term

3.4.1 Ridge

$$\hat{\beta}^{\text{ridge}} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p x_{ij}\beta_j)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\} \quad \lambda > 0$$

normally standardize the inputs before solving the above formula

$$\text{In matrix form: } RSS(\lambda) = (y - X\beta)^T(y - X\beta) + \lambda \beta^T \beta \quad -2X^T(y - X\beta) + 2\lambda \beta = 0$$

$$\text{Then } \hat{\beta}^{\text{ridge}} = (X^T X + \lambda I)^{-1} X^T y \quad -X^T y + X^T X \beta + \lambda \beta = 0 \quad (X^T X + \lambda I) \beta = X^T y$$

$$\text{For orthonormal inputs, } \hat{\beta}^{\text{ridge}} = \frac{\hat{\beta}}{1+\lambda}$$

$$df(\lambda) = \operatorname{tr}[X(X^T X + \lambda I)^{-1} X^T] = \operatorname{tr}(I/\lambda) = \sum_{j=1}^p \frac{\sigma_j^{-2}}{\sigma_j^2 + \lambda} \quad (3.50)$$

Exercise 3.6 prior dist. $\beta_j \stackrel{\text{i.i.d}}{\sim} N(0, \tau^2)$

$$f(\beta | y) = \frac{f(\beta, y)}{f(y)} = \frac{f(y|\beta)f(\beta)}{\int f(y|\beta)f(\beta)d\beta} \propto f(y|\beta)f(\beta) \quad \begin{matrix} \text{f}(y|\beta) \text{ is a function} \\ \text{of } f(y|\beta)f(\beta) \\ \text{or like} \\ \text{follow} \\ \text{the same} \\ (\text{scale different}) \\ \text{distribution} \end{matrix}$$

$$\int_{\beta} f(\beta | y) d\beta = 1$$

Thin SVD of X matrix.

$$X = UDV^T$$

$\downarrow \quad \downarrow \quad \downarrow$
 $N \times p \quad p \times p \quad p \times p$

only calculated the p column vectors which correspond to p row vectors of V^T

$$\text{Then } \hat{\beta} = X(X^T X)^{-1} X^T y = UV^T y$$

$$\begin{aligned} &UV^T (V D^T U^T U D V^T)^{-1} V D^T U^T y \\ &\quad \underbrace{UDV^T(V^T)^{-1}}_1 D^{-1} (D^T)^{-1} \underbrace{V^T V}_1 D^T U^T y \end{aligned}$$

Then ridge $\hat{X}\beta^{\text{ridge}} = X(X^T X + \lambda I)^{-1} X^T Y$

$$= X(V D^2 V^T + \lambda V V^T)^{-1} V D U^T Y$$

$$= X(V(D^2 + \lambda I)V^T)^{-1} V D U^T Y$$

$$= X(V^T)^{-1}(D^2 + \lambda I)^{-1} \underbrace{V^{-1} V D U^T Y}_{I}$$

$$= X(V^T)^{-1}(D^2 + \lambda I)^{-1} D U^T Y$$

$$= \underbrace{UDV^T}_{\mathbb{1}} (V^T)^{-1}(D^2 + \lambda I)^{-1} D U^T Y$$

$$= U D (D^2 + \lambda I)^{-1} D U^T Y$$

$$= \sum_{j=1}^p \lambda_j \frac{d_j^2}{d_j^2 + \lambda} u_j^T y$$

\downarrow
column of U

Sample covariance matrix: $S = X^T X / N$

where $X^T X = \underline{V D^2 V^T} \rightarrow \text{eigen decomposition}$

first principle component direction v_1 : $z_1 = X v_1$ has the largest ^{sample} variance.

$$\begin{aligned} \text{Var}(z_1) &= \text{Var}(X v_1) \\ &= v_1^T \text{Var}(X) v_1 \\ &= v_1^T \frac{X^T X}{N} v_1 \\ &= v_1^T \frac{V D^2 V^T}{N} v_1 = \frac{d_1^2}{N} \end{aligned}$$

Where $V_1^T V = [1 \ 0 \ 0 \ \dots]$

$$[1 \ 0 \ 0 \ \dots] \begin{bmatrix} d_1^2 & & 0 \\ & d_2^2 & \\ 0 & \ddots & d_p^2 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \end{bmatrix$$

$$\text{Var}(z_1) = \frac{d_1^2}{N}$$

3.4.2 Lasso

$$\hat{\beta}^{\text{Lasso}} = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^P x_{ij} \beta_j)^2 \text{ subject to } \sum_{j=1}^P |\beta_j| \leq t$$

equals

$$\hat{\beta}^{\text{Lasso}} = \underset{\beta}{\operatorname{argmin}} \left\{ \frac{1}{2} \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^P x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^P |\beta_j| \right\}$$

nonlinear in the y_i
no closed form expression as in ridge

e.g. $t_0 = \sum_j |\beta_j|$ if $t = \frac{t_0}{2}$, LSE coefficients are shrunk by about 50% on average.

3.4.3 Generalization of Regularization

TABLE 3.4. Estimators of β_j in the case of orthonormal columns of \mathbf{X} . M and λ are constants chosen by the corresponding techniques; sign denotes the sign of its argument (± 1), and x_+ denotes “positive part” of x . Below the table, estimators are shown by broken red lines. The 45° line in gray shows the unrestricted estimate for reference.

Estimator	Formula
Best subset (size M)	$\hat{\beta}_j \cdot I[\operatorname{rank}(\hat{\beta}_j) \leq M)]$ ← hard thresholding
Ridge	$\hat{\beta}_j / (1 + \lambda)$
Lasso	$\operatorname{sign}(\hat{\beta}_j) (\hat{\beta}_j - \lambda)_+$ ← soft thresholding

example: Section 5.9 ←

wavelet-based smoothing

↓
基于小波光滑

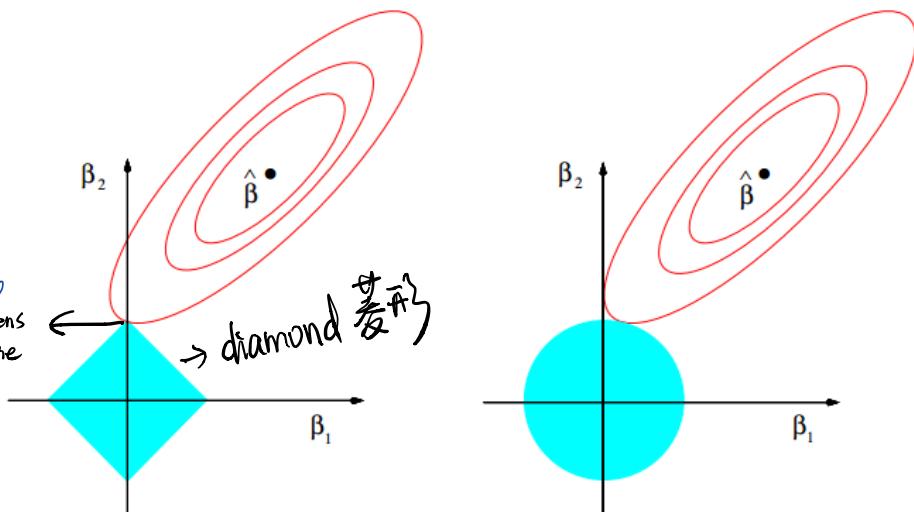
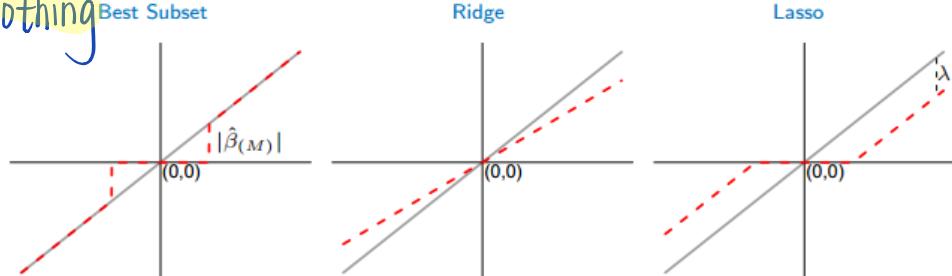


FIGURE 3.11. Estimation picture for the lasso (left) and ridge regression (right). Shown are contours of the error and constraint functions. The solid blue areas are the constraint regions $|\beta_1| + |\beta_2| \leq t$ and $\beta_1^2 + \beta_2^2 \leq t^2$, respectively, while the red ellipses are the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

Generalize them and view them as Bayes estimates.

The criterion:

$$\tilde{\beta} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{j=1}^N (y_j - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p |\beta_j|^q \right\} \text{ for } q \geq 0$$

log-prior density for β_j : $|\beta_j|^q$

非凸约束

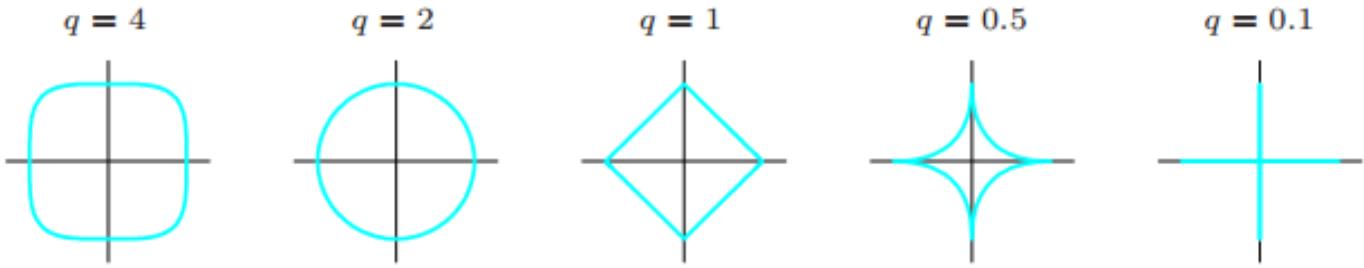


FIGURE 3.12. Contours of constant value of $\sum_j |\beta_j|^q$ for given values of q .

Notice that when $q \leq 1$, the prior is not uniform in direction, but concentrates more mass in the coordinate directions.

注意到 $q \leq 1$, 先验在各方向上不是均匀的, 而是更多地集中在坐标方向上. 对应 $q = 1$ 情形的先验分布是关于每个输入变量是的独立的二重指数分布 (或者 Laplace 分布), 概率密度为 $(1/2\tau)\exp(-|\beta|)/\tau$ 并且 $\tau = 1/\lambda$. $q=1$ 的情形 (lasso) 是使得约束区域为凸的最小 q 值; 非凸约束区域使得优化问题很困难.

★ Combination of Ridge & Lasso regression

Elastic Net $\lambda \sum_{j=1}^p (\alpha \beta_j^2 + (1-\alpha) |\beta_j|)$ which is a different compromise between ridge & lasso.
↓ further in Section 18.4

3.4.4 Least Angle Regression (LAR)

LAR can be viewed as a kind of "democratic" version of forward stepwise regression.

$\Gamma_k = y - X_{A_k} \beta_{A_k}$ is the current residual, then the direction for this step is

$$s_k = (X_{A_k}^T X_{A_k})^{-1} X_{A_k}^T \Gamma_k$$

↓ just like $\hat{\beta} = (X^T X)^{-1} X^T Y$

Steps: 1st direction $s_1 = (X_{A_1}^T X_{A_1})^{-1} X_{A_1}^T Y$

↳ (Just moving forward this direction, not reach to it)

2nd: get the direction for this step: $A_1 = (x_1)$

$$s_2 = (X_{A_2}^T X_{A_2})^{-1} X_{A_2}^T (Y - X_{A_1} \beta_{A_1})$$

coefficient profile evolves as

$$\beta_{A_2}(\alpha) = \beta_{A_2} + \alpha \cdot s_2$$

↳ This β_{A_2} just enters as 0

when x_2 is added

Ex. 3.23

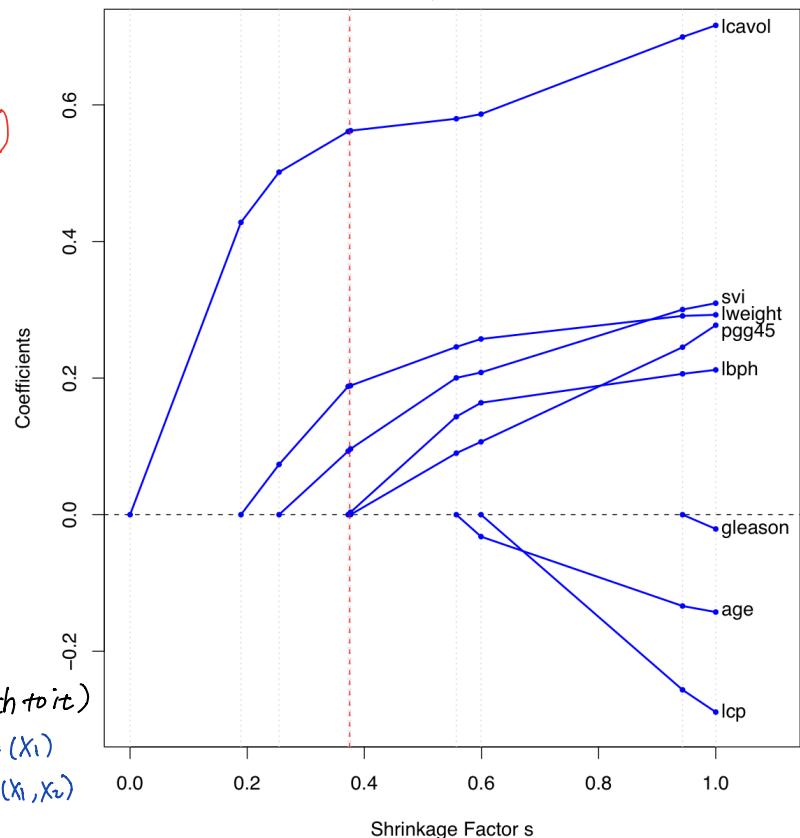


FIGURE 3.10. Profiles of lasso coefficients, as the tuning parameter t is varied. Coefficients are plotted versus $s = t / \sum_j |\beta_j|$. A vertical line is drawn at $s = 0.36$, the value chosen by cross-validation. Compare Figure 3.8 on page 65; the lasso profiles hit zero, while those for ridge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

Algorithm 3.2 Least Angle Regression.

1. Standardize the predictors to have mean zero and unit norm. Start with the residual $\mathbf{r} = \mathbf{y} - \bar{\mathbf{y}}$, $\beta_1, \beta_2, \dots, \beta_p = 0$.
2. Find the predictor \mathbf{x}_j most correlated with \mathbf{r} .

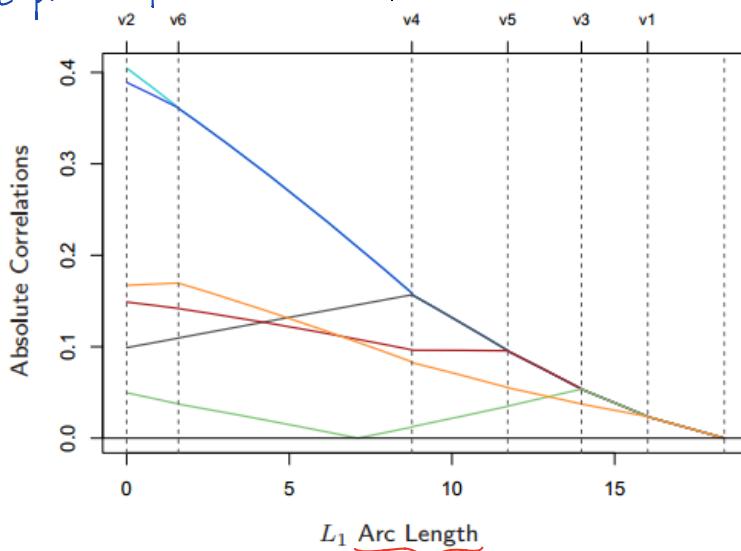
$$\hat{\beta}_j \leftarrow [0, \langle \mathbf{x}_j, \mathbf{r} \rangle]$$
3. Move β_j from 0 towards its least-squares coefficient $\langle \mathbf{x}_j, \mathbf{r} \rangle$, until some other competitor \mathbf{x}_k has as much correlation with the current residual as does \mathbf{x}_j .
4. Move β_j and β_k in the direction defined by their joint least squares coefficient of the current residual on $(\mathbf{x}_j, \mathbf{x}_k)$, until some other competitor \mathbf{x}_l has as much correlation with the current residual.
5. Continue in this way until all p predictors have been entered. After $\min(N - 1, p)$ steps, we arrive at the full least-squares solution.

手写注释：
 ① 与 \mathbf{r} 的相关性
 ② 方向不一致
 ③ 最小二乘法
 ④ 方向不一致

similar steps for fit vector \hat{f}_k , which would evolve as $\hat{f}_k(\alpha) = \hat{f}_k + \alpha \cdot \mathbf{u}_k$, where $\mathbf{u}_k = X_{Ak} \delta_k$

Ex 3.24

Whole path is piecewise-linear. 分段线性



$$\text{arc length} = 2\pi r \left(\frac{\theta}{360} \right)$$

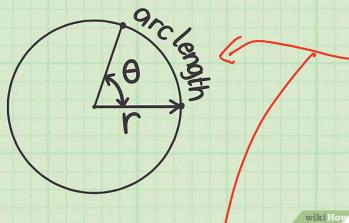


FIGURE 3.14. Progression of the absolute correlations during each step of the LAR procedure, using a simulated data set with six predictors. The labels at the top of the plot indicate which variables enter the active set at each step. The step length are measured in units of L_1 arc length.

²The L_1 arc-length of a differentiable curve $\beta(s)$ for $s \in [0, S]$ is given by $\text{TV}(\beta, S) = \int_0^S \|\dot{\beta}(s)\|_1 ds$, where $\dot{\beta}(s) = \partial \beta(s) / \partial s$. For the piecewise-linear LAR coefficient profile, this amounts to summing the L_1 norms of the changes in coefficients from step to step.

when $p \gg N$. Osborne et al. (2000a) also discovered a piecewise-linear path for computing the lasso, which they called a *homotopy* algorithm.

Effective degrees of freedom of the fitted vector $\hat{\mathbf{y}} = (\hat{y}_1, \hat{y}_2, \dots, \hat{y}_N)$ as

$$\text{Section 7.4-7.6 df}(\hat{\mathbf{y}}) = \frac{1}{\sigma^2} \sum_{i=1}^N \text{Cov}(\hat{y}_i, y_i) \quad (3.60)$$

$\text{predicted value } \hat{y}_i \leftarrow \rightarrow \text{corresponding outcome value } y_i$

3.5 Methods Using Derived Input Directions

3.5.1 Principal Components Regression

Components $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_M$ are **standardized X**

$$\hat{\mathbf{y}}^{PCR} = \bar{y}\mathbf{1} + \sum_{m=1}^M \hat{\theta}_m \mathbf{z}_m \quad \hat{\beta}^{PCR}(M) = \sum_{m=1}^M \hat{\theta}_m v_m \quad \text{where } \mathbf{z}_m = \mathbf{X}v_m$$

$\hat{\mathbf{y}} = \mathbf{X}\hat{\beta} = \mathbf{Z}\hat{\beta}^{PCR} = \mathbf{X}v\hat{\beta} \quad \text{if } \hat{\beta} = \begin{bmatrix} \hat{\theta}_1 \\ \vdots \\ \hat{\theta}_M \\ \vdots \\ \hat{\theta}_n \end{bmatrix}$

the m th principal component direction v_m solves:

$$\max_{\alpha} \text{Var}(\mathbf{X}\alpha) \quad (3.63)$$

subject to $\|\alpha\| = 1$, $\alpha^T \mathbf{S}v_\ell = 0$, $\ell = 1, \dots, m-1$,

where \mathbf{S} is the sample covariance matrix of the \mathbf{x}_j . The conditions $\alpha^T \mathbf{S}v_\ell = 0$ ensures that $\mathbf{z}_m = \mathbf{X}\alpha$ is uncorrelated with all the previous linear combinations $\mathbf{z}_\ell = \mathbf{X}v_\ell$.

3.5.2 Partial Least Squares **standardized X**

Algorithm 3.3 Partial Least Squares.

1. Standardize each \mathbf{x}_j to have mean zero and variance one. Set $\hat{\mathbf{y}}^{(0)} = \bar{y}\mathbf{1}$, and $\mathbf{x}_j^{(0)} = \mathbf{x}_j$, $j = 1, \dots, p$.
2. For $m = 1, 2, \dots, p$
 - (a) $\mathbf{z}_m = \sum_{j=1}^p \hat{\varphi}_{mj} \mathbf{x}_j^{(m-1)}$, where $\hat{\varphi}_{mj} = \langle \mathbf{x}_j^{(m-1)}, \mathbf{y} \rangle$.
 - (b) $\hat{\theta}_m = \langle \mathbf{z}_m, \mathbf{y} \rangle / \langle \mathbf{z}_m, \mathbf{z}_m \rangle$.
 - (c) $\hat{\mathbf{y}}^{(m)} = \hat{\mathbf{y}}^{(m-1)} + \hat{\theta}_m \mathbf{z}_m$.
 - (d) Orthogonalize each $\mathbf{x}_j^{(m-1)}$ with respect to \mathbf{z}_m : $\mathbf{x}_j^{(m)} = \mathbf{x}_j^{(m-1)} - [\langle \mathbf{z}_m, \mathbf{x}_j^{(m-1)} \rangle / \langle \mathbf{z}_m, \mathbf{z}_m \rangle] \mathbf{z}_m$, $j = 1, 2, \dots, p$.
3. Output the sequence of fitted vectors $\{\hat{\mathbf{y}}^{(m)}\}_1^p$. Since the $\{\mathbf{z}_\ell\}_1^m$ are linear in the original \mathbf{x}_j , so is $\hat{\mathbf{y}}^{(m)} = \mathbf{X}\hat{\beta}^{PLS}(m)$. These linear coefficients can be recovered from the sequence of PLS transformations.

PLS's solution path is a **nonlinear fcn of y**.

The m th PLS direction $\hat{\varphi}_m$ solves:

$$\max_{\alpha} \text{Corr}^2(\mathbf{y}, \mathbf{X}\alpha) \text{Var}(\mathbf{X}\alpha) \quad (3.64)$$

subject to $\|\alpha\| = 1$, $\alpha^T \mathbf{S}\hat{\varphi}_\ell = 0$, $\ell = 1, \dots, m-1$.

Ex. 3.14

Ex. 3.15

Ex. 3.18

3.6 Discussion: A Comparison of the Selection and Shrinkage Methods From Frank and Friedman (1993):

Method	Properties	
Best subset		
Ridge	shrinks all direction, but the lower-variance direction more	
Lasso	middle between ridge & Best subset, enjoy some of the properties of each	
PCA	leaves M high-variance directions only and discard the rest	
PLS	shrinks low-variance directions, & also inflate some high-var directions	a little unstable, a little higher prediction error compared to ridge

3.7 Multiple Outcome Shrinkage and Selection

I. Canonical correlation analysis (CCA)

Ex. 3.20

Algorithm 3.4 Incremental Forward Stagewise Regression— FS_ϵ .

- Start with the residual \mathbf{r} equal to \mathbf{y} and $\beta_1, \beta_2, \dots, \beta_p = 0$. All the predictors are standardized to have mean zero and unit norm.
 - Find the predictor \mathbf{x}_j most correlated with \mathbf{r}
 - Update $\beta_j \leftarrow \beta_j + \delta_j$, where $\delta_j = \epsilon \cdot \text{sign}[\langle \mathbf{x}_j, \mathbf{r} \rangle]$ and $\epsilon > 0$ is a small step size, and set $\mathbf{r} \leftarrow \mathbf{r} - \delta_j \mathbf{x}_j$.
(β的运动方向始终与最小二乘法方向一致)
 \downarrow
 $\text{sign}[\langle \mathbf{x}_j, \mathbf{r} \rangle]$
 - Repeat steps 2 and 3 many times, until the residuals are uncorrelated with all the predictors.
-

II. FS₀ ☆.

Algorithm 3.2b Least Angle Regression: FS₀ Modification.

4. Find the new direction by solving the constrained least squares problem

$$\min_b \|\mathbf{r} - \mathbf{X}_{\mathcal{A}} b\|_2^2 \text{ subject to } b_j s_j \geq 0, j \in \mathcal{A},$$

where s_j is the sign of $\langle \mathbf{x}_j, \mathbf{r} \rangle$.

III. Group Lasso

The grouped-lasso minimizes the convex criterion

$$\min_{\beta \in \mathbb{R}^p} \left(\|\mathbf{y} - \beta_0 \mathbf{1} - \sum_{\ell=1}^L \mathbf{X}_{\ell} \beta_{\ell}\|_2^2 + \lambda \sum_{\ell=1}^L \sqrt{p_{\ell}} \|\beta_{\ell}\|_2 \right), \quad (3.80)$$

where the $\sqrt{p_{\ell}}$ terms accounts for the varying group sizes, and $\|\cdot\|_2$ is the Euclidean norm (not squared). Since the Euclidean norm of a vector β_{ℓ} is zero only if all of its components are zero, this procedure encourages sparsity at both the group and individual levels. That is, for some values of λ , an entire group of predictors may drop out of the model.

Generalizations include more general L_2 norms $\|\eta\|_K = (\eta^T K \eta)^{1/2}$, as well as allowing overlapping groups of predictors (Zhao et al., 2008).

3.8.5 To be continued.