

# 3

## Linear Methods for Regression

### 3.1 Introduction

A linear regression model assumes that the regression function  $E(Y|X)$  is linear in the inputs  $X_1, \dots, X_p$ . Linear models were largely developed in the precomputer age of statistics, but even in today's computer era there are still good reasons to study and use them. They are simple and often provide an adequate and interpretable description of how the inputs affect the output. For prediction purposes they can sometimes outperform fancier nonlinear models, especially in situations with small numbers of training cases, low signal-to-noise ratio or sparse data. Finally, linear methods can be applied to transformations of the inputs and this considerably expands their scope. These generalizations are sometimes called basis-function methods, and are discussed in Chapter 5.

In this chapter we describe linear methods for regression, while in the next chapter we discuss linear methods for classification. On some topics we go into considerable detail, as it is our firm belief that an understanding of linear methods is essential for understanding nonlinear ones. In fact, many nonlinear techniques are direct generalizations of the linear methods discussed here.

## 3.2 Linear Regression Models and Least Squares

As introduced in Chapter 2, we have an input vector  $X^T = (X_1, X_2, \dots, X_p)$ , and want to predict a real-valued output  $Y$ . The linear regression model has the form

$$f(X) = \beta_0 + \sum_{j=1}^p X_j \beta_j. \quad (3.1)$$

The linear model either assumes that the regression function  $E(Y|X)$  is linear, or that the linear model is a reasonable approximation. Here the  $\beta_j$ 's are unknown parameters or coefficients, and the variables  $X_j$  can come from different sources:

- quantitative inputs;
- transformations of quantitative inputs, such as log, square-root or square;
- basis expansions, such as  $X_2 = X_1^2$ ,  $X_3 = X_1^3$ , leading to a polynomial representation;
- numeric or “dummy” coding of the levels of qualitative inputs. For example, if  $G$  is a five-level factor input, we might create  $X_j$ ,  $j = 1, \dots, 5$ , such that  $X_j = I(G = j)$ . Together this group of  $X_j$  represents the effect of  $G$  by a set of level-dependent constants, since in  $\sum_{j=1}^5 X_j \beta_j$ , one of the  $X_j$ 's is one, and the others are zero.
- interactions between variables, for example,  $X_3 = X_1 \cdot X_2$ .

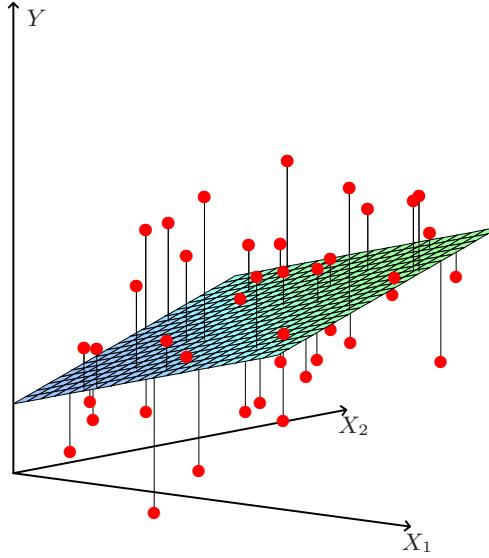
No matter the source of the  $X_j$ , the model is linear in the parameters.

Typically we have a set of training data  $(x_1, y_1), \dots, (x_N, y_N)$  from which to estimate the parameters  $\beta$ . Each  $x_i = (x_{i1}, x_{i2}, \dots, x_{ip})^T$  is a vector of feature measurements for the  $i$ th case. The most popular estimation method is *least squares*, in which we pick the coefficients  $\beta = (\beta_0, \beta_1, \dots, \beta_p)^T$  to minimize the residual sum of squares

$$\begin{aligned} \text{RSS}(\beta) &= \sum_{i=1}^N (y_i - f(x_i))^2 \\ &= \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2. \end{aligned} \quad (3.2)$$

From a statistical point of view, this criterion is reasonable if the training observations  $(x_i, y_i)$  represent independent random draws from their population. Even if the  $x_i$ 's were not drawn randomly, the criterion is still valid if the  $y_i$ 's are conditionally independent given the inputs  $x_i$ . Figure 3.1 illustrates the geometry of least-squares fitting in the  $\mathbb{R}^{p+1}$ -dimensional

cf econometrics  
 best linear approx to true f'n  
 identifiability constraints



**FIGURE 3.1.** Linear least squares fitting with  $X \in \mathbb{R}^2$ . We seek the linear function of  $X$  that minimizes the sum of squared residuals from  $Y$ .

space occupied by the pairs  $(X, Y)$ . Note that (3.2) makes no assumptions about the validity of model (3.1); it simply finds the best linear fit to the data. Least squares fitting is intuitively satisfying no matter how the data arise; the criterion measures the average lack of fit.

How do we minimize (3.2)? Denote by  $\mathbf{X}$  the  $N \times (p + 1)$  matrix with each row an input vector (with a 1 in the first position), and similarly let  $\mathbf{y}$  be the  $N$ -vector of outputs in the training set. Then we can write the residual sum-of-squares as

$$\text{RSS}(\beta) = (\mathbf{y} - \mathbf{X}\beta)^T(\mathbf{y} - \mathbf{X}\beta). \quad (3.3)$$

This is a quadratic function in the  $p + 1$  parameters. Differentiating with respect to  $\beta$  we obtain

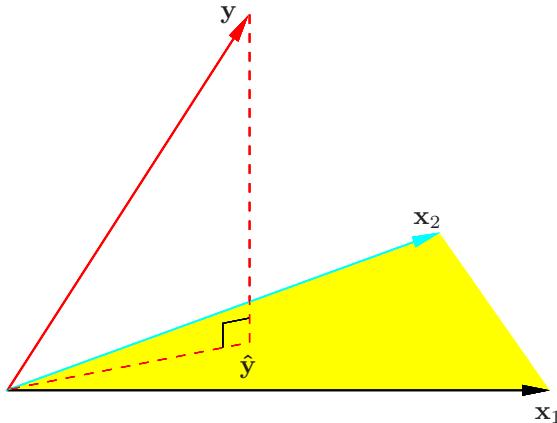
$$\begin{aligned} \frac{\partial \text{RSS}}{\partial \beta} &= -2\mathbf{X}^T(\mathbf{y} - \mathbf{X}\beta) \\ \frac{\partial^2 \text{RSS}}{\partial \beta \partial \beta^T} &= 2\mathbf{X}^T \mathbf{X}. \end{aligned} \quad (3.4)$$

Assuming (for the moment) that  $\mathbf{X}$  has full column rank, and hence  $\mathbf{X}^T \mathbf{X}$  is positive definite, we set the first derivative to zero

$$\mathbf{X}^T(\mathbf{y} - \mathbf{X}\beta) = 0 \quad (3.5)$$

to obtain the unique solution

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}. \quad (3.6)$$



**FIGURE 3.2.** The  $N$ -dimensional geometry of least squares regression with two predictors. The outcome vector  $\mathbf{y}$  is orthogonally projected onto the hyperplane spanned by the input vectors  $\mathbf{x}_1$  and  $\mathbf{x}_2$ . The projection  $\hat{\mathbf{y}}$  represents the vector of the least squares predictions

The predicted values at an input vector  $x_0$  are given by  $\hat{f}(x_0) = (1 : x_0)^T \hat{\beta}$ ; the fitted values at the training inputs are

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\beta} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}, \quad (3.7)$$

where  $\hat{y}_i = \hat{f}(x_i)$ . The matrix  $\mathbf{H} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$  appearing in equation (3.7) is sometimes called the “hat” matrix because it puts the hat on  $\mathbf{y}$ .

Figure 3.2 shows a different geometrical representation of the least squares estimate, this time in  $\mathbb{R}^N$ . We denote the column vectors of  $\mathbf{X}$  by  $\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_p$ , with  $\mathbf{x}_0 \equiv 1$ . For much of what follows, this first column is treated like any other. These vectors span a subspace of  $\mathbb{R}^N$ , also referred to as the column space of  $\mathbf{X}$ . We minimize  $\text{RSS}(\beta) = \|\mathbf{y} - \mathbf{X}\beta\|^2$  by choosing  $\hat{\beta}$  so that the residual vector  $\mathbf{y} - \hat{\mathbf{y}}$  is orthogonal to this subspace. This orthogonality is expressed in (3.5), and the resulting estimate  $\hat{\mathbf{y}}$  is hence the *orthogonal projection* of  $\mathbf{y}$  onto this subspace. The hat matrix  $\mathbf{H}$  computes the orthogonal projection, and hence it is also known as a projection matrix.

It might happen that the columns of  $\mathbf{X}$  are not linearly independent, so that  $\mathbf{X}$  is not of full rank. This would occur, for example, if two of the inputs were perfectly correlated, (e.g.,  $\mathbf{x}_2 = 3\mathbf{x}_1$ ). Then  $\mathbf{X}^T\mathbf{X}$  is singular and the least squares coefficients  $\hat{\beta}$  are not uniquely defined. However, the fitted values  $\hat{\mathbf{y}} = \mathbf{X}\hat{\beta}$  are still the projection of  $\mathbf{y}$  onto the column space of  $\mathbf{X}$ ; there is just more than one way to express that projection in terms of the column vectors of  $\mathbf{X}$ . The non-full-rank case occurs most often when one or more qualitative inputs are coded in a redundant fashion. There is usually a natural way to resolve the non-unique representation, by recoding and/or dropping redundant columns in  $\mathbf{X}$ . Most regression software packages detect these redundancies and automatically implement

L or pivoting (cf R)

→ also consider  
near singularity - when is it  
a problem?

(importance of collinearity is overstated)

proj  
matrix

← previous  
example

some strategy for removing them. Rank deficiencies can also occur in signal and image analysis, where the number of inputs  $p$  can exceed the number of training cases  $N$ . In this case, the features are typically reduced by filtering or else the fitting is controlled by regularization (Section 5.2.3 and Chapter 18).

Up to now we have made minimal assumptions about the true distribution of the data. In order to pin down the sampling properties of  $\hat{\beta}$ , we now assume that the observations  $y_i$  are uncorrelated and have constant variance  $\sigma^2$ , and that the  $x_i$  are fixed (non random). The variance–covariance matrix of the least squares parameter estimates is easily derived from (3.6) and is given by

$$\text{Var}(\hat{\beta}) = (\mathbf{X}^T \mathbf{X})^{-1} \sigma^2. \quad (3.8)$$

Typically one estimates the variance  $\sigma^2$  by

$$\hat{\sigma}^2 = \frac{1}{N - p - 1} \sum_{i=1}^N (y_i - \hat{y}_i)^2.$$

The  $N - p - 1$  rather than  $N$  in the denominator makes  $\hat{\sigma}^2$  an unbiased estimate of  $\sigma^2$ :  $E(\hat{\sigma}^2) = \sigma^2$ .

To draw inferences about the parameters and the model, additional assumptions are needed. We now assume that (3.1) is the correct model for the mean; that is, the conditional expectation of  $Y$  is linear in  $X_1, \dots, X_p$ . We also assume that the deviations of  $Y$  around its expectation are additive and Gaussian. Hence

$$\begin{aligned} Y &= E(Y|X_1, \dots, X_p) + \varepsilon \\ &= \beta_0 + \sum_{j=1}^p X_j \beta_j + \varepsilon, \end{aligned} \quad (3.9)$$

where the error  $\varepsilon$  is a Gaussian random variable with expectation zero and variance  $\sigma^2$ , written  $\varepsilon \sim N(0, \sigma^2)$ .

Under (3.9), it is easy to show that

$$\hat{\beta} \sim N(\beta, (\mathbf{X}^T \mathbf{X})^{-1} \sigma^2). \quad (3.10)$$

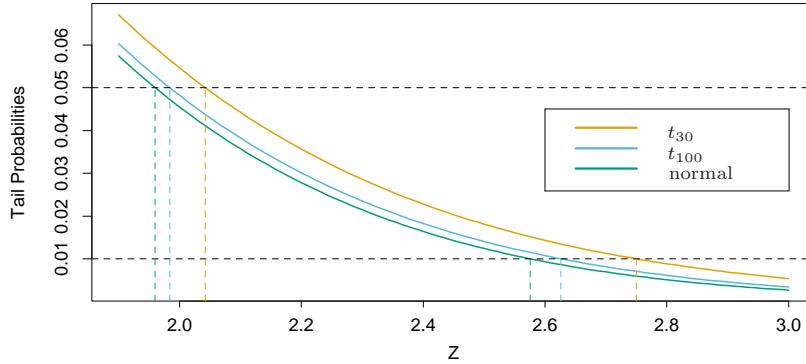
This is a multivariate normal distribution with mean vector and variance–covariance matrix as shown. Also

$$(N - p - 1) \hat{\sigma}^2 \sim \sigma^2 \chi_{N-p-1}^2, \quad (3.11)$$

a chi-squared distribution with  $N - p - 1$  degrees of freedom. In addition  $\hat{\beta}$  and  $\hat{\sigma}^2$  are statistically independent. We use these distributional properties to form tests of hypothesis and confidence intervals for the parameters  $\beta_j$ .

CONDITIONALLY

note Bessel correction



**FIGURE 3.3.** The tail probabilities  $\Pr(|Z| > z)$  for three distributions,  $t_{30}$ ,  $t_{100}$  and standard normal. Shown are the appropriate quantiles for testing significance at the  $p = 0.05$  and  $0.01$  levels. The difference between  $t$  and the standard normal becomes negligible for  $N$  bigger than about 100.

To test the hypothesis that a particular coefficient  $\beta_j = 0$ , we form the standardized coefficient or *Z-score*

$$z_j = \frac{\hat{\beta}_j}{\hat{\sigma} \sqrt{v_j}}, \quad (3.12)$$

where  $v_j$  is the  $j$ th diagonal element of  $(\mathbf{X}^T \mathbf{X})^{-1}$ . Under the null hypothesis that  $\beta_j = 0$ ,  $z_j$  is distributed as  $t_{N-p-1}$  (a  $t$  distribution with  $N - p - 1$  degrees of freedom), and hence a large (absolute) value of  $z_j$  will lead to rejection of this null hypothesis. If  $\hat{\sigma}$  is replaced by a known value  $\sigma$ , then  $z_j$  would have a standard normal distribution. The difference between the tail quantiles of a  $t$ -distribution and a standard normal become negligible as the sample size increases, and so we typically use the normal quantiles (see Figure 3.3).

Often we need to test for the significance of groups of coefficients simultaneously. For example, to test if a categorical variable with  $k$  levels can be excluded from a model, we need to test whether the coefficients of the dummy variables used to represent the levels can all be set to zero. Here we use the  $F$  statistic,

$$F = \frac{(\text{RSS}_0 - \text{RSS}_1)/(p_1 - p_0)}{\text{RSS}_1/(N - p_1 - 1)}, \quad (3.13)$$

where  $\text{RSS}_1$  is the residual sum-of-squares for the least squares fit of the bigger model with  $p_1 + 1$  parameters, and  $\text{RSS}_0$  the same for the nested smaller model with  $p_0 + 1$  parameters, having  $p_1 - p_0$  parameters constrained to be

poll:  
 what is  
 the  
 analogue  
 of  $t \rightarrow Z$   
 as  $N \rightarrow \infty$ ?  
 (i.e.  $F \rightarrow ?$ )  
 (see next page)

zero. The  $F$  statistic measures the change in residual sum-of-squares per additional parameter in the bigger model, and it is normalized by an estimate of  $\sigma^2$ . Under the Gaussian assumptions, and the null hypothesis that the smaller model is correct, the  $F$  statistic will have a  $F_{p_1-p_0, N-p_1-1}$  distribution. It can be shown (Exercise 3.1) that the  $z_j$  in (3.12) are equivalent to the  $F$  statistic for dropping the single coefficient  $\beta_j$  from the model. For large  $N$ , the quantiles of  $F_{p_1-p_0, N-p_1-1}$  approach those of  $\chi^2_{p_1-p_0}/(p_1-p_0)$ .

Similarly, we can isolate  $\beta_j$  in (3.10) to obtain a  $1-2\alpha$  confidence interval for  $\beta_j$ :

$$(\hat{\beta}_j - z^{(1-\alpha)} v_j^{\frac{1}{2}} \hat{\sigma}, \quad \hat{\beta}_j + z^{(1-\alpha)} v_j^{\frac{1}{2}} \hat{\sigma}). \quad (3.14)$$

Here  $z^{(1-\alpha)}$  is the  $1-\alpha$  percentile of the normal distribution:

$$\begin{aligned} z^{(1-0.025)} &= 1.96, \\ z^{(1-0.05)} &= 1.645, \text{ etc.} \end{aligned}$$

Hence the standard practice of reporting  $\hat{\beta} \pm 2 \cdot \text{se}(\hat{\beta})$  amounts to an approximate 95% confidence interval. Even if the Gaussian error assumption does not hold, this interval will be approximately correct, with its coverage approaching  $1-2\alpha$  as the sample size  $N \rightarrow \infty$ .

In a similar fashion we can obtain an approximate confidence set for the entire parameter vector  $\beta$ , namely

$$C_\beta = \{\beta | (\hat{\beta} - \beta)^T \mathbf{X}^T \mathbf{X} (\hat{\beta} - \beta) \leq \hat{\sigma}^2 \chi_{p+1}^{2(1-\alpha)}\}, \quad (3.15)$$

where  $\chi_\ell^{2(1-\alpha)}$  is the  $1-\alpha$  percentile of the chi-squared distribution on  $\ell$  degrees of freedom: for example,  $\chi_5^{2(1-0.05)} = 11.1$ ,  $\chi_5^{2(1-0.1)} = 9.2$ . This confidence set for  $\beta$  generates a corresponding confidence set for the true function  $f(x) = x^T \beta$ , namely  $\{x^T \beta | \beta \in C_\beta\}$  (Exercise 3.2; see also Figure 5.4 in Section 5.2.2 for examples of confidence bands for functions).

### 3.2.1 Example: Prostate Cancer

The data for this example come from a study by Stamey et al. (1989). They examined the correlation between the level of prostate-specific antigen and a number of clinical measures in men who were about to receive a radical prostatectomy. The variables are log cancer volume (`lcavol`), log prostate weight (`lweight`), `age`, log of the amount of benign prostatic hyperplasia (`lbph`), seminal vesicle invasion (`svi`), log of capsular penetration (`lcp`), Gleason score (`gleason`), and percent of Gleason scores 4 or 5 (`pgg45`). The correlation matrix of the predictors given in Table 3.1 shows many strong correlations. Figure 1.1 (page 3) of Chapter 1 is a scatterplot matrix showing every pairwise plot between the variables. We see that `svi` is a binary variable, and `gleason` is an ordered categorical variable. We see, for

why log-transform everything?

when is this useful?

complot/  
pairsplot/  
Gally ::  
ggpairs/  
car

**TABLE 3.1.** Correlations of predictors in the prostate cancer data.

	lcavol	lweight	age	lbph	svi	lcp	gleason
lweight	0.300						
age	0.286	0.317					
lbph	0.063	0.437	0.287				
svi	0.593	0.181	0.129	-0.139			
lcp	0.692	0.157	0.173	-0.089	0.671		
gleason	0.426	0.024	0.366	0.033	0.307	0.476	
pgg45	0.483	0.074	0.276	-0.030	0.481	0.663	0.757

**TABLE 3.2.** Linear model fit to the prostate cancer data. The Z score is the coefficient divided by its standard error (3.12). Roughly a Z score larger than two in absolute value is significantly nonzero at the  $p = 0.05$  level.

Term	Coefficient	Std. Error	Z Score
Intercept	2.46	0.09	27.60
lcavol	0.68	0.13	5.37
lweight	0.26	0.10	2.75
age	-0.14	0.10	-1.40
lbph	0.21	0.10	2.06
svi	0.31	0.12	2.47
lcp	-0.29	0.15	-1.87
gleason	-0.02	0.15	-0.15
pgg45	0.27	0.15	1.74

example, that both `lcavol` and `lcp` show a strong relationship with the response `lpsa`, and with each other. We need to fit the effects jointly to untangle the relationships between the predictors and the response.

We fit a linear model to the log of prostate-specific antigen, `lpsa`, after first standardizing the predictors to have unit variance. We randomly split the dataset into a training set of size 67 and a test set of size 30. We applied least squares estimation to the training set, producing the estimates, standard errors and Z-scores shown in Table 3.2. The Z-scores are defined in (3.12), and measure the effect of dropping that variable from the model. A Z-score greater than 2 in absolute value is approximately significant at the 5% level. (For our example, we have nine parameters, and the 0.025 tail quantiles of the  $t_{67-9}$  distribution are  $\pm 2.002!$ ) The predictor `lcavol` shows the strongest effect, with `lweight` and `svi` also strong. Notice that `lcp` is not significant, once `lcavol` is in the model (when used in a model without `lcavol`, `lcp` is strongly significant). We can also test for the exclusion of a number of terms at once, using the  $F$ -statistic (3.13). For example, we consider dropping all the non-significant terms in Table 3.2, namely `age`,

*I don't bother*

|| why?  
Schillzeth  
2010

*will this help?*

lcp, gleason, and pgg45. We get

$$F = \frac{(32.81 - 29.43)/(9 - 5)}{29.43/(67 - 9)} = 1.67, \quad (3.16)$$

which has a  $p$ -value of  $0.17$  ( $\Pr(F_{4,58} > 1.67) = 0.17$ ), and hence is not significant.

The mean prediction error on the test data is 0.521. In contrast, prediction using the mean training value of lpsa has a test error of 1.057, which is called the “base error rate.” Hence the linear model reduces the base error rate by about 50%. We will return to this example later to compare various selection and shrinkage methods.



### 3.2.2 The Gauss–Markov Theorem

One of the most famous results in statistics asserts that the least squares estimates of the parameters  $\beta$  have the smallest variance among all linear unbiased estimates. We will make this precise here, and also make clear that the restriction to unbiased estimates is not necessarily a wise one. This observation will lead us to consider biased estimates such as ridge regression later in the chapter. We focus on estimation of any linear combination of the parameters  $\theta = a^T \beta$ ; for example, predictions  $f(x_0) = x_0^T \beta$  are of this form. The least squares estimate of  $a^T \beta$  is

$$\hat{\theta} = a^T \hat{\beta} = a^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}. \quad (3.17)$$

Considering  $\mathbf{X}$  to be fixed, this is a linear function  $\mathbf{c}_0^T \mathbf{y}$  of the response vector  $\mathbf{y}$ . If we assume that the linear model is correct,  $a^T \hat{\beta}$  is unbiased since

$$\begin{aligned} E(a^T \hat{\beta}) &= E(a^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}) \\ &= a^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \beta \\ &= a^T \beta. \end{aligned} \quad (3.18)$$

The Gauss–Markov theorem states that if we have any other linear estimator  $\tilde{\theta} = \mathbf{c}^T \mathbf{y}$  that is unbiased for  $a^T \beta$ , that is,  $E(\mathbf{c}^T \mathbf{y}) = a^T \beta$ , then

$$\text{Var}(a^T \hat{\beta}) \leq \text{Var}(\mathbf{c}^T \mathbf{y}). \quad (3.19)$$

The proof (Exercise 3.3) uses the triangle inequality. For simplicity we have stated the result in terms of estimation of a single parameter  $a^T \beta$ , but with a few more definitions one can state it in terms of the entire parameter vector  $\beta$  (Exercise 3.3).

Consider the mean squared error of an estimator  $\tilde{\theta}$  in estimating  $\theta$ :

$$\begin{aligned} \underline{\text{MSE}}(\tilde{\theta}) &= E(\tilde{\theta} - \theta)^2 \\ &= \text{Var}(\tilde{\theta}) + [E(\tilde{\theta}) - \theta]^2. \end{aligned} \quad (3.20)$$

The first term is the variance, while the second term is the squared bias. The Gauss-Markov theorem implies that the least squares estimator has the smallest mean squared error of all linear estimators with no bias. However, there may well exist a biased estimator with smaller mean squared error. Such an estimator would trade a little bias for a larger reduction in variance. Biased estimates are commonly used. Any method that shrinks or sets to zero some of the least squares coefficients may result in a biased estimate. We discuss many examples, including variable subset selection and ridge regression, later in this chapter. From a more pragmatic point of view, most models are distortions of the truth, and hence are biased; picking the right model amounts to creating the right balance between bias and variance. We go into these issues in more detail in Chapter 7.

Mean squared error is intimately related to prediction accuracy, as discussed in Chapter 2. Consider the prediction of the new response at input  $x_0$ ,

$$Y_0 = f(x_0) + \varepsilon_0. \quad (3.21)$$

Then the expected prediction error of an estimate  $\tilde{f}(x_0) = x_0^T \tilde{\beta}$  is

$$\begin{aligned} E(Y_0 - \tilde{f}(x_0))^2 &= \sigma^2 + E(x_0^T \tilde{\beta} - f(x_0))^2 \\ &= \sigma^2 + \text{MSE}(\tilde{f}(x_0)). \end{aligned} \quad (3.22)$$

Therefore, expected prediction error and mean squared error differ only by the constant  $\sigma^2$ , representing the variance of the new observation  $y_0$ .

cf  
prediction  
vs  
confidence  
intervals

### 3.2.3 Multiple Regression from Simple Univariate Regression

The linear model (3.1) with  $p > 1$  inputs is called the *multiple linear regression model*. The least squares estimates (3.6) for this model are best understood in terms of the estimates for the *univariate* ( $p = 1$ ) linear model, as we indicate in this section.

Suppose first that we have a univariate model with no intercept, that is,

$$Y = X\beta + \varepsilon. \quad (3.23)$$

The least squares estimate and residuals are

$$\begin{aligned} \hat{\beta} &= \frac{\sum_1^N x_i y_i}{\sum_1^N x_i^2}, \\ r_i &= y_i - x_i \hat{\beta}. \end{aligned} \quad (3.24)$$

In convenient vector notation, we let  $\mathbf{y} = (y_1, \dots, y_N)^T$ ,  $\mathbf{x} = (x_1, \dots, x_N)^T$  and define

$$\begin{aligned} \langle \mathbf{x}, \mathbf{y} \rangle &= \sum_{i=1}^N x_i y_i, \\ &= \mathbf{x}^T \mathbf{y}, \end{aligned} \quad (3.25)$$

3 p doesn't  
count  
intercept?

the *inner product* between  $\mathbf{x}$  and  $\mathbf{y}$ <sup>1</sup>. Then we can write

$$\begin{aligned}\hat{\beta} &= \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle}, \\ \mathbf{r} &= \mathbf{y} - \mathbf{x}\hat{\beta}.\end{aligned}\tag{3.26}$$

As we will see, this simple univariate regression provides the building block for multiple linear regression. Suppose next that the inputs  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p$  (the columns of the data matrix  $\mathbf{X}$ ) are orthogonal; that is  $\langle \mathbf{x}_j, \mathbf{x}_k \rangle = 0$  for all  $j \neq k$ . Then it is easy to check that the multiple least squares estimates  $\hat{\beta}_j$  are equal to  $\langle \mathbf{x}_j, \mathbf{y} \rangle / \langle \mathbf{x}_j, \mathbf{x}_j \rangle$ —the univariate estimates. In other words, when the inputs are orthogonal, they have no effect on each other's parameter estimates in the model.

Orthogonal inputs occur most often with balanced, designed experiments (where orthogonality is enforced), but almost never with observational data. Hence we will have to orthogonalize them in order to carry this idea further. Suppose next that we have an intercept and a single input  $\mathbf{x}$ . Then the least squares coefficient of  $\mathbf{x}$  has the form

$$\hat{\beta}_1 = \frac{\langle \mathbf{x} - \bar{x}\mathbf{1}, \mathbf{y} \rangle}{\langle \mathbf{x} - \bar{x}\mathbf{1}, \mathbf{x} - \bar{x}\mathbf{1} \rangle},\tag{3.27}$$

where  $\bar{x} = \sum_i x_i / N$ , and  $\mathbf{1} = \mathbf{x}_0$ , the vector of  $N$  ones. We can view the estimate (3.27) as the result of two applications of the simple regression (3.26). The steps are:

1. regress  $\mathbf{x}$  on  $\mathbf{1}$  to produce the residual  $\mathbf{z} = \mathbf{x} - \bar{x}\mathbf{1}$ ;
2. regress  $\mathbf{y}$  on the residual  $\mathbf{z}$  to give the coefficient  $\hat{\beta}_1$ .

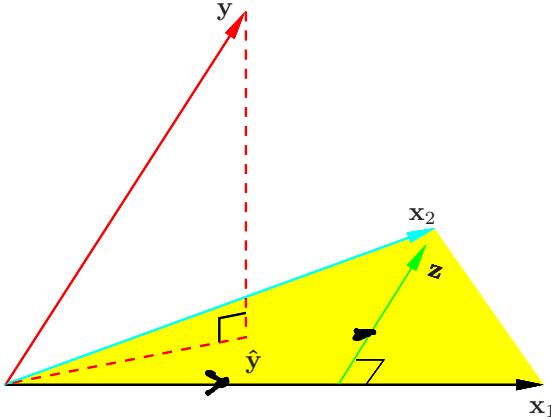
In this procedure, “regress  $\mathbf{b}$  on  $\mathbf{a}$ ” means a simple univariate regression of  $\mathbf{b}$  on  $\mathbf{a}$  with no intercept, producing coefficient  $\hat{\gamma} = \langle \mathbf{a}, \mathbf{b} \rangle / \langle \mathbf{a}, \mathbf{a} \rangle$  and residual vector  $\mathbf{b} - \hat{\gamma}\mathbf{a}$ . We say that  $\mathbf{b}$  is adjusted for  $\mathbf{a}$ , or is “orthogonalized” with respect to  $\mathbf{a}$ .

Step 1 orthogonalizes  $\mathbf{x}$  with respect to  $\mathbf{x}_0 = \mathbf{1}$ . Step 2 is just a simple univariate regression, using the orthogonal predictors  $\mathbf{1}$  and  $\mathbf{z}$ . Figure 3.4 shows this process for two general inputs  $\mathbf{x}_1$  and  $\mathbf{x}_2$ . The orthogonalization does not change the subspace spanned by  $\mathbf{x}_1$  and  $\mathbf{x}_2$ , it simply produces an orthogonal basis for representing it.

This recipe generalizes to the case of  $p$  inputs, as shown in Algorithm 3.1. Note that the inputs  $\mathbf{z}_0, \dots, \mathbf{z}_{j-1}$  in step 2 are orthogonal, hence the simple regression coefficients computed there are in fact also the multiple regression coefficients.

---

<sup>1</sup>The inner-product notation is suggestive of generalizations of linear regression to different metric spaces, as well as to probability spaces.



**FIGURE 3.4.** Least squares regression by orthogonalization of the inputs. The vector  $\mathbf{x}_2$  is regressed on the vector  $\mathbf{x}_1$ , leaving the residual vector  $\mathbf{z}$ . The regression of  $\mathbf{y}$  on  $\mathbf{z}$  gives the multiple regression coefficient of  $\mathbf{x}_2$ . Adding together the projections of  $\mathbf{y}$  on each of  $\mathbf{x}_1$  and  $\mathbf{z}$  gives the least squares fit  $\hat{\mathbf{y}}$ .

---

**Algorithm 3.1** Regression by Successive Orthogonalization.

---

1. Initialize  $\mathbf{z}_0 = \mathbf{x}_0 = \mathbf{1}$ .

2. For  $j = 1, 2, \dots, p$

Regress  $\mathbf{x}_j$  on  $\mathbf{z}_0, \mathbf{z}_1, \dots, \mathbf{z}_{j-1}$  to produce coefficients  $\hat{\gamma}_{\ell j} = \langle \mathbf{z}_\ell, \mathbf{x}_j \rangle / \langle \mathbf{z}_\ell, \mathbf{z}_\ell \rangle$ ,  $\ell = 0, \dots, j-1$  and residual vector  $\mathbf{z}_j = \mathbf{x}_j - \sum_{k=0}^{j-1} \hat{\gamma}_{kj} \mathbf{z}_k$ .

3. Regress  $\mathbf{y}$  on the residual  $\mathbf{z}_p$  to give the estimate  $\hat{\beta}_p$ .

---

The result of this algorithm is

$$\hat{\beta}_p = \frac{\langle \mathbf{z}_p, \mathbf{y} \rangle}{\langle \mathbf{z}_p, \mathbf{z}_p \rangle}. \quad (3.28)$$

Re-arranging the residual in step 2, we can see that each of the  $\mathbf{x}_j$  is a linear combination of the  $\mathbf{z}_k$ ,  $k \leq j$ . Since the  $\mathbf{z}_j$  are all orthogonal, they form a basis for the column space of  $\mathbf{X}$ , and hence the least squares projection onto this subspace is  $\hat{\mathbf{y}}$ . Since  $\mathbf{z}_p$  alone involves  $\mathbf{x}_p$  (with coefficient 1), we see that the coefficient (3.28) is indeed the multiple regression coefficient of  $\mathbf{y}$  on  $\mathbf{x}_p$ . This key result exposes the effect of correlated inputs in multiple regression. Note also that by rearranging the  $\mathbf{x}_j$ , any one of them could be in the last position, and a similar result holds. Hence stated more generally, we have shown that the  $j$ th multiple regression coefficient is the univariate regression coefficient of  $\mathbf{y}$  on  $\mathbf{x}_{j+1}, \dots, \mathbf{x}_p$ , the residual after regressing  $\mathbf{x}_j$  on  $\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_{j-1}, \mathbf{x}_{j+1}, \dots, \mathbf{x}_p$ :

To prove that it's independent of order

The multiple regression coefficient  $\hat{\beta}_j$  represents the additional contribution of  $\mathbf{x}_j$  on  $\mathbf{y}$ , after  $\mathbf{x}_j$  has been adjusted for  $\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_{j-1}, \mathbf{x}_{j+1}, \dots, \mathbf{x}_p$ .

If  $\mathbf{x}_p$  is highly correlated with some of the other  $\mathbf{x}_k$ 's, the residual vector  $\mathbf{z}_p$  will be close to zero, and from (3.28) the coefficient  $\hat{\beta}_p$  will be very unstable. This will be true for all the variables in the correlated set. In such situations, we might have all the Z-scores (as in Table 3.2) be small—any one of the set can be deleted—yet we cannot delete them all. From (3.28) we also obtain an alternate formula for the variance estimates (3.8),

$$\text{Var}(\hat{\beta}_p) = \frac{\sigma^2}{\langle \mathbf{z}_p, \mathbf{z}_p \rangle} = \frac{\sigma^2}{\|\mathbf{z}_p\|^2}. \quad (3.29)$$

} inferential  
stuff

In other words, the precision with which we can estimate  $\hat{\beta}_p$  depends on the length of the residual vector  $\mathbf{z}_p$ ; this represents how much of  $\mathbf{x}_p$  is unexplained by the other  $\mathbf{x}_k$ 's.

Algorithm 3.1 is known as the Gram-Schmidt procedure for multiple regression, and is also a useful numerical strategy for computing the estimates. We can obtain from it not just  $\hat{\beta}_p$ , but also the entire multiple least squares fit, as shown in Exercise 3.4.

We can represent step 2 of Algorithm 3.1 in matrix form:

$$\mathbf{X} = \mathbf{Z}\Gamma, \quad (3.30)$$

where  $\mathbf{Z}$  has as columns the  $\mathbf{z}_j$  (in order), and  $\Gamma$  is the upper triangular matrix with entries  $\hat{\gamma}_{kj}$ . Introducing the diagonal matrix  $\mathbf{D}$  with  $j$ th diagonal entry  $D_{jj} = \|\mathbf{z}_j\|$ , we get

$$\begin{aligned} \mathbf{X} &= \underbrace{\mathbf{Z}\mathbf{D}^{-1}\mathbf{D}\Gamma}_{\mathcal{Q}} \\ &= \mathbf{QR}, \end{aligned} \quad (3.31)$$

the so-called  $QR$  decomposition of  $\mathbf{X}$ . Here  $\mathbf{Q}$  is an  $N \times (p+1)$  orthogonal matrix,  $\mathbf{Q}^T \mathbf{Q} = \mathbf{I}$ , and  $\mathbf{R}$  is a  $(p+1) \times (p+1)$  upper triangular matrix.

? orthonormal?

The  $QR$  decomposition represents a convenient orthogonal basis for the column space of  $\mathbf{X}$ . It is easy to see, for example, that the least squares solution is given by

$$\hat{\beta} = \mathbf{R}^{-1} \mathbf{Q}^T \mathbf{y}, \quad (3.32)$$

$$\hat{\mathbf{y}} = \mathbf{QQ}^T \mathbf{y}. \quad (3.33)$$

Equation (3.32) is easy to solve because  $\mathbf{R}$  is upper triangular (Exercise 3.4).

$$\begin{aligned} (\mathbf{X}^T \mathbf{X}) \beta &= \mathbf{X}^T \mathbf{y} \\ (\mathbf{Q}\mathbf{R})^T (\mathbf{Q}\mathbf{R}) \beta &= \mathbf{R}^T \mathbf{Q}^T \mathbf{y} \\ \mathbf{R}^T (\mathbf{Q}^T \mathbf{Q}) \mathbf{R} \beta &= \mathbf{R}^T \mathbf{Q}^T \mathbf{y} \\ \mathbf{R}^T \mathbf{R} \beta &= \mathbf{R}^T \mathbf{Q}^T \mathbf{y} \end{aligned}$$

### 3.2.4 Multiple Outputs

Suppose we have multiple outputs  $Y_1, Y_2, \dots, Y_K$  that we wish to predict from our inputs  $X_0, X_1, X_2, \dots, X_p$ . We assume a linear model for each output

$$Y_k = \beta_{0k} + \sum_{j=1}^p X_j \beta_{jk} + \varepsilon_k \quad (3.34)$$

$$= f_k(X) + \varepsilon_k. \quad (3.35)$$

With  $N$  training cases we can write the model in matrix notation

$$\mathbf{Y} = \mathbf{XB} + \mathbf{E}. \quad (3.36)$$

Here  $\mathbf{Y}$  is the  $N \times K$  response matrix, with  $ik$  entry  $y_{ik}$ ,  $\mathbf{X}$  is the  $N \times (p+1)$  input matrix,  $\mathbf{B}$  is the  $(p+1) \times K$  matrix of parameters and  $\mathbf{E}$  is the  $N \times K$  matrix of errors. A straightforward generalization of the univariate loss function (3.2) is

$$\text{RSS}(\mathbf{B}) = \sum_{k=1}^K \sum_{i=1}^N (y_{ik} - f_k(x_i))^2 \quad (3.37)$$

$$= \text{tr}[(\mathbf{Y} - \mathbf{XB})^T (\mathbf{Y} - \mathbf{XB})]. \quad (3.38)$$

The least squares estimates have exactly the same form as before

$$\hat{\mathbf{B}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}. \quad (3.39)$$

Hence the coefficients for the  $k$ th outcome are just the least squares estimates in the regression of  $\mathbf{y}_k$  on  $\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_p$ . Multiple outputs do not affect one another's least squares estimates.

If the errors  $\varepsilon = (\varepsilon_1, \dots, \varepsilon_K)$  in (3.34) are correlated, then it might seem appropriate to modify (3.37) in favor of a multivariate version. Specifically, suppose  $\text{Cov}(\varepsilon) = \Sigma$ , then the multivariate weighted criterion

$$\text{RSS}(\mathbf{B}; \Sigma) = \sum_{i=1}^N (y_i - f(x_i))^T \Sigma^{-1} (y_i - f(x_i)) \quad (3.40)$$

arises naturally from multivariate Gaussian theory. Here  $f(x)$  is the vector function  $(f_1(x), \dots, f_K(x))$ , and  $y_i$  the vector of  $K$  responses for observation  $i$ . However, it can be shown that again the solution is given by (3.39);  $K$  separate regressions that ignore the correlations (Exercise 3.11). If the  $\Sigma_i$  vary among observations, then this is no longer the case, and the solution for  $\mathbf{B}$  no longer decouples.

In Section 3.7 we pursue the multiple outcome problem, and situations where it does pay to combine the regressions.

### 3.3 Subset Selection

There are two reasons why we are often not satisfied with the least squares estimates (3.6).

- The first is *prediction accuracy*: the least squares estimates often have low bias but large variance. Prediction accuracy can sometimes be improved by shrinking or setting some coefficients to zero. By doing so we sacrifice a little bit of bias to reduce the variance of the predicted values, and hence may improve the overall prediction accuracy.
- The second reason is *interpretation*. With a large number of predictors, we often would like to determine a smaller subset that exhibit the strongest effects. In order to get the “big picture,” we are willing to sacrifice some of the small details.

In this section we describe a number of approaches to variable subset selection with linear regression. In later sections we discuss shrinkage and hybrid approaches for controlling variance, as well as other dimension-reduction strategies. These all fall under the general heading *model selection*. Model selection is not restricted to linear models; Chapter 7 covers this topic in some detail.

With subset selection we retain only a subset of the variables, and eliminate the rest from the model. Least squares regression is used to estimate the coefficients of the inputs that are retained. There are a number of different strategies for choosing the subset.

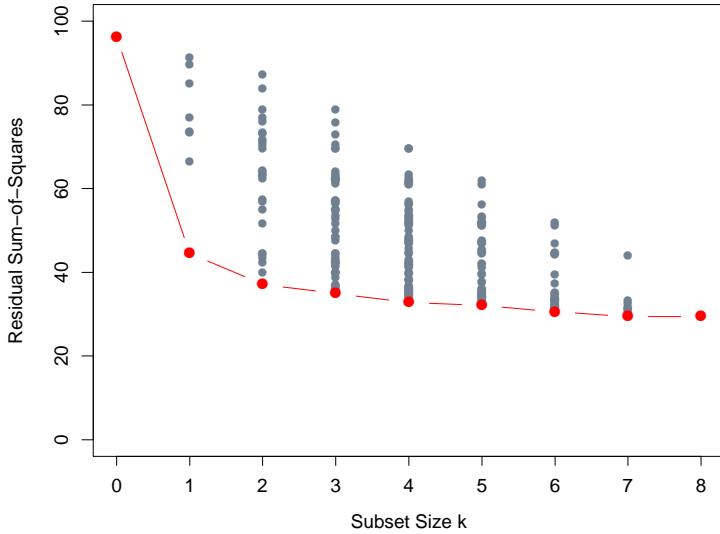
#### 3.3.1 Best-Subset Selection

Best subset regression finds for each  $k \in \{0, 1, 2, \dots, p\}$  the subset of size  $k$  that gives smallest residual sum of squares (3.2). An efficient algorithm—the *leaps and bounds* procedure (Furnival and Wilson, 1974)—makes this feasible for  $p$  as large as 30 or 40. Figure 3.5 shows all the subset models for the prostate cancer example. The lower boundary represents the models that are eligible for selection by the best-subsets approach. Note that the best subset of size 2, for example, need not include the variable that was in the best subset of size 1 (for this example all the subsets are nested). The best-subset curve (red lower boundary in Figure 3.5) is necessarily decreasing, so cannot be used to select the subset size  $k$ . The question of how to choose  $k$  involves the tradeoff between bias and variance, along with the more subjective desire for parsimony. There are a number of criteria that one may use; typically we choose the smallest model that minimizes an estimate of the expected prediction error.

Many of the other approaches that we discuss in this chapter are similar, in that they use the training data to produce a sequence of models varying in complexity and indexed by a single parameter. In the next section we use

cf  
disc  
fast?  
just  
discuss  
largest  
effects?

cf  
stepwise



**FIGURE 3.5.** All possible subset models for the prostate cancer example. At each subset size is shown the residual sum-of-squares for each model of that size.

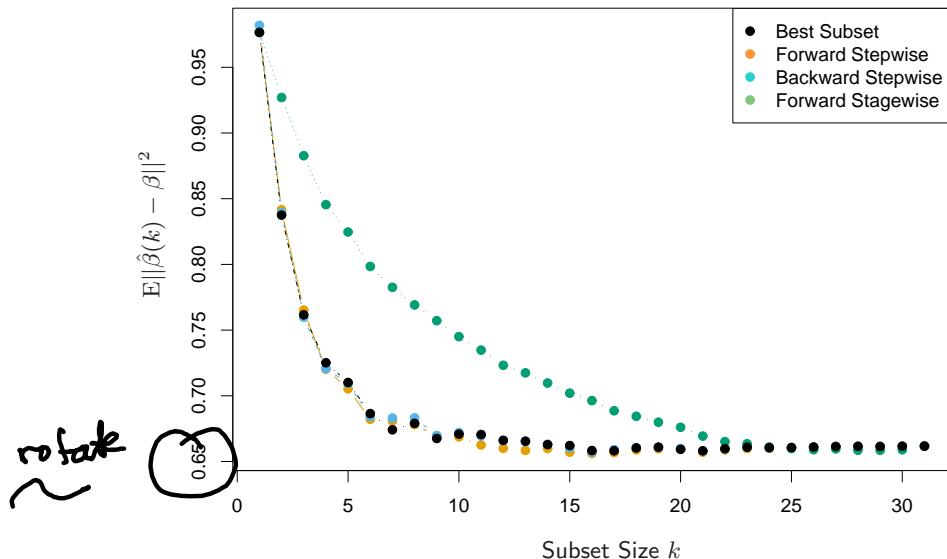
cross-validation to estimate prediction error and select  $k$ ; the AIC criterion is a popular alternative. We defer more detailed discussion of these and other approaches to Chapter 7.

### 3.3.2 Forward- and Backward-Stepwise Selection

Rather than search through all possible subsets (which becomes infeasible for  $p$  much larger than 40), we can seek a good path through them. *Forward-stepwise selection* starts with the intercept, and then sequentially adds into the model the predictor that most improves the fit. With many candidate predictors, this might seem like a lot of computation; however, clever updating algorithms can exploit the QR decomposition for the current fit to rapidly establish the next candidate (Exercise 3.9). Like best-subset regression, forward stepwise produces a sequence of models indexed by  $k$ , the subset size, which must be determined.

Forward-stepwise selection is a *greedy algorithm*, producing a nested sequence of models. In this sense it might seem sub-optimal compared to best-subset selection. However, there are several reasons why it might be preferred:

- *Computational*; for large  $p$  we cannot compute the best subset sequence, but we can always compute the forward stepwise sequence (even when  $p \gg N$ ).
- *Statistical*; a price is paid in variance for selecting the best subset of each size; forward stepwise is a more constrained search, and will have lower variance, but perhaps more bias.



**FIGURE 3.6.** Comparison of four subset-selection techniques on a simulated linear regression problem  $Y = X^T \beta + \varepsilon$ . There are  $N = 300$  observations on  $p = 31$  standard Gaussian variables, with pairwise correlations all equal to 0.85. For 10 of the variables, the coefficients are drawn at random from a  $N(0, 0.4)$  distribution; the rest are zero. The noise  $\varepsilon \sim N(0, 6.25)$ , resulting in a signal-to-noise ratio of 0.64. Results are averaged over 50 simulations. Shown is the mean-squared error of the estimated coefficient  $\hat{\beta}(k)$  at each step from the true  $\beta$ .

*Backward-stepwise selection* starts with the full model, and sequentially deletes the predictor that has the least impact on the fit. The candidate for dropping is the variable with the smallest Z-score (Exercise 3.10). Backward selection can only be used when  $N > p$ , while forward stepwise can always be used.

Figure 3.6 shows the results of a small simulation study to compare best-subset regression with the simpler alternatives forward and backward selection. Their performance is very similar, as is often the case. Included in the figure is forward stagewise regression (next section), which takes longer to reach minimum error.

→ cf Murtaugh

On the prostate cancer example, best-subset, forward and backward selection all gave exactly the same sequence of terms.

Some software packages implement hybrid stepwise-selection strategies that consider both forward and backward moves at each step, and select the “best” of the two. For example in the R package the `step` function uses the AIC criterion for weighing the choices, which takes proper account of the number of parameters fit; at each step an add or drop will be performed that minimizes the AIC score.

Other more traditional packages base the selection on  $F$ -statistics, adding “significant” terms, and dropping “non-significant” terms. These are out of fashion, since they do not take proper account of the multiple testing issues. It is also tempting after a model search to print out a summary of the chosen model, such as in Table 3.2; however, the standard errors are not valid, since they do not account for the search process. The bootstrap (Section 8.2) can be useful in such settings.

Finally, we note that often variables come in groups (such as the dummy variables that code a multi-level categorical predictor). Smart stepwise procedures (such as `step` in R) will add or drop whole groups at a time, taking proper account of their degrees-of-freedom.

\*\*

~ termwise  
- machinery  
(Julia?)

### 3.3.3 Forward-Stagewise Regression

Forward-stagewise regression (FS) is even more constrained than forward-stepwise regression. It starts like forward-stepwise regression, with an intercept equal to  $\bar{y}$ , and centered predictors with coefficients initially all 0. At each step the algorithm identifies the variable most correlated with the current residual. It then computes the simple linear regression coefficient of the residual on this chosen variable, and then adds it to the current coefficient for that variable. This is continued till none of the variables have correlation with the residuals—i.e. the least-squares fit when  $N > p$ .

Unlike forward-stepwise regression, none of the other variables are adjusted when a term is added to the model. As a consequence, forward stagewise can take many more than  $p$  steps to reach the least squares fit, and historically has been dismissed as being inefficient. It turns out that this “slow fitting” can pay dividends in high-dimensional problems. We see in Section 3.8.1 that both forward stagewise and a variant which is slowed down even further are quite competitive, especially in very high-dimensional problems.

Forward-stagewise regression is included in Figure 3.6. In this example it takes over 1000 steps to get all the correlations below  $10^{-4}$ . For subset size  $k$ , we plotted the error for the last step for which there were  $k$  nonzero coefficients. Although it catches up with the best fit, it takes longer to do so.

2

### 3.3.4 Prostate Cancer Data Example (Continued)

Table 3.3 shows the coefficients from a number of different selection and shrinkage methods. They are *best-subset selection* using an all-subsets search, *ridge regression*, the *lasso*, *principal components regression* and *partial least squares*. Each method has a complexity parameter, and this was chosen to minimize an estimate of prediction error based on tenfold cross-validation; full details are given in Section 7.10. Briefly, cross-validation works by dividing the training data randomly into ten equal parts. The learning method is fit—for a range of values of the complexity parameter—to nine-tenths of the data, and the prediction error is computed on the remaining one-tenth. This is done in turn for each one-tenth of the data, and the ten prediction error estimates are averaged. From this we obtain an estimated prediction error curve as a function of the complexity parameter.

Note that we have already divided these data into a training set of size 67 and a test set of size 30. Cross-validation is applied to the training set, since selecting the shrinkage parameter is part of the training process. The test set is there to judge the performance of the selected model.

The estimated prediction error curves are shown in Figure 3.7. Many of the curves are very flat over large ranges near their minimum. Included are estimated standard error bands for each estimated error rate, based on the ten error estimates computed by cross-validation. We have used the “one-standard-error” rule—we pick the most parsimonious model within one standard error of the minimum (Section 7.10, page 244). Such a rule acknowledges the fact that the tradeoff curve is estimated with error, and hence takes a conservative approach.

Best-subset selection chose to use the two predictors `lcvol` and `lweight`. The last two lines of the table give the average prediction error (and its estimated standard error) over the test set.

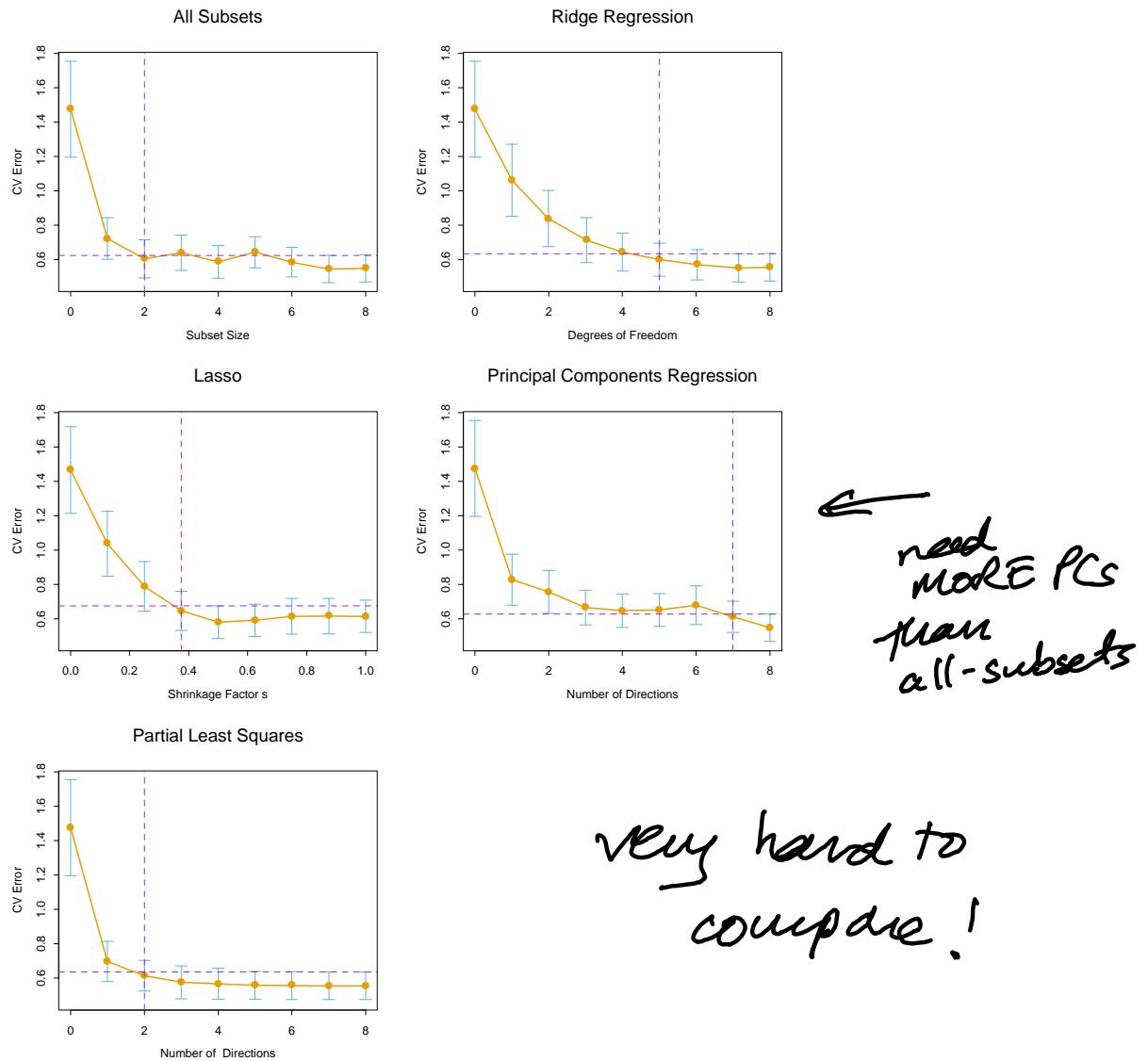


## 3.4 Shrinkage Methods

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.

### 3.4.1 Ridge Regression

Ridge regression shrinks the regression coefficients by imposing a penalty on their size. The ridge coefficients minimize a penalized residual sum of



**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.

**TABLE 3.3.** Estimated coefficients and test error results, for different subset and shrinkage methods applied to the prostate data. The blank entries correspond to variables omitted.

Term	LS	Best Subset	Ridge	Lasso	PCR	PLS
Intercept	2.465	2.477	2.452	2.468	2.497	2.452
lcavol	0.680	0.740	0.420	0.533	0.543	0.419
lweight	0.263	0.316	0.238	0.169	0.289	0.344
age	-0.141		-0.046		-0.152	-0.026
lbph	0.210		0.162	0.002	0.214	0.220
svi	0.305		0.227	0.094	0.315	0.243
lcp	-0.288		0.000		-0.051	0.079
gleason	-0.021		0.040		0.232	0.011
pgg45	0.267		0.133		-0.056	0.084
Test Error	0.521	0.492	0.492	0.479	0.449	0.528
Std Error	0.179	0.143	0.165	0.164	0.105	0.152

how  
etc?



squares,

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

Here  $\lambda \geq 0$  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). The idea of penalizing by the sum-of-squares of the parameters is also used in neural networks, where it is known as *weight decay* (Chapter 11).

An equivalent way to write the ridge problem is

$$\begin{aligned} \hat{\beta}^{\text{ridge}} &= \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p x_{ij}\beta_j)^2, \\ &\text{subject to } \sum_{j=1}^p \beta_j^2 \leq t, \end{aligned} \quad (3.42)$$

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). In addition,

~ ? stagewise?  
} not much help here...

HOW BIG  
IS THE  
DIFFERENCE  
?

does it  
matter?

notice that the intercept  $\beta_0$  has been left out of the penalty term. Penalization of the intercept would make the procedure depend on the origin chosen for  $Y$ ; that is, adding a constant  $c$  to each of the targets  $y_i$  would not simply result in a shift of the predictions by the same amount  $c$ . It can be shown (Exercise 3.5) that 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  $\bar{y} = \frac{1}{N} \sum_1^N y_i$ . The remaining coefficients get estimated by a ridge regression without intercept, using the centered  $x_{ij}$ . Henceforth we assume that this centering has been done, so that the input matrix  $\mathbf{X}$  has  $p$  (rather than  $p + 1$ ) columns.

Writing the criterion in (3.41) in matrix form,

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

the ridge regression solutions are easily seen to be

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

where  $\mathbf{I}$  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  $\mathbf{y}$ . The solution adds a positive constant to the diagonal of  $\mathbf{X}^T\mathbf{X}$  before inversion. This makes the problem nonsingular, even if  $\mathbf{X}^T\mathbf{X}$  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.

Figure 3.8 shows the ridge coefficient estimates for the prostate cancer example, plotted as functions of  $\text{df}(\lambda)$ , the effective degrees of freedom implied by the penalty  $\lambda$  (defined in (3.50) on page 68). In the case of orthonormal inputs, the ridge estimates are just a scaled version of the least squares estimates, that is,  $\hat{\beta}^{\text{ridge}} = \hat{\beta}/(1 + \lambda)$ .

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  $\lambda = \sigma^2/\tau^2$  (Exercise 3.6). Thus the ridge estimate is the mode of the posterior distribution; since the distribution is Gaussian, it is also the posterior mean.

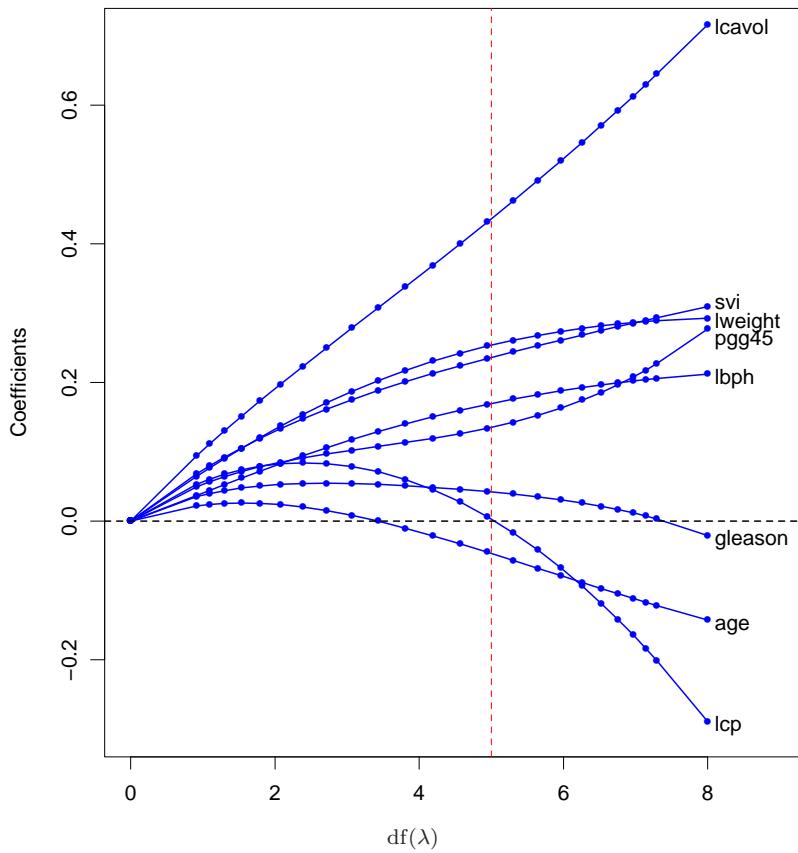
The *singular value decomposition* (SVD) of the centered input matrix  $\mathbf{X}$  gives us some additional insight into the nature of ridge regression. This decomposition is extremely useful in the analysis of many statistical methods. The SVD of the  $N \times p$  matrix  $\mathbf{X}$  has the form

$$\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T. \quad (3.45)$$

or solve

$$(\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})\beta = \mathbf{X}^T\mathbf{y}$$

(trace hat))



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

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

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

$$\begin{aligned}\mathbf{X}\hat{\beta}^{\text{ls}} &= \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y} \\ &= \mathbf{U}\mathbf{U}^T\mathbf{y},\end{aligned}\quad (3.46)$$

after some simplification. Note that  $\mathbf{U}^T\mathbf{y}$  are the coordinates of  $\mathbf{y}$  with respect to the orthonormal basis  $\mathbf{U}$ . Note also the similarity with (3.33);  $\mathbf{Q}$  and  $\mathbf{U}$  are generally different orthogonal bases for the column space of  $\mathbf{X}$  (Exercise 3.8).

Now the ridge solutions are

$$\begin{aligned}\mathbf{X}\hat{\beta}^{\text{ridge}} &= \mathbf{X}(\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})^{-1}\mathbf{X}^T\mathbf{y} \\ &= \mathbf{U}\mathbf{D}(\mathbf{D}^2 + \lambda\mathbf{I})^{-1}\mathbf{D}\mathbf{U}^T\mathbf{y} \\ &= \sum_{j=1}^p \mathbf{u}_j \frac{d_j^2}{d_j^2 + \lambda} \mathbf{u}_j^T \mathbf{y},\end{aligned}\quad (3.47)$$

where the  $\mathbf{u}_j$  are the columns of  $\mathbf{U}$ . Note that since  $\lambda \geq 0$ , we have  $d_j^2/(d_j^2 + \lambda) \leq 1$ . Like linear regression, ridge regression computes the coordinates of  $\mathbf{y}$  with respect to the orthonormal basis  $\mathbf{U}$ . 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 matrix  $\mathbf{X}$  is another way of expressing the *principal components* of the variables in  $\mathbf{X}$ . The sample covariance matrix is given by  $\mathbf{S} = \mathbf{X}^T\mathbf{X}/N$ , and from (3.45) we have

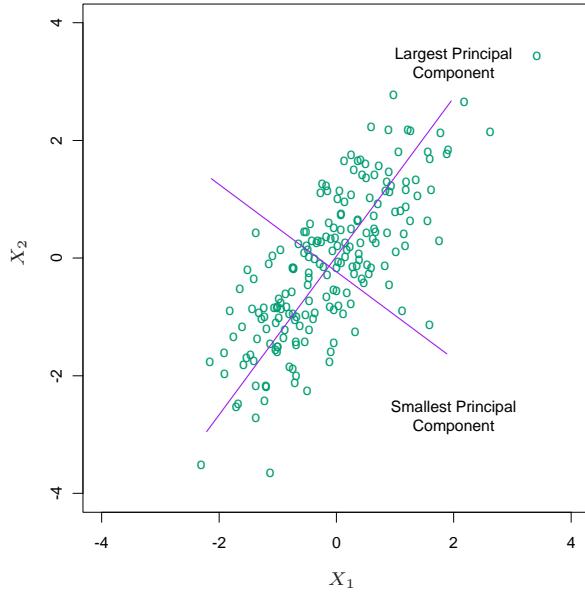
$$\mathbf{X}^T\mathbf{X} = \mathbf{V}\mathbf{D}^2\mathbf{V}^T, \quad (3.48)$$

which is the *eigen decomposition* of  $\mathbf{X}^T\mathbf{X}$  (and of  $\mathbf{S}$ , up to a factor  $N$ ). The eigenvectors  $v_j$  (columns of  $\mathbf{V}$ ) are also called the *principal components* (or Karhunen–Loeve) directions of  $\mathbf{X}$ . The first principal component direction  $v_1$  has the property that  $\mathbf{z}_1 = \mathbf{X}v_1$  has the largest sample variance amongst all normalized linear combinations of the columns of  $\mathbf{X}$ . This sample variance is easily seen to be

$$\text{Var}(\mathbf{z}_1) = \text{Var}(\mathbf{X}v_1) = \frac{d_1^2}{N}, \quad (3.49)$$

and in fact  $\mathbf{z}_1 = \mathbf{X}v_1 = \mathbf{u}_1 d_1$ . The derived variable  $\mathbf{z}_1$  is called the first principal component of  $\mathbf{X}$ , and hence  $\mathbf{u}_1$  is the normalized first principal

$$\mathbf{X}^T\mathbf{X} = \mathbf{V}\mathbf{D}^2\mathbf{V}^T$$



**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.

component. Subsequent principal components  $\mathbf{z}_j$  have maximum variance  $d_j^2/N$ , subject to being orthogonal to the earlier ones. Conversely the last principal component has *minimum* variance. Hence the small singular values  $d_j$  correspond to directions in the column space of  $\mathbf{X}$  having small variance, and ridge regression shrinks these directions the most.

Figure 3.9 illustrates the principal components of some data points in two dimensions. If we consider fitting a linear surface over this domain (the  $Y$ -axis is sticking out of the page), the configuration of the data allow us to determine its gradient more accurately in the long direction than the short. Ridge regression protects against the potentially high variance of gradients estimated in the short directions. The implicit assumption is that the response will tend to vary most in the directions of high variance of the inputs. This is often a reasonable assumption, since predictors are often chosen for study because they vary with the response variable, but need not hold in general.

} cf PCR

In Figure 3.7 we have plotted the estimated prediction error versus the quantity

$$\begin{aligned} \text{df}(\lambda) &= \text{tr}[\mathbf{X}(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T], \\ &= \text{tr}(\mathbf{H}_\lambda) \\ &= \sum_{j=1}^p \frac{d_j^2}{d_j^2 + \lambda}. \end{aligned} \quad (3.50)$$

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  $\text{df}(\lambda) = p$  when  $\lambda = 0$  (no regularization) and  $\text{df}(\lambda) \rightarrow 0$  as  $\lambda \rightarrow \infty$ . Of course there is always an additional one degree of freedom for the intercept, which was removed *a priori*. This definition is motivated in more detail in Section 3.4.4 and Sections 7.4–7.6. In Figure 3.7 the minimum occurs at  $\text{df}(\lambda) = 5.0$ . Table 3.3 shows that ridge regression reduces the test error of the full least squares estimates by a small amount.

### 3.4.2 The Lasso

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

$$\begin{aligned} \hat{\beta}^{\text{lasso}} &= \underset{\beta}{\text{argmin}} \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 \\ &\text{subject to } \sum_{j=1}^p |\beta_j| \leq t. \end{aligned} \quad (3.51)$$

Just as in ridge regression, we can re-parametrize the constant  $\beta_0$  by standardizing the predictors; the solution for  $\hat{\beta}_0$  is  $\bar{y}$ , and thereafter we fit a model without an intercept (Exercise 3.5). In the signal processing literature, the lasso is also known as *basis pursuit* (Chen et al., 1998).

We can also write the lasso problem in the equivalent *Lagrangian form*

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

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  $\sum_1^p |\beta_j|$ . This latter constraint makes the solutions nonlinear in the  $y_i$ , and there is no closed form expression as in ridge regression. Computing the lasso solution

is a quadratic programming problem, although we see in Section 3.4.4 that efficient algorithms are available for computing the entire path of solutions as  $\lambda$  is varied, with the same computational cost as for ridge regression. Because of the nature of the constraint, making  $t$  sufficiently small will cause some of the coefficients to be exactly zero. 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  $\hat{\beta}_j = \hat{\beta}_j^{\text{ls}}$ , 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. However, the nature of the shrinkage is not obvious, and we investigate it further in Section 3.4.4 below. Like the subset size in variable subset selection, or the penalty parameter in ridge regression,  $t$  should be adaptively chosen to minimize an estimate of expected prediction error.

In Figure 3.7, for ease of interpretation, we have plotted the lasso prediction error estimates versus the standardized parameter  $s = t / \sum_1^p |\hat{\beta}_j|$ . A value  $\hat{s} \approx 0.36$  was chosen by 10-fold cross-validation; this caused four coefficients to be set to zero (fifth column of Table 3.3). The resulting model has the second lowest test error, slightly lower than the full least squares model, but the standard errors of the test error estimates (last line of Table 3.3) are fairly large.

Figure 3.10 shows the lasso coefficients as the standardized tuning parameter  $s = t / \sum_1^p |\hat{\beta}_j|$  is varied. At  $s = 1.0$  these are the least squares estimates; they decrease to 0 as  $s \rightarrow 0$ . This decrease is not always strictly monotonic, although it is in this example. A vertical line is drawn at  $s = 0.36$ , the value chosen by cross-validation.

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

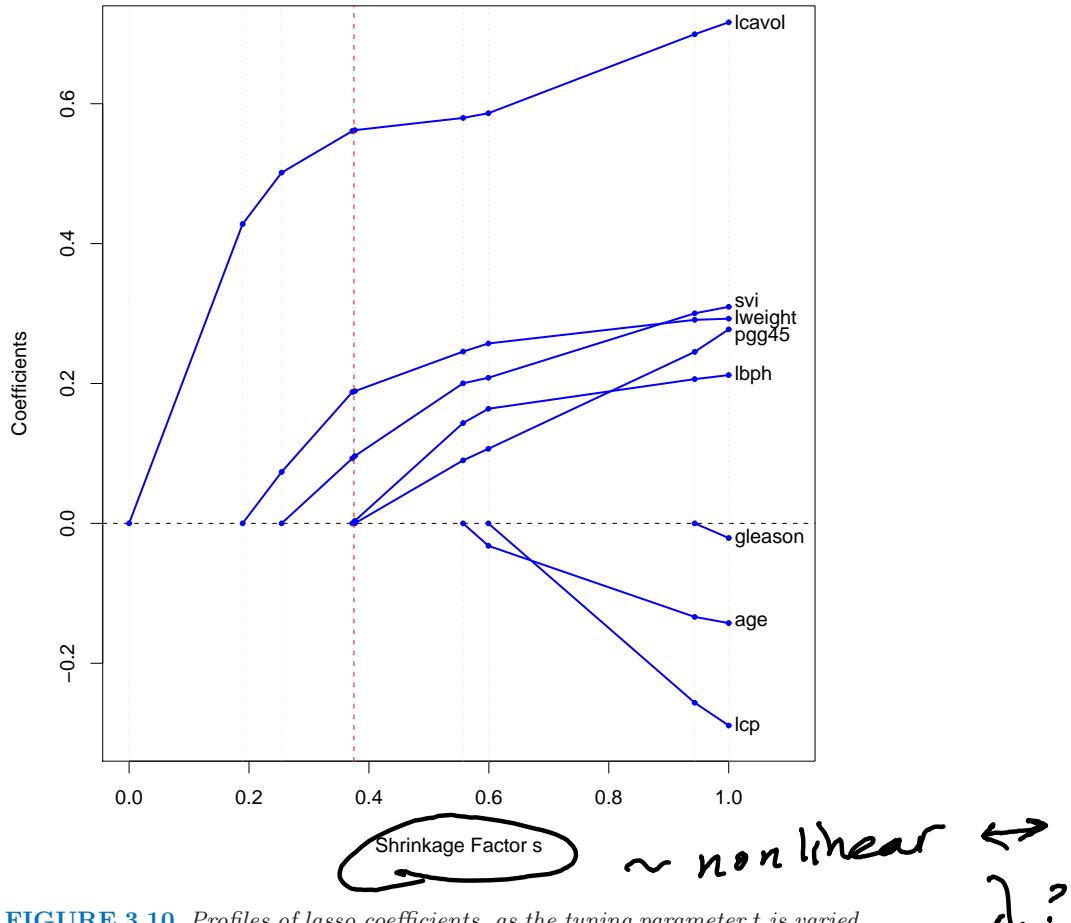
In this section we discuss and compare the three approaches discussed so far for restricting the linear regression model: subset selection, ridge regression and the lasso.

In the case of an orthonormal input matrix  $\mathbf{X}$  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  $\lambda$ , truncating at zero. This is called “soft thresholding,” and is used in the context of wavelet-based smoothing in Section 5.9. Best-subset selection drops all variables with coefficients smaller than the  $M$ th largest; this is a form of “hard-thresholding.”

Back to the nonorthogonal 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. The constraint

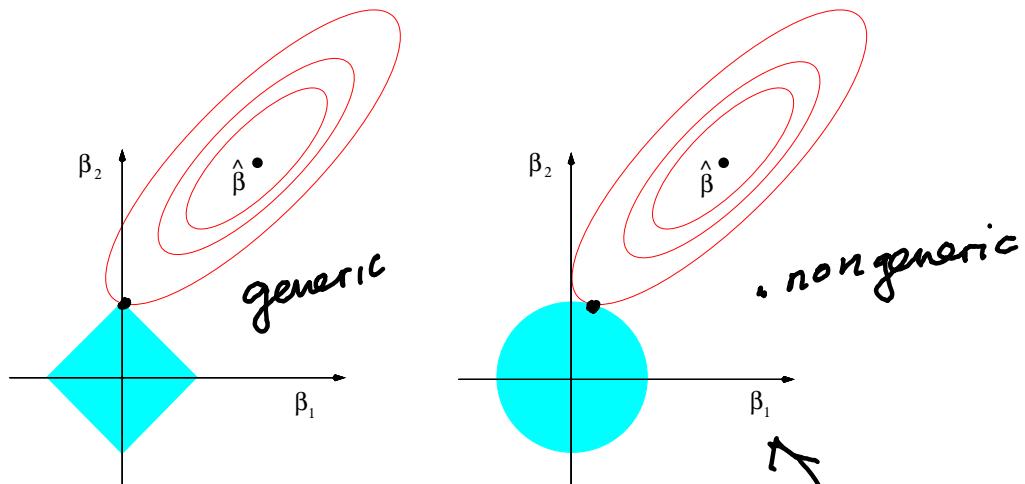
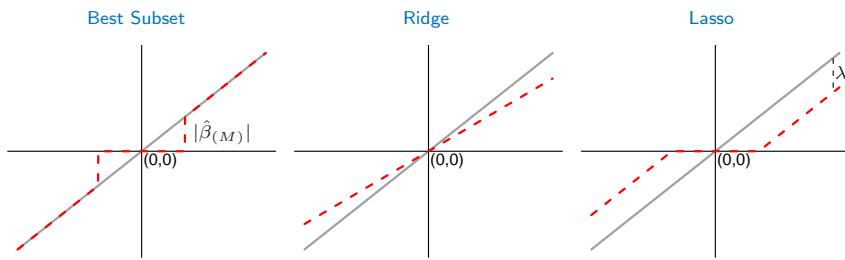
*other  
shrinkage  
(SCAD,  
spike/slab,  
...)*



**FIGURE 3.10.** Profiles of lasso coefficients, as the tuning parameter  $t$  is varied. Coefficients are plotted versus  $s = t / \sum_1^p |\hat{\beta}_j|$ . 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.

**TABLE 3.4.** Estimators of  $\beta_j$  in the case of orthonormal columns of  $\mathbf{X}$ .  $M$  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  $45^\circ$  line in gray shows the unrestricted estimate for reference.

Estimator	Formula
Best subset (size $M$ )	$\hat{\beta}_j \cdot I( \hat{\beta}_j  \geq  \hat{\beta}_{(M)} )$
Ridge	$\hat{\beta}_j / (1 + \lambda)$
Lasso	$\text{sign}(\hat{\beta}_j)( \hat{\beta}_j  - \lambda)_+$



**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.

HW

generate for  
a real  
example

region for ridge regression is the disk  $\beta_1^2 + \beta_2^2 \leq t$ , while that for lasso is the diamond  $|\beta_1| + |\beta_2| \leq t$ . 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.

We can generalize ridge regression and the lasso, and view them as Bayes estimates. Consider the criterion

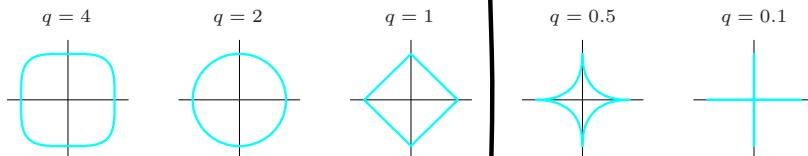
$$\tilde{\beta} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p x_{ij}\beta_j)^2 + \lambda \sum_{j=1}^p |\beta_j|^q \right\} \quad (3.53)$$

for  $q \geq 0$ . 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  $|\beta_j|^q$  as the log-prior density for  $\beta_j$ , these are also the equi-contours of the prior distribution of the parameters. The value  $q = 0$  corresponds to variable subset selection, as the penalty simply counts the number of nonzero parameters;  $q = 1$  corresponds to the lasso, while  $q = 2$  to ridge regression. Notice that for  $q \leq 1$ , the prior is not uniform in direction, but concentrates more mass in the coordinate directions. The prior corresponding to the  $q = 1$  case is an independent double exponential (or Laplace) distribution for each input, with density  $(1/2\tau) \exp(-|\beta|/\tau)$  and  $\tau = 1/\lambda$ . 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. Note, however, that they are derived as posterior modes, that is, maximizers of the posterior. It is more common to use the mean of the posterior as the Bayes estimate. Ridge regression is also the posterior mean, but the lasso and best subset selection are not.

Looking again at the criterion (3.53), we might try using other values of  $q$  besides 0, 1, or 2. Although one might consider estimating  $q$  from the data, our experience is that it is not worth the effort for the extra variance incurred. Values of  $q \in (1, 2)$  suggest a compromise between the lasso and ridge regression. Although this is the case, with  $q > 1$ ,  $|\beta_j|^q$  is differentiable at 0, and so does not share the ability of lasso ( $q = 1$ ) for



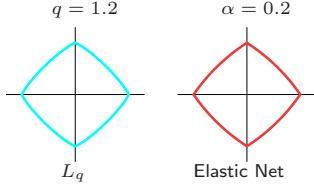
**FIGURE 3.12.** Contours of constant value of  $\sum_j |\beta_j|^q$  for given values of  $q$ .

$q \neq 2$

convexity  
differentiability

NOT what  
elastic net  
does ...

MAP vs  
Bayes  
estimate



**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-net has sharp (non-differentiable) corners, while the  $q = 1.2$  penalty does not.

setting coefficients exactly to zero. Partly for this reason as well as for computational tractability, Zou and Hastie (2005) introduced the elastic-net penalty

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

a different compromise between ridge and lasso. Figure 3.13 compares the  $L_q$  penalty with  $q = 1.2$  and the elastic-net penalty with  $\alpha = 0.2$ ; it is hard to detect the difference by eye. The elastic-net selects variables like the lasso, and shrinks together the coefficients of correlated predictors like ridge. It also has considerable computational advantages over the  $L_q$  penalties. We discuss the elastic-net further in Section 18.4.)

### 3.4.4 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 (Section 3.3.2). As we will see, LAR is intimately connected with the lasso, and in fact provides an extremely efficient algorithm for computing the entire lasso path as in 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 least-squares 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