## *Linear Methods for Regression - Hastie – Chap – III*

# (part – B ; PRML – CS5691)

# Shrinkage Methods

 By retaining a subset of the predictors and discarding the rest, subset selection produces a model that is interpretable and has possibly lower prediction error than the full model. However, because it is a discrete process—variables are either retained or discarded—it often exhibits high variance, and so doesn't reduce the prediction error of the full model. Shrinkage methods are more continuous, and don't suffer as much from high variability.

# Ridge Regression

• Ridge regression shrinks the regression coefficients by imposing a penalty on their size.



FIGURE 3.7. Estimated prediction error curves and their standard errors for the various selection and shrinkage methods. Each curve is plotted as a function of the corresponding complexity parameter for that method. The horizontal axis has been chosen so that the model complexity increases as we move from left to right. The estimates of prediction error and their standard errors were obtained by tenfold cross-validation; full details are given in Section 7.10. The least complex model within one standard error of the best is chosen, indicated by the purple vertical broken lines.

 The ridge coefficients minimize a penalized residual sum of squares,

$$\hat{\beta}^{ridge} = argmin_{\beta} \left\{ \sum_{i=1}^{N} (y_i - \beta_0 - \sum_{j=1}^{p} x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^{p} \beta_j^2 \right\}.$$
(3.41)

- Here  $\lambda \ge 0$  is a complexity parameter that controls the amount of shrinkage: the larger the value of  $\lambda$ , the greater the amount of shrinkage. The coefficients are shrunk toward zero (and each other).
- An equivalent way to write the ridge problem is

$$\hat{\beta}^{ridge} = argmin_{\beta} \sum_{i=1}^{N} (y_i - \beta_0 - \sum_{j=1}^{p} x_{ij}\beta_j)^2,$$

$$(3.42)$$

$$subject \ to \ \sum_{j=1}^{p} \beta_j^2 \le t,$$

- Which makes explicit the size constraint on the parameters. There is a one to-one correspondence between the parameters  $\lambda$  in (3.41) and t in (3.42). When there are many correlated variables in a linear regression model, their coefficients can become poorly determined and exhibit high variance. A wildly large positive coefficient on one variable can be canceled by a similarly large negative coefficient on its correlated cousin. By imposing a size constraint on the coefficients, as in (3.42), this problem is alleviated.
- The ridge solutions are not equivariant under scaling of the inputs, and so one normally standardizes the inputs before solving (3.41). The solution to (3.41) can be separated into two parts, after reparametrization using centered inputs: each  $x_{ij}$  gets replaced by  $x_{ij} \bar{x_j}$ . We estimate  $\beta_0$  by  $\bar{y} = \frac{1}{N} \sum_{i=1}^{N} y_i$ .
- Read about the process of whitening

- The remaining coefficients get estimated by a ridge regression without intercept, using the centered x<sub>ij</sub>. Henceforth we assume that this centering has been done, so that the input matrix X has p (rather than p + 1) columns.
- Writing the criterion in (3.41) in matrix form,

$$RSS(\lambda) = (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) + \lambda\beta^T \beta, \qquad (3.43)$$

• The ridge regression solutions are easily seen to be

$$\hat{\beta}^{ridge} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}, \qquad (3.44)$$

where, I is the  $p \times p$  identity matrix. Notice that with the choice of quadratic penalty  $\beta^T \beta$ , the ridge regression solution is again a linear function of y. The solution adds a positive constant to the diagonal of  $X^T X$  before inversion.

- This makes the problem nonsingular, even if X<sup>T</sup>X is not of full rank, and was the main motivation for ridge regression when it was first introduced in statistics (Hoerl and Kennard, 1970).
- Traditional descriptions of ridge regression start with definition (3.44). We choose to motivate it via (3.41) and (3.42), as these provide insight into how it works.

- Ridge regression can also be derived as the mean or mode of a posterior distribution, with a suitably chosen prior distribution. In detail, suppose  $y_i \sim N(\beta_0 + x_i^T \beta, \sigma^2)$ , and the parameters  $\beta_j$  are each distributed as  $N(0, \tau^2)$ , independently of one another. Then the (negative) log-posterior density of  $\beta$ , with  $\tau^2$  and  $\sigma^2$  assumed known, is equal to the expression in curly braces in (3.41), with  $\lambda = \sigma^2/\tau^2$  Thus the ridge estimate is the mode of the posterior distribution; since the distribution is Gaussian, it is also the posterior mean (Ex – 3.6 - Hastie).
- The singular value decomposition (SVD) of the centered input matrix **X** gives us some additional insight into the nature of ridge regression. The SVD of the  $N \times p$  matrix **X** has the form

$$\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T$$

(3.45)

Stack the (centered) observations into the rows of an  $N \times p$  matrix X. We construct the singular value decomposition of X:

 $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T$ .

Sec. 14.5 – PP 535

#### (14.54)

This is a standard decomposition in numerical analysis, and many algorithms exist for its computation (Golub and Van Loan, 1983, for example). Here U is an  $N \times p$  orthogonal matrix  $(\mathbf{U}^T \mathbf{U} = \mathbf{I}_p)$  whose columns  $\mathbf{u}_j$  are called the *left singular vectors*; V is a  $p \times p$  orthogonal matrix  $(\mathbf{V}^T \mathbf{V} = \mathbf{I}_p)$  with columns  $v_j$  called the *right singular vectors*, and D is a  $p \times p$  diagonal matrix, with diagonal elements  $d_1 \geq d_2 \geq \cdots \geq d_p \geq 0$  known as the *sin*-

Here U and V are  $N \times p$  and  $p \times p$  orthogonal matrices, with the columns of U spanning the column space of X, and the columns of V spanning the row space. D is a  $p \times p$  diagonal matrix, with diagonal entries  $d_1 \ge d_2 \ge$  $\dots \ge d_p \ge 0$  called the singular values of X. If one or more values  $d_j = 0$ , X is singular.

Using the singular value decomposition we can write the least squares fitted vector as

$$\begin{aligned} \mathbf{X} \hat{\beta}^{\text{ls}} &= \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \\ &= \mathbf{U} \mathbf{U}^T \mathbf{y}, \end{aligned}$$

 Using the singular value decomposition we can write the least squares fitted vector as

$$\begin{aligned} \boldsymbol{X}\hat{\beta}^{ls} &= \boldsymbol{X}(\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{X}^T \boldsymbol{y}, \\ &= \boldsymbol{y}, \end{aligned} \tag{3.46}$$

 Note that U<sup>T</sup>y are the coordinates of y with respect to the orthonormal basis U.

Now the ridge solutions are

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

$$=$$

$$=$$

$$y$$

$$(3.47)$$

• Where the  $u_j$  are the columns of U. Note that since  $\lambda \ge 0$ , we have  $d_j^2/(d_j^2 + \lambda) \le 1$ . Like linear regression, ridge regression computes the coordinates of y with respect to the orthonormal basis U. It then shrinks these coordinates by the factors  $d_j^2/(d_j^2 + \lambda)$ . This means that a greater amount of shrinkage is applied to the coordinates of basis vectors with smaller  $d_j^2$ .

• What does a *small value of*  $d_j^2$  mean? The *SVD* of the centered **X** is another way of expressing the *principal components* of the variables in **X**. The sample covariance matrix is given by  $S = X^T X/N$ , and from (3.45) we have

$$\boldsymbol{X}^T \boldsymbol{X} = \boldsymbol{V} \boldsymbol{D}^2 \boldsymbol{V}^T, \qquad (3.48)$$

- Which is the eigen decomposition of X<sup>T</sup>X (and of S, up to a factor N). The eigenvectors v<sub>j</sub> (columns of V) are also called the principal components (or Karhunen–Loeve) directions of X.
- The first principal component direction v<sub>1</sub> has the property that z<sub>1</sub> = Xv<sub>1</sub> has the largest sample variance amongst all normalized linear combinations of the columns of X. Sample variance is easily seen to be

$$Var(\mathbf{z}_{1}) = Var(\mathbf{X}v_{1}) = \frac{d_{1}^{2}}{N},$$
 (3.49)

and in fact

$$\boldsymbol{z}_1 = \boldsymbol{X}\boldsymbol{v}_1 = \boldsymbol{u}_1\boldsymbol{d}_1$$

 Subsequent principal components z<sub>j</sub> have maximum variance d<sup>2</sup><sub>j</sub>/N, subject to being orthogonal to the earlier ones.
 Conversely the last principal component has minimum variance. Hence the small singular values d<sub>j</sub> correspond to directions in the column space of X having small variance, and ridge regression shrinks these directions the most.  In Figure 3.7 we have plotted the estimated prediction error versus the quantity (DoF):

$$df(\lambda) = tr[X(X^TX + \lambda I)^{-1}X^T],$$

 $= tr(\mathbf{H}_{\lambda})$ 

$$= \sum_{j=1}^{p} \frac{d_j^2}{d_j^2 + \lambda} \,. \tag{3.50}$$

• This monotone decreasing function of  $\lambda$  is the *effective degrees of freedom* of the ridge regression fit. Usually in a linear-regression fit with p variables, the degrees-of-freedom of the fit is p, the number of free parameters. The idea is that although all p coefficients in a ridge fit will be non-zero, they are fit in a restricted fashion controlled by  $\lambda$ . Note that  $df(\lambda) = p$  when  $\lambda = 0$  (no regularization) and  $df(\lambda) \rightarrow 0$  as  $\lambda \rightarrow \infty$ .

## The Lasso

• The lasso is a shrinkage method like ridge, with subtle but important differences. The lasso estimate is defined by

$$\hat{\beta}^{lasso} = \arg \min_{\beta} \sum_{i=1}^{N} (y_i - \beta_0 - \sum_{j=1}^{p} x_{ij} \beta_j)^2$$

$$subject \ to \ \sum_{j=1}^{p} |\beta_j|. \tag{3.51}$$

• Write the lasso problem in the equivalent Lagrangian form

$$\hat{\beta}^{lasso} = argmin_{\beta} \left\{ \frac{1}{2} \sum_{i=1}^{N} (y_i - \beta_0 - \sum_{j=1}^{p} x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^{p} |\beta_j| \right\}.$$
(3.52)

• Notice the similarity to the ridge regression problem (3.42) or (3.41): the  $L_2$  ridge penalty  $\sum_{1}^{p} \beta_{j}^{2}$  is replaced by the  $L_1$  lasso penalty  $\sum_{1}^{p} |\beta_{j}|$ . Thus the lasso does a kind of continuous subset selection. If t is chosen larger than  $t_0 = \sum_{1}^{p} |\hat{\beta}_{j}|$  (where  $\hat{\beta}_{j} = \hat{\beta}_{j}^{ls}$ , the least squares estimates), then the lasso estimates are the  $\hat{\beta}_{j}$ 's. On the other hand, for  $t = t_0/2$  say, then the least squares coefficients are shrunk by about 50% on average.

## Discussion: Subset Selection, Ridge Regression and the Lasso

- In the case of an orthonormal input matrix **X** the three procedures have explicit solutions. Each method applies a simple transformation to the least squares estimate  $\hat{\beta}_j$ , as detailed in Table 3.4.
- Ridge regression does a proportional shrinkage. Lasso translates each coefficient by a constant factor λ, truncating at zero. This is called "soft thresholding,". Best-subset selection drops all variables with coefficients smaller than the M<sup>th</sup> largest; this is a form of "hard-thresholding."
- Back to the no orthogonal case; some pictures help understand their relationship. Figure 3.11 depicts the lasso (*left*) and ridge regression (*right*) when there are only two parameters. The residual sum of squares has elliptical contours, centered at the full least squares estimate.

**TABLE 3.4.** Estimators of  $\beta_j$  in the case of orthonormal columns of **X**. *M* and  $\lambda$  are constants chosen by the corresponding techniques; sign denotes the sign of its argument ( $\pm 1$ ), 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.





**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  $|\beta_1| + |\beta_2| \le t$  and  $\beta_1^2 + \beta_2^2 \le t^2$ , respectively, while the red ellipses are the contours of the least squares error function.

- Region for ridge regression is the disk  $\beta_1^2 + \beta_2^2 \le t$ , while that for lasso is the diamond  $|\beta_1| + |\beta_2| \le t$ . Both methods find the first point where the elliptical contours hit the constraint region. Unlike the disk, the diamond has corners; if the solution occurs at a corner, then it has one parameter  $\beta_j$  equal to zero. When p > 2, the diamond becomes a rhomboid, and has many corners, flat edges and faces; there are many more opportunities for the estimated parameters to be zero.
- Consider the criterion

$$\tilde{\beta} = argmin_{\beta} \left\{ \sum_{i=1}^{N} (yi - \beta_0 - \sum_{j=1}^{p} x_{ij}\beta_j)^2 + \lambda \sum_{j=1}^{p} |\beta_j|^q \right\} \quad (3.53)$$

for  $q \ge 0$ . The contours of constant value of  $\sum_{j} |\beta_{j}|^{q}$  are shown in Figure 3.12, for the case of two inputs.

Thinking of |β<sub>j</sub>|<sup>q</sup> as the log-prior density. The case
 q = 1 (lasso) is the smallest q such that the constraint region is convex; non-convex constraint regions make the optimization problem more difficult. In this view, the lasso, ridge regression and best subset selection are Bayes estimates with different priors. They are derived as posterior modes, that is, maximizers of the posterior.



**FIGURE 3.12**. Contours of constant value of  $\sum_{j} |\beta_{j}|^{q}$  for given values of q.



- **FIGURE 3.13.** Contours of constant value of  $\sum_{j} |\beta_{j}|^{q}$  for q = 1.2 (left plot), and the elastic-net penalty  $\sum_{j} (\alpha \beta_{j}^{2} + (1 \alpha) |\beta_{j}|)$  for  $\alpha = 0.2$  (right plot). Although visually very similar, the elastic
- *Zou and Hastie* (2005) *introduced the elastic net penalty*

$$\lambda \sum_{j=1}^{p} (\alpha \beta_j^2 + (1 - \alpha) |\beta_j|), \qquad (3.54)$$

# The Grouped Lasso

• In some problems, the predictors belong to pre-defined groups; In this situation it may be desirable to shrink and select the members of a group together. The grouped lasso is one way to achieve this. Suppose that the p predictors are divided into *L* groups, with p` the number in group  $\ell$ . For ease of notation, we use a matrix  $X_{\ell}$  to represent the predictors corresponding to the  $\ell th$  group, with corresponding coefficient vector  $\beta_{\ell}$ . The grouped-lasso minimizes the convex criterion

• 
$$\min_{\beta \in \mathbb{R}^p} \left( ||y - \beta_0 \mathbf{1} - \sum_{\ell=1}^L X_\ell \beta_\ell ||_2^2 + \lambda \sum_{\ell=1}^L \sqrt{p_\ell} ||\beta_\ell||_2 \right),$$
 (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<sub>ℓ</sub> 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. This procedure was proposed by Bakin (1999) and Lin and Zhang (2006), and studied and generalized by Yuan and Lin (2007).

# Further Properties of the Lasso

 A number of authors have studied the ability of the lasso and related procedures to recover the correct model, as *N* and *p* grow. Examples of this work include Knight and Fu (2000), Greenshtein and Ritov (2004), Tropp (2004), Donoho (2006b), Meinshausen (2007), Meinshausen and B¨uhlmann (2006), Tropp (2006), Zhao and Yu (2006), Wainwright (2006), and Bunea et al. (2007). Alternatively, one can modify the lasso penalty function so that larger coefficients are shrunken less severely; the smoothly clipped absolute deviation (SCAD) penalty of Fan and Li (2005) replaces  $\lambda|\beta|$  by  $J_a(\beta, \lambda)$ , where

$$\frac{dJ_a(\beta,\lambda)}{d\beta} = \lambda \cdot \operatorname{sign}(\beta) \left[ I(|\beta| \le \lambda) + \frac{(a\lambda - |\beta|)_+}{(a-1)\lambda} I(|\beta| > \lambda) \right] \quad (3.82)$$

for some  $a \ge 2$ . The second term in square-braces reduces the amount of shrinkage in the lasso for larger values of  $\beta$ , with ultimately no shrinkage as  $a \to \infty$ . Figure 3.20 shows the SCAD penalty, along with the lasso and



FIGURE 3.20. The lasso and two alternative non-convex penalties designed to penalize large coefficients less. For SCAD we use  $\lambda = 1$  and a = 4, and  $\nu = \frac{1}{2}$  in the last panel.

 $|\beta|^{1-\nu}$ . However this criterion is non-convex, which is a drawback since it makes the computation much more difficult. The *adaptive lasso* (Zou, 2006) uses a weighted penalty of the form  $\sum_{j=1}^{p} w_j |\beta_j|$  where  $w_j = 1/|\hat{\beta}_j|^{\nu}$ ,  $\hat{\beta}_j$  is the ordinary least squares estimate and  $\nu > 0$ . This is a practical approximation to the  $|\beta|^q$  penalties ( $q = 1 - \nu$  here) discussed in Section 3.4.3. The adaptive lasso yields consistent estimates of the parameters while retaining the attractive convexity property of the lasso.

#### **Other Competing methods:**

LAR PLS PCR FSW FS₀

## **Computational Considerations**

• Least squares fitting is usually done via the Cholesky decomposition of the matrix  $X^T X$  or a QR decomposition of X. With N observations and p features, the Cholesky decomposition requires  $p^3 + Np^2/2$  operations, while the QR decomposition requires  $Np^2$  operations. Depending on the relative size of N and p, the Cholesky can sometimes be faster; on the other hand, it can be less numerically stable (Lawson and Hansen, 1974). Computation of the lasso via the *LAR* algorithm has the same order of computation as a least squares fit.

# Discussion: A Comparison of the Selection and Shrinkage Methods

- To summarize, *PLS*, *PCR* and **ridge regression** tend to behave similarly.
- Ridge regression may be preferred because it shrinks smoothly, rather than in discrete steps. Lasso falls somewhere between ridge regression and best subset regression, and enjoys some of the properties of each.

# Least Angle Regression

- Least angle regression (*LAR*) is a relative newcomer (*Efron et al.*, 2004), and can be viewed as a kind of "*democratic*" version of forward stepwise regression
- Figure 3.10. Forward stepwise regression builds a model sequentially, adding one variable at a time. At each step, it identifies the best variable to include in the active set, and then updates the least squares fit to include all the active variables.
- Least angle regression uses a similar strategy, but only enters "as much" of a predictor as it deserves. At the first step it identifies the variable most correlated with the response. Rather than fit this variable completely, *LAR* moves the coefficient of this variable continuously toward its leastsquares value (causing its correlation with the evolving residual to decrease in absolute value).

#### As soon as another variable "catches up" in terms of correlation with the residual, the process is paused. The second variable then joins the active set, and their

coefficients are moved together in a way that keeps their correlations tied and decreasing. This process is continued until all the variables are in the model, and ends at the full least-squares fit. Algorithm 3.2 provides the details.

#### Algorithm 3.2 Least Angle Regression.

1. Standardize the predictors to have mean zero and unit norm. Start with the residual  $r = y - \overline{y}, \beta_1, \beta_2, \dots, \beta_p = 0$ .

2. Find the predictor  $\mathbf{x}_i$  most correlated with r.

3. Move  $\beta_j$  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  $\beta_j$  and 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.

• Suppose  $A_k$  is the active set of variables at the beginning of the  $k^{\text{th}}$ step, and let  $\beta_{A_k}$  be the coefficient vector for these variables at this step; there will be k - 1 nonzero values, and the one just entered will be zero. If  $\boldsymbol{r}_k = \boldsymbol{y} - \boldsymbol{X}_{A_k}\beta_{A_k}$  is the current residual, then the direction for this step is

$$\delta_k = \left( \boldsymbol{X}_{A_k}^T \boldsymbol{X} \right)^{-1} \boldsymbol{X}_{Ak}^T \boldsymbol{r}_k.$$
 (3.35)

- The coefficient profile then evolves as  $\beta_{A_k}(\alpha) = \beta_{A_k} + \alpha \cdot \delta_k$ .
- By construction the coefficients in LAR change in a piecewise linear fashion.

Algorithm 3.2a Least Angle Regression: Lasso Modification.

4*a*. If a non-zero coefficient hits zero, drop its variable from the active set of variables and recompute the current joint least squares direction.

• The LAR(lasso) algorithm is extremely efficient, requiring the same order of computation as that of a single least squares fit using the p predictors. Least angle regression always takes p steps to get to the full least squares estimates. The lasso path can have more than p steps, although the two are often quite similar. Algorithm 3.2 with the lasso modification 3.2a is an efficient way of computing the solution to any lasso problem, especially when  $p \gg N$ 

#### Methods Using Derived Input Directions

• In many situations we have a large number of inputs, often very correlated. The methods in this section produce a small number of linear combinations  $Z_m, m = 1, ..., M$  of the original inputs  $X_j$ , and the  $Z_m$  are then used in place of the  $X_j$  as inputs in the regression. The methods differ in how the linear combinations are constructed.

# Principal Components Regression

• In this approach the linear combinations  $Z_m$  used

# **Partial Least Squares**

 This technique also constructs a set of linear combinations of the inputs for regression, but unlike principal components regression it uses y (in addition to X) for this construction. Like principal component regression, partial least squares (*PLS*) is not scale invariant, so we assume that each  $\mathbf{x}_i$  is standardized to have mean 0 and variance 1. PLS begins by computing  $\hat{\varphi}_{1i} = \langle \mathbf{x}_i, \mathbf{y} \rangle$  for each *j*. From this we construct the derived input  $\mathbf{z}_1 = \sum_i \hat{\varphi}_{1i} \mathbf{x}_i$ , which is the first partial least squares direction. The outcome y is regressed on  $z_1$  giving coefficient  $\hat{\theta}_1$ , and then we orthogonalize  $\mathbf{x}_1, \ldots, \mathbf{x}_p$  with respect to  $z_1$ . We continue this process, until  $M \le p$  directions have been obtained. In this manner, partial least squares produces a sequence of derived, orthogonal inputs or directions  $z_1, z_2, \dots, z_M$ . As with principal-component regression, if we were to construct all M = p directions, we would get back a solution equivalent to the usual least squares estimates; using M < p directions produces a reduced regression.

#### 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  $\hat{\mathbf{y}}^{(m)} = \mathbf{X}\hat{\beta}^{\text{pls}}(m)$ . These linear coefficients can be recovered from the sequence of PLS transformations.

• What optimization problem is partial least squares solving? Since it uses the response y to construct its directions, its solution path is a nonlinear function of y. It can be shown (*Exercise* 3.15) that partial least squares seeks directions that have high variance and have high correlation with the response, in contrast to principal components regression which keys only on high variance the  $m^{th}$  principal component direction  $v_m$  solves:

$$\max_{\alpha} Var(\boldsymbol{X}_{\alpha})$$
(3.63)
$$subject \ to \ ||\alpha|| = 1, \alpha^{T} \boldsymbol{S} v_{\ell} = 0, \ \ell = 1, \dots, m - 1,$$

• here s is the sample covariance matrix of the  $\mathbf{x}_i$ .

• The *mth PLS* direction  $\hat{\varphi}_m$  solves:

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

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

• If the input matrix X is orthogonal, then partial least squares finds the least squares estimates after m = 1 steps.

#### Incremental Forward Stagewise Regression

Algorithm 3.4 Incremental Forward Stagewise Regression— $FS_{\epsilon}$ .

- 1. Start with the residual **r** equal to **y** and  $\beta_1, \beta_2, \ldots, \beta_p = 0$ . All the predictors are standardized to have mean zero and unit norm.
- 2. Find the predictor  $\mathbf{x}_j$  most correlated with  $\mathbf{r}$
- 3. 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  $\mathbf{r} \leftarrow \mathbf{r} \delta_j \mathbf{x}_j$ .
- 4. Repeat steps 2 and 3 many times, until the residuals are uncorrelated with all the predictors.

Algorithm 3.2b Least Angle Regression: FS<sub>0</sub> Modification.

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

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

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



**FIGURE 3.16**. Comparison of LAR and lasso with forward stepwise, forward stagewise (FS) and incremental forward stagewise (FS<sub>0</sub>) regression. The setup is the same as in Figure 3.6, except N = 100 here rather than 300. Here the slower FS regression ultimately outperforms forward stepwise. LAR and lasso show similar behavior to FS and FS<sub>0</sub>. Since the procedures take different numbers of steps (across simulation replicates and methods), we plot the MSE as a function of the fraction of total L1 arc-length toward the least-squares fit.