

3. Methods Based on Linear Regression

The methods discussed in this section are primarily intended for high-dimensional situations, where p is large relative to N , perhaps much larger than N .

We can write

$$f(\mathbf{x}) = \beta_0 + \mathbf{x}\boldsymbol{\beta} = \beta_0 + \sum_{j=1}^p \beta_j x_j, \quad (1)$$

treating the constant and other regressors separately. In the notation of ETM, p denotes $k - 1$, and N denotes n .

Of course, the x_j can include transformations of inputs, such as square roots, logs, squares, or higher powers, various types of dummy variables, and products of inputs or of transformations of them.

So a “linear” regression model can actually allow for very nonlinear relationships.

3.1. Regression by Successive Orthogonalization

When the \mathbf{x}_j are orthogonal, the $\hat{\beta}_j = \mathbf{x}_j^\top \mathbf{y} / \mathbf{x}_j^\top \mathbf{x}_j$ obtained by univariate regressions of \mathbf{y} on each of the \mathbf{x}_j are equal to the multiple regression least squares estimates.

We can accomplish something similar even when the data are not orthogonal. This is called **Gram-Schmidt Orthogonalization**. It is related to the FWL Theorem.

1. Set $\mathbf{z}_0 = \mathbf{x}_0 = \mathbf{1}$.
2. Regress \mathbf{x}_1 on \mathbf{z}_0 to obtain a residual vector \mathbf{z}_1 and coefficient $\hat{\gamma}_{01}$.
3. For $j = 2, \dots, p$, regress \mathbf{x}_j on $\mathbf{z}_0, \dots, \mathbf{z}_{j-1}$ to obtain a residual vector \mathbf{z}_j and coefficients $\hat{\gamma}_{\ell j}$ for $\ell = 0, \dots, j - 1$.
4. Regress \mathbf{y} on \mathbf{z}_p to obtain $\hat{\beta}_p = \mathbf{z}_p^\top \mathbf{y} / \mathbf{z}_p^\top \mathbf{z}_p$.

The scalar $\hat{\beta}_p$ is the coefficient on \mathbf{x}_p in the multiple regression of \mathbf{y} on a constant and all the \mathbf{x}_j .

Note that step 3 simply involves j univariate regressions.

If we did this $p + 1$ times, changing the ordering so that each of the \mathbf{x}_j came last, we could obtain all the OLS coefficients.

The **QR decomposition** of \mathbf{X} is

$$\mathbf{X} = \mathbf{Z}\mathbf{D}^{-1}\mathbf{D}\mathbf{\Gamma} = \mathbf{Q}\mathbf{R}, \quad (2)$$

where \mathbf{Z} contains the columns \mathbf{z}_0 to \mathbf{z}_p (in order), $\mathbf{\Gamma}$ is an upper triangular matrix containing the $\hat{\gamma}_{\ell j}$, and \mathbf{D} is a diagonal matrix with typical diagonal element $\|\mathbf{z}_j\|$. Thus $\mathbf{Q} = \mathbf{Z}\mathbf{D}^{-1}$ and $\mathbf{R} = \mathbf{D}\mathbf{\Gamma}$.

It can be seen that \mathbf{Q} is an $N \times (p+1)$ orthogonal matrix with the property that $\mathbf{Q}^\top \mathbf{Q} = \mathbf{I}$, and \mathbf{R} is a $(p+1) \times (p+1)$ upper triangular matrix.

The QR decomposition provides a convenient orthogonal basis for $\mathcal{S}(\mathbf{X})$.

It is easy to see that

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = \mathbf{R}^{-1} \mathbf{Q}^\top \mathbf{y}, \quad (3)$$

because

$$(\mathbf{X}^\top \mathbf{X})^{-1} = (\mathbf{R}^\top \mathbf{Q}^\top \mathbf{Q} \mathbf{R})^{-1} = (\mathbf{R}^\top \mathbf{R})^{-1} = \mathbf{R}^{-1} (\mathbf{R}^\top)^{-1},$$

and

$$\mathbf{X}^\top \mathbf{y} = \mathbf{R}^\top \mathbf{Q}^\top \mathbf{y}.$$

Since \mathbf{R} is triangular, it is very easy to invert. Also,

$$\hat{\mathbf{y}} = \mathbf{P}_X \mathbf{y} = \mathbf{Q}\mathbf{Q}^\top \mathbf{y}. \quad (4)$$

so that

$$\text{SSR}(\hat{\boldsymbol{\beta}}) = (\mathbf{y} - \mathbf{Q}\mathbf{Q}^\top \mathbf{y})^\top (\mathbf{y} - \mathbf{Q}\mathbf{Q}^\top \mathbf{y}) = \mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{Q}\mathbf{Q}^\top \mathbf{y}. \quad (5)$$

The diagonals of the hat matrix are the diagonals of $\mathbf{Q}\mathbf{Q}^\top$.

The QR decomposition is easy to implement and numerically stable. It requires $O(Np^2)$ operations.

An alternative is to form $\mathbf{X}^\top \mathbf{X}$ and $\mathbf{X}^\top \mathbf{y}$ efficiently, use the Cholesky decomposition to invert $\mathbf{X}^\top \mathbf{X}$, and then multiply $\mathbf{X}^\top \mathbf{y}$ by $(\mathbf{X}^\top \mathbf{X})^{-1}$.

The Cholesky approach requires $O(p^3 + Np^2/2)$ operations. For $p \ll N$, Cholesky wins, but for large p QR does. Never use a general matrix inversion routine, because such routines are not as stable as QR or Cholesky.

In extreme cases, QR can yield reliable results when Cholesky does not. If $\mathbf{X}^\top \mathbf{X}$ is not well-conditioned (i.e., if the ratio of the largest to the smallest eigenvalue is large), then it is much safer to use QR than Cholesky. Intelligent scaling helps.

3.2. Subset Selection

When p is large, OLS tends to have low bias but large variance.

By setting some coefficients to zero, we can often reduce variance substantially while increasing bias only modestly.

It can be hard to make sense of models with a great many coefficients.

Best-subset regression tries every subset of size $k = 1, \dots, p$ to find the one that, for each k , has the smallest SSR.

There is an efficient algorithm, but it becomes infeasible once p becomes even moderately large (say much over 40).

Forward-stepwise selection and **backward-stepwise selection** are feasible even when p is large. The former works even if $p > N$.

To begin the forward-stepwise selection procedure, we regress \mathbf{y} on each of the \mathbf{x}_j and pick the one that yields the smallest SSR.

Next, we add each of the remaining $p - 1$ regressors, one at a time, and pick the one that yields the smallest SSR.

Then we add each of the remaining $p - 2$ regressors, one at a time, and pick the one that yields the smallest SSR.

If we are using the QR decomposition, we can do this efficiently. We are just adding one more column to \mathbf{Z} and hence \mathbf{Q} , and one more row to \mathbf{I} and hence \mathbf{R} .

We do this for each of the candidate regressors and calculate the SSR from (5) for each of them. We can reuse most of the elements of the old \mathbf{Q} matrix.

For backward-stepwise regression, we start with the full model (assuming that $p < N$) and then successively drop regressors. At each step, we drop the one with the smallest t statistic.

Smart stepwise programs will recognize that regressors often come in groups. It generally does not make sense to

- Drop a regressor X if the model includes X^2 or XZ ;
- Drop some dummy variables that are part of a set (e.g. some province dummies but not all, or some categorical dummies but not all).

In cases like this, variables should be added, or subtracted, in groups.

If we use either forward-stepwise or backward-stepwise selection, how do we decide which model to choose?

The model with all p regressors will always fit best, in terms of SSR.

We either need to penalize the number of parameters (d) or use cross-validation.

$$C_p: \frac{1}{N}(\text{SSR} + 2d\hat{\sigma}^2) = (1 + 2d/N)\hat{\sigma}^2. \quad (6)$$

Thus C_p adds a penalty of $2d\hat{\sigma}^2$ to SSR.

The **Akaike information criterion**, or **AIC**, is defined for all models estimated by maximum likelihood. For linear regression models with Gaussian errors,

$$\text{AIC} = \frac{1}{N\hat{\sigma}^2}(\text{SSR} + 2d\hat{\sigma}^2), \quad (7)$$

which is proportional to C_p .

More generally, with the sign reversed, the AIC equals $2 \log L(\boldsymbol{\theta} | \mathbf{y}) - 2d$. Notice that the penalty for AIC does not increase with N .

The **Bayesian information criterion**, or **BIC**, for linear regression models is

$$\text{BIC} = \frac{1}{N} (\text{SSR} + \log(N)d\hat{\sigma}^2). \quad (8)$$

More generally, with the sign reversed, $\text{BIC} = 2 \log L(\boldsymbol{\theta} \mid \mathbf{y}) - d \log N$.

When N is large, BIC places a much heavier penalty on d than does AIC, because $\log(N) \gg 2$. Therefore, the chosen model is very likely to have lower complexity.

If the correct model is among the ones we estimate, BIC chooses it with probability 1 as $N \rightarrow \infty$. The other two criteria sometimes choose excessively complex models.

3.3. Ridge Regression

As we have seen, ridge regression minimizes the objective function

$$\sum_{i=1}^N \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2, \quad (9)$$

where λ is a complexity parameter that controls the amount of shrinkage. This is an example of ℓ_2 -regularization.

Minimizing (9) is equivalent to solving the constrained minimization problem

$$\min \sum_{i=1}^N \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \text{ subject to } \sum_{j=1}^p \beta_j^2 \leq t. \quad (10)$$

In (9), λ plays the role of a Lagrange multiplier. For any value of t , there is an associated value of λ , which might be 0 if t is sufficiently large.

Notice that the intercept is not included in the penalty term. The simplest way to accomplish this is to recenter all variables before running the regression.

It would seem very strange to impose the same penalty on each of the β_j^2 if the associated regressors were not scaled similarly. Therefore, it is usual to rescale all the regressors to have variance 1. They already have mean 0 from the recentering.

We have seen that

$$\hat{\boldsymbol{\beta}}_{\text{ridge}} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^\top \mathbf{y}, \quad (11)$$

where every column of \mathbf{X} now has sample mean 0 and sample variance 1.

A number that is often more interesting than λ is the **effective degrees of freedom** of the ridge regression fit. If the d_j are the **singular values** of \mathbf{X} (see below), then

$$\text{df}(\lambda) \equiv \text{Tr}(\mathbf{X}(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^\top) = \sum_{j=1}^p \frac{d_j^2}{d_j^2 + \lambda}. \quad (12)$$

As $\lambda \rightarrow 0$, $\text{df}(\lambda) \rightarrow p$.

Thus the effective degrees of freedom for the ridge regression is never greater than the actual degrees of freedom for the corresponding OLS regression.

The ridge regression estimator (11) can be obtained as the mode (and also the mean) of a Bayesian posterior. Suppose that

$$y_i \sim \text{N}(\beta_0 + \mathbf{x}_i^\top \boldsymbol{\beta}, \sigma^2), \quad (13)$$

and the β_j are believed to follow independent $\text{N}(0, \tau^2)$ distributions. Then expression (9) is minus the log of the posterior, with $\lambda = \sigma^2/\tau^2$.

Thus minimizing (9) yields the posterior mode, which is $\hat{\boldsymbol{\beta}}_{\text{ridge}}$.

Evidently, as τ increases, λ becomes smaller. The more diffuse our prior about the β_j , the less it affects the posterior mode.

Similarly, as σ^2 increases, so that the data become noisier, the more our prior influences the posterior mode.

A simple way to compute ridge regression for any given λ is to form the augmented dataset

$$\mathbf{X}_a = \begin{bmatrix} \mathbf{X} \\ \sqrt{\lambda} \mathbf{I}_p \end{bmatrix} \quad \text{and} \quad \mathbf{y}_a = \begin{bmatrix} \mathbf{y} \\ \mathbf{0} \end{bmatrix}, \quad (14)$$

which has $N + p$ observations. Then OLS estimation of the regression

$$\mathbf{y}_a = \mathbf{X}_a \boldsymbol{\beta} + \mathbf{u}_a \quad (15)$$

yields the ridge regression estimator (11), because

$$\mathbf{X}_a^\top \mathbf{X}_a = \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_p \quad \text{and} \quad \mathbf{X}_a^\top \mathbf{y}_a = \mathbf{X}^\top \mathbf{y}. \quad (16)$$

In this case, it is tempting to use Cholesky, at least when p is not too large, because it is easy to update $\mathbf{X}_a^\top \mathbf{X}_a$ and $\mathbf{X}_a^\top \mathbf{y}_a$ as we vary λ .

3.4. The Singular Value Decomposition

The **singular value decomposition** of \mathbf{X} is

$$\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^\top, \quad (17)$$

where \mathbf{U} and \mathbf{V} are $N \times p$ and $p \times p$ orthogonal matrices, with $\mathcal{S}(\mathbf{U}) = \mathcal{S}(\mathbf{X})$. Thus $\mathbf{U}^\top\mathbf{U} = \mathbf{V}^\top\mathbf{V} = \mathbf{I}_p$. The columns of \mathbf{V} are the **eigenvectors** of $\mathbf{X}^\top\mathbf{X}$.

Recall that λ is an **eigenvalue** of a matrix \mathbf{A} if there exists a nonzero vector $\boldsymbol{\xi}$, called an **eigenvector**, such that

$$\mathbf{A}\boldsymbol{\xi} = \lambda\boldsymbol{\xi}. \quad (18)$$

Thus the action of \mathbf{A} on $\boldsymbol{\xi}$ produces a vector with the same direction as $\boldsymbol{\xi}$, but a different length unless $\lambda = 1$.

In (17), \mathbf{D} is a $p \times p$ diagonal matrix, with diagonal elements d_j that have the property

$$d_1 \geq d_2 \geq \dots \geq d_p \geq 0. \quad (19)$$

These are the square roots of the eigenvalues of $\mathbf{X}^\top\mathbf{X}$.

The rank of the matrix \mathbf{D} , and also of the matrix $\mathbf{X}^\top \mathbf{X}$, is the number of d_j that are strictly positive.

Observe that

$$\begin{aligned}
 \mathbf{P}_\mathbf{X} \mathbf{y} &= \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \\
 &= \mathbf{U} \mathbf{D} \mathbf{V}^\top (\mathbf{V} \mathbf{D} \mathbf{U}^\top \mathbf{U} \mathbf{D} \mathbf{V}^\top)^{-1} \mathbf{V} \mathbf{D} \mathbf{U}^\top \mathbf{y} \\
 &= \mathbf{U} \mathbf{D} \mathbf{V}^\top (\mathbf{D} \mathbf{V}^\top)^{-1} (\mathbf{U}^\top \mathbf{U})^{-1} (\mathbf{V} \mathbf{D})^{-1} \mathbf{V} \mathbf{D} \mathbf{U}^\top \mathbf{y} \\
 &= \mathbf{U} \mathbf{U}^\top \mathbf{y}.
 \end{aligned} \tag{20}$$

We can now see that

$$\begin{aligned}
 \mathbf{X} \hat{\boldsymbol{\beta}}_{\text{ridge}} &= (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y} \\
 &= \mathbf{U} \mathbf{D} (\mathbf{D}^2 + \lambda \mathbf{I})^{-1} \mathbf{D} \mathbf{U}^\top \mathbf{y} \\
 &= \sum_{j=1}^p \mathbf{u}_j \frac{d_j^2}{d_j^2 + \lambda} \mathbf{u}_j^\top \mathbf{y},
 \end{aligned} \tag{21}$$

where \mathbf{u}_j is the j^{th} column of \mathbf{U} .

When $\lambda = 0$, the scalars in the last line of (21) collapse to 1, and we see that

$$\mathbf{P}_{\mathbf{X}}\mathbf{y} = \mathbf{X}\hat{\boldsymbol{\beta}}_{\text{OLS}} = \sum_{j=1}^p \mathbf{u}_j \mathbf{u}_j^\top \mathbf{y}, \quad (22)$$

which is another way of writing the last line of (20).

When $\lambda > 0$, the estimates are shrunk. There is more shrinkage applied to the coordinates of basis vectors with small values of d_j .

The **eigen decomposition** of $\mathbf{X}^\top \mathbf{X}$ is

$$\mathbf{X}^\top \mathbf{X} = \mathbf{V} \mathbf{D}^2 \mathbf{V}^\top. \quad (23)$$

The eigenvectors \mathbf{v}_j are the columns of \mathbf{V} .

The largest eigenvector \mathbf{v}_1 has the property that $\mathbf{z}_1 = \mathbf{X}\mathbf{v}_1$ has the largest sample variance among all normalized linear combinations of the columns of \mathbf{X} .

Since $\mathbf{z}_1 = \mathbf{u}_1 d_1$, this sample variance is

$$\text{Var}(\mathbf{X}\mathbf{v}_1) = \text{Var}(\mathbf{u}_1 d_1) = d_1^2/N. \quad (24)$$

The first **principal component** of \mathbf{X} is \mathbf{z}_1 . Subsequent principal components are orthogonal to \mathbf{z}_1 , and to each other. Their variance declines as j increases.

Thus small singular values d_j correspond to directions in $\mathcal{S}(\mathbf{X})$ that have small variance. Ridge regression shrinks these directions the most, because they provide relatively little information, which may easily be swamped by noise.

3.5. Principal Components Regression

Recall from (24) that the principal components of \mathbf{X} are $\mathbf{z}_m = \mathbf{X}\mathbf{v}_m$, where the \mathbf{v}_m are the eigenvectors, that is, the columns of \mathbf{V} . These are $p \times 1$.

If we regress \mathbf{y} on the first M principal components, we find that

$$\hat{\mathbf{y}}^{\text{pc}} = \bar{y}\mathbf{1} + \sum_{m=1}^M \hat{\theta}_m \mathbf{z}_m, \quad (25)$$

where, because the \mathbf{z}_m are orthogonal,

$$\hat{\theta}_m = \frac{\mathbf{z}_m^\top \mathbf{y}}{\mathbf{z}_m^\top \mathbf{z}_m}. \quad (26)$$

As with other methods, it is a good idea to standardize the \mathbf{x}_j so that they have mean 0 and variance 1.

The principal components are ordered in the same way as the eigenvalues:

$$||\mathbf{z}_1|| \geq ||\mathbf{z}_2|| \geq ||\mathbf{z}_3|| \geq \dots \geq ||\mathbf{z}_p||. \quad (27)$$

Thus \mathbf{z}_1 is the direction in which the columns of \mathbf{X} vary the most.

For principal components regression (PCR), we first regress \mathbf{y} on \mathbf{z}_1 , then on \mathbf{z}_1 and \mathbf{z}_2 , and so on.

Because the \mathbf{z}_m are orthogonal, adding another regressor, say \mathbf{z}_g , simply means regressing the current residual

$$\hat{\boldsymbol{\epsilon}}_{g-1} \equiv \mathbf{y} - \hat{\theta}_0 - \hat{\theta}_1 \mathbf{z}_1 - \hat{\theta}_2 \mathbf{z}_2 - \dots - \hat{\theta}_{g-1} \mathbf{z}_{g-1} \quad (28)$$

on \mathbf{z}_g . Then we set $\hat{\boldsymbol{\epsilon}}_g = \hat{\boldsymbol{\epsilon}}_{g-1} - \hat{\theta}_g \mathbf{z}_g$. We keep going in this way until $g = M$.

The raw fit will improve as M increases, but criteria like AIC may get better or worse. As usual, we can also use cross-validation to choose M .

Of course, one can recover the implied β coefficients from the $\hat{\theta}_m$:

$$\hat{\beta}^{\text{pc}} = \sum_{m=1}^M \hat{\theta}_m \mathbf{v}_m. \quad (29)$$

Especially when M is small, these are likely to be shrunk relative to the OLS estimates. When $M = p$, these will be equal to OLS estimates.

There is a close relationship between ridge regression and PCR. Both operate via the principal components of the input matrix.

In both cases, coefficients are shrunk but not set to zero.

For ridge, directions with small eigenvalues get shrunk more than ones with large eigenvalues; recall (21).

For PCR, directions with small eigenvalues are discarded, while ones with large eigenvalues are retained.

3.6. The Lasso

Lasso (Tibshirani, 1996) stands for **least absolute shrinkage and selection operator**.

The lasso can be obtained as the solution to

$$\min \frac{1}{2N} \sum_{i=1}^N \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \text{ subject to } \sum_{j=1}^p |\beta_j| \leq t. \quad (30)$$

This looks almost identical to (10). The main difference is that the constraint involves absolute values instead of squares. Thus the lasso involves ℓ_1 -regularization.

As with ridge regression, solving (30) is equivalent to minimizing

$$\frac{1}{2N} \sum_{i=1}^N \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p |\beta_j|, \quad (31)$$

where λ is a Lagrange multiplier; compare (9). Large values of λ are associated with small values of t .

Following SLS (but not ESL or ISLR), there is a factor of $1/(2N)$ in (30) and (31). This is not needed, but it makes the values of λ comparable across samples of different sizes.

If all variables are recentered to have mean zero, we can dispense with the β_0 parameter. We can recover it if necessary as

$$\hat{\beta}_0 = \bar{y} - \sum_{j=1}^p \hat{\beta}_j \bar{x}_j. \quad (32)$$

The obvious way to solve (30) is by quadratic programming, but there are better approaches, which we discuss below.

The key difference between ridge and lasso is that, for lasso, some of the coefficients will be 0 when t is sufficiently small (equivalently, when λ is sufficiently large).

If t is sufficiently large — greater than $t_0 = \sum_{j=1}^p |\hat{\beta}_j|$ — then $\hat{\beta}^{\text{lasso}} = \hat{\beta}^{\text{ols}}$. In this case, the constraints do not bind.

For $t < t_0$, the lasso coefficients are shrunk, and some of them may well be 0. Ultimately, as $t \rightarrow 0$, they will all be 0.

For the tuning parameter, we can use λ , t , or the **standardized tuning parameter**

$$s \equiv t / \sum_{j=1}^p |\hat{\beta}_j|. \quad (33)$$

Whichever one we use, it is typically chosen by cross-validation.

The above statement is true when lasso is used for prediction, but not when it is used for inference. Other methods are proposed in papers by Chernozhukov et al.

ESL and SLS recommend using K -fold cross-validation, where $K = 5$ or $K = 10$. We will discuss these methods in detail later.

For K -fold cross-validation, we divide the training sample into K folds (subsamples) of equal or roughly equal size.

Then we sequentially omit one fold at a time, applying the lasso to the other $K - 1$ folds. This gives us K sets of estimates.

For $k = 1, \dots, K$, the estimates using all folds except the k^{th} are then used to compute fitted values for observations in fold k .

We use these fitted values to compute the mean squared prediction error for all observations. This is then plotted against t , s , or λ to choose the optimal value.

As usual, when the bound is too tight, the bias is large, and when the bound is too loose, the variance is large.

In the former case, there will be few non-zero coefficients, and in the latter case there will be many.

See SLS-fig-2.3.pdf for cross-validated estimate of prediction error.

If \mathbf{X} is orthonormal, there is an explicit solution for lasso:

$$\hat{\beta}_j^{\text{lasso}} = \text{sign}(\hat{\beta}_j)(|\hat{\beta}_j| - \lambda)_+, \quad (34)$$

where $(z)_+ = z$ if $z > 0$ and $(z)_+ = 0$ if $z \leq 0$.

Thus $\hat{\beta}_j^{\text{lasso}}$ either has the same sign as $\hat{\beta}_j$ or equals 0.

Like ridge, lasso is defined even when $p + 1 > N$.

The equivalent solution for ridge in the orthonormal case is

$$\hat{\beta}_j^{\text{ridge}} = \hat{\beta}_j / (1 + \lambda). \quad (35)$$

So it always has the same sign as $\hat{\beta}_j$ and never equals 0.

Like ridge regression, the lasso has a Bayesian interpretation. The prior for β_j is a Laplace (double exponential) distribution:

$$p(\beta_j) = \frac{\lambda}{2} \exp(-\lambda|\beta_j|). \quad (36)$$

The lasso has been generalized in many ways. One important one is the **elastic net**, which is discussed below.

ESL-fig3.10.pdf plots the profiles of lasso coefficients against s — see (33) — for the prostate cancer data.

This may be compared with ESL-fig3.08.pdf, which plots the profiles of ridge coefficients against $\text{df}(\lambda)$ for the same data.

The lasso shrinkage causes the estimates of the non-zero coefficients to be biased towards zero and inconsistent.

This bias can be reduced in at least three ways.

- Use the lasso to identify the set of predictors with non-zero coefficients. Then run another regression on just that set of predictors. This is called **post-lasso** estimation. Its properties were studied in various Chernozhukov et al. papers.
- Use the lasso twice. This is called the **relaxed lasso**. In the second step, include only the variables with non-zero coefficients in the first step.

If cross-validation is used at each step, the penalty parameter λ in the first step is likely to be larger than in the second step, because there are more noise variables to eliminate in the first step.

Since bias increases with λ , there should be less bias after the second step.

- Use the **adaptive lasso** (Zou, 2006). It replaces the penalty in (31) with

$$\lambda \sum_{j=1}^p w_j |\beta_j|, \quad w_j = 1/|\tilde{\beta}_j|^\nu, \quad \nu > 0. \quad (37)$$

Here the $\tilde{\beta}_j$ are **pilot estimates**, perhaps OLS estimates if $p \ll N$ or univariate regression coefficients otherwise.

The penalty on small coefficients is larger than the penalty on large ones. Thus the former tend to get driven to zero, and the latter get shrunk less than ordinary lasso would do.

3.7. Properties of the Lasso

The lasso objective function is a special case of

$$\frac{1}{2N} \sum_{i=1}^N \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p |\beta_j|^q. \quad (38)$$

For lasso, $q = 1$, and for ridge, $q = 2$. When $q = 0$, we have best-subset regression.

The smallest value of q for which the constraint set is convex is $q = 1$. Convexity makes optimization much easier.

The maximum number of non-zero coefficients is $\min(N, p)$.

Various consistency results have been proved for the lasso. The true parameter vector needs to be sparse relative to $N/\log(p)$, so it must contain a lot of zeros. See SLS, Chapter 11.

The obvious way to obtain standard errors for lasso is to use the pairs bootstrap. Resample from (y_i, \mathbf{x}_i) pairs and apply the procedure to each bootstrap sample.

We can use the standard deviations of the bootstrap lasso estimates around their mean as standard errors.

We can also calculate the proportion of the time that each parameter estimate is non-zero.

Better inference methods have been developed in papers by Chernozhukov and others. These include **double selection** and **double machine learning**. However, they work only for coefficients on regressors that are *known* to be in the model.

3.8. Least Angle Regression

Least angle regression (Efron, Hastie, Johnstone, and Tibshirani, 2004) can be thought of as an alternative to stepwise regression, but it also provides an efficient way to compute the lasso.

Before doing anything else, standardize all predictors to have mean 0 and unit norm, so that $\|\mathbf{x}_j\| = 1$.

1. Let $\mathbf{r} = \mathbf{y} - \bar{\mathbf{y}}$.
2. Find the predictor \mathbf{x}_j most correlated with \mathbf{r} .
3. Move β_j from 0 towards $\mathbf{x}_j^\top \mathbf{r} = (\mathbf{x}_j^\top \mathbf{x}_j)^{-1} \mathbf{x}_j^\top \mathbf{r}$ until some other \mathbf{x}_k has as much correlation with the current residual as does \mathbf{x}_j .
4. Move both β_j and β_k together in the direction of the least squares coefficient of the current residual on $[\mathbf{x}_j, \mathbf{x}_k]$ until some other \mathbf{x}_ℓ has as much correlation with the current residual.
5. Continue in this way until all p predictors have been included in the regression. After $\min(N - 1, p)$ steps, we arrive at the least squares solution.

The direction of the active step after variable k is entered is

$$\boldsymbol{\delta}_k = (\mathbf{X}_{[k]}^\top \mathbf{X}_{[k]})^{-1} \mathbf{X}_{[k]} (\mathbf{y} - \mathbf{X}_{[k]} \boldsymbol{\beta}_{[k]}). \quad (39)$$

Here $\mathbf{X}_{[k]}$ denotes the columns of \mathbf{X} that are now included, and $\boldsymbol{\beta}_{[k]}$ denotes the coefficients, one of which is 0.

The coefficients then evolve as

$$\boldsymbol{\beta}_{[k]}(\alpha) = \boldsymbol{\beta}_{[k]} + \alpha \boldsymbol{\delta}_k, \quad (40)$$

and the fit vector evolves as

$$\hat{\mathbf{f}}_k(\alpha) = \hat{\mathbf{f}}_k + \mathbf{X}_{[k]} \boldsymbol{\delta}_k. \quad (41)$$

The exact step length and the identity of the next variable to add can be calculated analytically, but ESL does not give details.

If we start at $t = 0$, the lasso coefficient(s) look just like the LAR ones until one of the non-zero lasso coefficients hits zero again.

Therefore, a modified version of the LAR algorithm can be used to compute lasso for all values of t or λ .

Add an additional step between 4. and 5.

- 4a. If a non-zero coefficient hits 0, drop its variable from the active set and recompute the current least-squares directions.

Note that dropped variables may reappear later.

LAR always takes p steps to obtain the OLS estimates. The lasso version may take more than that, but it tends to be similar.

This procedure works even when $p \gg N$. It provides a whole sequence of lasso estimates for $0 \leq \lambda < \infty$.

Chapter 5 of SLS discusses optimization methods for lasso and other convex optimization problems in much greater detail. It suggests that, although LAR is theoretically attractive, it is not the fastest method.

3.9. Pathwise Coordinate Optimization

We can rewrite (31) as

$$\frac{1}{2N} \sum_{i=1}^N \left(y_i - \beta_0 - \sum_{k \neq j} \tilde{\beta}_k(\lambda) x_{ik} - \beta_j x_{ij} \right)^2 + \lambda \sum_{k \neq j} |\tilde{\beta}_k(\lambda)| + \lambda |\beta_j|, \quad (42)$$

where $\tilde{\beta}_k(\lambda)$ denotes the current estimate of β_k at penalty parameter λ .

Minimizing (42) is a one-dimensional problem. It is effectively a univariate lasso problem with response variable

$$\tilde{u}_i^{(j)} = y_i - \sum_{k \neq j} \tilde{\beta}_k(\lambda) x_{ik}. \quad (43)$$

The solution to this problem is

$$\tilde{\beta}_j(\lambda) = \text{sign}(t) (|t| - \lambda)_+, \quad t = \sum_{i=1}^N x_{ij} \tilde{u}_i^{(j)}; \quad (44)$$

recall (34).

Repeated iteration of (43) and (44) over all j yields the lasso estimate $\hat{\beta}(\lambda)$; see Friedman, Hastie, Hoefling, and Tibshirani (2007). This approach is an example of **coordinate descent**.

Pathwise coordinate optimization is extremely simple. But one disadvantage, relative to LAR, is that it just yields estimates for a single value of λ .

However, estimates for a large value of λ can be used as starting points for a smaller value, and so on. When λ is sufficiently large, $\hat{\beta}^{\text{lasso}} = \mathbf{0}$. We can start at the smallest value of λ for which that is true.

We still only get results for a grid of values, rather than all values as with LAR, but that is often sufficient.

Observe that the actual computation of the \tilde{u}_i^j in (43) and of $\tilde{\beta}_j(\lambda)$ in (44) are extremely simple. Moreover, in many cases, when $\tilde{\beta}_j(\lambda) = 0$, it just stays at 0.

SLS reports simulation results which suggest that coordinate descent is the fastest way to obtain lasso estimates, both when $N \gg p$ and when $p \gg N$.

3.10. Elastic-Net Penalization

The elastic-net penalty (Zou and Hastie, 2005) has the form

$$\sum_{j=1}^p (\alpha |\beta_j| + (1 - \alpha) \beta_j^2). \quad (45)$$

This reduces to the lasso penalty when $\alpha = 1$ and the ridge penalty when $\alpha = 0$.

As usual, expression (45) is either multiplied by a tuning parameter λ and added to the sum of squares, or it is required to be less than an upper limit t .

Since there are two tuning parameters, cross-validation is a lot more work than for lasso or ridge. To reduce effort, people often just choose α and condition on it.

Elastic-net tends to yield more non-zero coefficients than lasso, but they are shrunk more towards zero.

Unlike lasso, it can yield more than N non-zero coefficients when $p > N$.

We can use any program for lasso to compute the elastic net estimator by using the augmentation trick discussed above to compute ridge regression.

Conversely, if we have a program for elastic net (such as `elasticnet` in Stata 16), we can use it to obtain ridge estimates by setting $\alpha = 0$.

If we multiply (45) by λ and then separate the two penalties, the ridge one becomes

$$(1 - \alpha)\lambda \sum_{j=1}^p \beta_j^2. \quad (46)$$

As in (14), we can form the augmented dataset

$$\mathbf{X}_a = \begin{bmatrix} \mathbf{X} \\ \sqrt{(1-\alpha)\lambda} \mathbf{I}_p \end{bmatrix} \quad \text{and} \quad \mathbf{y}_a = \begin{bmatrix} \mathbf{y} \\ \mathbf{0} \end{bmatrix}. \quad (47)$$

Then we can obtain the elastic-net estimates by solving the lasso problem

$$\operatorname{argmin}_{\boldsymbol{\beta}} \left(\frac{1}{2} (\mathbf{y}_a - \mathbf{X}_a \boldsymbol{\beta})^\top (\mathbf{y}_a - \mathbf{X}_a \boldsymbol{\beta}) + \alpha \lambda \sum_{j=1}^p |\beta_j| \right). \quad (48)$$

This just depends on one tuning parameter, $\alpha\lambda$. However, since \mathbf{X}_a depends on $(1-\alpha)\lambda$, we ultimately have to choose two tuning parameters.