# Chapter 7
# MLR Variable Selection and Lasso

This chapter considers MLR variable selection and prediction intervals. Prediction regions and prediction intervals applied to a bootstrap sample can result in confidence regions and confidence intervals. The bootstrap confidence regions will be used for inference after variable selection.

Some shrinkage methods do variable selection: the MLR method, such as a OLS, uses the predictors that had nonzero shrinkage estimator coefficients. These methods include least angle regression, lasso, relaxed lasso, and elastic net. Least angle regression variable selection is the LARS-OLS hybrid estimator of Efron et al. (2004, p. 421). Lasso variable selection is called relaxed lasso by Hastie et al. (2015, p. 12), and the relaxed lasso estimator with $\phi = 0$ by Meinshausen (2007, p. 376). Also see Fan and Li (2001), Tibshirani (1996), and Zou and Hastie (2005). The Meinshausen (2007) relaxed lasso estimator fits lasso with penalty $\lambda_n$ to get a subset of variables with nonzero coefficients, and then fits lasso with a smaller penalty $\phi_n$ to this subset of variables where $n$ is the sample size.

## 7.1 Introduction

*Variable selection*, also called subset or model selection, is the search for a subset of predictor variables that can be deleted without important loss of information if $n/p$ is large (and the search for a useful subset of predictors if $n/p$ is not large). Consider the 1D regression model where $Y \perp\!\!\!\perp \boldsymbol{x}|SP$ where $SP = \boldsymbol{x}^T\boldsymbol{\beta}$. See Chapters 1 and 10. A *model for variable selection* can be described by

$$\boldsymbol{x}^T\boldsymbol{\beta} = \boldsymbol{x}_S^T\boldsymbol{\beta}_S + \boldsymbol{x}_E^T\boldsymbol{\beta}_E = \boldsymbol{x}_S^T\boldsymbol{\beta}_S \qquad (7.1)$$

where $\boldsymbol{x} = (\boldsymbol{x}_S^T, \boldsymbol{x}_E^T)^T$ is a $p \times 1$ vector of predictors, $\boldsymbol{x}_S$ is an $a_S \times 1$ vector, and $\boldsymbol{x}_E$ is a $(p - a_S) \times 1$ vector. Given that $\boldsymbol{x}_S$ is in the model, $\boldsymbol{\beta}_E = \boldsymbol{0}$ and

$E$ denotes the subset of terms that can be eliminated given that the subset $S$ is in the model.

Since $S$ is unknown, candidate subsets will be examined. Let $\boldsymbol{x}_I$ be the vector of $a$ terms from a candidate subset indexed by $I$, and let $\boldsymbol{x}_O$ be the vector of the remaining predictors (out of the candidate submodel). Then

$$\boldsymbol{x}^T\boldsymbol{\beta} = \boldsymbol{x}_I^T\boldsymbol{\beta}_I + \boldsymbol{x}_O^T\boldsymbol{\beta}_O.$$

Suppose that $S$ is a subset of $I$ and that model (7.1) holds. Then

$$\boldsymbol{x}^T\boldsymbol{\beta} = \boldsymbol{x}_S^T\boldsymbol{\beta}_S = \boldsymbol{x}_S^T\boldsymbol{\beta}_S + \boldsymbol{x}_{I/S}^T\boldsymbol{\beta}_{(I/S)} + \boldsymbol{x}_O^T\boldsymbol{0} = \boldsymbol{x}_I^T\boldsymbol{\beta}_I$$

where $\boldsymbol{x}_{I/S}$ denotes the predictors in $I$ that are not in $S$. Since this is true regardless of the values of the predictors, $\boldsymbol{\beta}_O = \boldsymbol{0}$ and the sample correlation $\mathrm{corr}(\boldsymbol{x}_i^T\boldsymbol{\beta}, \boldsymbol{x}_{I,i}^T\boldsymbol{\beta}_I) = 1.0$ for the population model if $S \subseteq I$. The estimated sufficient predictor (ESP) is $\boldsymbol{x}^T\hat{\boldsymbol{\beta}}$, and *a submodel $I$ is worth considering if the correlation* $\mathrm{corr}(ESP, ESP(I)) \geq 0.95$.

To clarify notation, suppose $p = 4$, a constant $x_1 = 1$ corresponding to $\beta_1$ is always in the model, and $\boldsymbol{\beta} = (\beta_1, \beta_2, 0, 0)^T$. Then the $J = 2^{p-1} = 8$ possible subsets of $\{1, 2, ..., p\}$ that always contain 1 are $I_1 = \{1\}$, $S = I_2 = \{1, 2\}$, $I_3 = \{1, 3\}$, $I_4 = \{1, 4\}$, $I_5 = \{1, 2, 3\}$, $I_6 = \{1, 2, 4\}$, $I_7 = \{1, 3, 4\}$, and $I_8 = \{1, 2, 3, 4\}$. There are $2^{p-a_S} = 4$ subsets $I_2, I_5, I_6$, and $I_8$ such that $S \subseteq I_j$. Let $\hat{\boldsymbol{\beta}}_{I_7} = (\hat{\beta}_1, \hat{\beta}_3, \hat{\beta}_4)^T$ and $\boldsymbol{x}_{I_7} = (x_1, x_3, x_4)^T$.

Let $I_{min}$ correspond to the set of predictors selected by a variable selection method such as forward selection or lasso variable selection. If $\hat{\boldsymbol{\beta}}_I$ is $a \times 1$, use zero padding to form the $p \times 1$ vector $\hat{\boldsymbol{\beta}}_{I,0}$ from $\hat{\boldsymbol{\beta}}_I$ by adding 0s corresponding to the omitted variables. For example, if $p = 4$ and $\hat{\boldsymbol{\beta}}_{I_{min}} = (\hat{\beta}_1, \hat{\beta}_3)^T$, then the observed variable selection estimator $\hat{\boldsymbol{\beta}}_{VS} = \hat{\boldsymbol{\beta}}_{I_{min},0} = (\hat{\beta}_1, 0, \hat{\beta}_3, 0)^T$. As a statistic, $\hat{\boldsymbol{\beta}}_{VS} = \hat{\boldsymbol{\beta}}_{I_k,0}$ with probabilities $\pi_{kn} = P(I_{min} = I_k)$ for $k = 1, ..., J$ where there are $J$ subsets, e.g. $J = 2^p - 1$.

**Definition 7.1.** The model $Y \perp\!\!\!\perp \boldsymbol{x} | \boldsymbol{x}^T\boldsymbol{\beta}$ that uses all of the predictors is called the *full model*. A model $Y \perp\!\!\!\perp \boldsymbol{x}_I | \boldsymbol{x}_I^T\boldsymbol{\beta}_I$ that uses a subset $\boldsymbol{x}_I$ of the predictors is called a *submodel*. The **full model is always a submodel**. The full model has *sufficient predictor* $SP = \boldsymbol{x}^T\boldsymbol{\beta}$ and the submodel has $SP = \boldsymbol{x}_I^T\boldsymbol{\beta}_I$.

Forward selection or backward elimination with the Akaike (1973) AIC criterion or Schwarz (1978) BIC criterion are often used for variable selection. Lasso variable selection or elastic net variable selection fits OLS to the predictors than had nonzero lasso or elastic net coefficients. .

Underfitting occurs if submodel $I$ does not contain $S$. Following, for example, Pelawa Watagoda (2019), let $\boldsymbol{X} = [\boldsymbol{X}_I \ \boldsymbol{X}_O]$ and $\boldsymbol{\beta} = (\boldsymbol{\beta}_I^T, \boldsymbol{\beta}_O^T)^T$. Then $\boldsymbol{X}\boldsymbol{\beta} = \boldsymbol{X}_I\boldsymbol{\beta}_I + \boldsymbol{X}_O\boldsymbol{\beta}_O$, and $\hat{\boldsymbol{\beta}}_I = (\boldsymbol{X}_I\boldsymbol{X}_I)^{-1}\boldsymbol{X}_I^T\boldsymbol{Y} = \boldsymbol{A}\boldsymbol{Y}$. Assuming the usual MLR model, $\mathrm{Cov}(\hat{\boldsymbol{\beta}}_I) = \mathrm{Cov}(\boldsymbol{A}\boldsymbol{Y}) = \boldsymbol{A}\sigma^2\boldsymbol{I}\boldsymbol{A}^T = \sigma^2(\boldsymbol{X}_I^T\boldsymbol{X}_I)^{-1}$.

Now $E(\hat{\boldsymbol{\beta}}_I) = E(\boldsymbol{AY}) = \boldsymbol{AX}\boldsymbol{\beta} = (\boldsymbol{X}_I\boldsymbol{X}_I)^{-1}\boldsymbol{X}_I^T(\boldsymbol{X}_I\boldsymbol{\beta}_I + \boldsymbol{X}_O\boldsymbol{\beta}_O) =$

$$\boldsymbol{\beta}_I + (\boldsymbol{X}_I\boldsymbol{X}_I)^{-1}\boldsymbol{X}_I^T\boldsymbol{X}_O\boldsymbol{\beta}_O = \boldsymbol{\beta}_I + \boldsymbol{AX}_O\boldsymbol{\beta}_O.$$

If $S \subseteq I$, then $\boldsymbol{\beta}_O = \boldsymbol{0}$, but if underfitting occurs then the bias vector $\boldsymbol{AX}_O\boldsymbol{\beta}_O$ can be large.

## 7.2 OLS Variable Selection

Simpler models are easier to explain and use than more complicated models, and there are several other important reasons to perform variable selection. For example, an OLS MLR model with unnecessary predictors has $\sum_{i=1}^{n} V(\hat{Y}_i)$ that is too large. If (7.1) holds, $S \subseteq I$, $\boldsymbol{\beta}_S$ is an $a_S \times 1$ vector, and $\boldsymbol{\beta}_I$ is a $j \times 1$ vector with $j > a_S$, then

$$\frac{1}{n}\sum_{i=1}^{n} V(\hat{Y}_{Ii}) = \frac{\sigma^2 j}{n} > \frac{\sigma^2 a_S}{n} = \frac{1}{n}\sum_{i=1}^{n} V(\hat{Y}_{Si}). \tag{7.2}$$

In particular, the full model has $j = p$. Hence having unnecessary predictors decreases the precision for prediction. Fitting unnecessary predictors is sometimes called *fitting noise* or *overfitting*. As an extreme case, suppose that the full model contains $p = n$ predictors, including a constant, so that the hat matrix $\boldsymbol{H} = \boldsymbol{I}_n$, the $n \times n$ identity matrix. Then $\hat{Y} = Y$ so that $\text{VAR}(\hat{Y}|\boldsymbol{x}) = \text{VAR}(Y)$. A model $I$ underfits if it does not include all of the predictors in $S$. A model $I$ does not underfit if $S \subseteq I$.

To see that (7.2) holds, assume that the full model includes all $p$ possible terms so the full model may overfit but does not underfit. Then $\hat{\boldsymbol{Y}} = \boldsymbol{HY}$ and $\text{Cov}(\hat{\boldsymbol{Y}}) = \sigma^2 \boldsymbol{HIH}^T = \sigma^2 \boldsymbol{H}$. Thus

$$\frac{1}{n}\sum_{i=1}^{n} V(\hat{Y}_i) = \frac{1}{n}tr(\sigma^2 \boldsymbol{H}) = \frac{\sigma^2}{n}tr((\boldsymbol{X}^T\boldsymbol{X})^{-1}\boldsymbol{X}^T\boldsymbol{X}) = \frac{\sigma^2 p}{n}$$

where $tr(\boldsymbol{A})$ is the trace operation. Replacing $p$ by $j$ and $a_S$ and replacing $\boldsymbol{H}$ by $\boldsymbol{H}_I$ and $\boldsymbol{H}_S$ implies Equation (7.2). Hence if only $a_S$ parameters are needed and $p >> a_S$, then serious overfitting occurs and increases $\frac{1}{n}\sum_{i=1}^{n} V(\hat{Y}_i)$.

Two important summaries for submodel $I$ are $R^2(I)$, the proportion of the variability of $Y$ explained by the nontrivial predictors in the model, and $MSE(I) = \hat{\sigma}_I^2$, the estimated error variance. See Definitions 5.17 and 5.18. Suppose that model $I$ contains $k$ predictors, including a constant. Since adding predictors does not decrease $R^2$, the adjusted $R_A^2(I)$ is often used, where

$$R_A^2(I) = 1 - (1 - R^2(I))\frac{n}{n-k} = 1 - MSE(I)\frac{n}{SST}.$$

See Seber and Lee (2003, pp. 400-401). Hence the model with the maximum $R_A^2(I)$ is also the model with the minimum $MSE(I)$.

For multiple linear regression, recall that if the candidate model of $\boldsymbol{x}_I$ has $k$ terms (including the constant), then the partial $F$ statistic for testing whether the $p - k$ predictor variables in $\boldsymbol{x}_O$ can be deleted is

$$F_I = \frac{SSE(I) - SSE}{(n-k) - (n-p)} \Big/ \frac{SSE}{n-p} = \frac{n-p}{p-k}\left[\frac{SSE(I)}{SSE} - 1\right]$$

where SSE is the error sum of squares from the full model, and SSE(I) is the error sum of squares from the candidate submodel. An extremely important criterion for variable selection is the $C_p$ criterion.

**Definition 7.2.**

$$C_p(I) = \frac{SSE(I)}{MSE} + 2k - n = (p-k)(F_I - 1) + k$$

where MSE is the error mean square for the full model.

Note that when $H_0$ is true, $(p-k)(F_I - 1) + k \xrightarrow{D} \chi^2_{p-k} + 2k - p$ for a large class of iid error distributions. Minimizing $C_p(I)$ is equivalent to minimizing $MSE\,[C_p(I)] = SSE(I) + (2k - n)MSE = \boldsymbol{r}^T(I)\boldsymbol{r}(I) + (2k - n)MSE$. The following theorem helps explain why $C_p$ is a useful criterion and suggests that for subsets $I$ with $k$ terms, submodels with $C_p(I) \le \min(2k, p)$ are especially interesting. Olive and Hawkins (2005) show that this interpretation of $C_p$ can be generalized to 1D regression models with a linear predictor $\boldsymbol{\beta}^T\boldsymbol{x} = \boldsymbol{x}^T\boldsymbol{\beta}$, such as generalized linear models. Denote the residuals and fitted values from the *full model* by $r_i = Y_i - \boldsymbol{x}_i^T\hat{\boldsymbol{\beta}} = Y_i - \hat{Y}_i$ and $\hat{Y}_i = \boldsymbol{x}_i^T\hat{\boldsymbol{\beta}}$ respectively. Similarly, let $\hat{\boldsymbol{\beta}}_I$ be the estimate of $\boldsymbol{\beta}_I$ obtained from the regression of $Y$ on $\boldsymbol{x}_I$ and denote the corresponding residuals and fitted values by $r_{I,i} = Y_i - \boldsymbol{x}_{I,i}^T\hat{\boldsymbol{\beta}}_I$ and $\hat{Y}_{I,i} = \boldsymbol{x}_{I,i}^T\hat{\boldsymbol{\beta}}_I$ where $i = 1, ..., n$.

**Theorem 7.1.** Suppose that a numerical variable selection method suggests several submodels with $k$ predictors, including a constant, where $2 \le k \le p$.

a) The model $I$ that minimizes $C_p(I)$ maximizes $\mathrm{corr}(r, r_I)$.

b) $C_p(I) \le 2k$ implies that $\mathrm{corr}(r, r_I) \ge \sqrt{1 - \frac{p}{n}}$.

c) As $\mathrm{corr}(r, r_I) \to 1$,

$$\mathrm{corr}(\boldsymbol{x}^T\hat{\boldsymbol{\beta}}, \boldsymbol{x}_I^T\hat{\boldsymbol{\beta}}_I) = \mathrm{corr}(\mathrm{ESP}, \mathrm{ESP}(I)) = \mathrm{corr}(\hat{Y}, \hat{Y}_I) \to 1.$$

**Proof.** These results are a corollary of Theorem 7.2 below. $\square$

**Remark 7.1.** Consider the model $I_i$ that deletes the predictor $x_i$. Then the model has $k = p - 1$ predictors including the constant, and the test statistic is $t_i$ where

$$t_i^2 = F_{I_i}.$$

Using Definition 7.2 and $C_p(I_{full}) = p$, it can be shown that

$$C_p(I_i) = C_p(I_{full}) + (t_i^2 - 2).$$

Using the screen $C_p(I) \leq \min(2k, p)$ suggests that the predictor $x_i$ should not be deleted if

$$|t_i| > \sqrt{2} \approx 1.414.$$

If $|t_i| < \sqrt{2}$ then the predictor can probably be deleted since $C_p$ decreases. The literature suggests using the $C_p(I) \leq k$ screen, but this screen eliminates too many potentially useful submodels.

More generally, it can be shown that $C_p(I) \leq 2k$ iff

$$F_I \leq \frac{p}{p - k}.$$

Now $k$ is the number of terms in the model $I$ including a constant while $p - k$ is the number of terms set to 0. As $k \to 0$, the partial $F$ test will reject Ho: $\boldsymbol{\beta}_O = \mathbf{0}$ (i.e. say that the full model should be used instead of the submodel $I$) unless $F_I$ is not much larger than 1. If $p$ is very large and $p - k$ is very small, then the partial $F$ test will tend to suggest that there is a model $I$ that is about as good as the full model even though model $I$ deletes $p - k$ predictors.

**Definition 7.3.** The "fit–fit" or *FF plot* is a plot of $\hat{Y}_{I,i}$ versus $\hat{Y}_i$ while a "residual–residual" or *RR plot* is a plot $r_{I,i}$ versus $r_i$. A *response plot* is a plot of $\hat{Y}_{I,i}$ versus $Y_i$. An *EE plot* is a plot of ESP(I) versus ESP. For MLR, the EE and FF plots are equivalent.

Six graphs will be used to compare the full model and the candidate submodel: the FF plot, RR plot, the response plots from the full and submodel, and the residual plots from the full and submodel. These six plots will contain a great deal of information about the candidate subset provided that Equation (7.1) holds and that a good estimator (such as OLS) for $\hat{\boldsymbol{\beta}}$ and $\hat{\boldsymbol{\beta}}_I$ is used.

**Application 7.1.** To visualize whether a candidate submodel using predictors $\boldsymbol{x}_I$ is good, use the fitted values and residuals from the submodel and full model to make an RR plot of the $r_{I,i}$ versus the $r_i$ and an FF plot of $\hat{Y}_{I,i}$ versus $\hat{Y}_i$. Add the OLS line to the RR plot and identity line to both plots as visual aids. The subset $I$ is good if the plotted points cluster tightly about

the identity line in *both plots*. In particular, the OLS line and the identity line should "nearly coincide" so that it is difficult to tell that the two lines intersect at the origin in the RR plot.

To verify that the six plots are useful for assessing variable selection, the following notation will be useful. Suppose that all submodels include a constant and that $X$ is the full rank $n \times p$ design matrix for the full model. Let the corresponding vectors of OLS fitted values and residuals be $\hat{Y} = X(X^T X)^{-1} X^T Y = HY$ and $r = (I - H)Y$, respectively. Suppose that $X_I$ is the $n \times k$ design matrix for the candidate submodel and that the corresponding vectors of OLS fitted values and residuals are $\hat{Y}_I = X_I (X_I^T X_I)^{-1} X_I^T Y = H_I Y$ and $r_I = (I - H_I)Y$, respectively.

A plot can be very useful if the OLS line can be compared to a reference line and if the OLS slope is related to some quantity of interest. Suppose that a plot of $w$ versus $z$ places $w$ on the horizontal axis and $z$ on the vertical axis. Then denote the OLS line by $\hat{z} = a + bw$. The following theorem shows that the plotted points in the FF, RR, and response plots will cluster about the identity line. Notice that the theorem is a property of OLS and holds even if the data does not follow an MLR model. Let $\mathrm{corr}(x, y)$ denote the correlation between $x$ and $y$.

**Theorem 7.2.** Suppose that every submodel contains a constant and that $X$ is a full rank matrix.
**Response Plot:** i) If $w = \hat{Y}_I$ and $z = Y$ then the OLS line is the identity line.
ii) If $w = Y$ and $z = \hat{Y}_I$ then the OLS line has slope $b = [\mathrm{corr}(Y, \hat{Y}_I)]^2 = R^2(I)$ and intercept $a = \overline{Y}(1 - R^2(I))$ where $\overline{Y} = \sum_{i=1}^{n} Y_i/n$ and $R^2(I)$ is the coefficient of multiple determination from the candidate model.
**FF or EE Plot:** iii) If $w = \hat{Y}_I$ and $z = \hat{Y}$ then the OLS line is the identity line. Note that $ESP(I) = \hat{Y}_I$ and $ESP = \hat{Y}$.
iv) If $w = \hat{Y}$ and $z = \hat{Y}_I$ then the OLS line has slope $b = [\mathrm{corr}(\hat{Y}, \hat{Y}_I)]^2 = SSR(I)/SSR$ and intercept $a = \overline{Y}[1 - (SSR(I)/SSR)]$ where SSR is the regression sum of squares.
**RR Plot:** v) If $w = r$ and $z = r_I$ then the OLS line is the identity line.
vi) If $w = r_I$ and $z = r$ then $a = 0$ and the OLS slope $b = [\mathrm{corr}(r, r_I)]^2$ and

$$\mathrm{corr}(r, r_I) = \sqrt{\frac{SSE}{SSE(I)}} = \sqrt{\frac{n - p}{C_p(I) + n - 2k}} = \sqrt{\frac{n - p}{(p - k)F_I + n - p}}.$$

**Proof:** Recall that $H$ and $H_I$ are symmetric idempotent matrices and that $HH_I = H_I$. The mean of OLS fitted values is equal to $\overline{Y}$ and the mean of OLS residuals is equal to 0. If the OLS line from regressing $z$ on $w$ is $\hat{z} = a + bw$, then $a = \overline{z} - b\overline{w}$ and

$$b = \frac{\sum(w_i - \overline{w})(z_i - \overline{z})}{\sum(w_i - \overline{w})^2} = \frac{SD(z)}{SD(w)}\mathrm{corr}(z, w).$$

Also recall that the OLS line passes through the means of the two variables $(\overline{w}, \overline{z})$.

(*) Notice that the OLS slope from regressing $z$ on $w$ is equal to one if and only if the OLS slope from regressing $w$ on $z$ is equal to $[\mathrm{corr}(z, w)]^2$.

i) The slope $b = 1$ if $\sum \hat{Y}_{I,i}Y_i = \sum \hat{Y}_{I,i}^2$. This equality holds since $\hat{\boldsymbol{Y}}_I^T\boldsymbol{Y} = \boldsymbol{Y}^T\boldsymbol{H}_I\boldsymbol{Y} = \boldsymbol{Y}^T\boldsymbol{H}_I\boldsymbol{H}_I\boldsymbol{Y} = \hat{\boldsymbol{Y}}_I^T\hat{\boldsymbol{Y}}_I$. Since $b = 1$, $a = \overline{Y} - \overline{Y} = 0$.

ii) By (*), the slope

$$b = [\mathrm{corr}(Y, \hat{Y}_I)]^2 = R^2(I) = \frac{\sum(\hat{Y}_{I,i} - \overline{Y})^2}{\sum(Y_i - \overline{Y})^2} = SSR(I)/SSTO.$$

The result follows since $a = \overline{Y} - b\overline{Y}$.

iii) The slope $b = 1$ if $\sum \hat{Y}_{I,i}\hat{Y}_i = \sum \hat{Y}_{I,i}^2$. This equality holds since $\hat{\boldsymbol{Y}}^T\hat{\boldsymbol{Y}}_I = \boldsymbol{Y}^T\boldsymbol{H}\boldsymbol{H}_I\boldsymbol{Y} = \boldsymbol{Y}^T\boldsymbol{H}_I\boldsymbol{Y} = \hat{\boldsymbol{Y}}_I^T\hat{\boldsymbol{Y}}_I$. Since $b = 1$, $a = \overline{Y} - \overline{Y} = 0$.

iv) From iii),

$$1 = \frac{SD(\hat{Y})}{SD(\hat{Y}_I)}[\mathrm{corr}(\hat{Y}, \hat{Y}_I)].$$

Hence

$$\mathrm{corr}(\hat{Y}, \hat{Y}_I) = \frac{SD(\hat{Y}_I)}{SD(\hat{Y})}$$

and the slope

$$b = \frac{SD(\hat{Y}_I)}{SD(\hat{Y})}\mathrm{corr}(\hat{Y}, \hat{Y}_I) = [\mathrm{corr}(\hat{Y}, \hat{Y}_I)]^2.$$

Also the slope

$$b = \frac{\sum(\hat{Y}_{I,i} - \overline{Y})^2}{\sum(\hat{Y}_i - \overline{Y})^2} = SSR(I)/SSR.$$

The result follows since $a = \overline{Y} - b\overline{Y}$.

v) The OLS line passes through the origin. Hence $a = 0$. The slope $b = \boldsymbol{r}^T\boldsymbol{r}_I/\boldsymbol{r}^T\boldsymbol{r}$. Since $\boldsymbol{r}^T\boldsymbol{r}_I = \boldsymbol{Y}^T(\boldsymbol{I} - \boldsymbol{H})(\boldsymbol{I} - \boldsymbol{H}_I)\boldsymbol{Y}$ and $(\boldsymbol{I} - \boldsymbol{H})(\boldsymbol{I} - \boldsymbol{H}_I) = \boldsymbol{I} - \boldsymbol{H}$, the numerator $\boldsymbol{r}^T\boldsymbol{r}_I = \boldsymbol{r}^T\boldsymbol{r}$ and $b = 1$.

vi) Again $a = 0$ since the OLS line passes through the origin. From v),

$$1 = \sqrt{\frac{SSE(I)}{SSE}}[\mathrm{corr}(r, r_I)].$$

Hence

$$\text{corr}(r, r_I) = \sqrt{\frac{SSE}{SSE(I)}}$$

and the slope

$$b = \sqrt{\frac{SSE}{SSE(I)}}[\text{corr}(r, r_I)] = [\text{corr}(r, r_I)]^2.$$

Algebra shows that

$$\text{corr}(r, r_I) = \sqrt{\frac{n-p}{C_p(I) + n - 2k}} = \sqrt{\frac{n-p}{(p-k)F_I + n - p}}. \quad \square$$

**Remark 7.2.** Let $I_{min}$ be the model than minimizes $C_p(I)$ among the models $I$ generated from the variable selection method such as forward selection. Assuming the the full model $I_p$ is one of the models generated, then $C_p(I_{min}) \leq C_p(I_p) = p$, and $\text{corr}(r, r_{I_{min}}) \to 1$ as $n \to \infty$ by Theorem 7.2 vi). Referring to Equation (7.1), if $P(S \subseteq I_{min})$ does not go to 1 as $n \to \infty$, then the above correlation would not go to one. Hence $P(S \subseteq I_{min}) \to 1$ as $n \to \infty$.

A standard model selection procedure will often be needed to suggest models. For example, forward selection or backward elimination could be used. If $p < 30$, Furnival and Wilson (1974) provide a technique for selecting a few candidate subsets after examining all possible subsets.

**Remark 7.3.** Daniel and Wood (1980, p. 85) suggest using Mallows' graphical method for screening subsets by plotting $k$ versus $C_p(I)$ for models close to or under the $C_p = k$ line. Theorem 7.2 vi) implies that if $C_p(I) \leq k$ or $F_I < 1$, then $\text{corr}(r, r_I)$ and $\text{corr}(ESP, ESP(I))$ both go to 1.0 as $n \to \infty$. Hence models $I$ that satisfy the $C_p(I) \leq k$ screen will contain the true model $S$ with high probability when $n$ is large. This result does not guarantee that the true model $S$ will satisfy the screen, but overfit is likely. Let $d$ be a lower bound on $\text{corr}(r, r_I)$. Theorem 7.2 vi) implies that if

$$C_p(I) \leq 2k + n\left[\frac{1}{d^2} - 1\right] - \frac{p}{d^2},$$

then $\text{corr}(r, r_I) \geq d$. The simple screen $C_p(I) \leq 2k$ corresponds to

$$d \equiv d_n = \sqrt{1 - \frac{p}{n}}.$$

To avoid excluding too many good submodels, consider models $I$ with $C_p(I) \leq \min(2k, p)$. Models under both the $C_p = k$ line and the $C_p = 2k$ line are of interest.

**Rule of thumb 7.1.** a) After using a numerical method such as forward selection or backward elimination, let $I_{min}$ correspond to the submodel with the smallest $C_p$. Find the submodel $I_I$ with the fewest number of predictors such that $C_p(I_I) \leq C_p(I_{min}) + 1$. Then $I_I$ is the initial submodel that should be examined. It is possible that $I_I = I_{min}$ or that $I_I$ is the full model. Do not use more predictors than model $I_I$ to avoid overfitting.

b) Models $I$ with fewer predictors than $I_I$ such that $C_p(I) \leq C_p(I_{min}) + 4$ are interesting and should also be examined.

c) Models $I$ with $k$ predictors, including a constant and with fewer predictors than $I_I$ such that $C_p(I_{min}) + 4 < C_p(I) \leq \min(2k, p)$ should be checked but often underfit: important predictors are deleted from the model. Underfit is especially likely to occur if a predictor with one degree of freedom is deleted (if the $c - 1$ indicator variables corresponding to a factor are deleted, then the factor has $c - 1$ degrees of freedom) and the jump in $C_p$ is large, greater than 4, say.

d) If there are no models $I$ with fewer predictors than $I_I$ such that $C_p(I) \leq \min(2k, p)$, then model $I_I$ is a good candidate for the best subset found by the numerical procedure.

Forward selection forms a sequence of submodels $I_1, ..., I_p$ where $I_j$ uses $j$ predictors including the constant. Let $I_1$ use $x_1^* = x_1 \equiv 1$: the model has a constant but no nontrivial predictors. To form $I_2$, consider all models $I$ with two predictors including $x_1^*$. Compute $SSE(I) = RSS(I) = \boldsymbol{r}^T(I)\boldsymbol{r}(I) = \sum_{i=1}^n r_i^2(I) = \sum_{i=1}^n (Y_i - \hat{Y}_i(I))^2$. Let $I_2$ minimize $SSE(I)$ for the $p-1$ models $I$ that contain $x_1^*$ and one other predictor. Denote the predictors in $I_2$ by $x_1^*, x_2^*$. In general, to form $I_j$ consider all models $I$ with $j$ predictors including variables $x_1^*, ..., x_{j-1}^*$. Compute $SSE(I)$, and let $I_j$ minimize $SSE(I)$ for the $p-j+1$ models $I$ that contain $x_1^*, ..., x_{j-1}^*$ and one other predictor not already selected. Denote the predictors in $I_j$ by $x_1^*, ..., x_j^*$. Continue in this manner for $j = 2, ..., M = p$.

Backward elimination also forms a sequence of submodels $I_1, ..., I_p$ where $I_j$ uses $j$ predictors including the constant. Let $I_p$ be the full model. To form $I_{p-1}$ consider all models $I$ with $p-1$ predictors including the constant. Compute $SSE(I)$ and let $I_{p-1}$ minimize $SSE(I)$ for the $p-1$ models $I$ that exclude one of the predictors $x_2, ..., x_p$. Denote the predictors in $I_{p-1}$ by $x_1^*, x_2^*, ..., x_{p-1}^*$. In general, to form $I_j$ consider all models $I$ with $j$ predictors including variables $x_1^*, ..., x_{j+1}^*$. Compute $SSE(I)$, and let $I_j$ minimize $SSE(I)$ for the $p-j+1$ models $I$ that exclude one of the predictors $x_2^*, ..., x_{j+1}^*$. Denote the predictors in $I_j$ by $x_1^*, ..., x_j^*$. Continue in this manner for $j = p = M, p-1, ..., 2, 1$ where $I_1$ uses $x_1^* = x_1 \equiv 1$.

Several criterion produce the same sequence of models if forward selection or backward elimination are used, including $MSE(I), C_p(I), R_A^2(I), AIC(I),$

$BIC(I)$, and $EBIC(I)$. This result holds since if the number of predictors $k$ in the model $I$ is fixed, the criterion is equivalent to minimizing $SSE(I)$ plus a constant. The constants differ so the model $I_{min}$ that minimizes the criterion often differ. Heuristically, backward elimination tries to delete the variable that will increase $C_p$ the least while forward selection tries to add the variable that will decrease $C_p$ the most.

When there is a sequence of $M$ submodels, the final submodel $I_d$ needs to be selected with $a_d$ terms, including a constant. Let the candidate model $I$ contain $a$ terms, including a constant, and let $\boldsymbol{x}_I$ and $\hat{\boldsymbol{\beta}}_I$ be $a \times 1$ vectors. Then there are many criteria used to select the final submodel $I_d$. For a given data set, the quantities $p, n$, and $\hat{\sigma}^2$ act as constants, and a criterion below may add a constant or be divided by a positive constant without changing the subset $I_{min}$ that minimizes the criterion.

Let criteria $C_S(I)$ have the form

$$C_S(I) = SSE(I) + aK_n\hat{\sigma}^2.$$

These criteria need a good estimator of $\sigma^2$ and $n/p$ large. See Shibata (1984). The criterion $C_p(I) = AIC_S(I)$ uses $K_n = 2$ while the $BIC_S(I)$ criterion uses $K_n = \log(n)$. See Jones (1946) and Mallows (1973) for $C_p$. It can be shown that $C_p(I) = AIC_S(I)$ is equivalent to the $C_P(I)$ criterion of Definition 7.2. Typically $\hat{\sigma}^2$ is the OLS full model $MSE$ when $n/p$ is large.

The following criteria also need $n/p$ large. $AIC$ is due to Akaike (1973), $AIC_C$ is due to Hurvich and Tsai (1989), and $BIC$ to Schwarz (1978) and Akaike (1977, 1978). Also see Burnham and Anderson (2004).

$$AIC(I) = n\log\left(\frac{SSE(I)}{n}\right) + 2a,$$

$$AIC_C(I) = n\log\left(\frac{SSE(I)}{n}\right) + \frac{2a(a+1)}{n-a-1},$$

$$\text{and } BIC(I) = n\log\left(\frac{SSE(I)}{n}\right) + a\log(n).$$

Forward selection with $C_p$ and $AIC$ often gives useful results if $n \geq 5p$ and if the final model has $n \geq 10a_d$. For $p < n < 5p$, forward selection with $C_p$ and AIC tends to pick the full model (which overfits since $n < 5p$) too often, especially if $\hat{\sigma}^2 = MSE$. The Hurvich and Tsai (1989, 1991) $AIC_C$ criterion can be useful if $n \geq \max(2p, 10a_d)$.

The EBIC criterion given in Luo and Chen (2013) may be useful when $n/p$ is not large. Let $0 \leq \gamma \leq 1$ and $|I| = a \leq \min(n, p)$ if $\hat{\boldsymbol{\beta}}_I$ is $a \times 1$. We may use $a \leq \min(n/5, p)$. Then $EBIC(I) =$

$$n\log\left(\frac{SSE(I)}{n}\right) + a\log(n) + 2\gamma\log\left[\binom{p}{a}\right] = BIC(I) + 2\gamma\log\left[\binom{p}{a}\right].$$

This criterion can give good results if $p = p_n = O(n^k)$ and $\gamma > 1 - 1/(2k)$. Hence we will use $\gamma = 1$. Then minimizing $EBIC(I)$ is equivalent to minimizing $BIC(I) - 2\log[(p-a)!] - 2\log(a!)$ since $\log(p!)$ is a constant.

The above criteria can be applied to forward selection and relaxed lasso. The $C_p$ criterion can also be applied to lasso. See Efron and Hastie (2016, pp. 221, 231).

Now suppose $p = 6$ and $S$ in Equation (7.1) corresponds to $x_1 \equiv 1, x_2$, and $x_3$. Suppose the data set is such that underfitting (omitting a predictor in $S$) does not occur. Then there are eight possible submodels that contain $S$: i) $x_1, x_2, x_3$; ii) $x_1, x_2, x_3, x_4$; iii) $x_1, x_2, x_3, x_5$; iv) $x_1, x_2, x_3, x_6$; v) $x_1, x_2, x_3, x_4, x_5$; vi) $x_1, x_2, x_3, x_4, x_6$; vii) $x_1, x_2, x_3, x_5, x_6$; and the full model viii) $x_1, x_2, x_3, x_4, x_5, x_6$. The possible submodel sizes are $k = 3, 4, 5$, or 6. Since the variable selection criteria for forward selection described above minimize the MSE given that $x_1^*, ..., x_{k-1}^*$ are in the model, the $MSE(I_k)$ are too small and underestimate $\sigma^2$. Also the model $I_{min}$ fits the data a bit too well. Suppose $I_{min} = I_d$. Compared to selecting a model $I_k$ before examining the data, the residuals $r_i(I_{min})$ are too small in magnitude, the $|\hat{Y}_{I_{min},i} - Y_i|$ are too small, and $MSE(I_{min})$ is too small. Hence using $I_{min} = I_d$ as the full model for inference does not work. In particular, the partial $F$ test statistic $F_R$ in Theorem 5.7, using $I_d$ as the full model, is too large since the $MSE$ is too small. Thus the partial $F$ test rejects $H_0$ too often. Similarly, the confidence intervals for $\beta_i$ are too short, and hypothesis tests reject $H_0 : \beta_i = 0$ too often when $H_0$ is true. The fact that the selected model $I_{min}$ from variable selection cannot be used as the full model for classical inference is known as **selection bias**. Also see Hurvich and Tsai (1990).

This chapter offers two remedies: i) use the large sample theory of $\hat{\boldsymbol{\beta}}_{VS} = \hat{\boldsymbol{\beta}}_{I_{min},0}$ from Definition 7.3 and the bootstrap for inference after variable selection, and ii) use data splitting for inference after variable selection.

## 7.3 Large Sample Theory for Some Variable Selection Estimators

Large sample theory is often tractable if the optimization problem is convex. The optimization problem for variable selection is not convex, so new tools are needed. Tibshirani et al. (2018) and Leeb and Pötscher (2006, 2008) note that we can not find the limiting distribution of $\boldsymbol{Z}_n = \sqrt{n}\boldsymbol{A}(\hat{\boldsymbol{\beta}}_{I_{min}} - \boldsymbol{\beta}_I)$ after variable selection. One reason is that with positive probability, $\hat{\boldsymbol{\beta}}_{I_{min}}$ does not have the same dimension as $\boldsymbol{\beta}_I$ if AIC or $C_p$ is used. Hence $\boldsymbol{Z}_n$ is not defined with positive probability.

The large sample theory for OLS variable selection estimators, such as forward selection and lasso variable selection, in this section is due to Pelawa Watagoda and Olive (2019, 2020). Rathnayake and Olive (2020) extend this

theory to many other variable selection estimators such as generalized linear models. Charkhi and Claeskens (2018) have a related result for forward selection with AIC when the iid errors are $N(0, \sigma^2)$. Assume $p$ is fixed, and $n \to \infty$. Suppose that model (7.1) holds. Assume the maximum leverage

$$\max_{i=1,\ldots,n} \boldsymbol{x}_{iI_j}^T (\boldsymbol{X}_{I_j}^T \boldsymbol{X}_{I_j})^{-1} \boldsymbol{x}_{iI_j} \to 0$$

in probability as $n \to \infty$ for each $I_j$ with $S \subseteq I_j$ where the dimension of $I_j$ is $a_j$. For the OLS model with $S \subseteq I_j$, $\sqrt{n}(\hat{\boldsymbol{\beta}}_{I_j} - \boldsymbol{\beta}_{I_j}) \xrightarrow{D} N_{a_j}(\boldsymbol{0}, \boldsymbol{V}_j)$ where $\boldsymbol{V}_j = \sigma^2 \boldsymbol{W}_j$ and $(\boldsymbol{X}_{I_j}^T \boldsymbol{X}_{I_j})/n \xrightarrow{P} \boldsymbol{W}_j^{-1}$ by the OLS CLT Theorem 5.9. Then

$$\boldsymbol{u}_{jn} = \sqrt{n}(\hat{\boldsymbol{\beta}}_{I_j,0} - \boldsymbol{\beta}) \xrightarrow{D} \boldsymbol{u}_j \sim N_p(\boldsymbol{0}, \boldsymbol{V}_{j,0}) \tag{7.3}$$

where $\boldsymbol{V}_{j,0}$ adds columns and rows of zeros corresponding to the $x_i$ not in $I_j$, and $\boldsymbol{V}_{j,0}$ is singular unless $I_j$ corresponds to the full model.

For MLR, $\boldsymbol{V}_{j,0} = \sigma^2 \boldsymbol{W}_{j,0}$. For example, if $p = 3$ and model $I_j$ uses a constant $x_1 \equiv 1$ and $x_3$ with

$$\boldsymbol{V}_j = \begin{bmatrix} V_{11} & V_{12} \\ V_{21} & V_{22} \end{bmatrix}, \quad \text{then} \quad \boldsymbol{V}_{j,0} = \begin{bmatrix} V_{11} & 0 & V_{12} \\ 0 & 0 & 0 \\ V_{21} & 0 & V_{22} \end{bmatrix}.$$

Let $I_{min}$ correspond to the set of predictors selected by a variable selection method such as forward selection or lasso variable selection. Use zero padding to form the $p \times 1$ variable selection estimator $\hat{\boldsymbol{\beta}}_{VS}$. For example, if $p = 4$ and $\hat{\boldsymbol{\beta}}_{I_{min}} = (\hat{\beta}_1, \hat{\beta}_3)^T$, then $\hat{\boldsymbol{\beta}}_{VS} = \hat{\boldsymbol{\beta}}_{I_{min},0} = (\hat{\beta}_1, 0, \hat{\beta}_3, 0)^T$. In the following definition, if each subset contains at least one variable, then there are $J = 2^p - 1$ subsets.

**Definition 7.4.** The *variable selection estimator* $\hat{\boldsymbol{\beta}}_{VS} = \hat{\boldsymbol{\beta}}_{I_{min},0}$, and $\hat{\boldsymbol{\beta}}_{VS} = \hat{\boldsymbol{\beta}}_{I_k,0}$ with probabilities $\pi_{kn} = P(I_{min} = I_k)$ for $k = 1, \ldots, J$ where there are $J$ subsets.

**Definition 7.5.** Let $\hat{\boldsymbol{\beta}}_{MIX}$ be a random vector with a mixture distribution of the $\hat{\boldsymbol{\beta}}_{I_k,0}$ with probabilities equal to $\pi_{kn}$. Hence $\hat{\boldsymbol{\beta}}_{MIX} = \hat{\boldsymbol{\beta}}_{I_k,0}$ with same probabilities $\pi_{kn}$ of the variable selection estimator $\hat{\boldsymbol{\beta}}_{VS}$, but the $I_k$ are randomly selected.

The large sample distribution of $\hat{\boldsymbol{\beta}}_{MIX}$ is simpler than that of $\hat{\boldsymbol{\beta}}_{VS}$, and is useful for explaining the large sample distribution of $\hat{\boldsymbol{\beta}}_{VS}$. For how to bootstrap $\hat{\boldsymbol{\beta}}_{MIX}$, see Rathnayake and Olive (2020). For mixture distributions, see Section 11.7.

The first assumption in Theorem 7.3 is $P(S \subseteq I_{min}) \to 1$ as $n \to \infty$. Then the variable selection estimator corresponding to $I_{min}$ underfits with probability going to zero, and the assumption holds under regularity conditions

if BIC or AIC is used. See Charkhi and Claeskens (2018) and Claeskens and Hjort (2008, pp. 70, 101, 102, 114, 232). For multiple linear regression with Mallows (1973) $C_p$ or AIC, see Li (1987), Nishii (1984), and Shao (1993). For a shrinkage estimator that does variable selection, let $\hat{\boldsymbol{\beta}}_{I_{min}}$ be the OLS estimator applied to a constant and the variables with nonzero shrinkage estimator coefficients. If the shrinkage estimator is a consistent estimator of $\boldsymbol{\beta}$, then $P(S \subseteq I_{min}) \to 1$ as $n \to \infty$. See Zhao and Yu (2006, p. 2554). Hence Theorem 7.3c) proves that the lasso variable selection and elastic net variable selection estimators are $\sqrt{n}$ consistent estimators of $\boldsymbol{\beta}$ if lasso and elastic net are consistent. Also see Theorem 7.4 and Remark 7.5. The assumption on $\boldsymbol{u}_{jn}$ in Theorem 7.3 is reasonable by (7.3) since $S \subseteq I_j$ for each $\pi_j$, and since $\hat{\boldsymbol{\beta}}_{MIX}$ uses random selection.

**Theorem 7.3.** Assume $P(S \subseteq I_{min}) \to 1$ as $n \to \infty$, and let $\hat{\boldsymbol{\beta}}_{MIX} = \hat{\boldsymbol{\beta}}_{I_k,0}$ with probabilities $\pi_{kn}$ where $\pi_{kn} \to \pi_k$ as $n \to \infty$. Denote the positive $\pi_k$ by $\pi_j$. Assume $\boldsymbol{u}_{jn} = \sqrt{n}(\hat{\boldsymbol{\beta}}_{I_j,0} - \boldsymbol{\beta}) \xrightarrow{D} \boldsymbol{u}_j \sim N_p(\boldsymbol{0}, \boldsymbol{V}_{j,0})$. a) Then

$$\boldsymbol{u}_n = \sqrt{n}(\hat{\boldsymbol{\beta}}_{MIX} - \boldsymbol{\beta}) \xrightarrow{D} \boldsymbol{u} \qquad (7.4)$$

where the cdf of $\boldsymbol{u}$ is $F_{\boldsymbol{u}}(\boldsymbol{t}) = \sum_j \pi_j F_{\boldsymbol{u}_j}(\boldsymbol{t})$. Thus $\boldsymbol{u}$ has a mixture distribution of the $\boldsymbol{u}_j$ with probabilities $\pi_j$, $E(\boldsymbol{u}) = \boldsymbol{0}$, and $\text{Cov}(\boldsymbol{u}) = \boldsymbol{\Sigma}_{\boldsymbol{u}} = \sum_j \pi_j \boldsymbol{V}_{j,0}$.

b) Let $\boldsymbol{A}$ be a $g \times p$ full rank matrix with $1 \leq g \leq p$. Then

$$\boldsymbol{v}_n = \boldsymbol{A}\boldsymbol{u}_n = \sqrt{n}(\boldsymbol{A}\hat{\boldsymbol{\beta}}_{MIX} - \boldsymbol{A}\boldsymbol{\beta}) \xrightarrow{D} \boldsymbol{A}\boldsymbol{u} = \boldsymbol{v} \qquad (7.5)$$

where $\boldsymbol{v}$ has a mixture distribution of the $\boldsymbol{v}_j = \boldsymbol{A}\boldsymbol{u}_j \sim N_g(\boldsymbol{0}, \boldsymbol{A}\boldsymbol{V}_{j,0}\boldsymbol{A}^T)$ with probabilities $\pi_j$.

c) The estimator $\hat{\boldsymbol{\beta}}_{VS}$ is a $\sqrt{n}$ consistent estimator of $\boldsymbol{\beta}$. Hence $\sqrt{n}(\hat{\boldsymbol{\beta}}_{VS} - \boldsymbol{\beta}) = O_P(1)$.

d) If $\pi_d = 1$, then $\sqrt{n}(\hat{\boldsymbol{\beta}}_{SEL} - \boldsymbol{\beta}) \xrightarrow{D} \boldsymbol{u} \sim N_p(\boldsymbol{0}, \boldsymbol{V}_{d,0})$ where $SEL$ is $VS$ or $MIX$.

**Proof.** a) Since $\boldsymbol{u}_n$ has a mixture distribution of the $\boldsymbol{u}_{kn}$ with probabilities $\pi_{kn}$, the cdf of $\boldsymbol{u}_n$ is $F_{\boldsymbol{u}_n}(\boldsymbol{t}) = \sum_k \pi_{kn} F_{\boldsymbol{u}_{kn}}(\boldsymbol{t}) \to F_{\boldsymbol{u}}(\boldsymbol{t}) = \sum_j \pi_j F_{\boldsymbol{u}_j}(\boldsymbol{t})$ at continuity points of the $F_{\boldsymbol{u}_j}(\boldsymbol{t})$ as $n \to \infty$.

b) Since $\boldsymbol{u}_n \xrightarrow{D} \boldsymbol{u}$, then $\boldsymbol{A}\boldsymbol{u}_n \xrightarrow{D} \boldsymbol{A}\boldsymbol{u}$.

c) The result follows since selecting from a finite number $J$ of $\sqrt{n}$ consistent estimators (even on a set that goes to one in probability) results in a $\sqrt{n}$ consistent estimator by Pratt (1959).

d) If $\pi_d = 1$, there is no selection bias, asymptotically. The result also follows by Pötscher (1991, Lemma 1). $\square$

The following subscript notation is useful. Subscripts before the $MIX$ are used for subsets of $\hat{\boldsymbol{\beta}}_{MIX} = (\hat{\beta}_1, ..., \hat{\beta}_p)^T$. Let $\hat{\boldsymbol{\beta}}_{i,MIX} = \hat{\beta}_i$. Similarly, if $I = \{i_1, ..., i_a\}$, then $\hat{\boldsymbol{\beta}}_{I,MIX} = (\hat{\beta}_{i_1}, ..., \hat{\beta}_{i_a})^T$. Subscripts after $MIX$ denote

the $i$th vector from a sample $\hat{\boldsymbol{\beta}}_{MIX,1}, ..., \hat{\boldsymbol{\beta}}_{MIX,B}$. Similar notation is used for other estimators such as $\hat{\boldsymbol{\beta}}_{VS}$. The subscript 0 is still used for zero padding. We may use $FULL$ to denote the full model $\hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}}_{FULL}$.

Typically the mixture distribution is not asymptotically normal unless a $\pi_d = 1$ (e.g. if $S$ is the full model), or if for each $\pi_j$, $\boldsymbol{Au}_j \sim N_g(\boldsymbol{0}, \boldsymbol{AV}_{j,0}\boldsymbol{A}^T) = N_g(\boldsymbol{0}, \boldsymbol{A\Sigma A}^T)$. Then $\sqrt{n}(\boldsymbol{A}\hat{\boldsymbol{\beta}}_{MIX} - \boldsymbol{A\beta}) \xrightarrow{D} \boldsymbol{Au} \sim N_g(\boldsymbol{0}, \boldsymbol{A\Sigma A}^T)$. This special case occurs for $\hat{\boldsymbol{\beta}}_{S,MIX}$ if $\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \xrightarrow{D} N_p(\boldsymbol{0}, \boldsymbol{V})$ where the asymptotic covariance matrix $\boldsymbol{V}$ is diagonal and nonsingular. Then $\hat{\boldsymbol{\beta}}_{S,MIX}$ and $\hat{\boldsymbol{\beta}}_{S,FULL}$ have the same multivariate normal limiting distribution. For several criteria, this result should hold for $\hat{\boldsymbol{\beta}}_{VS}$ since asymptotically, $\sqrt{n}(\boldsymbol{A}\hat{\boldsymbol{\beta}}_{VS} - \boldsymbol{A\beta})$ is selecting from the $\boldsymbol{Au}_j$ which have the same distribution. Then the confidence regions applied to $\boldsymbol{A}\hat{\boldsymbol{\beta}}_{SEL}^* = \boldsymbol{B}\hat{\boldsymbol{\beta}}_{S,SEL}^*$ should have similar volume and cutoffs where $SEL$ is $MIX$, $VS$, or $FULL$.

Theorem 7.3 can be used to justify prediction intervals after variable selection. See Pelawa Watagoda and Olive (2020). Theorem 7.3d) is useful for *variable selection consistency* and the *oracle property* where $\pi_d = \pi_S = 1$ if $P(I_{min} = S) \to 1$ as $n \to \infty$. See Claeskens and Hjort (2008, pp. 101-114) and Fan and Li (2001) for references. A necessary condition for $P(I_{min} = S) \to 1$ is that $S$ is one of the models considered with probability going to one. This condition holds under strong regularity conditions for fast methods. See Wieczorek (2018) for forward selection and Hastie et al. (2015, pp. 295-302) for lasso, where the predictors need a "near orthogonality" condition.

**Remark 7.4.** If $A_1, A_2, ..., A_k$ are pairwise disjoint and if $\cup_{i=1}^k A_i = S$, then the collection of sets $A_1, A_2, ..., A_k$ is a *partition* of $S$. Then the *Law of Total Probability* states that if $A_1, A_2, ..., A_k$ form a partition of $S$ such that $P(A_i) > 0$ for $i = 1, ..., k$, then

$$P(B) = \sum_{j=1}^k P(B \cap A_j) = \sum_{j=1}^k P(B|A_j)P(A_j).$$

Let sets $A_{k+1}, ..., A_m$ satisfy $P(A_i) = 0$ for $i = k+1, ..., m$. Define $P(B|A_j) = 0$ if $P(A_j = 0$. Then a Generalized Law of Total Probability is

$$P(B) = \sum_{j=1}^m P(B \cap A_j) = \sum_{j=1}^m P(B|A_j)P(A_j),$$

and will be used in the following paragraph.

Pötscher (1991) used the conditional distribution of $\hat{\boldsymbol{\beta}}_{VS}|(\hat{\boldsymbol{\beta}}_{VS} = \hat{\boldsymbol{\beta}}_{I_k,0})$ to find the distribution of $\boldsymbol{w}_n = \sqrt{n}(\hat{\boldsymbol{\beta}}_{VS} - \boldsymbol{\beta})$. Let $W = W_{VS} = k$ if $\hat{\boldsymbol{\beta}}_{VS} = \hat{\boldsymbol{\beta}}_{I_k,0}$ where $P(W_{VS} = k) = \pi_{kn}$ for $k = 1, ..., J$. Then $(\hat{\boldsymbol{\beta}}_{VS:n}, W_{VS:n}) = (\hat{\boldsymbol{\beta}}_{VS}, W_{VS})$ has a joint distribution where the sample size $n$ is usually suppressed. Note that $\hat{\boldsymbol{\beta}}_{VS} = \hat{\boldsymbol{\beta}}_{I_W,0}$. Define $P(B|A_k)P(A_k) = 0$ if $P(A_k) = 0$.

Let $\hat{\boldsymbol{\beta}}_{I_k,0}^C$ be a random vector from the conditional distribution $\hat{\boldsymbol{\beta}}_{I_k,0}|(W_{VS} = k)$. Let $\boldsymbol{w}_{kn} = \sqrt{n}(\hat{\boldsymbol{\beta}}_{I_k,0} - \boldsymbol{\beta})|(W_{VS} = k) \sim \sqrt{n}(\hat{\boldsymbol{\beta}}_{I_k,0}^C - \boldsymbol{\beta})$. Denote $F_{\boldsymbol{z}}(\boldsymbol{t}) = P(z_1 \leq t_1, ..., z_p \leq t_p)$ by $P(\boldsymbol{z} \leq \boldsymbol{t})$. Then

$$F_{\boldsymbol{w}_n}(\boldsymbol{t}) = P[n^{1/2}(\hat{\boldsymbol{\beta}}_{VS} - \boldsymbol{\beta}) \leq \boldsymbol{t}] =$$

$$\sum_{k=1}^J P[n^{1/2}(\hat{\boldsymbol{\beta}}_{VS} - \boldsymbol{\beta}) \leq \boldsymbol{t}|(\hat{\boldsymbol{\beta}}_{VS} = \hat{\boldsymbol{\beta}}_{I_k,0})]P(\hat{\boldsymbol{\beta}}_{VS} = \hat{\boldsymbol{\beta}}_{I_k,0}) =$$

$$\sum_{k=1}^J P[n^{1/2}(\hat{\boldsymbol{\beta}}_{I_k,0} - \boldsymbol{\beta}) \leq \boldsymbol{t}|(\hat{\boldsymbol{\beta}}_{VS} = \hat{\boldsymbol{\beta}}_{I_k,0})]\pi_{kn}$$

$$= \sum_{k=1}^J P[n^{1/2}(\hat{\boldsymbol{\beta}}_{I_k,0}^C - \boldsymbol{\beta}) \leq \boldsymbol{t}]\pi_{kn} = \sum_{k=1}^J F_{\boldsymbol{w}_{kn}}(\boldsymbol{t})\pi_{kn}.$$

Hence $\hat{\boldsymbol{\beta}}_{VS}$ has a mixture distribution of the $\hat{\boldsymbol{\beta}}_{I_k,0}^C$ with probabilities $\pi_{kn}$, and $\boldsymbol{w}_n$ has a mixture distribution of the $\boldsymbol{w}_{kn}$ with probabilities $\pi_{kn}$.

Charkhi and Claeskens (2018) showed that $\boldsymbol{w}_{jn} = \sqrt{n}(\hat{\boldsymbol{\beta}}_{I_j,0}^C - \boldsymbol{\beta}) \xrightarrow{D} \boldsymbol{w}_j$ if $S \subseteq I_j$ for the MLE with AIC. Here $\boldsymbol{w}_j$ is a multivariate truncated normal distribution (where no truncation is possible) that is symmetric about $\boldsymbol{0}$. Hence $E(\boldsymbol{w}_j) = 0$, and $\text{Cov}(\boldsymbol{w}_j) = \boldsymbol{\Sigma}_j$ exits. Referring to Definitions 7.3 and 7.4, note that both $\sqrt{n}(\hat{\boldsymbol{\beta}}_{MIX} - \boldsymbol{\beta})$ and $\sqrt{n}(\hat{\boldsymbol{\beta}}_{VS} - \boldsymbol{\beta})$ are selecting from the $\boldsymbol{u}_{kn} = \sqrt{n}(\hat{\boldsymbol{\beta}}_{I_k,0} - \boldsymbol{\beta})$ and asymptotically from the $\boldsymbol{u}_j$ of Equation (7.3). The random selection for $\hat{\boldsymbol{\beta}}_{MIX}$ does not change the distribution of $\boldsymbol{u}_{jn}$, but selection bias does change the distribution of the selected $\boldsymbol{u}_{jn}$ to that of $\boldsymbol{w}_{jn}$. Similarly, selection bias does change the distribution of the selected $\boldsymbol{u}_j$ to that of $\boldsymbol{w}_j$. The reasonable Theorem 7.4 assumption that $\boldsymbol{w}_{jn} \xrightarrow{D} \boldsymbol{w}_j$ may not be mild.

**Theorem 7.4, Variable Selection CLT.** Assume $P(S \subseteq I_{min}) \to 1$ as $n \to \infty$, and let $\hat{\boldsymbol{\beta}}_{VS} = \hat{\boldsymbol{\beta}}_{I_k,0}$ with probabilities $\pi_{kn}$ where $\pi_{kn} \to \pi_k$ as $n \to \infty$. Denote the positive $\pi_k$ by $\pi_j$. Assume $\boldsymbol{w}_{jn} = \sqrt{n}(\hat{\boldsymbol{\beta}}_{I_j,0}^C - \boldsymbol{\beta}) \xrightarrow{D} \boldsymbol{w}_j$. Then

$$\boldsymbol{w}_n = \sqrt{n}(\hat{\boldsymbol{\beta}}_{VS} - \boldsymbol{\beta}) \xrightarrow{D} \boldsymbol{w} \qquad (7.6)$$

where the cdf of $\boldsymbol{w}$ is $F_{\boldsymbol{w}}(\boldsymbol{t}) = \sum_j \pi_j F_{\boldsymbol{w}_j}(\boldsymbol{t})$. Thus $\boldsymbol{w}$ is a mixture distribution of the $\boldsymbol{w}_j$ with probabilities $\pi_j$.

**Proof.** Since $\boldsymbol{w}_n$ has a mixture distribution of the $\boldsymbol{w}_{kn}$ with probabilities $\pi_{kn}$, the cdf of $\boldsymbol{w}_n$ is $F_{\boldsymbol{w}_n}(\boldsymbol{t}) = \sum_k \pi_{kn} F_{\boldsymbol{w}_{kn}}(\boldsymbol{t}) \to F_{\boldsymbol{w}}(\boldsymbol{t}) = \sum_j \pi_j F_{\boldsymbol{w}_j}(\boldsymbol{t})$ at continuity points of the $F_{\boldsymbol{w}_j}(\boldsymbol{t})$ as $n \to \infty$. $\square$

**Remark 7.5.** If $P(S \subseteq I_{min}) \to 1$ as $n \to \infty$, then $\hat{\boldsymbol{\beta}}_{VS}$ is a $\sqrt{n}$ consistent estimator of $\boldsymbol{\beta}$ since selecting from a finite number $J$ of $\sqrt{n}$ consistent estima-

tors (even on a set that goes to one in probability) results in a $\sqrt{n}$ consistent estimator by Pratt (1959). By both this result and Theorems 7.3 and 7.4, the lasso variable selection and elastic net variable selection estimators are $\sqrt{n}$ consistent if lasso and elastic net are consistent.

## 7.4 Bootstrapping Variable Selection

This section considers bootstrapping the MLR variable selection model. Rathnayake and Olive (2020) shows how to bootstrap variable selection for many other regression models. This section will explain why the bootstrap confidence regions (4.13), (4.14), and (4.15) give useful results. Much of the theory in Section 4.3 does not apply to the variable selection estimator $T_n = \boldsymbol{A}\hat{\boldsymbol{\beta}}_{I_{min},0}$ with $\boldsymbol{\theta} = \boldsymbol{A}\boldsymbol{\beta}$, because $T_n$ is not smooth since $T_n$ is equal to the estimator $T_{jn}$ with probability $\pi_{jn}$ for $j = 1, ..., J$. Here $\boldsymbol{A}$ is a known full rank $g \times p$ matrix with $1 \leq g \leq p$.

Obtaining the bootstrap samples for $\hat{\boldsymbol{\beta}}_{VS}$ and $\hat{\boldsymbol{\beta}}_{MIX}$ is simple. Generate $\boldsymbol{Y}^*$ and $\boldsymbol{X}^*$ that would be used to produce $\hat{\boldsymbol{\beta}}^*$ if the full model estimator $\hat{\boldsymbol{\beta}}$ was being bootstrapped. Instead of computing $\hat{\boldsymbol{\beta}}^*$, compute the variable selection estimator $\hat{\boldsymbol{\beta}}_{VS,1}^* = \hat{\boldsymbol{\beta}}_{I_{k_1},0}^{*C}$. Then generate another $\boldsymbol{Y}^*$ and $\boldsymbol{X}^*$ and compute $\hat{\boldsymbol{\beta}}_{MIX,1}^* = \hat{\boldsymbol{\beta}}_{I_{k_1},0}^*$ (using the same subset $I_{k_1}$). This process is repeated $B$ times to get the two bootstrap samples for $i = 1, ..., B$. Let the selection probabilities for the bootstrap variable selection estimator be $\rho_{kn}$. Then this bootstrap procedure bootstraps both $\hat{\boldsymbol{\beta}}_{VS}$ and $\hat{\boldsymbol{\beta}}_{MIX}$ with $\pi_{kn} = \rho_{kn}$.

The key idea is to show that the bootstrap data cloud is slightly more variable than the iid data cloud, so confidence region (4.14) applied to the bootstrap data cloud has coverage bounded below by $(1 - \delta)$ for large enough $n$ and $B$.

For the bootstrap, suppose that $T_i^*$ is equal to $T_{ij}^*$ with probability $\rho_{jn}$ for $j = 1, ..., J$ where $\sum_j \rho_{jn} = 1$, and $\rho_{jn} \to \pi_j$ as $n \to \infty$. Let $B_{jn}$ count the number of times $T_i^* = T_{ij}^*$ in the bootstrap sample. Then the bootstrap sample $T_1^*, ..., T_B^*$ can be written as

$$T_{1,1}^*, ..., T_{B_{1n},1}^*, ..., T_{1,J}^*, ..., T_{B_{Jn},J}^*$$

where the $B_{jn}$ follow a multinomial distribution and $B_{jn}/B \xrightarrow{P} \rho_{jn}$ as $B \to \infty$. Denote $T_{1j}^*, ..., T_{B_{jn},j}^*$ as the $j$th bootstrap component of the bootstrap sample with sample mean $\overline{T}_j^*$ and sample covariance matrix $\boldsymbol{S}_{T,j}^*$. Then

$$\overline{T}^* = \frac{1}{B} \sum_{i=1}^B T_i^* = \sum_j \frac{B_{jn}}{B} \frac{1}{B_{jn}} \sum_{i=1}^{B_{jn}} T_{ij}^* = \sum_j \hat{\rho}_{jn} \overline{T}_j^*.$$

Similarly, we can define the $j$th component of the iid sample $T_1, ..., T_B$ to have sample mean $\overline{T}_j$ and sample covariance matrix $\boldsymbol{S}_{T,j}$.

Let $T_n = \hat{\boldsymbol{\beta}}_{MIX}$ and $T_{ij} = \hat{\boldsymbol{\beta}}_{I_j,0}$. If $S \subseteq I_j$, assume $\sqrt{n}(\hat{\boldsymbol{\beta}}_{I_j} - \boldsymbol{\beta}_{I_j}) \xrightarrow{D} N_{a_j}(\boldsymbol{0}, \boldsymbol{V}_j)$ and $\sqrt{n}(\hat{\boldsymbol{\beta}}_{I_j}^* - \hat{\boldsymbol{\beta}}_{I_j}) \xrightarrow{D} N_{a_j}(\boldsymbol{0}, \boldsymbol{V}_j)$. Then by Equation (7.3),

$$\sqrt{n}(\hat{\boldsymbol{\beta}}_{I_j,0} - \boldsymbol{\beta}) \xrightarrow{D} N_p(\boldsymbol{0}, \boldsymbol{V}_{j,0}) \text{ and } \sqrt{\text{n}}(\hat{\boldsymbol{\beta}}_{I_j,0}^* - \hat{\boldsymbol{\beta}}_{I_j,0}) \xrightarrow{D} N_p(\boldsymbol{0}, \boldsymbol{V}_{j,0}). \quad (7.7)$$

This result means that the component clouds have the same variability asymptotically. The iid data component clouds are all centered at $\boldsymbol{\beta}$. If the bootstrap data component clouds were all centered at the same value $\tilde{\boldsymbol{\beta}}$, then the bootstrap cloud would be like an iid data cloud shifted to be centered at $\tilde{\boldsymbol{\beta}}$, and (4.14) would be a confidence region for $\boldsymbol{\theta} = \boldsymbol{\beta}$. Instead, the bootstrap data component clouds are shifted slightly from a common center, and are each centered at a $\hat{\boldsymbol{\beta}}_{I_j,0}$. Geometrically, the shifting of the bootstrap component data clouds makes the bootstrap data cloud similar but more variable than the iid data cloud asymptotically (we want $n \geq 20p$), and centering the bootstrap data cloud at $T_n$ results in the confidence region (4.14) having slightly higher asymptotic coverage than applying (4.14) to the iid data cloud. Also, (4.14) tends to have higher coverage than (4.15) since the cutoff for (4.14) tends to be larger than the cutoff for (4.15). Region (4.13) has the same volume as region (4.15), but tends to have higher coverage since empirically, the bagging estimator $\overline{T}^*$ tends to estimate $\boldsymbol{\theta}$ at least as well as $T_n$ for a mixture distribution. A similar argument holds if $T_n = \boldsymbol{A}\hat{\boldsymbol{\beta}}_{MIX}$, $T_{ij} = \boldsymbol{A}\hat{\boldsymbol{\beta}}_{I_j,0}$, and $\boldsymbol{\theta} = \boldsymbol{A}\boldsymbol{\beta}$.

To see that $T^*$ has more variability than $T_n$, asymptotically, look at Figure 3.1. Imagine that $n$ is huge and the $J = 6$ ellipsoids are 99.9% covering regions for the component data clouds corresponding to $T_{jn}$ for $j = 1, ..., J$. Separating the clouds slightly, without rotation, increases the variability of the overall data cloud. The bootstrap distribution of $T^*$ corresponds to the separated clouds. The shape of the overall data cloud does not change much, but the volume does increase.

**Remark 7.6.** Note that there are several important variable selection models, including the model given by Equation (7.1) where $\boldsymbol{x}^T\boldsymbol{\beta} = \boldsymbol{x}_S^T\boldsymbol{\beta}_S$. Another model is $\boldsymbol{x}^T\boldsymbol{\beta} = \boldsymbol{x}_{S_i}^T\boldsymbol{\beta}_{S_i}$ for $i = 1, ..., K$. Then there are $K \geq 2$ competing "true" nonnested submodels where $\boldsymbol{\beta}_{S_i}$ is $a_{S_i} \times 1$. For example, suppose the $K = 2$ models have predictors $x_1, x_2, x_3$ for $S_1$ and $x_1, x_2, x_4$ for $S_2$. Then $x_3$ and $x_4$ are likely to be selected and omitted often by forward selection for the $B$ bootstrap samples. Hence omitting all predictors $x_i$ that have a $\beta_{ij}^* = 0$ for at least one of the bootstrap samples $j = 1, ..., B$ could result in underfitting, e.g. using just $x_1$ and $x_2$ in the above $K = 2$ example. Theorems 7.3 and 7.4 still hold if "$P(S \subseteq I_{min}) \to 1$" is replaced by "$P(S_i \subseteq I_{min}$ for some $i) \to 1$," and the bootstrap sample is still more variable than the iid sample.

In the simulations for $H_0 : \boldsymbol{A\beta} = \boldsymbol{B\beta}_S = \boldsymbol{\theta}_0$ with $n \geq 20p$, the coverage tended to get close to $1 - \delta$ for $B \geq \max(200, 50p)$ so that $\boldsymbol{S}_T^*$ is a good estimator of $\text{Cov}(T^*)$. In the simulations where $S$ is not the full model, inference with backward elimination with $I_{min}$ using $AIC$ was often more precise than inference with the full model if $n \geq 20p$ and $B \geq 50p$.

The matrix $\boldsymbol{S}_T^*$ can be singular due to one or more columns of zeros in the bootstrap sample for $\beta_1, ..., \beta_p$. The variables corresponding to these columns are likely not needed in the model given that the other predictors are in the model. A simple remedy is to add $d$ bootstrap samples of the full model estimator $\hat{\boldsymbol{\beta}}^* = \hat{\boldsymbol{\beta}}_{FULL}^*$ to the bootstrap sample. For example, take $d = \lceil cB \rceil$ with $c = 0.01$. A confidence interval $[L_n, U_n]$ can be computed without $\boldsymbol{S}_T^*$ for (4.13), (4.14), and (4.15). Using the confidence interval $[\max(L_n, T_{(1)}^*), \min(U_n, T_{(B)}^*)]$ can give a shorter covering region.

Undercoverage can occur if bootstrap sample data cloud is less variable than the iid data cloud, e.g., if $(n - p)/n$ is not close to one. Coverage can be higher than the nominal coverage for two reasons: i) the bootstrap data cloud is more variable than the iid data cloud of $T_1, ..., T_B$, and ii) zero padding.

The bootstrap component clouds for $\hat{\boldsymbol{\beta}}_{VS}^*$ are again separated compared to the iid clouds for $\hat{\boldsymbol{\beta}}_{VS}$, which are centered about $\boldsymbol{\beta}$. Heuristically, most of the selection bias is due to predictors in $E$, not to the predictors in $S$. Hence $\hat{\boldsymbol{\beta}}_{S,VS}^*$ is roughly similar to $\hat{\boldsymbol{\beta}}_{S,MIX}^*$. Typically the distributions of $\hat{\boldsymbol{\beta}}_{E,VS}^*$ and $\hat{\boldsymbol{\beta}}_{E,MIX}^*$ are not similar, but use the same zero padding. In simulations, confidence regions for $\hat{\boldsymbol{\beta}}_{VS}$ tended to have less undercoverage than confidence regions for $\hat{\boldsymbol{\beta}}_{MIX}^*$.

### *7.4.1* **The Parametric Bootstrap**

For the multiple linear regression model, $\boldsymbol{Y} = \boldsymbol{X\beta} + \boldsymbol{e}$, assume a constant $x_1$ is in the model, and the zero mean $e_i$ are iid with variance $V(e_i) = \sigma^2$. Let $\boldsymbol{H} = \boldsymbol{X}(\boldsymbol{X}^T\boldsymbol{X})^{-1}\boldsymbol{X}^T$. For each $I$ with $S \subseteq I$, assume the maximum leverage $\max_{i=1,...,n} \boldsymbol{x}_{iI}^T(\boldsymbol{X}_I^T\boldsymbol{X}_I)^{-1}\boldsymbol{x}_{iI} \rightarrow 0$ in probability as $n \rightarrow \infty$. For OLS with $S \subseteq I$, $\sqrt{n}(\hat{\boldsymbol{\beta}}_I - \boldsymbol{\beta}_I) \xrightarrow{D} N_{a_I}(\boldsymbol{0}, \boldsymbol{V}_I)$ by Equation (7.3).

The parametric bootstrap generates $\boldsymbol{Y}_j^* = (Y_i^*)$ from a parametric distribution. Then regress $\boldsymbol{Y}_j^*$ on $\boldsymbol{X}$ to get $\hat{\boldsymbol{\beta}}_j^*$ for $j = 1, ..., B$. Consider the parametric bootstrap for the MLR model with $\boldsymbol{Y}^* \sim N_n(\boldsymbol{X}\hat{\boldsymbol{\beta}}, \hat{\sigma}_n^2\boldsymbol{I}) \sim N_n(\boldsymbol{HY}, \hat{\sigma}_n^2\boldsymbol{I})$ where **we are not assuming** that the $e_i \sim N(0, \sigma^2)$, and

$$\hat{\sigma}_n^2 = MSE = \frac{1}{n - p} \sum_{i=1}^{n} r_i^2$$

where the residuals are from the full OLS model. Then $MSE$ is a $\sqrt{n}$ consistent estimator of $\sigma^2$ under mild conditions by Su and Cook (2012). Hence

$$\boldsymbol{Y}^* = \boldsymbol{X}\hat{\boldsymbol{\beta}}_{OLS} + \boldsymbol{e}^*$$

where the $e_i^*$ are iid $N(0, MSE)$ and $\hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}}_{OLS}$.

Thus $\hat{\boldsymbol{\beta}}_I^* = (\boldsymbol{X}_I^T \boldsymbol{X}_I)^{-1} \boldsymbol{X}_I^T \boldsymbol{Y}^* \sim N_{a_I}(\hat{\boldsymbol{\beta}}_I, \hat{\sigma}_n^2 (\boldsymbol{X}_I^T \boldsymbol{X}_I)^{-1})$ since $E(\hat{\boldsymbol{\beta}}_I^*) = (\boldsymbol{X}_I^T \boldsymbol{X}_I)^{-1} \boldsymbol{X}_I^T \boldsymbol{H} \boldsymbol{Y} = \hat{\boldsymbol{\beta}}_I$ because $\boldsymbol{H} \boldsymbol{X}_I = \boldsymbol{X}_I$, and $\text{Cov}(\hat{\boldsymbol{\beta}}_I^*) = \hat{\sigma}_n^2 (\boldsymbol{X}_I^T \boldsymbol{X}_I)^{-1}$. Hence

$$\sqrt{n}(\hat{\boldsymbol{\beta}}_I^* - \hat{\boldsymbol{\beta}}_I) \sim N_{a_I}(\boldsymbol{0}, n\hat{\sigma}_n^2 (\boldsymbol{X}_I^T \boldsymbol{X}_I)^{-1}) \xrightarrow{D} N_{a_I}(\boldsymbol{0}, \boldsymbol{V}_I)$$

as $n, B \to \infty$ if $S \subseteq I$.

## 7.4.2 The Residual Bootstrap

The *residual bootstrap* is often useful for additive error regression models of the form $Y_i = m(\boldsymbol{x}_i) + e_i = \hat{m}(\boldsymbol{x}_i) + r_i = \hat{Y}_i + r_i$ for $i = 1, ..., n$ where the $i$th residual $r_i = Y_i - \hat{Y}_i$. Let $\boldsymbol{Y} = (Y_1, ..., Y_n)^T$, $\boldsymbol{r} = (r_1, ..., r_n)^T$, and let $\boldsymbol{X}$ be an $n \times p$ matrix with $i$th row $\boldsymbol{x}_i^T$. Then the fitted values $\hat{Y}_i = \hat{m}(\boldsymbol{x}_i)$, and the residuals are obtained by regressing $\boldsymbol{Y}$ on $\boldsymbol{X}$. Here the errors $e_i$ are iid, and it would be useful to be able to generate $B$ iid samples $e_{1j}, ..., e_{nj}$ from the distribution of $e_i$ where $j = 1, ..., B$. If the $m(\boldsymbol{x}_i)$ were known, then we could form a vector $\boldsymbol{Y}_j$ where the $i$th element $Y_{ij} = m(\boldsymbol{x}_i) + e_{ij}$ for $i = 1, ..., n$. Then regress $\boldsymbol{Y}_j$ on $\boldsymbol{X}$. Instead, draw samples $r_{1j}^*, ..., r_{nj}^*$ with replacement from the residuals, then form a vector $\boldsymbol{Y}_j^*$ where the $i$th element $Y_{ij}^* = \hat{m}(\boldsymbol{x}_i) + r_{ij}^*$ for $i = 1, ..., n$. Then regress $\boldsymbol{Y}_j^*$ on $\boldsymbol{X}$. If the residuals do not sum to 0, replace $r_i$ by $\epsilon_i = r_i - \overline{r}$, and $r_{ij}^*$ by $\epsilon_{ij}^*$.

**Example 7.1.** For multiple linear regression, $Y_i = \boldsymbol{x}_i^T \boldsymbol{\beta} + e_i$ is written in matrix form as $\boldsymbol{Y} = \boldsymbol{X}\boldsymbol{\beta} + \boldsymbol{e}$. Regress $\boldsymbol{Y}$ on $\boldsymbol{X}$ to obtain $\hat{\boldsymbol{\beta}}$, $\boldsymbol{r}$, and $\hat{\boldsymbol{Y}}$ with $i$th element $\hat{Y}_i = \hat{m}(\boldsymbol{x}_i) = \boldsymbol{x}_i^T \hat{\boldsymbol{\beta}}$. For $j = 1, ..., B$, regress $\boldsymbol{Y}_j^*$ on $\boldsymbol{X}$ to form $\hat{\boldsymbol{\beta}}_{1,n}^*, ..., \hat{\boldsymbol{\beta}}_{B,n}^*$ using the residual bootstrap.

Now examine the OLS model. Let $\hat{\boldsymbol{Y}} = \hat{\boldsymbol{Y}}_{OLS} = \boldsymbol{X}\hat{\boldsymbol{\beta}}_{OLS} = \boldsymbol{H}\boldsymbol{Y}$ be the fitted values from the OLS full model. Let $\boldsymbol{r}^W$ denote an $n \times 1$ random vector of elements selected with replacement from the OLS full model residuals. Following Freedman (1981) and Efron (1982, p. 36),

$$\boldsymbol{Y}^* = \boldsymbol{X}\hat{\boldsymbol{\beta}}_{OLS} + \boldsymbol{r}^W$$

follows a standard linear model where the elements $r_i^W$ of $\boldsymbol{r}^W$ are iid from the empirical distribution of the OLS full model residuals $r_i$. Hence

$$E(r_i^W) = \frac{1}{n}\sum_{i=1}^{n} r_i = 0, \quad V(r_i^W) = \sigma_n^2 = \frac{1}{n}\sum_{i=1}^{n} r_i^2 = \frac{n-p}{n}MSE,$$

$$E(\boldsymbol{r}^W) = \boldsymbol{0}, \text{ and } \text{Cov}(\boldsymbol{Y}^*) = \text{Cov}(\boldsymbol{r}^W) = \sigma_n^2 \boldsymbol{I}_n.$$

Let $\hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}}_{OLS}$. Then $\hat{\boldsymbol{\beta}}^* = (\boldsymbol{X}^T\boldsymbol{X})^{-1}\boldsymbol{X}^T\boldsymbol{Y}^*$ with $\text{Cov}(\hat{\boldsymbol{\beta}}^*) = \sigma_n^2(\boldsymbol{X}^T\boldsymbol{X})^{-1} = \frac{n-p}{n}MSE(\boldsymbol{X}^T\boldsymbol{X})^{-1}$, and $E(\hat{\boldsymbol{\beta}}^*) = (\boldsymbol{X}^T\boldsymbol{X})^{-1}\boldsymbol{X}^T E(\boldsymbol{Y}^*) = (\boldsymbol{X}^T\boldsymbol{X})^{-1}\boldsymbol{X}^T\boldsymbol{H}\boldsymbol{Y} = \hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}}_n$ since $\boldsymbol{H}\boldsymbol{X} = \boldsymbol{X}$. The expectations are with respect to the bootstrap distribution where $\hat{\boldsymbol{Y}}$ acts as a constant.

For the OLS estimator $\hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}}_{OLS}$, the estimated covariance matrix of $\hat{\boldsymbol{\beta}}_{OLS}$ is $\widehat{\text{Cov}}(\hat{\boldsymbol{\beta}}_{OLS}) = MSE(\boldsymbol{X}^T\boldsymbol{X})^{-1}$. The sample covariance matrix of the $\hat{\boldsymbol{\beta}}^*$ is estimating $\text{Cov}(\hat{\boldsymbol{\beta}}^*)$ as $B \to \infty$. Hence the residual bootstrap standard error $SE(\hat{\beta}_i^*) \approx \sqrt{\frac{n-p}{n}} \, SE(\hat{\beta}_i)$ for $i = 1,...,p$ where $\hat{\boldsymbol{\beta}}_{OLS} = \hat{\boldsymbol{\beta}} = (\hat{\beta}_1,...,\hat{\beta}_p)^T$. The OLS CLT Theorem 5.9 says

$$\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \xrightarrow{D} N_p(\boldsymbol{0}, \lim_{n\to\infty} n\widehat{\text{Cov}}(\hat{\boldsymbol{\beta}}_{OLS})) \sim N_p(\boldsymbol{0}, \sigma^2\boldsymbol{W})$$

where $n(\boldsymbol{X}^T\boldsymbol{X})^{-1} \to \boldsymbol{W}$. Since $\boldsymbol{Y}^* = \boldsymbol{X}\hat{\boldsymbol{\beta}}_{OLS} + \boldsymbol{r}^W$ follows a standard linear model, it may not be surprising that

$$\sqrt{n}(\hat{\boldsymbol{\beta}}^* - \hat{\boldsymbol{\beta}}_{OLS}) \xrightarrow{D} N_p(\boldsymbol{0}, \lim_{n\to\infty} n\widehat{\text{Cov}}(\hat{\boldsymbol{\beta}}^*)) \sim N_p(\boldsymbol{0}, \sigma^2\boldsymbol{W}).$$

See Freedman (1981).

For the above residual bootstrap, $\hat{\boldsymbol{\beta}}_{I_j}^* = (\boldsymbol{X}_{I_j}^T\boldsymbol{X}_{I_j})^{-1}\boldsymbol{X}_{I_j}^T\boldsymbol{Y}^* = \boldsymbol{D}_j\boldsymbol{Y}^*$ with $\text{Cov}(\hat{\boldsymbol{\beta}}_{I_j}^*) = \sigma_n^2(\boldsymbol{X}_{I_j}^T\boldsymbol{X}_{I_j})^{-1}$ and $E(\hat{\boldsymbol{\beta}}_{I_j}^*) = (\boldsymbol{X}_{I_j}^T\boldsymbol{X}_{I_j})^{-1}\boldsymbol{X}_{I_j}^T E(\boldsymbol{Y}^*) = (\boldsymbol{X}_{I_j}^T\boldsymbol{X}_{I_j})^{-1}\boldsymbol{X}_{I_j}^T\boldsymbol{H}\boldsymbol{Y} = \hat{\boldsymbol{\beta}}_{I_j}$ since $\boldsymbol{H}\boldsymbol{X}_{I_j} = \boldsymbol{X}_{I_j}$. The expectations are with respect to the bootstrap distribution where $\hat{\boldsymbol{Y}}$ acts as a constant.

Thus for $S \subseteq I$ and the residual bootstrap using residuals from the full OLS model, $E(\hat{\boldsymbol{\beta}}_I^*) = \hat{\boldsymbol{\beta}}_I$ and $n\text{Cov}(\hat{\boldsymbol{\beta}}_I^*) = n[(n-p)/n]\hat{\sigma}_n^2(\boldsymbol{X}_I^T\boldsymbol{X}_I)^{-1} \xrightarrow{P} \boldsymbol{V}_I$ as $n \to \infty$ with $\hat{\sigma}_n^2 = MSE$. Hence $\hat{\boldsymbol{\beta}}_I^* - \hat{\boldsymbol{\beta}}_I \xrightarrow{P} \boldsymbol{0}$ as $n \to \infty$ by Lai et al (1979). Note that $\hat{\boldsymbol{\beta}}_I^* = \hat{\boldsymbol{\beta}}_{I,n}^*$ and $\hat{\boldsymbol{\beta}}_I = \hat{\boldsymbol{\beta}}_{I,n}$ depend on $n$.

**Remark 7.7.** The Cauchy Schwartz inequality says $|\boldsymbol{a}^T\boldsymbol{b}| \leq \|\boldsymbol{a}\| \, \|\boldsymbol{b}\|$. Suppose $\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) = O_P(1)$ is bounded in probability. This will occur if $\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \xrightarrow{D} N_p(\boldsymbol{0}, \boldsymbol{\Sigma})$, e.g. if $\hat{\boldsymbol{\beta}}$ is the OLS estimator. Then

$$|r_i - e_i| = |Y_i - \boldsymbol{x}_i^T\hat{\boldsymbol{\beta}} - (Y_i - \boldsymbol{x}_i^T\boldsymbol{\beta})| = |\boldsymbol{x}_i^T(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})|.$$

Hence

$$\sqrt{n} \max_{i=1,\dots,n} |r_i - e_i| \leq (\max_{i=1,\dots,n} \|\boldsymbol{x}_i\|) \ \|\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})\| = O_P(1)$$

since $\max \|\boldsymbol{x}_i\| = O_P(1)$ or there is extrapolation. Hence OLS residuals behave well if the zero mean error distribution of the iid $e_i$ has a finite variance $\sigma^2$.

**Remark 7.8.** Note that both the residual bootstrap and parametric bootstrap for OLS are robust to the unknown error distribution of the iid $e_i$. For the residual bootstrap with $S \subseteq I$ where $I$ is not the full model, it may not be true that $\sqrt{n}(\hat{\boldsymbol{\beta}}_I^* - \hat{\boldsymbol{\beta}}_I) \xrightarrow{D} N_{a_I}(\boldsymbol{0}, \boldsymbol{V}_I)$ as $n, B \to \infty$. For the model $\boldsymbol{Y} = \boldsymbol{X}\boldsymbol{\beta} + \boldsymbol{e}$, the $e_i$ are iid from a distribution that does not depend on $n$, and $\boldsymbol{\beta}_E = \boldsymbol{0}$. For $\boldsymbol{Y}^* = \boldsymbol{X}\hat{\boldsymbol{\beta}} + \boldsymbol{r}^W$, the distribution of the $r_i^W$ depends on $n$ and $\hat{\boldsymbol{\beta}}_E \neq \boldsymbol{0}$ although $\sqrt{n}\hat{\boldsymbol{\beta}}_E = O_P(1)$.

## *7.4.3* The Nonparametric Bootstrap

The nonparametric bootstrap (also called the empirical bootstrap, naive bootstrap, the pairwise bootstrap, and the pairs bootstrap) draws a sample of $n$ cases $(Y_i^*, \boldsymbol{x}_i^*)$ with replacement from the $n$ cases $(Y_i, \boldsymbol{x}_i)$, and regresses the $Y_i^*$ on the $\boldsymbol{x}_i^*$ to get $\hat{\boldsymbol{\beta}}_{VS,1}^*$, and then draws another sample to get $\hat{\boldsymbol{\beta}}_{MIX,1}^*$. This process is repeated $B$ times to get the two bootstrap samples for $i = 1, \dots, B$.

Then for the full model,

$$\boldsymbol{Y}^* = \boldsymbol{X}^*\hat{\boldsymbol{\beta}}_{OLS} + \boldsymbol{r}^W$$

and for a submodel $I$,

$$\boldsymbol{Y}^* = \boldsymbol{X}_I^*\hat{\boldsymbol{\beta}}_{I,OLS} + \boldsymbol{r}_I^W.$$

Freedman (1981) showed that under regularity conditions for the OLS MLR model, $\sqrt{n}(\hat{\boldsymbol{\beta}}^* - \hat{\boldsymbol{\beta}}) \xrightarrow{D} N_p(\boldsymbol{0}, \sigma^2 \boldsymbol{W}) \sim N_p(\boldsymbol{0}, \boldsymbol{V})$. Hence if $S \subseteq I_j$,

$$\sqrt{n}(\hat{\boldsymbol{\beta}}_I^* - \hat{\boldsymbol{\beta}}_I) \xrightarrow{D} N_{a_I}(\boldsymbol{0}, \boldsymbol{V}_I)$$

as $n, B \to \infty$. (Treat $I_j$ as if $I_j$ is the full model.)

One set of regularity conditions is that the MLR model holds, and if $\boldsymbol{x}_i = (1 \ \boldsymbol{u}_i^T)^T$, then the $\boldsymbol{w}_i = (Y_i \ \boldsymbol{u}_i^T)^T$ are iid from some population with a nonsingular covariance matrix.

The nonparametric bootstrap uses $\boldsymbol{w}_1^*, \dots, \boldsymbol{w}_n^*$ where the $\boldsymbol{w}_i^*$ are sampled with replacement from $\boldsymbol{w}_1, \dots, \boldsymbol{w}_n$. By Example 4.3, $E(\boldsymbol{w}^*) = \overline{\boldsymbol{w}}$, and

$$\text{Cov}(\boldsymbol{w}^*) = \frac{1}{n} \sum_{i=1}^{n} (\boldsymbol{w}_i - \overline{\boldsymbol{w}})(\boldsymbol{w}_i - \overline{\boldsymbol{w}})^T = \widetilde{\boldsymbol{\Sigma}}\boldsymbol{w} = \begin{bmatrix} \tilde{S}_Y^2 & \tilde{\boldsymbol{\Sigma}}_Y \boldsymbol{u} \\ \tilde{\boldsymbol{\Sigma}}_{\boldsymbol{u} Y} & \tilde{\boldsymbol{\Sigma}}_{\boldsymbol{u}} \end{bmatrix}.$$

Note that $\hat{\boldsymbol{\beta}}$ is a constant with respect to the bootstrap distribution. Assume all inverse matrices exist. Then by Section 6.1.1,

$$\hat{\boldsymbol{\beta}}^* = \begin{bmatrix} \hat{\beta}_1^* \\ \hat{\boldsymbol{\beta}}_{\boldsymbol{u}}^* \end{bmatrix} = \begin{bmatrix} \overline{Y}^* - \hat{\boldsymbol{\beta}}_{\boldsymbol{u}}^{*T} \overline{\boldsymbol{u}}^* \\ \tilde{\boldsymbol{\Sigma}}_{\boldsymbol{u}}^{-1*} \tilde{\boldsymbol{\Sigma}}_{\boldsymbol{u} Y}^* \end{bmatrix} \xrightarrow{P} \begin{bmatrix} \overline{Y} - \hat{\boldsymbol{\beta}}_{\boldsymbol{u}}^{T} \overline{\boldsymbol{u}} \\ \tilde{\boldsymbol{\Sigma}}_{\boldsymbol{u}}^{-1} \tilde{\boldsymbol{\Sigma}}_{\boldsymbol{u} Y} \end{bmatrix} = \begin{bmatrix} \hat{\beta}_1 \\ \hat{\boldsymbol{\beta}}_{\boldsymbol{u}} \end{bmatrix} = \hat{\boldsymbol{\beta}}$$

as $B \to \infty$. This result suggests that the nonparametric bootstrap for OLS MLR might work under milder regularity conditions than the $\boldsymbol{w}_i$ being iid from some population with a nonsingular covariance matrix.

## 7.4.4 Bootstrapping OLS Variable Selection

Undercoverage can occur if the bootstrap sample data cloud is less variable than the iid data cloud, e.g., if $(n-p)/n$ is not close to one. Coverage can be higher than the nominal coverage for two reasons: i) the bootstrap data cloud is more variable than the iid data cloud of $T_1, ..., T_B$, and ii) zero padding.

To see the effect of zero padding, consider $H_0 : \boldsymbol{A}\boldsymbol{\beta} = \boldsymbol{\beta}_O = \boldsymbol{0}$ where $\boldsymbol{\beta}_O = (\beta_{i_1}, ...., \beta_{i_g})^T$ and $O \subseteq E$ in (7.1) so that $H_0$ is true. Suppose a nominal 95% confidence region is used and $U_B = 0.96$. Hence the confidence region (4.13) or (4.14) covers at least 96% of the bootstrap sample. If $\hat{\boldsymbol{\beta}}_{O,j}^* = \boldsymbol{0}$ for more than 4% of the $\hat{\boldsymbol{\beta}}_{O,1}^*, ..., \hat{\boldsymbol{\beta}}_{O,B}^*$, then $\boldsymbol{0}$ is in the confidence region and the bootstrap test fails to reject $H_0$. If this occurs for each run in the simulation, then the observed coverage will be 100%.

Now suppose $\hat{\boldsymbol{\beta}}_{O,j}^* = \boldsymbol{0}$ for $j = 1, ..., B$. Then $\boldsymbol{S}_T^*$ is singular, but the singleton set $\{\boldsymbol{0}\}$ is the large sample $100(1 - \delta)\%$ confidence region (4.13), (4.14), or (4.15) for $\boldsymbol{\beta}_O$ and $\delta \in (0, 1)$, and the pvalue for $H_0 : \boldsymbol{\beta}_O = \boldsymbol{0}$ is one. (This result holds since $\{\boldsymbol{0}\}$ contains 100% of the $\hat{\boldsymbol{\beta}}_{O,j}^*$ in the bootstrap sample.) For large sample theory tests, the pvalue estimates the population pvalue. Let $I$ denote the other predictors in the model so $\boldsymbol{\beta} = (\boldsymbol{\beta}_I^T, \boldsymbol{\beta}_O^T)^T$. For the $I_{min}$ model from forward selection, there may be strong evidence that $\boldsymbol{x}_O$ is not needed in the model given $\boldsymbol{x}_I$ is in the model if the "100%" confidence region is $\{\boldsymbol{0}\}$, $n \geq 20p$, $B \geq 50p$, and the error distribution is unimodal and not highly skewed. (Since the pvalue is one, this technique may be useful for data snooping: applying OLS theory to submodel $I$ may have negligible selection bias.)

**Remark 7.9.** Note that there are several important variable selection models, including the model given by Equation (7.1) where $\boldsymbol{x}^T\boldsymbol{\beta} = \boldsymbol{x}_S^T\boldsymbol{\beta}_S$. Another model is $\boldsymbol{x}^T\boldsymbol{\beta} = \boldsymbol{x}_{S_i}^T\boldsymbol{\beta}_{S_i}$ for $i = 1, ..., K$. Then there are $K \geq 2$ competing "true" nonnested submodels where $\boldsymbol{\beta}_{S_i}$ is $a_{S_i} \times 1$. For example,

suppose the $K = 2$ models have predictors $x_1, x_2, x_3$ for $S_1$ and $x_1, x_2, x_4$ for $S_2$. Then $x_3$ and $x_4$ are likely to be selected and omitted often by forward selection for the $B$ bootstrap samples. Hence omitting all predictors $x_i$ that have a $\beta_{ij}^* = 0$ for at least one of the bootstrap samples $j = 1, ..., B$ could result in underfitting, e.g. using just $x_1$ and $x_2$ in the above $K = 2$ example. If $n$ and $B$ are large enough, the singleton set $\{\mathbf{0}\}$ could still be the "100%" confidence region for a vector $\boldsymbol{\beta}_O$. See Remark 7.7.

Suppose the predictors $x_i$ have been standardized. Then another important regression model has the $\beta_i$ taper off rapidly, but no coefficients are equal to zero. For example, $\beta_i = e^{-i}$ for $i = 1, ..., p$.

**Example 7.2.** Cook and Weisberg (1999a, pp. 351, 433, 447) gives a data set on 82 mussels sampled off the coast of New Zealand. Let the response variable be the logarithm $\log(M)$ of the *muscle mass*, and the predictors are the *length L* and *height H* of the shell in mm, the logarithm $\log(W)$ of the *shell width W*, the logarithm $\log(S)$ of the *shell mass S*, and a constant. Inference for the full model is shown below along with the shorth($c$) nominal 95% confidence intervals for $\beta_i$ computed using the nonparametric and residual bootstraps. As expected, the residual bootstrap intervals are close to the classical least squares confidence intervals $\approx \hat{\beta}_i \pm 1.96 SE(\hat{\beta}_i)$.

```
      large sample full model inference
      Est.    SE   t   Pr(>|t|)   nparboot         resboot
 int -1.249 0.838 -1.49 0.14 [-2.93,-0.093][-3.045,0.473]
 L   -0.001 0.002 -0.28 0.78 [-0.005,0.003][-0.005,0.004]
 logW 0.130 0.374  0.35 0.73 [-0.457,0.829][-0.703,0.890]
 H    0.008 0.005  1.50 0.14 [-0.002,0.018][-0.003,0.016]
 logS 0.640 0.169  3.80 0.00 [ 0.244,1.040][ 0.336,1.012]
 output and shorth intervals for the min Cp submodel FS
      Est.     SE     95% shorth CI    95% shorth CI
 int  -0.9573 0.1519 [-3.294, 0.495] [-2.769, 0.460]
 L     0              [-0.005, 0.004] [-0.004, 0.004]
 logW  0              [ 0.000, 1.024] [-0.595, 0.869]
 H     0.0072 0.0047 [ 0.000, 0.016] [ 0.000, 0.016]
 logS  0.6530 0.1160 [ 0.322, 0.901] [ 0.324, 0.913]
                for forward selection for all subsets
```

The minimum $C_p$ model from all subsets variable selection and forward selection both used a constant, $H$, and $\log(S)$. The shorth($c$) nominal 95% confidence intervals for $\beta_i$ using the residual bootstrap are shown. Note that the intervals for $H$ are right skewed and contain 0 when closed intervals are used instead of open intervals. Some least squares output is shown, but should only be used for inference if the model was selected before looking at the data.

It was expected that $\log(S)$ may be the only predictor needed, along with a constant, since $\log(S)$ and $\log(M)$ are both $\log(\text{mass})$ measurements and likely highly correlated. Hence we want to test $H_0 : \beta_2 = \beta_3 = \beta_4 = 0$ with

the $I_{min}$ model selected by all subsets variable selection. (Of course this test would be easy to do with the full model using least squares theory.) Then $H_0 : \boldsymbol{A\beta} = (\beta_2, \beta_3, \beta_4)^T = \boldsymbol{0}$. Using the prediction region method with the full model gave an interval [0,2.930] with $D_{\boldsymbol{0}} = 1.641$. Note that $\sqrt{\chi^2_{3,0.95}} = 2.795$. So fail to reject $H_0$. Using the prediction region method with the $I_{min}$ variable selection model had $[0, D_{(U_B)}] = [0, 3.293]$ while $D_{\boldsymbol{0}} = 1.134$. So fail to reject $H_0$.

Then we redid the bootstrap with the full model and forward selection. The full model had $[0, D_{(U_B)}] = [0, 2.908]$ with $D_{\boldsymbol{0}} = 1.577$. So fail to reject $H_0$. Using the prediction region method with the $I_{min}$ forward selection model had $[0, D_{(U_B)}] = [0, 3.258]$ while $D_{\boldsymbol{0}} = 1.245$. So fail to reject $H_0$. The ratio of the volumes of the bootstrap confidence regions for this test was 0.392. (Use (4.16) with $\boldsymbol{S}_T^*$ and $D$ from forward selection for the numerator, and from the full model for the denominator.) Hence the forward selection bootstrap test was more precise than the full model bootstrap test. Some $R$ code used to produce the above output is shown below.

```
library(leaps)
y <- log(mussels[,5]); x <- mussels[,1:4]
x[,4] <- log(x[,4]); x[,2] <- log(x[,2])
out <- regboot(x,y,B=1000)
tem <- rowboot(x,y,B=1000)
outvs <- vselboot(x,y,B=1000) #get bootstrap CIs
outfs <- fselboot(x,y,B=1000) #get bootstrap CIs
apply(out$betas,2,shorth3);
apply(tem$betas,2,shorth3);
apply(outvs$betas,2,shorth3)  #for all subsets
apply(outfs$betas,2,shorth3) #for forward selection
ls.print(outvs$full)
ls.print(outvs$sub)
ls.print(outfs$sub)
#test if beta_2 = beta_3 = beta_4 = 0
Abeta <- out$betas[,2:4]  #full model
#prediction region method with residual bootstrap
out<-predreg(Abeta)
Abeta <- outvs$betas[,2:4]
#prediction region method with Imin all subsets
outvs <- predreg(Abeta)
Abeta <- outfs$betas[,2:4]
#prediction region method with Imin forward sel.
outfs<-predreg(Abeta)
#ratio of volumes for forward selection and full model
(sqrt(det(outfs$cov))*outfs$D0^3)/(sqrt(det(out$cov))*out$D0^3)
```

**Example 7.3.** Consider the Gladstone (1905) data set that has 12 variables on 267 persons after death. The response variable was *brain weight*.

Head measurements were *breadth, circumference, head height, length,* and *size* as well as *cephalic index* and *brain weight. Age*, *height*, and two categorical variables *ageclass* (0: under 20, 1: 20-45, 2: over 45) and *sex* were also given. The eight predictor variables shown in the output were used.

Output is shown below for the full model and the bootstrapped minimum $C_p$ forward selection estimator. Note that the shorth intervals for *length* and *sex* are quite long. These variables are often in and often deleted from the bootstrap forward selection. Model $I_I$ is the model with the fewest predictors such that $C_P(I_I) \leq C_P(I_{min}) + 1$. For this data set, $I_I = I_{min}$. The bootstrap CIs differ due to different random seeds.

```
large sample full model inference for Ex. 7.3
        Estimate    SE       t     Pr(>|t|) 95% shorth CI
Int    -3021.255 1701.070 -1.77 0.077 [-6549.8,322.79]
age        -1.656    0.314 -5.27 0.000 [ -2.304,-1.050]
breadth  -8.717    12.025 -0.72 0.469 [-34.229,14.458]
cephalic 21.876    22.029  0.99 0.322 [-20.911,67.705]
circum     0.852     0.529  1.61 0.109 [ -0.065, 1.879]
headht     7.385     1.225  6.03 0.000 [  5.138, 9.794]
height    -0.407     0.942 -0.43 0.666 [ -2.211, 1.565]
len       13.475     9.422  1.43 0.154 [ -5.519,32.605]
sex       25.130    10.015  2.51 0.013 [  6.717,44.19]
output and shorth intervals for the min Cp submodel
        Estimate    SE       t     Pr(>|t|) 95% shorth CI
Int    -1764.516  186.046 -9.48 0.000 [-6151.6,-415.4]
age        -1.708    0.285 -5.99 0.000 [ -2.299,-1.068]
breadth    0                            [-32.992, 8.148]
cephalic   5.958     2.089  2.85 0.005 [-10.859,62.679]
circum     0.757     0.512  1.48 0.140 [  0.000, 1.817]
headht     7.424     1.161  6.39 0.000 [  5.028, 9.732]
height     0                            [ -2.859, 0.000]
len        6.716     1.466  4.58 0.000 [  0.000,30.508]
sex       25.313     9.920  2.55 0.011 [  0.000,42.144]
output and shorth for I_I model
        Estimate  Std.Err t-val Pr(>|t|) 95% shorth CI
Int    -1764.516  186.046 -9.48 0.000 [-6104.9,-778.2]
age        -1.708    0.285 -5.99 0.000 [ -2.259,-1.003]
breadth    0                            [-31.012, 6.567]
cephalic   5.958     2.089  2.85 0.005 [ -6.700,61.265]
circum     0.757     0.512  1.48 0.140 [  0.000, 1.866]
headht     7.424     1.161  6.39 0.000 [  5.221,10.090]
height     0                            [ -2.173, 0.000]
len        6.716     1.466  4.58 0.000 [  0.000,28.819]
sex       25.313     9.920  2.55 0.011 [  0.000,42.847]
```

The $R$ code used to produce the above output is shown below. The last four commands are useful for examining the variable selection output.

```
x<-cbrainx[,c(1,3,5,6,7,8,9,10)]
y<-cbrainy
library(leaps)
out <- regboot(x,y,B=1000)
outvs <- fselboot(x,cbrainy) #get bootstrap CIs,
apply(out$betas,2,shorth3)
apply(outvs$betas,2,shorth3)
ls.print(outvs$full)
ls.print(outvs$sub)
outvs <- modIboot(x,cbrainy) #get bootstrap CIs,
apply(outvs$betas,2,shorth3)
ls.print(outvs$sub)
tem<-regsubsets(x,y,method="forward")
tem2<-summary(tem)
tem2$which
tem2$cp
```

## 7.4.5 Simulations

For variable selection with the $p \times 1$ vector $\hat{\boldsymbol{\beta}}_{I_{min},0}$, consider testing $H_0 :$ $\boldsymbol{A\beta} = \boldsymbol{\theta}_0$ versus $H_1 : \boldsymbol{A\beta} \neq \boldsymbol{\theta}_0$ with $\boldsymbol{\theta} = \boldsymbol{A\beta}$ where often $\boldsymbol{\theta}_0 = \mathbf{0}$. Then let $T_n = \boldsymbol{A}\hat{\boldsymbol{\beta}}_{I_{min},0}$ and let $T_i^* = \boldsymbol{A}\hat{\boldsymbol{\beta}}_{I_{min},0,i}^*$ for $i = 1, ..., B$. The shorth estimator can be applied to a bootstrap sample $\hat{\beta}_{i1}^*, ..., \hat{\beta}_{iB}^*$ to get a confidence interval for $\beta_i$. Here $T_n = \hat{\beta}_i$ and $\theta = \beta_i$.

Assume $p$ is fixed, $n \geq 20p$, and that the error distribution is unimodal and not highly skewed. Then the plotted points in the response and residual plots should scatter in roughly even bands about the identity line (with unit slope and zero intercept) and the $r = 0$ line, respectively. See Figure 5.8. If the error distribution is skewed or multimodal, then much larger sample sizes may be needed.

Next, we describe a small simulation study that was done using $B = \max(1000, n/25, 50p)$ and 5000 runs. The simulation used $p = 4, 6, 7, 8$, and 10; $n = 25p$ and $50p$; $\psi = 0, 1/\sqrt{p}$, and 0.9; and $k = 1$ and $p - 2$ where $k$ and $\psi$ are defined in the following paragraph. In the simulations, we use $\theta = \boldsymbol{A\beta} = \beta_i$, $\theta = \boldsymbol{A\beta} = \boldsymbol{\beta}_S = \mathbf{1}$ and $\theta = \boldsymbol{A\beta} = \boldsymbol{\beta}_E = \mathbf{0}$.

Let $\boldsymbol{x} = (1 \ \boldsymbol{u}^T)^T$ where $\boldsymbol{u}$ is the $(p-1) \times 1$ vector of nontrivial predictors. In the simulations, for $i = 1, ..., n$, we generated $\boldsymbol{w}_i \sim N_{p-1}(\mathbf{0}, \boldsymbol{I})$ where the $m = p - 1$ elements of the vector $\boldsymbol{w}_i$ are iid N(0,1). Let the $m \times m$ matrix $\boldsymbol{A} = (a_{ij})$ with $a_{ii} = 1$ and $a_{ij} = \psi$ where $0 \leq \psi < 1$ for $i \neq j$. Then the vector $\boldsymbol{u}_i = \boldsymbol{A}\boldsymbol{w}_i$ so that $Cov(\boldsymbol{u}_i) = \boldsymbol{\Sigma_u} = \boldsymbol{AA}^T = (\sigma_{ij})$ where the diagonal

entries $\sigma_{ii} = [1+(m-1)\psi^2]$ and the off diagonal entries $\sigma_{ij} = [2\psi+(m-2)\psi^2]$. Hence the correlations are $Cor(x_i, x_j) = \rho = (2\psi + (m-2)\psi^2)/(1 + (m-1)\psi^2)$ for $i \neq j$ where $x_i$ and $x_j$ are nontrivial predictors. If $\psi = 1/\sqrt{cp}$, then $\rho \to 1/(c+1)$ as $p \to \infty$ where $c > 0$. As $\psi$ gets close to 1, the predictor vectors cluster about the line in the direction of $(1, ..., 1)^T$. Let $Y_i = 1 + 1x_{i,2} + \cdots + 1x_{i,k+1} + e_i$ for $i = 1, ..., n$. Hence $\boldsymbol{\beta} = (1, .., 1, 0, ..., 0)^T$ with $k + 1$ ones and $p - k - 1$ zeros. The zero mean errors $e_i$ were iid from five distributions: i) N(0,1), ii) $t_3$, iii) EXP(1) - 1, iv) uniform$(-1, 1)$, and v) 0.9 N(0,1) + 0.1 N(0,100). Only distribution iii) is not symmetric.

When $\psi = 0$, the full model least squares confidence intervals for $\beta_i$ should have length near $2t_{96,0.975}\sigma/\sqrt{n} \approx 2(1.96)\sigma/10 = 0.392\sigma$ when $n = 100$ and the iid zero mean errors have variance $\sigma^2$. The simulation computed the Frey shorth($c$) interval for each $\beta_i$ and used bootstrap confidence regions to test $H_0 : \boldsymbol{\beta}_S = \mathbf{1}$ (whether first $k + 1$ $\beta_i = 1$) and $H_0 : \boldsymbol{\beta}_E = \mathbf{0}$ (whether the last $p - k - 1$ $\beta_i = 0$). The nominal coverage was 0.95 with $\delta = 0.05$. Observed coverage between 0.94 and 0.96 suggests coverage is close to the nominal value.

The regression models used the residual bootstrap on the forward selection estimator $\hat{\boldsymbol{\beta}}_{I_{min},0}$. Table 7.1 gives results for when the iid errors $e_i \sim N(0, 1)$ with $n = 100$, $p = 4$, and $k = 1$. Table 7.1 shows two rows for each model giving the observed confidence interval coverages and average lengths of the confidence intervals. The term "reg" is for the full model regression, and the term "vs" is for forward selection. The last six columns give results for the tests. The terms pr, hyb, and br are for the prediction region method (4.13), hybrid region (4.15), and Bickel and Ren region (4.14). The 0 indicates the test was $H_0 : \boldsymbol{\beta}_E = \mathbf{0}$, while the 1 indicates that the test was $H_0 : \boldsymbol{\beta}_S = \mathbf{1}$. The length and coverage = P(fail to reject $H_0$) for the interval $[0, D_{(U_B)}]$ or $[0, D_{(U_B,T)}]$ where $D_{(U_B)}$ or $D_{(U_B,T)}$ is the cutoff for the confidence region. The cutoff will often be near $\sqrt{\chi^2_{g,0.95}}$ if the statistic $T$ is asymptotically normal. Note that $\sqrt{\chi^2_{2,0.95}} = 2.448$ is close to 2.45 for the full model regression bootstrap tests.

Volume ratios of the three confidence regions can be compared using (4.16), but there is not enough information in Table 7.1 to compare the volume of the confidence region for the full model regression versus that for the forward selection regression since the two methods have different determinants $|\boldsymbol{S}_T^*|$.

The inference for forward selection was often as precise or more precise than the inference for the full model. The coverages were near 0.95 for the regression bootstrap on the full model, although there was slight undercoverage for the tests since $(n - p)/n = 0.96$ when $n = 25p$. Suppose $\psi = 0$. Then from Section 7.2, $\hat{\boldsymbol{\beta}}_S$ may have the same limiting distribution for $I_{min}$ and the full model. Note that the average lengths and coverages were similar for the full model and forward selection $I_{min}$ for $\beta_1$, $\beta_2$, and $\boldsymbol{\beta}_S = (\beta_1, \beta_2)^T$. Forward selection inference was more precise for $\boldsymbol{\beta}_E = (\beta_3, \beta_4)^T$. The Bickel

**Table 7.1** Bootstrapping OLS Forward Selection with $C_p$, $e_i \sim N(0,1)$

| $\psi$ | $\beta_1$ | $\beta_2$ | $\beta_{p-1}$ | $\beta_p$ | pr0 | hyb0 | br0 | pr1 | hyb1 | br1 |
|---|---|---|---|---|---|---|---|---|---|---|
| reg,0 | 0.946 | 0.950 | 0.947 | 0.948 | 0.940 | 0.941 | 0.941 | 0.937 | 0.936 | 0.937 |
| len | 0.396 | 0.399 | 0.399 | 0.398 | 2.451 | 2.451 | 2.452 | 2.450 | 2.450 | 2.451 |
| vs,0 | 0.948 | 0.950 | 0.997 | 0.996 | 0.991 | 0.979 | 0.991 | 0.938 | 0.939 | 0.940 |
| len | 0.395 | 0.398 | 0.323 | 0.323 | 2.699 | 2.699 | 3.002 | 2.450 | 2.450 | 2.457 |
| reg,0.5 | 0.946 | 0.944 | 0.946 | 0.945 | 0.938 | 0.938 | 0.938 | 0.934 | 0.936 | 0.936 |
| len | 0.396 | 0.661 | 0.661 | 0.661 | 2.451 | 2.451 | 2.452 | 2.451 | 2.451 | 2.452 |
| vs,0.5 | 0.947 | 0.968 | 0.997 | 0.998 | 0.993 | 0.984 | 0.993 | 0.955 | 0.955 | 0.963 |
| len | 0.395 | 0.658 | 0.537 | 0.539 | 2.703 | 2.703 | 2.994 | 2.461 | 2.461 | 2.577 |
| reg,0.9 | 0.946 | 0.941 | 0.944 | 0.950 | 0.940 | 0.940 | 0.940 | 0.935 | 0.935 | 0.935 |
| len | 0.396 | 3.257 | 3.253 | 3.259 | 2.451 | 2.451 | 2.452 | 2.451 | 2.451 | 2.452 |
| vs,0.9 | 0.947 | 0.968 | 0.994 | 0.996 | 0.992 | 0.981 | 0.992 | 0.962 | 0.959 | 0.970 |
| len | 0.395 | 2.751 | 2.725 | 2.735 | 2.716 | 2.716 | 2.971 | 2.497 | 2.497 | 2.599 |

and Ren (4.14) cutoffs and coverages were at least as high as those of the hybrid region (4.15).

For $\psi > 0$ and $I_{min}$, the coverages for the $\beta_i$ corresponding to $\boldsymbol{\beta}_S$ were near 0.95, but the average length could be shorter since $I_{min}$ tends to have less multicorrelation than the full model. For $\psi \geq 0$, the $I_{min}$ coverages were higher than 0.95 for $\beta_3$ and $\beta_4$ and for testing $H_0 : \boldsymbol{\beta}_E = \mathbf{0}$ since zeros often occurred for $\hat{\beta}_j^*$ for $j = 3, 4$. The average CI lengths were shorter for $I_{min}$ than for the OLS full model for $\beta_3$ and $\beta_4$. Note that for $I_{min}$, the coverage for testing $H_0 : \boldsymbol{\beta}_S = \mathbf{1}$ was higher than that for the OLS full model.

**Table 7.2** Bootstrap CIs with $C_p$, $p = 10, k = 8, \psi = 0.9$, error type v)

| $n$ | $\beta_1$ | $\beta_2$ | $\beta_3$ | $\beta_4$ | $\beta_5$ | $\beta_6$ | $\beta_7$ | $\beta_8$ | $\beta_9$ | $\beta_{10}$ |
|---|---|---|---|---|---|---|---|---|---|---|
| 250 | 0.945 | 0.824 | 0.822 | 0.827 | 0.827 | 0.824 | 0.826 | 0.817 | 0.827 | 0.999 |
| shlen | 0.825 | 6.490 | 6.490 | 6.482 | 6.485 | 6.479 | 6.512 | 6.496 | 6.493 | 6.445 |
| 250 | 0.946 | 0.979 | 0.980 | 0.985 | 0.981 | 0.983 | 0.983 | 0.977 | 0.983 | 0.998 |
| prlen | 0.807 | 7.836 | 7.850 | 7.842 | 7.830 | 7.830 | 7.851 | 7.840 | 7.839 | 7.802 |
| 250 | 0.947 | 0.976 | 0.978 | 0.984 | 0.978 | 0.978 | 0.979 | 0.973 | 0.980 | 0.996 |
| brlen | 0.811 | 8.723 | 8.760 | 8.765 | 8.736 | 8.764 | 8.745 | 8.747 | 8.753 | 8.756 |
| 2500 | 0.951 | 0.947 | 0.948 | 0.948 | 0.948 | 0.947 | 0.949 | 0.944 | 0.951 | 0.999 |
| shlen | 0.263 | 2.268 | 2.271 | 2.271 | 2.273 | 2.262 | 2.632 | 2.277 | 2.272 | 2.047 |
| 2500 | 0.945 | 0.961 | 0.959 | 0.955 | 0.960 | 0.960 | 0.961 | 0.958 | 0.961 | 0.998 |
| prlen | 0.258 | 2.630 | 2.639 | 2.640 | 2.632 | 2.632 | 2.641 | 2.638 | 2.642 | 2.517 |
| 2500 | 0.946 | 0.958 | 0.954 | 0.960 | 0.956 | 0.960 | 0.962 | 0.955 | 0.961 | 0.997 |
| brlen | 0.258 | 2.865 | 2.875 | 2.882 | 2.866 | 2.871 | 2.887 | 2.868 | 2.875 | 2.830 |
| 25000 | 0.952 | 0.940 | 0.939 | 0.935 | 0.940 | 0.942 | 0.938 | 0.937 | 0.942 | 1.000 |
| shlen | 0.083 | 0.809 | 0.808 | 0.806 | 0.805 | 0.807 | 0.808 | 0.808 | 0.809 | 0.224 |
| 25000 | 0.948 | 0.964 | 0.968 | 0.962 | 0.964 | 0.966 | 0.964 | 0.964 | 0.967 | 0.991 |
| prlen | 0.082 | 0.806 | 0.805 | 0.801 | 0.800 | 0.805 | 0.805 | 0.803 | 0.806 | 0.340 |
| 25000 | 0.949 | 0.969 | 0.972 | 0.968 | 0.967 | 0.971 | 0.969 | 0.969 | 0.973 | 0.999 |
| brlen | 0.082 | 0.810 | 0.810 | 0.805 | 0.804 | 0.809 | 0.810 | 0.808 | 0.810 | 0.317 |

Results for other values of $n$, $p$, $k$, and distributions of $e_i$ were similar. For forward selection with $\psi = 0.9$ and $C_p$, the hybrid region (4.15) and shorth confidence intervals occasionally had coverage less than 0.93. It was also rare for the bootstrap to have one or more columns of zeroes so $\boldsymbol{S}_T^*$ was singular. For error distributions i)-iv) and $\psi = 0.9$, sometimes the shorth CIs needed $n \geq 100p$ for all $p$ CIs to have good coverage. For error distribution v) and $\psi = 0.9$, even larger values of $n$ were needed. Confidence intervals based on (4.13) and (4.14) worked for much smaller $n$, but tended to be longer than the shorth CIs.

See Table 7.2 for one of the worst scenarios for the shorth, where shlen, prlen, and brlen are for the average CI lengths based on the shorth, (4.13), and (4.14), respectively. In Table 4.3, $k = 8$ and the two nonzero $\pi_j$ correspond to the full model $\hat{\boldsymbol{\beta}}$ and $\hat{\boldsymbol{\beta}}_{S,0}$. Hence $\beta_i = 1$ for $i = 1, ..., 9$ and $\beta_{10} = 0$. Hence confidence intervals for $\beta_{10}$ had the highest coverage and usually the shortest average length (for $i \neq 1$) due to zero padding. Theory in Section 7.2 showed that the CI lengths are proportional to $1/\sqrt{n}$. When $n = 25000$, the shorth CI uses the 95.16th percentile while CI (4.13) uses the 95.00th percentile, allowing the average CI length of (4.13) to be shorter than that of the shorth CI, but the distribution for $\hat{\beta}_i^*$ is likely approximately symmetric for $i \neq 10$ since the average lengths of the three confidence intervals were about the same for each $i \neq 10$.

When BIC was used, undercoverage was a bit more common and severe, and undercoverage occasionally occurred with regions (4.13) and (4.14). BIC also occasionally had 100% coverage since BIC produces more zeroes than $C_p$.

Some $R$ code for the simulation is shown below.

```
record coverages and ``lengths" for
b1, b2, bp-1, bp, pm0, hyb0, br0, pm1, hyb1, br1


regbootsim3(n=100,p=4,k=1,nruns=5000,type=1,psi=0)
$cicov
[1] 0.9458 0.9500 0.9474 0.9484 0.9400 0.9408 0.9410
0.9368 0.9362 0.9370
$avelen
[1] 0.3955 0.3990 0.3987 0.3982 2.4508 2.4508 2.4521
[8] 2.4496 2.4496 2.4508
$beta
[1] 1 1 0 0
$k
[1] 1
library(leaps)
vsbootsim4(n=100,p=4,k=1,nruns=5000,type=1,psi=0)
$cicov
[1] 0.9480 0.9496 0.9972 0.9958 0.9910 0.9786 0.9914
```

```
0.9384 0.9394 0.9402
$avelen
[1]  0.3954 0.3987 0.3233 0.3231 2.6987 2.6987 3.0020
[8]  2.4497 2.4497 2.4570
$beta
[1] 1 1 0 0
$k
[1] 1
```

## 7.5 Data Splitting

Data splitting is used for inference after model selection. Use a training set to select a full model, and a validation set for inference with the selected full model. Here $p >> n$ is possible. See Hurvich and Tsai (1990, p. 216) and Rinaldo et al. (2019). Typically when training and validation sets are used, the training set is bigger than the validation set or half sets are used, often causing large efficiency loss.

Let $J$ be a positive integer and let $\lfloor x \rfloor$ be the integer part of $x$, e.g., $\lfloor 7.7 \rfloor = 7$. Initially divide the data into two sets $H_1$ with $n_1 = \lfloor n/(2J) \rfloor$ cases and $V_1$ with $n - n_1$ cases. If the fitted model from $H_1$ is not good enough, randomly select $n_1$ cases from $V_1$ to add to $H_1$ to form $H_2$. Let $V_2$ have the remaining cases from $V_1$. Continue in this manner, possibly forming sets $(H_1, V_1), (H_2, V_2), ..., (H_J, V_J)$ where $H_i$ has $n_i = in_1$ cases. Stop when $H_d$ gives a reasonable model $I_d$ with $a_d$ predictors if $d < J$. Use $d = J$, otherwise. Use the model $I_d$ as the full model for inference with the data in $V_d$.

This procedure is simple for a fixed data set, but it would be good to automate the procedure. Forward selection with the Chen and Chen (2008) EBIC criterion and lasso are useful for finding a reasonable fitted model. BIC and the Hurvich and Tsai (1989) $AIC_C$ criterion can be useful if $n \geq \max(2p, 10a_d)$. For example, if $n = 500000$ and $p = 90$, using $n_1 = 900$ would result in a much smaller loss of efficiency than $n_1 = 250000$.

## 7.6 Some Alternative MLR Estimators

From Definition 5.11, the multiple linear regression (MLR) model is

$$Y_i = \beta_1 + x_{i,2}\beta_2 + \cdots + x_{i,p}\beta_p + e_i = \boldsymbol{x}_i^T \boldsymbol{\beta} + e_i \qquad (7.8)$$

for $i = 1, ..., n$. This model is also called the **full model**. Here $n$ is the sample size and the random variable $e_i$ is the $i$th error. Assume that the $e_i$ are iid

with variance $V(e_i) = \sigma^2$. In matrix notation, these $n$ equations become $\boldsymbol{Y} = \boldsymbol{X\beta} + \boldsymbol{e}$ where $\boldsymbol{Y}$ is an $n \times 1$ vector of dependent variables, $\boldsymbol{X}$ is an $n \times p$ matrix of predictors, $\boldsymbol{\beta}$ is a $p \times 1$ vector of unknown coefficients, and $\boldsymbol{e}$ is an $n \times 1$ vector of unknown errors.

There are many methods for estimating $\boldsymbol{\beta}$, including (ordinary) least squares (OLS) for the full model, forward selection with OLS, elastic net, principal components regression (PCR), partial least squares (PLS), lasso, lasso variable selection, and ridge regression (RR). For the last six methods, it is convenient to use centered or scaled data. Suppose $U$ has observed values $U_1, ..., U_n$. For example, if $U_i = Y_i$ then $U$ corresponds to the response variable $Y$. The observed values of a random variable $V$ are *centered* if their sample mean is 0. The centered values of $U$ are $V_i = U_i - \overline{U}$ for $i = 1, ..., n$. Let $g$ be an integer near 0. If the sample variance of the $U_i$ is

$$\hat{\sigma}_g^2 = \frac{1}{n-g} \sum_{i=1}^n (U_i - \overline{U})^2,$$

then the sample standard deviation of $U_i$ is $\hat{\sigma}_g$. If the values of $U_i$ are not all the same, then $\hat{\sigma}_g > 0$, and the standardized values of the $U_i$ are

$$W_i = \frac{U_i - \overline{U}}{\hat{\sigma}_g}.$$

Typically $g = 1$ or $g = 0$ are used: $g = 1$ gives an unbiased estimator of $\sigma^2$ while $g = 0$ gives the method of moments estimator. Note that the standardized values are centered, $\overline{W} = 0$, and the sample variance of the standardized values

$$\frac{1}{n-g} \sum_{i=1}^n W_i^2 = 1. \tag{7.9}$$

**Remark 7.10.** Let the nontrivial predictors $\boldsymbol{u}_i^T = (x_{i,2}, ..., x_{i,p}) = (u_{i,1}, ..., u_{i,p-1})$. Then $\boldsymbol{x}_i = (1, \boldsymbol{u}_i^T)^T$. Let the $n \times (p-1)$ matrix of standardized nontrivial predictors $\boldsymbol{W}_g = (W_{ij})$ when the predictors are standardized using $\hat{\sigma}_g$. Thus, $\sum_{i=1}^n W_{ij} = 0$ and $\sum_{i=1}^n W_{ij}^2 = n - g$ for $j = 1, ..., p-1$. Hence

$$W_{ij} = \frac{x_{i,j+1} - \overline{x}_{j+1}}{\hat{\sigma}_{j+1}} \quad \text{where} \quad \hat{\sigma}_{j+1}^2 = \frac{1}{n-g} \sum_{i=1}^n (x_{i,j+1} - \overline{x}_{j+1})^2$$

is $\hat{\sigma}_g$ for the $(j+1)$th variable $x_{j+1}$. Let $\boldsymbol{w}_i^T = (w_{i,1}, ..., w_{i,p-1})$ be the standardized vector of nontrivial predictors for the $i$th case. Since the standardized data are also centered, $\overline{\boldsymbol{w}} = \boldsymbol{0}$. Then the sample covariance matrix of the $\boldsymbol{w}_i$ is the sample correlation matrix of the $\boldsymbol{u}_i$:

$$\hat{\boldsymbol{\rho}}_{\boldsymbol{u}} = \boldsymbol{R}_{\boldsymbol{u}} = (r_{ij}) = \frac{\boldsymbol{W}_g^T \boldsymbol{W}_g}{n-g}$$

where $r_{ij}$ is the sample correlation of $u_i = x_{i+1}$ and $u_j = x_{j+1}$. Thus the sample correlation matrix $\boldsymbol{R_u}$ does not depend on $g$. Let $\boldsymbol{Z} = \boldsymbol{Y} - \overline{\boldsymbol{Y}}$ where $\overline{\boldsymbol{Y}} = \overline{Y}\boldsymbol{1}$. Since the $R$ software tends to use $g = 0$, let $\boldsymbol{W} = \boldsymbol{W}_0$. Note that $n \times (p-1)$ matrix $\boldsymbol{W}$ does not include a vector $\boldsymbol{1}$ of ones. Then regression through the origin is used for the model

$$\boldsymbol{Z} = \boldsymbol{W}\boldsymbol{\eta} + \boldsymbol{e} \tag{7.10}$$

where $\boldsymbol{Z} = (Z_1, ..., Z_n)^T$ and $\boldsymbol{\eta} = (\eta_1, ..., \eta_{p-1})^T$. The vector of fitted values $\hat{\boldsymbol{Y}} = \overline{\boldsymbol{Y}} + \hat{\boldsymbol{Z}}$.

**Remark 7.11.** i) Interest is in model (7.8): estimate $\hat{Y}_f$ and $\hat{\boldsymbol{\beta}}$. For many regression estimators, a method is needed so that everyone who uses the same units of measurements for the predictors and $Y$ gets the same $(\hat{\boldsymbol{Y}}, \hat{\boldsymbol{\beta}})$. Also, see Remark 5.3. Equation (7.10) $\boldsymbol{Z} = \boldsymbol{W}\boldsymbol{\eta} + \boldsymbol{e}$ is a commonly used method for achieving this goal. Suppose $g = 0$. The method of moments estimator of the variance $\sigma_w^2$ is

$$\hat{\sigma}_{g=0}^2 = S_M^2 = \frac{1}{n}\sum_{i=1}^{n}(w_i - \overline{w})^2.$$

When data $x_i$ are standardized to have $\overline{w} = 0$ and $S_M^2 = 1$, the standardized data $w_i$ has no units. ii) Hence the estimators $\hat{\boldsymbol{Z}}$ and $\hat{\boldsymbol{\eta}}$ do not depend on the units of measurement of the $x_i$ if standardized data and Equation (7.10) are used. Linear combinations of the $\boldsymbol{w}_i$ are linear combinations of the $\boldsymbol{u}_i$, which are linear combinations of the $\boldsymbol{x}_i$. (Note that $\boldsymbol{\gamma}^T\boldsymbol{u} = (0 \ \boldsymbol{\gamma}^T) \ \boldsymbol{x}$.) Thus the estimators $\hat{\boldsymbol{Y}}$ and $\hat{\boldsymbol{\beta}}$ are obtained using $\hat{\boldsymbol{Z}}$, $\hat{\boldsymbol{\eta}}$, and $\overline{\boldsymbol{Y}}$. The linear transformation to obtain $(\hat{\boldsymbol{Y}}, \hat{\boldsymbol{\beta}})$ from $(\hat{\boldsymbol{Z}}, \hat{\boldsymbol{\eta}})$ is unique for a given set of units of measurements for the $x_i$ and $Y$. Hence everyone using the same units of measurements gets the same $(\hat{\boldsymbol{Y}}, \hat{\boldsymbol{\beta}})$. iii) Also, since $\overline{W}_j = 0$ and $S_{M,j}^2 = 1$, the standardized predictor variables have similar spread, and the magnitude of $\hat{\eta}_i$ is a measure of the importance of the predictor variable $W_j$ for predicting $Y$.

**Remark 7.12.** Let $\hat{\sigma}_j$ be the sample standard deviation of variable $x_j$ (often with $g = 0$) for $j = 2, ...., p$. Let $\hat{Y}_i = \hat{\beta}_1 + x_{i,2}\hat{\beta}_2 + \cdots + x_{i,p}\hat{\beta}_p = \boldsymbol{x}_i^T\hat{\boldsymbol{\beta}}$. If standardized nontrivial predictors are used, then

$$\hat{Y}_i = \hat{\gamma} + w_{i,1}\hat{\eta}_1 + \cdots + w_{i,p-1}\hat{\eta}_{p-1} = \hat{\gamma} + \frac{x_{i,2} - \overline{x}_2}{\hat{\sigma}_2}\hat{\eta}_1 + \cdots + \frac{x_{i,p} - \overline{x}_p}{\hat{\sigma}_p}\hat{\eta}_{p-1}$$

$$= \hat{\gamma} + \boldsymbol{w}_i^T\hat{\boldsymbol{\eta}} = \hat{\gamma} + \hat{Z}_i \tag{7.11}$$

where

$$\hat{\eta}_j = \hat{\sigma}_{j+1}\hat{\beta}_{j+1} \tag{7.12}$$

for $j = 1, ..., p-1$. Often $\hat{\gamma} = \overline{Y}$ so that $\hat{Y}_i = \overline{Y}$ if $x_{i,j} = \overline{x}_j$ for $j = 2, ..., p$. Then $\hat{\boldsymbol{Y}} = \overline{\boldsymbol{Y}} + \hat{\boldsymbol{Z}}$ where $\overline{\boldsymbol{Y}} = \overline{Y}\boldsymbol{1}$. Note that

$$\hat{\gamma} = \hat{\beta}_1 + \frac{\overline{x}_2}{\hat{\sigma}_2}\hat{\eta}_1 + \cdots + \frac{\overline{x}_p}{\hat{\sigma}_p}\hat{\eta}_{p-1}.$$

**Notation.** The symbol $A \equiv B = f(c)$ means that $A$ and $B$ are equivalent and equal, and that $f(c)$ is the formula used to compute $A$ and $B$.

Most regression methods attempt to find an estimate $\hat{\boldsymbol{\beta}}$ of $\boldsymbol{\beta}$ which minimizes some criterion function $Q(\boldsymbol{b})$ of the residuals. As in Definition 5.1, given an estimate $\boldsymbol{b}$ of $\boldsymbol{\beta}$, the corresponding vector of *fitted values* is $\widehat{\boldsymbol{Y}} \equiv \widehat{\boldsymbol{Y}}(\boldsymbol{b}) = \boldsymbol{X}\boldsymbol{b}$, and the vector of *residuals* is $\boldsymbol{r} \equiv \boldsymbol{r}(\boldsymbol{b}) = \boldsymbol{Y} - \widehat{\boldsymbol{Y}}(\boldsymbol{b})$. See Definition 5.2 for the OLS model for $\boldsymbol{Y} = \boldsymbol{X}\boldsymbol{\beta} + \boldsymbol{e}$. The following model is useful for the centered response and standardized nontrivial predictors, or if $\boldsymbol{Z} = \boldsymbol{Y}$, $\boldsymbol{W} = \boldsymbol{X}_I$, and $\boldsymbol{\eta} = \boldsymbol{\beta}_I$ corresponds to a submodel $I$.

**Definition 7.6.** If $\boldsymbol{Z} = \boldsymbol{W}\boldsymbol{\eta} + \boldsymbol{e}$, where the $n \times q$ matrix $\boldsymbol{W}$ has full rank $q = p - 1$, then the *OLS estimator*

$$\hat{\boldsymbol{\eta}}_{OLS} = (\boldsymbol{W}^T\boldsymbol{W})^{-1}\boldsymbol{W}^T\boldsymbol{Z}$$

minimizes the OLS criterion $Q_{OLS}(\boldsymbol{\eta}) = \boldsymbol{r}(\boldsymbol{\eta})^T\boldsymbol{r}(\boldsymbol{\eta})$ over all vectors $\boldsymbol{\eta} \in \mathbb{R}^{p-1}$. The vector of *predicted* or *fitted values* $\widehat{\boldsymbol{Z}}_{OLS} = \boldsymbol{W}\hat{\boldsymbol{\eta}}_{OLS} = \boldsymbol{H}\boldsymbol{Z}$ where $\boldsymbol{H} = \boldsymbol{W}(\boldsymbol{W}^T\boldsymbol{W})^{-1}\boldsymbol{W}^T$. The vector of residuals $\boldsymbol{r} = \boldsymbol{r}(\boldsymbol{Z}, \boldsymbol{W}) = \boldsymbol{Z} - \widehat{\boldsymbol{Z}} = (\boldsymbol{I} - \boldsymbol{H})\boldsymbol{Z}$.

Assume that the sample correlation matrix

$$\boldsymbol{R_u} = \frac{\boldsymbol{W}^T\boldsymbol{W}}{n} \xrightarrow{P} \boldsymbol{V}^{-1}. \tag{7.13}$$

Note that $\boldsymbol{V}^{-1} = \boldsymbol{\rho_u}$, the population correlation matrix of the nontrivial predictors $\boldsymbol{u}_i$, if the $\boldsymbol{u}_i$ are a random sample from a population. Let $\boldsymbol{H} = \boldsymbol{W}(\boldsymbol{W}^T\boldsymbol{W})^{-1}\boldsymbol{W}^T = (h_{ij})$, and assume that $\max_{i=1,\ldots,n} h_{ii} \xrightarrow{P} 0$ as $n \to \infty$. Then by Theorem 5.9 (the OLS CLT), the OLS estimator satisfies

$$\sqrt{n}(\hat{\boldsymbol{\eta}}_{OLS} - \boldsymbol{\eta}) \xrightarrow{D} N_{p-1}(\boldsymbol{0}, \sigma^2\boldsymbol{V}). \tag{7.14}$$

**Remark 7.13:** Variable selection is the search for a subset of predictor variables that can be deleted without important loss of information if $n/p$ is large (and the search for a useful subset of predictors if $n/p$ is not large). Refer to Equation (7.1) where $\boldsymbol{x}^T\boldsymbol{\beta} = \boldsymbol{x}_S^T\boldsymbol{\beta}_S + \boldsymbol{x}_E^T\boldsymbol{\beta}_E = \boldsymbol{x}_S^T\boldsymbol{\beta}_S$. Let $p$ be the number of predictors in the full model, including a constant. Let $q = p - 1$ be the number of nontrivial predictors in the full model. Let $a = a_I$ be the number of predictors in the submodel $I$, including a constant. Let $k = k_I = a_I - 1$ be the number of nontrivial predictors in the submodel. For submodel $I$, think of $I$ as indexing the predictors in the model, including the constant. Let $A$ index the nontrivial predictors in the model. Hence $I$ adds the constant

(trivial predictor) to the collection of nontrivial predictors in $A$. In Equation (7.1), there is a "true submodel" $\boldsymbol{Y} = \boldsymbol{X}_S \boldsymbol{\beta}_S + \boldsymbol{e}$ where all of the elements of $\boldsymbol{\beta}_S$ are nonzero but all of the elements of $\boldsymbol{\beta}$ that are not elements of $\boldsymbol{\beta}_S$ are zero. Then $a = a_S$ is the number of predictors in that submodel, including a constant, and $k = k_S$ is the number of active predictors = number of nonnoise variables = number of nontrivial predictors in the true model $S = I_S$. Then there are $p - a$ noise variables ($x_i$ that have coefficient $\beta_i = 0$) in the full model. The true model is generally only known in simulations. For Equation (7.1), we also assume that if $\boldsymbol{x}^T \boldsymbol{\beta} = \boldsymbol{x}_I^T \boldsymbol{\beta}_I$, then $S \subseteq I$. Hence $S$ is the unique smallest subset of predictors such that $\boldsymbol{x}^T \boldsymbol{\beta} = \boldsymbol{x}_S^T \boldsymbol{\beta}_S$. An alternative variable selection model was given by Remark 7.6.

Model selection generates $M$ models. Then a hopefully good model is selected from these $M$ models. Variable selection is a special case of model selection. Many methods for variable and model selection have been suggested for the MLR model. We will consider several $R$ functions including i) forward selection computed with the regsubsets function from the leaps library, ii) principal components regression (PCR) with the pcr function from the pls library, iii) partial least squares (PLS) with the plsr function from the pls library, iv) ridge regression with the cv.glmnet or glmnet function from the glmnet library, v) lasso with the cv.glmnet or glmnet function from the glmnet library, and vi) relaxed lasso which is OLS applied to the lasso active set (nontrivial predictors with nonzero coefficients) and a constant. See Sections 7.7–7.11, Olive (2020: ch. 3, 2021a: ch. 4), and James et al. (2013, ch. 6). For this chapter, PLS and PCR are MLR alternative MLR methods, but will not be discussed in detail.

These six methods produce $M$ models and use a criterion to select the final model (e.g. $C_p$ or 10-fold cross validation (CV)). The number of models $M$ depends on the method. Often one of the models is the full model (7.8) that uses all $p - 1$ nontrivial predictors. The full model is (approximately) fit with (ordinary) least squares. For one of the $M$ models, some of the methods use $\hat{\boldsymbol{\eta}} = \boldsymbol{0}$ and fit the model $Y_i = \beta_1 + e_i$ with $\hat{Y}_i \equiv \overline{Y}$ that uses none of the nontrivial predictors. Forward selection, PCR, and PLS use variables $v_1 = 1$ (the constant or trivial predictor) and $v_j = \boldsymbol{\gamma}_j^T \boldsymbol{x}$ that are linear combinations of the predictors for $j = 2, ..., p$. Model $I_i$ uses variables $v_1, v_2, ..., v_i$ for $i = 1, ..., M$ where $M \le p$ and often $M \le \min(p, n/10)$. Then $M$ models $I_i$ are used. (For forward selection and PCR, OLS is used to regress $Y$ (or $Z$) on $v_1, ..., v_i$.) Then a criterion chooses the final submodel $I_d$ from candidates $I_1, ..., I_M$.

**Remark 7.14.** Prediction interval (7.34) used a number $d$ that was often the number of predictors in the selected model. For forward selection, PCR, PLS, lasso, and lasso variable selection, let $d$ be the number of predictors $v_j = \boldsymbol{\gamma}_j^T \boldsymbol{x}$ in the final model (with nonzero coefficients), including a constant $v_1$. For forward selection, lasso, and lasso variable selection, $v_j$ corresponds to a single nontrivial predictor, say $v_j = x_j^* = x_{k_j}$. Another method for

obtaining $d$ is to let $d = j$ if $j$ is the degrees of freedom of the selected model if that model was chosen in advance without model or variable selection. Hence $d = j$ is not the model degrees of freedom if model selection was used.

Overfitting or "fitting noise" occurs when there is not enough data to estimate the $p \times 1$ vector $\boldsymbol{\beta}$ well with the estimation method, such as OLS. The OLS model is overfitting if $n < 5p$. When $n > p$, $\boldsymbol{X}$ is not invertible, but if $n = p$, then $\hat{\boldsymbol{Y}} = \boldsymbol{HY} = \boldsymbol{X}(\boldsymbol{X}^T\boldsymbol{X})^{-1}\boldsymbol{X}^T\boldsymbol{Y} = \boldsymbol{I}_n\boldsymbol{Y} = \boldsymbol{Y}$ regardless of how bad the predictors are. If $n < p$, then the OLS program fails or $\hat{\boldsymbol{Y}} = \boldsymbol{Y}$: the fitted regression plane interpolates the training data response variables $Y_1, ..., Y_n$. The following rule of thumb is useful for many regression methods. Note that $d = p$ for the full OLS model.

**Rule of thumb 7.2.** We want $n \geq 10d$ to avoid overfitting. Occasionally $n$ as low as $5d$ is used, but models with $n < 5d$ are overfitting.

**Remark 7.15.** Use $\boldsymbol{Z}_n \sim AN_r(\boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n)$ to indicate that a normal approximation is used: $\boldsymbol{Z}_n \approx N_r(\boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n)$. Let $a$ be a constant, let $\boldsymbol{A}$ be a $k \times r$ constant matrix (often with full rank $k \leq r$), and let $\boldsymbol{c}$ be a $k \times 1$ constant vector. If $\sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}) \xrightarrow{D} N_r(\boldsymbol{0}, \boldsymbol{V})$, then $a\boldsymbol{Z}_n = a\boldsymbol{I}_r\boldsymbol{Z}_n$ with $\boldsymbol{A} = a\boldsymbol{I}_r$,

$$a\boldsymbol{Z}_n \sim AN_r\left(a\boldsymbol{\mu}_n, a^2\boldsymbol{\Sigma}_n\right), \quad \text{and} \quad \boldsymbol{A}\boldsymbol{Z}_n + \boldsymbol{c} \sim AN_k\left(\boldsymbol{A}\boldsymbol{\mu}_n + \boldsymbol{c}, \boldsymbol{A}\boldsymbol{\Sigma}_n\boldsymbol{A}^T\right),$$

$$\hat{\boldsymbol{\theta}}_n \sim AN_r\left(\boldsymbol{\theta}, \frac{\boldsymbol{V}}{n}\right), \quad \text{and} \quad \boldsymbol{A}\hat{\boldsymbol{\theta}}_n + \boldsymbol{c} \sim AN_k\left(\boldsymbol{A}\boldsymbol{\theta} + \boldsymbol{c}, \frac{\boldsymbol{A}\boldsymbol{V}\boldsymbol{A}^T}{n}\right).$$

Theorem 5.9 gives the large sample theory for the OLS full model. Then $\hat{\boldsymbol{\beta}} \approx N_p(\boldsymbol{\beta}, \sigma^2(\boldsymbol{X}^T\boldsymbol{X})^{-1}))$ or $\hat{\boldsymbol{\beta}} \sim AN_p(\boldsymbol{\beta}, MSE(\boldsymbol{X}^T\boldsymbol{X})^{-1}))$.

When minimizing or maximizing a real valued function $Q(\boldsymbol{\eta})$ of the $k \times 1$ vector $\boldsymbol{\eta}$, the solution $\hat{\boldsymbol{\eta}}$ is found by setting the gradient of $Q(\boldsymbol{\eta})$ equal to $\boldsymbol{0}$. The following definition and lemma follow Graybill (1983, pp. 351-352) closely. Maximum likelihood estimators are examples of estimating equations. There is a vector of parameters $\boldsymbol{\eta}$, and the gradient of the log likelihood function $\log L(\boldsymbol{\eta})$ is set to zero. The solution $\hat{\boldsymbol{\eta}}$ is the MLE, an estimator of the parameter vector $\boldsymbol{\eta}$, but in the log likelihood, $\boldsymbol{\eta}$ is a dummy variable vector, not the fixed unknown parameter vector.

**Definition 7.7.** Let $Q(\boldsymbol{\eta})$ be a real valued function of the $k \times 1$ vector $\boldsymbol{\eta}$. The gradient of $Q(\boldsymbol{\eta})$ is the $k \times 1$ vector

$$\bigtriangledown Q = \bigtriangledown Q(\boldsymbol{\eta}) = \frac{\partial Q}{\partial \boldsymbol{\eta}} = \frac{\partial Q(\boldsymbol{\eta})}{\partial \boldsymbol{\eta}} = \begin{bmatrix} \frac{\partial}{\partial \eta_1} Q(\boldsymbol{\eta}) \\ \frac{\partial}{\partial \eta_2} Q(\boldsymbol{\eta}) \\ \vdots \\ \frac{\partial}{\partial \eta_k} Q(\boldsymbol{\eta}) \end{bmatrix}.$$

Suppose there is a model with unknown parameter vector $\boldsymbol{\eta}$. A set of *estimating equations* $f(\boldsymbol{\eta})$ is used to maximize or minimize $Q(\boldsymbol{\eta})$ where $\boldsymbol{\eta}$ is a dummy variable vector.

Often $f(\boldsymbol{\eta}) = \bigtriangledown Q$, and we solve $f(\boldsymbol{\eta}) = \bigtriangledown Q \overset{set}{=} \boldsymbol{0}$ for the solution $\hat{\boldsymbol{\eta}}$, and $f : \mathbb{R}^k \rightarrow \mathbb{R}^k$. Note that $\hat{\boldsymbol{\eta}}$ is an estimator of the unknown parameter vector $\boldsymbol{\eta}$ in the model, but $\boldsymbol{\eta}$ is a dummy variable in $Q(\boldsymbol{\eta})$. Hence we could use $Q(\boldsymbol{b})$ instead of $Q(\boldsymbol{\eta})$, but the solution of the estimating equations would still be $\hat{\boldsymbol{b}} = \hat{\boldsymbol{\eta}}$.

As a mnemonic (memory aid) for the following theorem, note that the derivative $\frac{d}{dx} ax = \frac{d}{dx} xa = a$ and $\frac{d}{dx} ax^2 = \frac{d}{dx} xax = 2ax$.

**Theorem 7.5.** a) If $Q(\boldsymbol{\eta}) = \boldsymbol{a}^T \boldsymbol{\eta} = \boldsymbol{\eta}^T \boldsymbol{a}$ for some $k \times 1$ constant vector $\boldsymbol{a}$, then $\bigtriangledown Q = \boldsymbol{a}$.

b) If $Q(\boldsymbol{\eta}) = \boldsymbol{\eta}^T \boldsymbol{A} \boldsymbol{\eta}$ for some $k \times k$ constant matrix $\boldsymbol{A}$, then $\bigtriangledown Q = 2\boldsymbol{A}\boldsymbol{\eta}$.

c) If $Q(\boldsymbol{\eta}) = \sum_{i=1}^{k} |\eta_i| = \|\boldsymbol{\eta}\|_1$, then $\bigtriangledown Q = \boldsymbol{s} = \boldsymbol{s}_{\boldsymbol{\eta}}$ where $s_i = \text{sign}(\eta_i)$ where $\text{sign}(\eta_i) = 1$ if $\eta_i > 0$ and $\text{sign}(\eta_i) = -1$ if $\eta_i < 0$. This gradient is only defined for $\boldsymbol{\eta}$ where none of the $k$ values of $\eta_i$ are equal to 0.

**Example 7.4.** If $\boldsymbol{Z} = \boldsymbol{W}\boldsymbol{\eta} + \boldsymbol{e}$, then the OLS estimator minimizes $Q(\boldsymbol{\eta}) = \|\boldsymbol{Z} - \boldsymbol{W}\boldsymbol{\eta}\|_2^2 = (\boldsymbol{Z} - \boldsymbol{W}\boldsymbol{\eta})^T (\boldsymbol{Z} - \boldsymbol{W}\boldsymbol{\eta}) = \boldsymbol{Z}^T \boldsymbol{Z} - 2\boldsymbol{Z}^T \boldsymbol{W}\boldsymbol{\eta} + \boldsymbol{\eta}^T (\boldsymbol{W}^T \boldsymbol{W})\boldsymbol{\eta}$. Using Theorem 7.5 with $\boldsymbol{a}^T = \boldsymbol{Z}^T \boldsymbol{W}$ and $\boldsymbol{A} = \boldsymbol{W}^T \boldsymbol{W}$ shows that $\bigtriangledown Q = -2\boldsymbol{W}^T \boldsymbol{Z} + 2(\boldsymbol{W}^T \boldsymbol{W})\boldsymbol{\eta}$. Let $\bigtriangledown Q(\hat{\boldsymbol{\eta}})$ denote the gradient evaluated at $\hat{\boldsymbol{\eta}}$. Then the OLS estimator satisfies the normal equations $(\boldsymbol{W}^T \boldsymbol{W})\hat{\boldsymbol{\eta}} = \boldsymbol{W}^T \boldsymbol{Z}$.

**Example 7.5.** The Hebbler (1847) data was collected from $n = 26$ districts in Prussia in 1843. We will study the relationship between $Y = $ the *number of women married to civilians* in the district with the predictors $x_1$ $= $ constant, $x_2 = pop = $ the *population of the district in 1843*, $x_3 = mmen$ $= $ the *number of married civilian men* in the district, $x_4 = mmilmen = $ the *number of married men in the military* in the district, and $x_5 = milwmn = $ the *number of women married to husbands in the military* in the district. Sometimes the person conducting the survey would not count a spouse if the spouse was not at home. Hence $Y$ is highly correlated but not equal to $x_3$. Similarly, $x_4$ and $x_5$ are highly correlated but not equal. We expect that $Y = x_3 + e$ is a good model, but $n/p = 5.2$ is small. See the following output.

```
ls.print(out)
Residual Standard Error=392.8709
```

```
R-Square=0.9999, p-value=0
F-statistic (df=4, 21)=67863.03
          Estimate  Std.Err t-value Pr(>|t|)
Intercept 242.3910 263.7263  0.9191   0.3685
pop         0.0004   0.0031  0.1130   0.9111
mmen        0.9995   0.0173 57.6490   0.0000
mmilmen    -0.2328   2.6928 -0.0864   0.9319
milwmn      0.1531   2.8231  0.0542   0.9572
res<-out$res
yhat<-Y-res #d = 5 predictors used including x_1
AERplot2(yhat,Y,res=res,d=5)
#response plot with 90% pointwise PIs
$respi #90% PI for a future residual
[1] -950.4811 1445.2584 #90% PI length = 2395.74
```

## 7.7 Forward Selection

Variable selection methods such as forward selection were covered in Sections 7.2–7.4 where model $I_j$ uses $j$ predictors $x_1^*, ..., x_j^*$ including the constant $x_1^* \equiv 1$. If $n/p$ is not large, forward selection can be done as in Section 7.2 except instead of forming $p$ submodels $I_1, ..., I_p$, form the sequence of $M$ submodels $I_1, ..., I_M$ where $M = \min(\lceil n/J \rceil, p)$ for some positive integer $J$ such as $J = 5, 10$, or 20. Here $\lceil x \rceil$ is the smallest integer $\geq x$, e.g., $\lceil 7.7 \rceil = 8$. Then for each submodel $I_j$, OLS is used to regress $Y$ on $1, x_2^*, ..., x_j^*$. Then a criterion chooses which model $I_d$ from candidates $I_1, ..., I_M$ is to be used as the final submodel.

**Remark 7.16.** Suppose $n/J$ is an integer. If $p \leq n/J$, then forward selection fits $(p-1) + (p-2) + \cdots + 2 + 1 = p(p-1)/2 \approx p^2/2$ models, where $p - i$ models are fit at step $i$ for $i = 1, ..., (p-1)$. If $n/J < p$, then forward selection uses $(n/J) - 1$ steps and fits $\approx (p-1) + (p-2) + \cdots + (p-(n/J)+1) = p((n/J) - 1) - (1 + 2 + \cdots + ((n/J) - 1)) =$

$$p\left(\frac{n}{J} - 1\right) - \frac{\frac{n}{J}\left(\frac{n}{J} - 1\right)}{2} \approx \frac{n}{J} \; \frac{\left(2p - \frac{n}{J}\right)}{2}$$

models. Thus forward selection can be slow if $n$ and $p$ are both large, although the $R$ package leaps uses a branch and bound algorithm that likely eliminates many of the possible fits. Note that after step $i$, the model has $i + 1$ predictors, including the constant.

The $R$ function regsubsets can be used for forward selection if $p < n$, and if $p \geq n$ if the maximum number of variables is less than $n$. Then warning messages are common. Some $R$ code is shown below.

```
#regsubsets works if p < n, e.g. p = n-1, and works
#if p > n with warnings if nvmax is small enough
set.seed(13)
n<-100
p<-200
k<-19 #the first 19 nontrivial predictors are active
J<-5
q <- p-1
b <- 0 * 1:q
b[1:k] <- 1 #beta = (1, 1, ..., 1, 0, 0, ..., 0)^T
x <- matrix(rnorm(n * q), nrow = n, ncol = q)
y <- 1 + x %*% b + rnorm(n)
nc <- ceiling(n/J)-1 #the constant will also be used
nc <- min(nc,q)
nc <- max(nc,1) #nc is the maximum number of
#nontrivial predictors used by forward selection
pp <- nc+1  #d = pp is used for PI (4.14)
vars <- as.vector(1:(p-1))
temp<-regsubsets(x,y,nvmax=nc,method="forward")
out<-summary(temp)
num <- length(out$cp)
mod <- out$which[num,] #use the last model
#do not need the constant in vin
vin <- vars[mod[-1]]

out$rss
 [1] 1496.49625 1342.95915 1214.93174 1068.56668
     973.36395  855.15436  745.35007  690.03901
     638.40677  590.97644  542.89273  503.68666
     467.69423  420.94132  391.41961  328.62016
     242.66311  178.77573   79.91771
out$bic
 [1]    -9.4032  -15.6232  -21.0367  -29.2685
       -33.9949  -42.3374  -51.4750  -54.5804
       -57.7525  -60.8673  -64.7485  -67.6391
       -70.4479  -76.3748  -79.0410  -91.9236
     -117.6413 -143.5903 -219.498595
tem <- lsfit(x[,1:19],y) #last model used the
sum(tem$resid^2)        #first 19 predictors
[1] 79.91771            #SSE(I) = RSS(I)
n*log(out$rss[19]/n) + 20*log(n)
[1] 69.68613            #BIC(I)
for(i in 1:19)   #a formula for BIC(I)
print( n*log(out$rss[i]/n) + (i+1)*log(n) )
bic <- c(279.7815, 273.5616, 268.1480, 259.9162,
255.1898, 246.8474, 237.7097, 234.6043, 231.4322,
```

```
228.3175, 224.4362, 221.5456, 218.7368, 212.8099,
210.1437, 197.2611, 171.5435, 145.5944,  69.6861)
tem<-lsfit(bic,out$bic)
tem$coef
    Intercept            X
-289.1846831    0.9999998 #bic - 289.1847 = out$bic
xx <- 1:min(length(out$bic),p-1)+1
ebic <- out$bic+2*log(dbinom(x=xx,size=p,prob=0.5))
#actually EBIC(I) - 2 p log(2).
```

**Example 7.5**, continued. The output below shows results from forward selection for the marry data. The minimum $C_p$ model $I_{min}$ uses a constant and *mmem*. The forward selection PIs are shorter than the OLS full model PIs.

```
library(leaps);Y <- marry[,3]; X <- marry[,-3]
temp<-regsubsets(X,Y,method="forward")
out<-summary(temp)
Selection Algorithm: forward
        pop mmen mmilmen milwmn
1  ( 1 ) " " "*"   " "     " "
2  ( 1 ) " " "*"   "*"     " "
3  ( 1 ) "*" "*"   "*"     " "
4  ( 1 ) "*" "*"   "*"     "*"
out$cp
[1] -0.8268967  1.0151462  3.0029429  5.0000000
#mmen and a constant = Imin
mincp <- out$which[out$cp==min(out$cp),]
#do not need the constant in vin
vin <- vars[mincp[-1]]
sub <- lsfit(X[,vin],Y)
ls.print(sub)
Residual Standard Error=369.0087
R-Square=0.9999
F-statistic (df=1, 24)=307694.4
          Estimate  Std.Err   t-value Pr(>|t|)
Intercept 241.5445 190.7426   1.2663   0.2175
X           1.0010   0.0018 554.7021   0.0000
res<-sub$res
yhat<-Y-res #d = 2 predictors used including x_1
AERplot2(yhat,Y,res=res,d=2)
#response plot with 90% pointwise PIs
$respi   #90% PI for a future residual
[1] -778.2763 1336.4416 #length 2114.72
```

Consider forward selection where $\boldsymbol{x}_I$ is $a \times 1$. Underfitting occurs if $S$ is not a subset of $I$ so $\boldsymbol{x}_I$ is missing important predictors. A special case

of underfitting is $d = a < a_S$. Overfitting for forward selection occurs if i) $n < 5a$ so there is not enough data to estimate the $a$ parameters in $\boldsymbol{\beta}_I$ well, or ii) $S \subseteq I$ but $S \neq I$. Overfitting is serious if $n < 5a$, but "not much of a problem" if $n > Jp$ where $J = 10$ or 20 for many data sets. Underfitting is a serious problem. Let $Y_i = \boldsymbol{x}_{I,i}^T \boldsymbol{\beta}_I + e_{I,i}$. Then $V(e_{I,i})$ may not be a constant $\sigma^2$: $V(e_{I,i})$ could depend on case $i$, and the model may no longer be linear. Check model $I$ with response and residual plots.

Forward selection is a *shrinkage* method: $p$ models are produced and except for the full model, some $|\hat{\beta}_i|$ are shrunk to 0. Lasso and ridge regression are also shrinkage methods. Ridge regression is a shrinkage method, but $|\hat{\beta}_i|$ is not shrunk to 0. Shrinkage methods that shrink $\hat{\beta}_i$ to 0 are also variable selection methods. See Sections 7.8, 7.9, and 7.11.

**Definition 7.8.** Suppose the population MLR model has $\boldsymbol{\beta}_S$ an $a_S \times 1$ vector. The population MLR model is *sparse* if $a_S$ is small. The population MLR model is *dense* or abundant if $n/a_S < J$ where $J = 5$ or $J = 10$, say. The fitted model $\hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}}_{I_{min},0}$ is *sparse* if $d =$ number of nonzero coefficients is small. The fitted model is *dense* if $n/d < J$ where $J = 5$ or $J = 10$.

## 7.8 Ridge Regression

Consider the MLR model $\boldsymbol{Y} = \boldsymbol{X}\boldsymbol{\beta} + \boldsymbol{e}$. Ridge regression uses the centered response $Z_i = Y_i - \overline{Y}$ and standardized nontrivial predictors in the model $\boldsymbol{Z} = \boldsymbol{W}\boldsymbol{\eta} + \boldsymbol{e}$. Then $\hat{Y}_i = \hat{Z}_i + \overline{Y}$. Note that in Definition 7.10, $\lambda_{1,n}$ is a tuning parameter, not an eigenvalue. The residuals $\boldsymbol{r} = \boldsymbol{r}(\hat{\boldsymbol{\beta}}_R) = \boldsymbol{Y} - \hat{\boldsymbol{Y}}$. Refer to Definition 7.6 for the OLS estimator $\hat{\boldsymbol{\eta}}_{OLS} = (\boldsymbol{W}^T\boldsymbol{W})^{-1}\boldsymbol{W}^T\boldsymbol{Z}$.

**Definition 7.9.** Consider the MLR model $\boldsymbol{Z} = \boldsymbol{W}\boldsymbol{\eta} + \boldsymbol{e}$. Let $\boldsymbol{b}$ be a $(p-1) \times 1$ vector. Then the fitted value $\hat{Z}_i(\boldsymbol{b}) = \boldsymbol{w}_i^T\boldsymbol{b}$ and the residual $r_i(\boldsymbol{b}) = Z_i - \hat{Z}_i(\boldsymbol{b})$. The vector of fitted values $\hat{\boldsymbol{Z}}(\boldsymbol{b}) = \boldsymbol{W}\boldsymbol{b}$ and the vector of residuals $\boldsymbol{r}(\boldsymbol{b}) = \boldsymbol{Z} - \hat{\boldsymbol{Z}}(\boldsymbol{b})$.

**Definition 7.10.** Consider fitting the MLR model $\boldsymbol{Y} = \boldsymbol{X}\boldsymbol{\beta} + \boldsymbol{e}$ using $\boldsymbol{Z} = \boldsymbol{W}\boldsymbol{\eta} + \boldsymbol{e}$. Let $\lambda \geq 0$ be a constant. The *ridge regression estimator* $\hat{\boldsymbol{\eta}}_R$ minimizes the *ridge regression criterion*

$$Q_R(\boldsymbol{\eta}) = \frac{1}{a}(\boldsymbol{Z} - \boldsymbol{W}\boldsymbol{\eta})^T(\boldsymbol{Z} - \boldsymbol{W}\boldsymbol{\eta}) + \frac{\lambda_{1,n}}{a}\sum_{i=1}^{p-1}\eta_i^2 \qquad (7.15)$$

over all vectors $\boldsymbol{\eta} \in \mathbb{R}^{p-1}$ where $\lambda_{1,n} \geq 0$ and $a > 0$ are known constants with $a = 1, 2, n$, and $2n$ common. Then

$$\hat{\boldsymbol{\eta}}_R = (\boldsymbol{W}^T\boldsymbol{W} + \lambda_{1,n}\boldsymbol{I}_{p-1})^{-1}\boldsymbol{W}^T\boldsymbol{Z}. \qquad (7.16)$$

The residual sum of squares $RSS(\boldsymbol{\eta}) = (\boldsymbol{Z} - \boldsymbol{W}\boldsymbol{\eta})^T(\boldsymbol{Z} - \boldsymbol{W}\boldsymbol{\eta})$, and $\lambda_{1,n} = 0$ corresponds to the OLS estimator $\hat{\boldsymbol{\eta}}_{OLS}$. The ridge regression vector of fitted values is $\hat{\boldsymbol{Z}} = \hat{\boldsymbol{Z}}_R = \boldsymbol{W}\hat{\boldsymbol{\eta}}_R$, and the ridge regression vector of residuals $\boldsymbol{r}_R = \boldsymbol{r}(\hat{\boldsymbol{\eta}}_R) = \boldsymbol{Z} - \hat{\boldsymbol{Z}}_R$. The estimator is said to be *regularized* if $\lambda_{1,n} > 0$. Obtain $\hat{\boldsymbol{Y}}$ and $\hat{\boldsymbol{\beta}}_R$ using $\hat{\boldsymbol{\eta}}_R$, $\hat{\boldsymbol{Z}}$, and $\overline{\boldsymbol{Y}}$.

Using a vector of parameters $\boldsymbol{\eta}$ and a dummy vector $\boldsymbol{\eta}$ in $Q_R$ is common for minimizing a criterion $Q(\boldsymbol{\eta})$, often with estimating equations. See the paragraphs above and below Definition 7.7. We could also write

$$Q_R(\boldsymbol{b}) = \frac{1}{a}\boldsymbol{r}(\boldsymbol{b})^T\boldsymbol{r}(\boldsymbol{b}) + \frac{\lambda_{1,n}}{a}\boldsymbol{b}^T\boldsymbol{b}$$

where the minimization is over all vectors $\boldsymbol{b} \in \mathbb{R}^{p-1}$. Note that $\sum_{i=1}^{p-1}\eta_i^2 = \boldsymbol{\eta}^T\boldsymbol{\eta} = \|\boldsymbol{\eta}\|_2^2$. The literature often uses $\lambda_a = \lambda = \lambda_{1,n}/a$.

Note that $\lambda_{1,n}\boldsymbol{b}^T\boldsymbol{b} = \lambda_{1,n}\sum_{i=1}^{p-1}b_i^2$. Each coefficient $b_i$ is penalized equally by $\lambda_{1,n}$. Hence using standardized nontrivial predictors makes sense so that if $\eta_i$ is large in magnitude, then the standardized variable $w_i$ is important.

**Remark 7.17.** i) If $\lambda_{1,n} = 0$, the ridge regression estimator becomes the OLS full model estimator: $\hat{\boldsymbol{\eta}}_R = \hat{\boldsymbol{\eta}}_{OLS}$.

ii) If $\lambda_{1,n} > 0$, then $\boldsymbol{W}^T\boldsymbol{W} + \lambda_{1,n}\boldsymbol{I}_{p-1}$ is nonsingular. Hence $\hat{\boldsymbol{\eta}}_R$ exists even if $\boldsymbol{X}$ and $\boldsymbol{W}$ are singular or ill conditioned, or if $p > n$.

iii) Following Hastie et al. (2009, p. 96), let the augmented matrix $\boldsymbol{W}_A$ and the augmented response vector $\boldsymbol{Z}_A$ be defined by

$$\boldsymbol{W}_A = \begin{pmatrix} \boldsymbol{W} \\ \sqrt{\lambda_{1,n}} \ \boldsymbol{I}_{p-1} \end{pmatrix}, \quad \text{and} \quad \boldsymbol{Z}_A = \begin{pmatrix} \boldsymbol{Z} \\ \boldsymbol{0} \end{pmatrix},$$

where $\boldsymbol{0}$ is the $(p-1) \times 1$ zero vector. For $\lambda_{1,n} > 0$, the OLS estimator from regressing $\boldsymbol{Z}_A$ on $\boldsymbol{W}_A$ is

$$\hat{\boldsymbol{\eta}}_A = (\boldsymbol{W}_A^T\boldsymbol{W}_A)^{-1}\boldsymbol{W}_A^T\boldsymbol{Z}_A = \hat{\boldsymbol{\eta}}_R$$

since $\boldsymbol{W}_A^T\boldsymbol{Z}_A = \boldsymbol{W}^T\boldsymbol{Z}$ and

$$\boldsymbol{W}_A^T\boldsymbol{W}_A = \begin{pmatrix} \boldsymbol{W}^T & \sqrt{\lambda_{1,n}} \ \boldsymbol{I}_{p-1} \end{pmatrix} \begin{pmatrix} \boldsymbol{W} \\ \sqrt{\lambda_{1,n}} \ \boldsymbol{I}_{p-1} \end{pmatrix} = \boldsymbol{W}^T\boldsymbol{W} + \lambda_{1,n} \ \boldsymbol{I}_{p-1}.$$

iv) A simple way to regularize a regression estimator, such as the $L_1$ estimator, is to compute that estimator from regressing $\boldsymbol{Z}_A$ on $\boldsymbol{W}_A$.

Remark 7.17 iii) is interesting. Note that for $\lambda_{1,n} > 0$, the $(n+p-1) \times (p-1)$ matrix $\boldsymbol{W}_A$ has full rank $p-1$. The augmented OLS model consists of adding $p-1$ pseudo-cases $(\boldsymbol{w}_{n+1}^T, Z_{n+1})^T, ..., (\boldsymbol{w}_{n+p-1}^T, Z_{n+p-1})^T$ where $Z_j = 0$ and $\boldsymbol{w}_j = (0, ..., \sqrt{\lambda_{1,n}}, 0, ..., 0)^T$ for $j = n+1, ..., n+p-1$ where the nonzero entry

is in the $k$th position if $j = n + k$. For centered response and standardized nontrivial predictors, the population OLS regression fit runs through the origin $(\boldsymbol{w}^T, Z)^T = (\boldsymbol{0}^T, 0)^T$. Hence for $\lambda_{1,n} = 0$, the augmented OLS model adds $p - 1$ typical cases at the origin. If $\lambda_{1,n}$ is not large, then the pseudo-data can still be regarded as typical cases. If $\lambda_{1,n}$ is large, the pseudo-data act as $w$–outliers (outliers in the standardized predictor variables), and the OLS slopes go to zero as $\lambda_{1,n}$ gets large, making $\hat{\boldsymbol{Z}} \approx \boldsymbol{0}$ so $\hat{\boldsymbol{Y}} \approx \overline{\boldsymbol{Y}}$.

To prove Remark 7.17 ii), let $(\psi, \boldsymbol{g})$ be an eigenvalue eigenvector pair of $\boldsymbol{W}^T \boldsymbol{W} = n \boldsymbol{R_u}$. Then $[\boldsymbol{W}^T \boldsymbol{W} + \lambda_{1,n} \boldsymbol{I}_{p-1}] \boldsymbol{g} = (\psi + \lambda_{1,n}) \boldsymbol{g}$, and $(\psi + \lambda_{1,n}, \boldsymbol{g})$ is an eigenvalue eigenvector pair of $\boldsymbol{W}^T \boldsymbol{W} + \lambda_{1,n} \boldsymbol{I}_{p-1} > 0$ provided $\lambda_{1,n} > 0$.

The degrees of freedom for a ridge regression with known $\lambda_{1,n}$ is also interesting and will be found in the next paragraph. The sample correlation matrix of the nontrivial predictors

$$\boldsymbol{R_u} = \frac{1}{n - g} \boldsymbol{W}_g^T \boldsymbol{W}_g$$

where we will use $g = 0$ and $\boldsymbol{W} = \boldsymbol{W}_0$. Then $\boldsymbol{W}^T \boldsymbol{W} = n \boldsymbol{R_u}$. By singular value decomposition (SVD) theory, the SVD of $\boldsymbol{W}$ is $\boldsymbol{W} = \boldsymbol{U} \boldsymbol{\Lambda} \boldsymbol{V}^T$ where the positive singular values $\sigma_i$ are square roots of the positive eigenvalues of both $\boldsymbol{W}^T \boldsymbol{W}$ and of $\boldsymbol{W} \boldsymbol{W}^T$. Also $\boldsymbol{V} = (\hat{\boldsymbol{e}}_1 \ \hat{\boldsymbol{e}}_2 \ \cdots \ \hat{\boldsymbol{e}}_p)$, and $\boldsymbol{W}^T \boldsymbol{W} \hat{\boldsymbol{e}}_i = \sigma_i^2 \hat{\boldsymbol{e}}_i$. Hence $\hat{\lambda}_i = \sigma_i^2$ where $\hat{\lambda}_i = \hat{\lambda}_i(\boldsymbol{W}^T \boldsymbol{W})$ is the $i$th eigenvalue of $\boldsymbol{W}^T \boldsymbol{W}$, and $\hat{\boldsymbol{e}}_i$ is the $i$th orthonormal eigenvector of $\boldsymbol{R_u}$ and of $\boldsymbol{W}^T \boldsymbol{W}$. The SVD of $\boldsymbol{W}^T$ is $\boldsymbol{W}^T = \boldsymbol{V} \boldsymbol{\Lambda}^T \boldsymbol{U}^T$, and the *Gram matrix*

$$\boldsymbol{W} \boldsymbol{W}^T = \begin{bmatrix} \boldsymbol{w}_1^T \boldsymbol{w}_1 \ \boldsymbol{w}_1^T \boldsymbol{w}_2 \ \ldots \ \boldsymbol{w}_1^T \boldsymbol{w}_n \\ \vdots \quad \vdots \quad \ddots \quad \vdots \\ \boldsymbol{w}_n^T \boldsymbol{w}_1 \ \boldsymbol{w}_n^T \boldsymbol{w}_2 \ \ldots \ \boldsymbol{w}_n^T \boldsymbol{w}_n \end{bmatrix}$$

which is the matrix of scalar products. **Warning:** Note that $\sigma_i$ is the $i$th singular value of $\boldsymbol{W}$, not the standard deviation of $w_i$.

Following Hastie et al. (2009, p. 68), if $\hat{\lambda}_i = \hat{\lambda}_i(\boldsymbol{W}^T \boldsymbol{W})$ is the $i$th eigenvalue of $\boldsymbol{W}^T \boldsymbol{W}$ where $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \cdots \geq \hat{\lambda}_{p-1}$, then the (effective) degrees of freedom for the ridge regression of $\boldsymbol{Z}$ on $\boldsymbol{W}$ with known $\lambda_{1,n}$ is $df(\lambda_{1,n}) =$

$$tr[\boldsymbol{W} (\boldsymbol{W}^T \boldsymbol{W} + \lambda_{1,n} \boldsymbol{I}_{p-1})^{-1} \boldsymbol{W}^T] = \sum_{i=1}^{p-1} \frac{\sigma_i^2}{\sigma_i^2 + \lambda_{1,n}} = \sum_{i=1}^{p-1} \frac{\hat{\lambda}_i}{\hat{\lambda}_i + \lambda_{1,n}} \quad (7.17)$$

where the trace of a square $(p - 1) \times (p - 1)$ matrix $\boldsymbol{A} = (a_{ij})$ is $tr(\boldsymbol{A}) = \sum_{i=1}^{p-1} a_{ii} = \sum_{i=1}^{p-1} \hat{\lambda}_i(\boldsymbol{A})$. Note that the trace of $\boldsymbol{A}$ is the sum of the diagonal elements of $\boldsymbol{A}$ = the sum of the eigenvalues of $\boldsymbol{A}$.

Note that $0 \le df(\lambda_{1,n}) \le p-1$ where $df(\lambda_{1,n}) = p-1$ if $\lambda_{1,n} = 0$ and $df(\lambda_{1,n}) \to 0$ as $\lambda_{1,n} \to \infty$. The $R$ code below illustrates how to compute ridge regression degrees of freedom.

```
set.seed(13)
n<-100; q<-3   #q = p-1
b <- 0 * 1:q + 1
u <- matrix(rnorm(n * q), nrow = n, ncol = q)
y <- 1 + u %*% b + rnorm(n) #make MLR model
w1 <- scale(u) #t(w1) %*% w1 = (n-1) R = (n-1)*cor(u)
w <- sqrt(n/(n-1))*w1    #t(w) %*% w = n R = n cor(u)
t(w) %*% w/n
              [,1]         [,2]         [,3]
[1,]   1.00000000 -0.04826094 -0.06726636
[2,]  -0.04826094  1.00000000 -0.12426268
[3,]  -0.06726636 -0.12426268  1.00000000
cor(u) #same as above
rs <- t(w)%*%w #scaled correlation matrix n R
svs <-svd(w)$d  #singular values of w
lambda <- 0
d <- sum(svs^2/(svs^2+lambda))
#effective df for ridge regression using w
d
[1] 3   #= q = p-1
112.60792 103.88089  83.51119
svs^2 #as above
uu<-scale(u,scale=F) #centered but not scaled
svs <-svd(uu)$d #singular values of uu
svs^2
[1] 135.78205 108.85903  85.83395
d <- sum(svs^2/(svs^2+lambda))
#effective df for ridge regression using uu
#d is again 3 if lambda = 0
```

In general, if $\hat{\boldsymbol{Z}} = \boldsymbol{H}_\lambda \boldsymbol{Z}$, then $df(\hat{\boldsymbol{Z}}) = tr(\boldsymbol{H}_\lambda)$ where $\boldsymbol{H}_\lambda$ is a $(p-1) \times (p-1)$ "hat matrix." For computing $\hat{\boldsymbol{Y}}$, $df(\hat{\boldsymbol{Y}}) = df(\hat{\boldsymbol{Z}}) + 1$ since a constant $\hat{\boldsymbol{\beta}}_1$ also needs to be estimated. These formulas for degrees of freedom assume that $\lambda$ is known before fitting the model. The formulas do not give the model degrees of freedom if $\hat{\lambda}$ is selected from $M$ values $\lambda_1, ..., \lambda_M$ using a criterion such as $k$-fold cross validation.

Suppose the ridge regression criterion is written, using $a = 2n$, as

$$Q_{R,n}(\boldsymbol{b}) = \frac{1}{2n}\boldsymbol{r}(\boldsymbol{b})^T \boldsymbol{r}(\boldsymbol{b}) + \lambda_{2n}\boldsymbol{b}^T \boldsymbol{b}, \qquad (7.18)$$

as in Hastie et al. (2015, p. 10). Then $\lambda_{2n} = \lambda_{1,n}/(2n)$ using the $\lambda_{1,n}$ from (7.15).

The following remark is interesting if $\lambda_{1,n}$ and $p$ are fixed. However, $\hat{\lambda}_{1,n}$ is usually used, for example, after 10-fold cross validation. The fact that $\hat{\boldsymbol{\eta}}_R = \boldsymbol{A}_{n,\lambda}\hat{\boldsymbol{\eta}}_{OLS}$ appears in Efron and Hastie (2016, p. 98), and Marquardt and Snee (1975). See Theorem 7.6 for the ridge regression central limit theorem.

**Remark 7.18.** Ridge regression has a simple relationship with OLS if $n > p$ and $(\boldsymbol{W}^T\boldsymbol{W})^{-1}$ exists. Then $\hat{\boldsymbol{\eta}}_R = (\boldsymbol{W}^T\boldsymbol{W} + \lambda_{1,n}\boldsymbol{I}_{p-1})^{-1}\boldsymbol{W}^T\boldsymbol{Z} = (\boldsymbol{W}^T\boldsymbol{W}+\lambda_{1,n}\boldsymbol{I}_{p-1})^{-1}(\boldsymbol{W}^T\boldsymbol{W})(\boldsymbol{W}^T\boldsymbol{W})^{-1}\boldsymbol{W}^T\boldsymbol{Z} = \boldsymbol{A}_{n,\lambda}\hat{\boldsymbol{\eta}}_{OLS}$ where $\boldsymbol{A}_{n,\lambda} \equiv \boldsymbol{A}_n = (\boldsymbol{W}^T\boldsymbol{W} + \lambda_{1,n}\boldsymbol{I}_{p-1})^{-1}\boldsymbol{W}^T\boldsymbol{W}$. By the OLS CLT Equation (7.14) with $\hat{\boldsymbol{V}}/n = (\boldsymbol{W}^T\boldsymbol{W})^{-1}$, a normal approximation for OLS is

$$\hat{\boldsymbol{\eta}}_{OLS} \sim AN_{n-p}(\boldsymbol{\eta}, MSE\ (\boldsymbol{W}^T\boldsymbol{W})^{-1}).$$

Hence a normal approximation for ridge regression is

$$\hat{\boldsymbol{\eta}}_R \sim AN_{p-1}(\boldsymbol{A}_n\boldsymbol{\eta}, MSE\ \boldsymbol{A}_n(\boldsymbol{W}^T\boldsymbol{W})^{-1}\boldsymbol{A}_n^T) \sim$$

$$AN_{p-1}[\boldsymbol{A}_n\boldsymbol{\eta}, MSE\ (\boldsymbol{W}^T\boldsymbol{W} + \lambda_{1,n}\boldsymbol{I}_{p-1})^{-1}(\boldsymbol{W}^T\boldsymbol{W})(\boldsymbol{W}^T\boldsymbol{W} + \lambda_{1,n}\boldsymbol{I}_{p-1})^{-1}].$$

If Equation (7.14) holds and $\lambda_{1,n}/n \to 0$ as $n \to \infty$, then $\boldsymbol{A}_n \xrightarrow{P} \boldsymbol{I}_{p-1}$.

**Remark 7.19.** The ridge regression criterion from Definition 7.10 can also be defined by

$$Q_R(\boldsymbol{\eta}) = \|\boldsymbol{Z} - \boldsymbol{W}\boldsymbol{\eta}\|_2^2 + \lambda_{1,n}\boldsymbol{\eta}^T\boldsymbol{\eta}. \tag{7.19}$$

Then by Theorem 7.5, the gradient $\bigtriangledown Q_R = -2\boldsymbol{W}^T\boldsymbol{Z} + 2(\boldsymbol{W}^T\boldsymbol{W})\boldsymbol{\eta} + 2\lambda_{1,n}\boldsymbol{\eta}$. Cancelling constants and evaluating the gradient at $\hat{\boldsymbol{\eta}}_R$ gives the score equations

$$-\boldsymbol{W}^T(\boldsymbol{Z} - \boldsymbol{W}\hat{\boldsymbol{\eta}}_R) + \lambda_{1,n}\hat{\boldsymbol{\eta}}_R = \boldsymbol{0}. \tag{7.20}$$

Following Hastie and Efron (2016, pp. 381-382, 392), this means $\hat{\boldsymbol{\eta}}_R = \boldsymbol{W}^T\boldsymbol{a}$ for some $n \times 1$ vector $\boldsymbol{a}$. Hence $-\boldsymbol{W}^T(\boldsymbol{Z} - \boldsymbol{W}\boldsymbol{W}^T\boldsymbol{a}) + \lambda_{1,n}\boldsymbol{W}^T\boldsymbol{a} = \boldsymbol{0}$, or

$$\boldsymbol{W}^T(\boldsymbol{W}\boldsymbol{W}^T + \lambda_{1,n}\boldsymbol{I}_n)]\boldsymbol{a} = \boldsymbol{W}^T\boldsymbol{Z}$$

which has solution $\boldsymbol{a} = (\boldsymbol{W}\boldsymbol{W}^T + \lambda_{1,n}\boldsymbol{I}_n)^{-1}\boldsymbol{Z}$. Hence

$$\hat{\boldsymbol{\eta}}_R = \boldsymbol{W}^T\boldsymbol{a} = \boldsymbol{W}^T(\boldsymbol{W}\boldsymbol{W}^T + \lambda_{1,n}\boldsymbol{I}_n)^{-1}\boldsymbol{Z} = (\boldsymbol{W}^T\boldsymbol{W} + \lambda_{1,n}\boldsymbol{I}_{p-1})^{-1}\boldsymbol{W}^T\boldsymbol{Z}.$$

Using the $n \times n$ matrix $\boldsymbol{W}\boldsymbol{W}^T$ is computationally efficient if $p > n$ while using the $p \times p$ matrix $\boldsymbol{W}^T\boldsymbol{W}$ is computationally efficient if $n > p$. If $\boldsymbol{A}$ is $k \times k$, then computing $\boldsymbol{A}^{-1}$ has $O(k^3)$ complexity.

The following identity from Gunst and Mason (1980, p. 342) is useful for ridge regression inference: $\hat{\boldsymbol{\eta}}_R = (\boldsymbol{W}^T\boldsymbol{W} + \lambda_{1,n}\boldsymbol{I}_{p-1})^{-1}\boldsymbol{W}^T\boldsymbol{Z}$

$$= (\boldsymbol{W}^T\boldsymbol{W} + \lambda_{1,n}\boldsymbol{I}_{p-1})^{-1}\boldsymbol{W}^T\boldsymbol{W}(\boldsymbol{W}^T\boldsymbol{W})^{-1}\boldsymbol{W}^T\boldsymbol{Z}$$

$$= (\boldsymbol{W}^T\boldsymbol{W} + \lambda_{1,n}\boldsymbol{I}_{p-1})^{-1}\boldsymbol{W}^T\boldsymbol{W}\hat{\boldsymbol{\eta}}_{OLS} = \boldsymbol{A}_n\hat{\boldsymbol{\eta}}_{OLS} =$$

$$[\boldsymbol{I}_{p-1} - \lambda_{1,n}(\boldsymbol{W}^T\boldsymbol{W} + \lambda_{1,n}\boldsymbol{I}_{p-1})^{-1}]\hat{\boldsymbol{\eta}}_{OLS} = \boldsymbol{B}_n\hat{\boldsymbol{\eta}}_{OLS} =$$

$$\hat{\boldsymbol{\eta}}_{OLS} - \frac{\lambda_{1n}}{n}n(\boldsymbol{W}^T\boldsymbol{W} + \lambda_{1,n}\boldsymbol{I}_{p-1})^{-1}\hat{\boldsymbol{\eta}}_{OLS}$$

since $\boldsymbol{A}_n - \boldsymbol{B}_n = \boldsymbol{0}$. See Problem 7.7. Assume Equation (7.13) holds. If $\lambda_{1,n}/n \to 0$ then

$$\frac{\boldsymbol{W}^T\boldsymbol{W} + \lambda_{1,n}\boldsymbol{I}_{p-1}}{n} \xrightarrow{P} \boldsymbol{V}^{-1}, \quad \text{and} \quad n(\boldsymbol{W}^{\mathrm{T}}\boldsymbol{W} + \lambda_{1,\mathrm{n}}\boldsymbol{I}_{\mathrm{p}-1})^{-1} \xrightarrow{\mathrm{P}} \boldsymbol{V}.$$

Note that

$$\boldsymbol{A}_n = \boldsymbol{A}_{n,\lambda} = \left(\frac{\boldsymbol{W}^T\boldsymbol{W} + \lambda_{1,n}\boldsymbol{I}_{p-1}}{n}\right)^{-1}\frac{\boldsymbol{W}^T\boldsymbol{W}}{n} \xrightarrow{P} \boldsymbol{V}\,\boldsymbol{V}^{-1} = \boldsymbol{I}_{p-1}$$

if $\lambda_{1,n}/n \to 0$ since matrix inversion is a continuous function of a positive definite matrix. See, for example, Bhatia et al. (1990), Stewart (1969), and Severini (2005, pp. 348-349).

For model selection, the $M$ values of $\lambda = \lambda_{1,n}$ are denoted by $\lambda_1, \lambda_2, ..., \lambda_M$ where $\lambda_i = \lambda_{1,n,i}$ depends on $n$ for $i = 1, ..., M$. If $\lambda_s$ corresponds to the model selected, then $\hat{\lambda}_{1,n} = \lambda_s$. The following theorem shows that ridge regression and the OLS full model are asymptotically equivalent if $\hat{\lambda}_{1,n} = o_P(n^{1/2})$ so $\hat{\lambda}_{1,n}/\sqrt{n} \xrightarrow{P} 0$.

**Theorem 7.6, RR CLT (Ridge Regression Central Limit Theorem.** Assume $p$ is fixed and that the conditions of the OLS CLT Theorem Equation (7.14) hold for the model $\boldsymbol{Z} = \boldsymbol{W}\boldsymbol{\eta} + \boldsymbol{e}$.

a) If $\hat{\lambda}_{1,n}/\sqrt{n} \xrightarrow{P} 0$, then

$$\sqrt{n}(\hat{\boldsymbol{\eta}}_R - \boldsymbol{\eta}) \xrightarrow{D} N_{p-1}(\boldsymbol{0}, \sigma^2\boldsymbol{V}).$$

b) If $\hat{\lambda}_{1,n}/\sqrt{n} \xrightarrow{P} \tau \geq 0$ then

$$\sqrt{n}(\hat{\boldsymbol{\eta}}_R - \boldsymbol{\eta}) \xrightarrow{D} N_{p-1}(-\tau\boldsymbol{V}\boldsymbol{\eta}, \sigma^2\boldsymbol{V}).$$

**Proof:** If $\hat{\lambda}_{1,n}/\sqrt{n} \xrightarrow{P} \tau \geq 0$, then by the above Gunst and Mason (1980) identity,

$$\hat{\boldsymbol{\eta}}_R = [\boldsymbol{I}_{p-1} - \hat{\lambda}_{1,n}(\boldsymbol{W}^T\boldsymbol{W} + \hat{\lambda}_{1,n}\boldsymbol{I}_{p-1})^{-1}]\hat{\boldsymbol{\eta}}_{OLS}.$$

Hence

$$\sqrt{n}(\hat{\boldsymbol{\eta}}_R - \boldsymbol{\eta}) = \sqrt{n}(\hat{\boldsymbol{\eta}}_R - \hat{\boldsymbol{\eta}}_{OLS} + \hat{\boldsymbol{\eta}}_{OLS} - \boldsymbol{\eta}) =$$

$$\sqrt{n}(\hat{\boldsymbol{\eta}}_{OLS} - \boldsymbol{\eta}) - \sqrt{n}\frac{\hat{\lambda}_{1,n}}{n}n(\boldsymbol{W}^T\boldsymbol{W} + \hat{\lambda}_{1,n}\boldsymbol{I}_{p-1})^{-1}\hat{\boldsymbol{\eta}}_{OLS}$$

$$\overset{D}{\to} N_{p-1}(\mathbf{0}, \sigma^2 \boldsymbol{V}) - \tau \boldsymbol{V} \boldsymbol{\eta} \sim N_{p-1}(-\tau \boldsymbol{V} \boldsymbol{\eta}, \sigma^2 \boldsymbol{V}). \quad \square$$

For $p$ fixed, Knight and Fu (2000) note i) that $\hat{\boldsymbol{\eta}}_R$ is a consistent estimator of $\boldsymbol{\eta}$ if $\lambda_{1,n} = o(n)$ so $\lambda_{1,n}/n \to 0$ as $n \to \infty$, ii) OLS and ridge regression are asymptotically equivalent if $\lambda_{1,n}/\sqrt{n} \to 0$ as $n \to \infty$, iii) ridge regression is a $\sqrt{n}$ consistent estimator of $\boldsymbol{\eta}$ if $\lambda_{1,n} = O(\sqrt{n})$ (so $\lambda_{1,n}/\sqrt{n}$ is bounded), and iv) if $\lambda_{1,n}/\sqrt{n} \to \tau \geq 0$, then

$$\sqrt{n}(\hat{\boldsymbol{\eta}}_R - \boldsymbol{\eta}) \overset{D}{\to} N_{p-1}(-\tau \boldsymbol{V} \boldsymbol{\eta}, \sigma^2 \boldsymbol{V}).$$

Hence the bias can be considerable if $\tau \neq 0$. If $\tau = 0$, then OLS and ridge regression have the same limiting distribution.

Even if $p$ is fixed, there are several problems with ridge regression inference if $\hat{\lambda}_{1,n}$ is selected, e.g. after 10-fold cross validation. For OLS forward selection, the probability that the model $I_{min}$ underfits goes to zero, and each model with $S \subseteq I$ produced a $\sqrt{n}$ consistent estimator $\hat{\boldsymbol{\beta}}_{I,0}$ of $\boldsymbol{\beta}$. Ridge regression with 10-fold CV often shrinks $\hat{\boldsymbol{\beta}}_R$ too much if both i) the number of population active predictors $k_S = a_S - 1$ in Equation (7.1) and Remark 7.13 is greater than about 20, and ii) the predictors are highly correlated. If $p$ is fixed and $\lambda_{1,n} = o_P(\sqrt{n})$, then the OLS full model and ridge regression are asymptotically equivalent, but much larger sample sizes may be needed for the normal approximation to be good for ridge regression since the ridge regression estimator can have large bias for moderate $n$. Ten fold CV does not appear to guarantee that $\hat{\lambda}_{1,n}/\sqrt{n} \overset{P}{\to} 0$ or $\hat{\lambda}_{1,n}/n \overset{P}{\to} 0$.

Ridge regression can be a lot better than the OLS full model if i) $\boldsymbol{X}^T \boldsymbol{X}$ is singular or ill conditioned or ii) $n/p$ is small. Ridge regression can be much faster than forward selection if $M = 100$ and $n$ and $p$ are large.

Roughly speaking, the biased estimation of the ridge regression estimator can make the MSE of $\hat{\boldsymbol{\beta}}_R$ or $\hat{\boldsymbol{\eta}}_R$ less than that of $\hat{\boldsymbol{\beta}}_{OLS}$ or $\hat{\boldsymbol{\eta}}_{OLS}$, but the large sample inference may need larger $n$ for ridge regression than for OLS. However, the large sample theory has $n >> p$. We will try to use prediction intervals to compare OLS, forward selection, ridge regression, and lasso for data sets where $p > n$. See Section 7.12.

**Warning.** Although the $R$ functions `glmnet` and `cv.glmnet` appear to do ridge regression, getting the fitted values, $\hat{\lambda}_{1,n}$, and degrees of freedom to match up with the formulas of this section can be difficult.

**Example 7.5**, continued. The ridge regression output below shows results for the marry data where 10-fold CV was used. A grid of 100 $\lambda$ values was used, and $\lambda_0 > 0$ was selected. A problem with getting the false degrees of freedom $d$ for ridge regression is that it is not clear that $\lambda = \lambda_{1,n}/(2n)$. We need to know the relationship between $\lambda$ and $\lambda_{1,n}$ in order to compute $d$. It seems unlikely that $d \approx 1$ if $\lambda_0$ is selected.

```
library(glmnet); y <- marry[,3]; x <- marry[,-3]
```

```
out<-cv.glmnet(x,y,alpha=0)
lam <- out$lambda.min #value of lambda that minimizes
#the 10-fold CV criterion
yhat <- predict(out,s=lam,newx=x)
res <- y - yhat
n <- length(y)
w1 <- scale(x)
w <- sqrt(n/(n-1))*w1    #t(w) %*% w = n R_u, u = x
diag(t(w)%*%w)
     pop    mmen mmilmen  milwmn
      26      26      26      26
#sum w_i^2 = n = 26 for i = 1, 2, 3, and 4
svs <- svd(w)$d  #singular values of w,
pp <- 1 + sum(svs^2/(svs^2+2*n*lam))  #approx 1
# d for ridge regression if lam = lam_{1,n}/(2n)
AERplot2(yhat,y,res=res,d=pp)
$respi #90% PI for a future residual
[1] -5482.316 14854.268 #length = 20336.584
#try to reproduce the fitted values
z <- y - mean(y)
q<-dim(w)[2]
I <- diag(q)
M<- w%*%solve(t(w)%*%w + lam*I/(2*n))%*%t(w)
fit <- M%*%z + mean(y)
plot(fit,yhat) #they are not the same
max(abs(fit-yhat))
[1] 46789.11
M<- w%*%solve(t(w)%*%w + lam*I/(1547.1741))%*%t(w)
fit <- M%*%z + mean(y)
max(abs(fit-yhat)) #close
[1] 8.484979
```

## 7.9 Lasso

Consider the MLR model $\boldsymbol{Y} = \boldsymbol{X}\boldsymbol{\beta} + \boldsymbol{e}$. Lasso uses the centered response $Z_i = Y_i - \overline{Y}$ and standardized nontrivial predictors in the model $\boldsymbol{Z} = \boldsymbol{W}\boldsymbol{\eta} + \boldsymbol{e}$ as described in Remark 7.9. Then $\hat{Y}_i = \hat{Z}_i + \overline{Y}$. The residuals $\boldsymbol{r} = \boldsymbol{r}(\hat{\boldsymbol{\beta}}_L) = \boldsymbol{Y} - \hat{\boldsymbol{Y}}$. Recall that $\overline{\boldsymbol{Y}} = \overline{Y}\boldsymbol{1}$.

**Definition 7.11.** Consider fitting the MLR model $\boldsymbol{Y} = \boldsymbol{X}\boldsymbol{\beta} + \boldsymbol{e}$ using $\boldsymbol{Z} = \boldsymbol{W}\boldsymbol{\eta} + \boldsymbol{e}$. The *lasso estimator* $\hat{\boldsymbol{\eta}}_L$ minimizes the *lasso criterion*

$$Q_L(\boldsymbol{\eta}) = \frac{1}{a}(\boldsymbol{Z} - \boldsymbol{W}\boldsymbol{\eta})^T(\boldsymbol{Z} - \boldsymbol{W}\boldsymbol{\eta}) + \frac{\lambda_{1,n}}{a}\sum_{i=1}^{p-1}|\eta_i| \qquad (7.21)$$

over all vectors $\boldsymbol{\eta} \in \mathbb{R}^{p-1}$ where $\lambda_{1,n} \geq 0$ and $a > 0$ are known constants with $a = 1, 2, n$, and $2n$ are common. The residual sum of squares $RSS(\boldsymbol{\eta}) = (\boldsymbol{Z} - \boldsymbol{W}\boldsymbol{\eta})^T(\boldsymbol{Z} - \boldsymbol{W}\boldsymbol{\eta})$, and $\lambda_{1,n} = 0$ corresponds to the OLS estimator $\hat{\boldsymbol{\eta}}_{OLS} = (\boldsymbol{W}^T\boldsymbol{W})^{-1}\boldsymbol{W}^T\boldsymbol{Z}$ if $\boldsymbol{W}$ has full rank $p-1$. The lasso vector of fitted values is $\hat{\boldsymbol{Z}} = \hat{\boldsymbol{Z}}_L = \boldsymbol{W}\hat{\boldsymbol{\eta}}_L$, and the lasso vector of residuals $\boldsymbol{r}(\hat{\boldsymbol{\eta}}_L) = \boldsymbol{Z} - \hat{\boldsymbol{Z}}_L$. The estimator is said to be *regularized* if $\lambda_{1,n} > 0$. Obtain $\hat{\boldsymbol{Y}}$ and $\hat{\boldsymbol{\beta}}_L$ using $\hat{\boldsymbol{\eta}}_L$, $\hat{\boldsymbol{Z}}$, and $\overline{\boldsymbol{Y}}$.

Using a vector of parameters $\boldsymbol{\eta}$ and a dummy vector $\boldsymbol{\eta}$ in $Q_L$ is common for minimizing a criterion $Q(\boldsymbol{\eta})$, often with estimating equations. See the paragraphs above and below Definition 7.7. We could also write

$$Q_L(\boldsymbol{b}) = \frac{1}{a}\boldsymbol{r}(\boldsymbol{b})^T\boldsymbol{r}(\boldsymbol{b}) + \frac{\lambda_{1,n}}{a}\sum_{j=1}^{p-1}|b_j|, \qquad (7.22)$$

where the minimization is over all vectors $\boldsymbol{b} \in \mathbb{R}^{p-1}$. The literature often uses $\lambda_a = \lambda = \lambda_{1,n}/a$.

For fixed $\lambda_{1,n}$, the lasso optimization problem is convex. Hence fast algorithms exist. As $\lambda_{1,n}$ increases, some of the $\hat{\eta}_i = 0$. If $\lambda_{1,n}$ is large enough, then $\hat{\boldsymbol{\eta}}_L = \boldsymbol{0}$ and $\hat{Y}_i = \overline{Y}$ for $i = 1, ..., n$. If none of the elements $\hat{\eta}_i$ of $\hat{\boldsymbol{\eta}}_L$ are zero, then $\hat{\boldsymbol{\eta}}_L$ can be found, in principle, by setting the partial derivatives of $Q_L(\boldsymbol{\eta})$ to 0. Potential minimizers also occur at values of $\boldsymbol{\eta}$ where not all of the partial derivatives exist. An analogy is finding the minimizer of a real valued function of one variable $h(x)$. Possible values for the minimizer include values of $x_c$ satisfying $h'(x_c) = 0$, and values $x_c$ where the derivative does not exist. Typically some of the elements $\hat{\eta}_i$ of $\hat{\boldsymbol{\eta}}_L$ that minimizes $Q_L(\boldsymbol{\eta})$ are zero, and differentiating does not work.

The following identity from Efron and Hastie (2016, p. 308), for example, is useful for inference for the lasso estimator $\hat{\boldsymbol{\eta}}_L$:

$$\frac{-1}{n}\boldsymbol{W}^T(\boldsymbol{Z} - \boldsymbol{W}\hat{\boldsymbol{\eta}}_L) + \frac{\lambda_{1,n}}{2n}\boldsymbol{s}_n = \boldsymbol{0} \text{ or } -\boldsymbol{W}^T(\boldsymbol{Z} - \boldsymbol{W}\hat{\boldsymbol{\eta}}_L) + \frac{\lambda_{1,n}}{2}\boldsymbol{s}_n = \boldsymbol{0}$$

where $s_{in} \in [-1, 1]$ and $s_{in} = \text{sign}(\hat{\eta}_{i,L})$ if $\hat{\eta}_{i,L} \neq 0$. Here $\text{sign}(\eta_i) = 1$ if $\eta_i > 1$ and $\text{sign}(\eta_i) = -1$ if $\eta_i < 1$. Note that $\boldsymbol{s}_n = \boldsymbol{s}_{n,\hat{\boldsymbol{\eta}}_L}$ depends on $\hat{\boldsymbol{\eta}}_L$. Thus $\hat{\boldsymbol{\eta}}_L$

$$= (\boldsymbol{W}^T\boldsymbol{W})^{-1}\boldsymbol{W}^T\boldsymbol{Z} - \frac{\lambda_{1,n}}{2n}n(\boldsymbol{W}^T\boldsymbol{W})^{-1}\boldsymbol{s}_n = \hat{\boldsymbol{\eta}}_{OLS} - \frac{\lambda_{1,n}}{2n}n(\boldsymbol{W}^T\boldsymbol{W})^{-1}\boldsymbol{s}_n.$$

If none of the elements of $\boldsymbol{\eta}$ are zero, and if $\hat{\boldsymbol{\eta}}_L$ is a consistent estimator of $\boldsymbol{\eta}$, then $\boldsymbol{s}_n \xrightarrow{P} \boldsymbol{s} = \boldsymbol{s}_{\boldsymbol{\eta}}$. If $\lambda_{1,n}/\sqrt{n} \to 0$, then OLS and lasso are asymptotically equivalent even if $\boldsymbol{s}_n$ does not converge to a vector $\boldsymbol{s}$ as $n \to \infty$ since $\boldsymbol{s}_n$ is bounded. For model selection, the $M$ values of $\lambda$ are denoted by $0 \leq \lambda_1 < \lambda_2 < \cdots < \lambda_M$ where $\lambda_i = \lambda_{1,n,i}$ depends on $n$ for $i = 1, ..., M$. Also, $\lambda_M$ is the smallest value of $\lambda$ such that $\hat{\boldsymbol{\eta}}_{\lambda_M} = \boldsymbol{0}$. Hence $\hat{\boldsymbol{\eta}}_{\lambda_i} \neq \boldsymbol{0}$ for $i < M$. If $\lambda_s$ corresponds to the model selected, then $\hat{\lambda}_{1,n} = \lambda_s$. The following theorem shows that lasso and the OLS full model are asymptotically equivalent if $\hat{\lambda}_{1,n} = o_P(n^{1/2})$ so $\hat{\lambda}_{1,n}/\sqrt{n} \xrightarrow{P} 0$: thus $\sqrt{n}(\hat{\boldsymbol{\eta}}_L - \hat{\boldsymbol{\eta}}_{OLS}) = o_p(1)$.

**Theorem 7.7, Lasso CLT.** Assume $p$ is fixed and that the conditions of the OLS CLT Theorem Equation (7.14) hold for the model $\boldsymbol{Z} = \boldsymbol{W}\boldsymbol{\eta} + \boldsymbol{e}$.
a) If $\hat{\lambda}_{1,n}/\sqrt{n} \xrightarrow{P} 0$, then

$$\sqrt{n}(\hat{\boldsymbol{\eta}}_L - \boldsymbol{\eta}) \xrightarrow{D} N_{p-1}(\boldsymbol{0}, \sigma^2 \boldsymbol{V}).$$

b) If $\hat{\lambda}_{1,n}/\sqrt{n} \xrightarrow{P} \tau \geq 0$ and $\boldsymbol{s}_n \xrightarrow{P} \boldsymbol{s} = \boldsymbol{s}_{\boldsymbol{\eta}}$, then

$$\sqrt{n}(\hat{\boldsymbol{\eta}}_L - \boldsymbol{\eta}) \xrightarrow{D} N_{p-1}\left(\frac{-\tau}{2}\boldsymbol{V}\boldsymbol{s}, \sigma^2 \boldsymbol{V}\right).$$

**Proof.** If $\hat{\lambda}_{1,n}/\sqrt{n} \xrightarrow{P} \tau \geq 0$ and $\boldsymbol{s}_n \xrightarrow{P} \boldsymbol{s} = \boldsymbol{s}_{\boldsymbol{\eta}}$, then

$$\sqrt{n}(\hat{\boldsymbol{\eta}}_L - \boldsymbol{\eta}) = \sqrt{n}(\hat{\boldsymbol{\eta}}_L - \hat{\boldsymbol{\eta}}_{OLS} + \hat{\boldsymbol{\eta}}_{OLS} - \boldsymbol{\eta}) =$$

$$\sqrt{n}(\hat{\boldsymbol{\eta}}_{OLS} - \boldsymbol{\eta}) - \sqrt{n}\frac{\lambda_{1,n}}{2n}n(\boldsymbol{W}^T\boldsymbol{W})^{-1}\boldsymbol{s}_n \xrightarrow{D} N_{p-1}(\boldsymbol{0}, \sigma^2 \boldsymbol{V}) - \frac{\tau}{2}\boldsymbol{V}\boldsymbol{s}$$

$$\sim N_{p-1}\left(\frac{-\tau}{2}\boldsymbol{V}\boldsymbol{s}, \sigma^2 \boldsymbol{V}\right)$$

since under the LS CLT, $n(\boldsymbol{W}^T\boldsymbol{W})^{-1} \xrightarrow{P} \boldsymbol{V}$.

Part a) does not need $\boldsymbol{s}_n \xrightarrow{P} \boldsymbol{s}$ as $n \to \infty$, since $\boldsymbol{s}_n$ is bounded. $\square$

Suppose $p$ is fixed. Knight and Fu (2000) note i) that $\hat{\boldsymbol{\eta}}_L$ is a consistent estimator of $\boldsymbol{\eta}$ if $\lambda_{1,n} = o(n)$ so $\lambda_{1,n}/n \to 0$ as $n \to \infty$, ii) OLS and lasso are asymptotically equivalent if $\lambda_{1,n} \to \infty$ too slowly as $n \to \infty$ (e.g. if $\lambda_{1,n} = \lambda$ is fixed), iii) lasso is a $\sqrt{n}$ consistent estimator of $\boldsymbol{\eta}$ if $\lambda_{1,n} = O(\sqrt{n})$ (so $\lambda_{1,n}/\sqrt{n}$ is bounded). Note that Theorem 7.7 shows that OLS and lasso are asymptotically equivalent if $\lambda_{1,n}/\sqrt{n} \to 0$ as $n \to 0$.

In the literature, the criterion often uses $\lambda_a = \lambda_{1,n}/a$:

$$Q_{L,a}(\boldsymbol{b}) = \frac{1}{a}\boldsymbol{r}(\boldsymbol{b})^T\boldsymbol{r}(\boldsymbol{b}) + \lambda_a \sum_{j=1}^{p-1} |b_j|.$$

The values $a = 1,\ 2,$ and $2n$ are common. Following Hastie et al. (2015, pp. 9, 17, 19) for the next two paragraphs, it is convenient to use $a = 2n$:

$$Q_{L,2n}(\boldsymbol{b}) = \frac{1}{2n}\boldsymbol{r}(\boldsymbol{b})^T\boldsymbol{r}(\boldsymbol{b}) + \lambda_{2n}\sum_{j=1}^{p-1}|b_j|, \qquad (7.23)$$

where the $Z_i$ are centered and the $w_j$ are standardized using $g = 0$ so $\overline{w}_j = 0$ and $n\hat{\sigma}_j^2 = \sum_{i=1}^n w_{i,j}^2 = n$. Then $\lambda = \lambda_{2n} = \lambda_{1,n}/(2n)$ in Equation (7.21). For model selection, the $M$ values of $\lambda$ are denoted by $0 \leq \lambda_{2n,1} < \lambda_{2n,2} < \cdots < \lambda_{2n,M}$ where $\hat{\boldsymbol{\eta}}_\lambda = \boldsymbol{0}$ iff $\lambda \geq \lambda_{2n,M}$ and

$$\lambda_{2n,max} = \lambda_{2n,M} = \max_j \left|\frac{1}{n}\boldsymbol{s}_j^T\boldsymbol{Z}\right|$$

and $\boldsymbol{s}_j$ is the $j$th column of $\boldsymbol{W}$ corresponding to the $j$th standardized non-trivial predictor $W_j$. In terms of the $0 \leq \lambda_1 < \lambda_2 < \cdots < \lambda_M$, used above Theorem 7.7, we have $\lambda_i = \lambda_{1,n,i} = 2n\lambda_{2n,i}$ and

$$\lambda_M = 2n\lambda_{2n,M} = 2\max_j\left|\boldsymbol{s}_j^T\boldsymbol{Z}\right|.$$

For model selection we let $I$ denote the index set of the predictors in the fitted model including the constant. The set $A$ defined below is the index set without the constant.

**Definition 7.12.** The *active set* $A$ is the index set of the nontrivial predictors in the fitted model: the predictors with nonzero $\hat{\eta}_i$.

Suppose that there are $k$ active nontrivial predictors. Then for lasso, $k \leq n$. Let the $n \times k$ matrix $\boldsymbol{W}_A$ correspond to the standardized active predictors. If the columns of $\boldsymbol{W}_A$ are in general position, then the lasso vector of fitted values

$$\hat{\boldsymbol{Z}}_L = \boldsymbol{W}_A(\boldsymbol{W}_A^T\boldsymbol{W}_A)^{-1}\boldsymbol{W}_A^T\boldsymbol{Z} - n\lambda_{2n}\boldsymbol{W}_A(\boldsymbol{W}_A^T\boldsymbol{W}_A)^{-1}\boldsymbol{s}_A$$

where $\boldsymbol{s}_A$ is the vector of signs of the active lasso coefficients. Here we are using the $\lambda_{2n}$ of (7.23), and $n\lambda_{2n} = \lambda_{1,n}/2$. We could replace $n\ \lambda_{2n}$ by $\lambda_2$ if we used $a = 2$ in the criterion

$$Q_{L,2}(\boldsymbol{b}) = \frac{1}{2}\boldsymbol{r}(\boldsymbol{b})^T\boldsymbol{r}(\boldsymbol{b}) + \lambda_2\sum_{j=1}^{p-1}|b_j|. \qquad (7.24)$$

See, for example, Tibshirani (2015). Note that $\boldsymbol{W}_A(\boldsymbol{W}_A^T\boldsymbol{W}_A)^{-1}\boldsymbol{W}_A^T\boldsymbol{Z}$ is the vector of OLS fitted values from regressing $\boldsymbol{Z}$ on $\boldsymbol{W}_A$ without an intercept.

**Example 7.5**, continued. The lasso output below shows results for the marry data where 10-fold CV was used. A grid of 38 $\lambda$ values was used, and $\lambda_0 > 0$ was selected.

```
library(glmnet); y <- marry[,3]; x <- marry[,-3]
out<-cv.glmnet(x,y)
lam <- out$lambda.min #value of lambda that minimizes
#the 10-fold CV criterion
yhat <- predict(out,s=lam,newx=x)
res <- y - yhat
pp <- out$nzero[out$lambda==lam] + 1 #d for lasso
AERplot2(yhat,y,res=res,d=pp)
$respi #90% PI for a future residual
-4102.672  4379.951  #length = 8482.62
```

There are some problems with lasso. i) Lasso large sample theory is worse or as good as that of the OLS full model if $n/p$ is large. ii) Ten fold CV does not appear to guarantee that $\hat{\lambda}_{1,n}/\sqrt{n} \xrightarrow{P} 0$ or $\hat{\lambda}_{1,n}/n \xrightarrow{P} 0$. iii) Lasso often shrinks $\hat{\boldsymbol{\beta}}$ too much if $a_S \geq 20$ and the predictors are highly correlated. iv) Ridge regression can be better than lasso if $a_S > n$.

Lasso can be a lot better than the OLS full model if i) $\boldsymbol{X}^T\boldsymbol{X}$ is singular or ill conditioned or ii) $n/p$ is small. iii) For lasso, $M = M(lasso)$ is often near 100. Let $J \geq 5$. If $n/J$ and $p$ are both a lot larger than $M(lasso)$, then lasso can be considerably faster than forward selection, PLS, and PCR if $M = M(lasso) = 100$ and $M = M(F) = \min(\lceil n/J \rceil, p)$ where $F$ stands for forward selection, PLS, or PCR. iv) The number of nonzero coefficients in $\hat{\boldsymbol{\eta}}_L \leq n$ even if $p > n$. This property of lasso can be useful if $p >> n$ and the population model is sparse.

## 7.10 Lasso Variable Selection

Lasso variable selection applies OLS on a constant and the active predictors that have nonzero lasso $\hat{\eta}_i$. The method is called relaxed lasso by Hastie et al. (2015, p. 12), and the relaxed lasso ($\phi = 0$) estimator by Meinshausen (2007). The method is also called OLS-post lasso and post model selection OLS. Let $\boldsymbol{X}_A$ denote the matrix with a column of ones and the unstandardized active nontrivial predictors. Hence the lasso variable selection estimator is $\hat{\boldsymbol{\beta}}_{LVS} = (\boldsymbol{X}_A^T\boldsymbol{X}_A)^{-1}\boldsymbol{X}_A^T\boldsymbol{Y}$, and lasso variable selection is an alternative to forward selection. Let $k$ be the number of active (nontrivial) predictors so $\hat{\boldsymbol{\beta}}_{LVS}$ is $(k+1) \times 1$.

Let $I_{min}$ correspond to the lasso variable selection estimator and $\hat{\boldsymbol{\beta}}_{VS} = \hat{\boldsymbol{\beta}}_{LVS,0} = \hat{\boldsymbol{\beta}}_{I_{min},0}$ to the zero padded lasso variable selection estimator. Then by Remark 7.5 where $p$ is fixed, $\hat{\boldsymbol{\beta}}_{LVS,0}$ is $\sqrt{n}$ consistent when lasso is consis-

tent, with the limiting distribution for $\hat{\boldsymbol{\beta}}_{VS} = \hat{\boldsymbol{\beta}}_{LVS,0}$ given by Theorem 7.4. Hence lasso variable selection can be bootstrapped as in Section 7.4. Lasso vaiable selection will often be better than lasso when the model is sparse or if $n \geq 10(k+1)$. Lasso can be better than lasso variable selection if $(\boldsymbol{X}_A^T \boldsymbol{X}_A)$ is ill conditioned or if $n/(k+1) < 10$. Also see Pelawa Watagoda and Olive (2020) and Rathnayake and Olive (2020).

Suppose the $n \times q$ matrix $x$ has the $q = p - 1$ nontrivial predictors. The following $R$ code gives some output for a lasso estimator and then the corresponding lasso variable selection estimator.

```
library(glmnet)
y <- marry[,3]
x <- marry[,-3]
out<-glmnet(x,y,dfmax=2)  #Use 2 for illustration:
#often dfmax approx min(n/J,p) for some J >= 5.
lam<-out$lambda[length(out$lambda)]
yhat <- predict(out,s=lam,newx=x)
#lasso with smallest lambda in grid such that df = 2
lcoef <- predict(out,type="coefficients",s=lam)
as.vector(lcoef) #first term is the intercept
#3.000397e+03 1.800342e-03 9.618035e-01 0.0 0.0
res <- y - yhat
AERplot(yhat,y,res,d=3,alph=1) #lasso response plot
##relaxed lasso =
#OLS on lasso active predictors and a constant
vars <- 1:dim(x)[2]
lcoef<-as.vector(lcoef)[-1] #don't need an intercept
vin <- vars[lcoef>0] #the lasso active set
vin
#1  2  since predictors 1 and 2 are active
sub <- lsfit(x[,vin],y) # lasso variable selection
sub$coef
#  Intercept           pop          mmen
#2.380912e+02 6.556895e-05 1.000603e+00
# 238.091     6.556895e-05 1.0006
res <- sub$resid
yhat <- y - res
AERplot(yhat,y,res,d=3,alph=1) #response plot
```

**Example 7.5**, continued. The lasso variable selection output below shows results for the marry data where 10-fold CV was used to choose the lasso estimator. Then lasso variable selection is OLS applied to the active variables with nonzero lasso coefficients and a constant. A grid of 38 $\lambda$ values was used, and $\lambda_0 > 0$ was selected. The OLS SE, t statistic and pvalue are generally not valid for lasso variable selection by Theorem 7.4.

```
library(glmnet); y <- marry[,3]; x <- marry[,-3]
```

```
out<-cv.glmnet(x,y)
lam <- out$lambda.min #value of lambda that minimizes
#the 10-fold CV criterion
pp <- out$nzero[out$lambda==lam] + 1
#d for lasso variable selection
#get lasso variable selection
lcoef <- predict(out,type="coefficients",s=lam)
lcoef<-as.vector(lcoef)[-1]
vin <- vars[lcoef!=0]
sub <- lsfit(x[,vin],y)
ls.print(sub)
Residual Standard Error=376.9412
R-Square=0.9999
F-statistic (df=2, 23)=147440.1
          Estimate  Std.Err t-value Pr(>|t|)58
Intercept 238.0912 248.8616  0.9567   0.3487
pop         0.0001   0.0029  0.0223   0.9824
mmen        1.0006   0.0164 60.9878   0.0000
res <- sub$resid
yhat <- y - res
AERplot2(yhat,y,res=res,d=pp)
$respi #90% PI for a future residual
-822.759 1403.771  #length = 2226.53
```

To summarize Example 7.5, forward selection selected the model with the minimum $C_p$ while the other methods used 10-fold CV. PLS and PCR used the OLS full model with PI length 2395.74, forward selection used a constant and *mmen* with PI length 2114.72, ridge regression had PI length 20336.58, lasso and lasso variable selection used a constant, *mmen*, and *pop* with lasso PI length 8482.62 and relaxed lasso PI length 2226.53. PI (4.14) was used. Figure 7.1 shows the response plots for forward selection, ridge regression, lasso, and lasso variable selection. The plots for PLS=PCR=OLS full model were similar to those of forward selection and lasso variable selection. The plots suggest that the MLR model is appropriate since the plotted points scatter about the identity line. The 90% pointwise prediction bands are also shown, and consist of two lines parallel to the identity line. These bands are very narrow in Figure 7.1 a) and d).

## 7.11 The Elastic Net

Following Hastie et al. (2015, p. 57), let $\boldsymbol{\beta} = (\beta_1, \boldsymbol{\beta}_S^T)^T$, let $\lambda_{1,n} \geq 0$, and let $\alpha \in [0,1]$. Let

$$RSS(\boldsymbol{\beta}) = (\boldsymbol{Y} - \boldsymbol{X}\boldsymbol{\beta})^T(\boldsymbol{Y} - \boldsymbol{X}\boldsymbol{\beta}) = \|\boldsymbol{Y} - \boldsymbol{X}\boldsymbol{\beta}\|_2^2.$$

**a) Forward Selection**



**b) Ridge Regression**
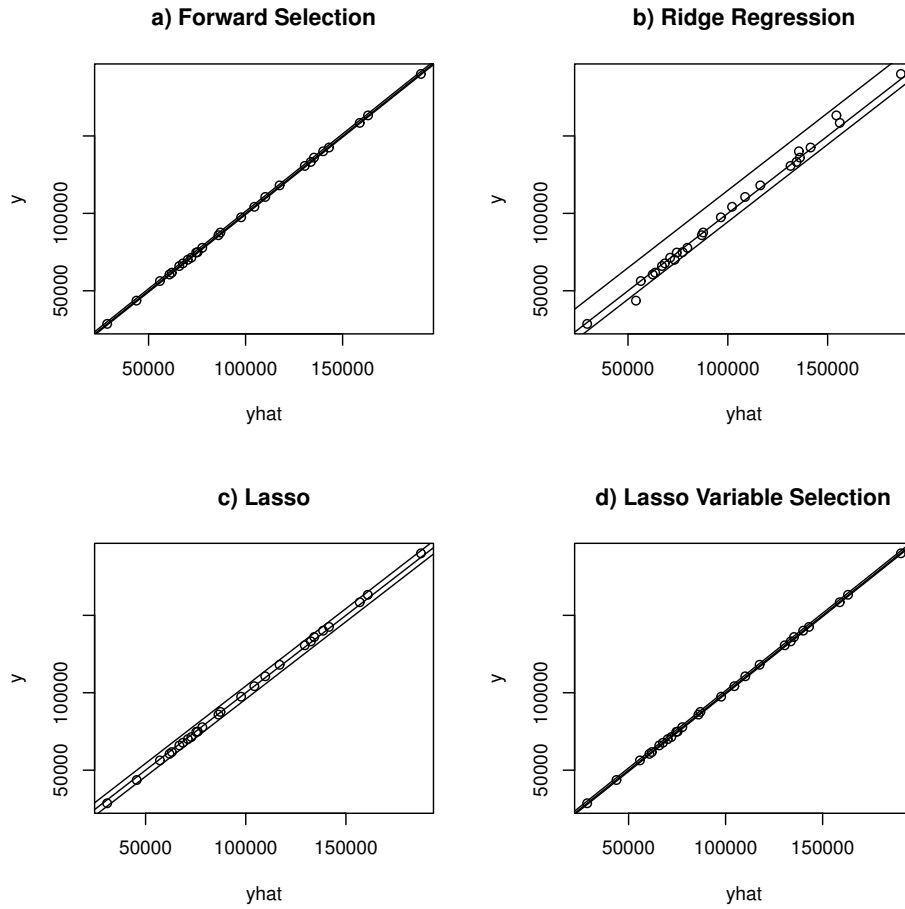


**c) Lasso**



**d) Lasso Variable Selection**



**Fig. 7.1** Marry Data Response Plots

For a $k \times 1$ vector $\boldsymbol{\eta}$, the squared (Euclidean) $L_2$ norm $\|\boldsymbol{\eta}\|_2^2 = \boldsymbol{\eta}^T \boldsymbol{\eta} = \sum_{i=1}^{k} \eta_i^2$ and the $L_1$ norm $\|\boldsymbol{\eta}\|_1 = \sum_{i=1}^{k} |\eta_i|$.

**Definition 7.13.** The *elastic net* estimator $\hat{\boldsymbol{\beta}}_{EN}$ minimizes the criterion

$$Q_{EN}(\boldsymbol{\beta}) = \frac{1}{2}RSS(\boldsymbol{\beta}) + \lambda_{1,n} \left[ \frac{1}{2}(1-\alpha)\|\boldsymbol{\beta}_S\|_2^2 + \alpha\|\boldsymbol{\beta}_S\|_1 \right], \text{ or} \qquad (7.25)$$

$$Q_2(\boldsymbol{\beta}) = RSS(\boldsymbol{\beta}) + \lambda_1\|\boldsymbol{\beta}_S\|_2^2 + \lambda_2\|\boldsymbol{\beta}_S\|_1 \qquad (7.26)$$

where $0 \le \alpha \le 1$, $\lambda_1 = (1-\alpha)\lambda_{1,n}$ and $\lambda_2 = 2\alpha\lambda_{1,n}$.

Note that $\alpha = 1$ corresponds to lasso (using $\lambda_{a=0.5}$), and $\alpha = 0$ corresponds to ridge regression. For $\alpha < 1$ and $\lambda_{1,n} > 0$, the optimization problem is *strictly convex* with a unique solution. The elastic net is due to Zou and Hastie (2005). It has been observed that the elastic net can have much better prediction accuracy than lasso when the predictors are highly correlated.

As with lasso, it is often convenient to use the centered response $\boldsymbol{Z} = \boldsymbol{Y} - \overline{\boldsymbol{Y}}$ where $\overline{\boldsymbol{Y}} = \overline{Y}\mathbf{1}$, and the $n \times (p-1)$ matrix of standardized nontrivial predictors $\boldsymbol{W}$. Then regression through the origin is used for the model

$$\boldsymbol{Z} = \boldsymbol{W}\boldsymbol{\eta} + \boldsymbol{e} \tag{7.27}$$

where the vector of fitted values $\hat{\boldsymbol{Y}} = \overline{\boldsymbol{Y}} + \hat{\boldsymbol{Z}}$.

Ridge regression can be computed using OLS on augmented matrices. Similarly, the elastic net can be computed using lasso on augmented matrices. Let the elastic net estimator $\hat{\boldsymbol{\eta}}_{EN}$ minimize

$$Q_{EN}(\boldsymbol{\eta}) = RSS_W(\boldsymbol{\eta}) + \lambda_1 \|\boldsymbol{\eta}\|_2^2 + \lambda_2 \|\boldsymbol{\eta}\|_1 \tag{7.28}$$

where $\lambda_1 = (1-\alpha)\lambda_{1,n}$ and $\lambda_2 = 2\alpha\lambda_{1,n}$. Let the $(n+p-1) \times (p-1)$ augmented matrix $\boldsymbol{W}_A$ and the $(n+p-1) \times 1$ augmented response vector $\boldsymbol{Z}_A$ be defined by

$$\boldsymbol{W}_A = \begin{pmatrix} \boldsymbol{W} \\ \sqrt{\lambda_1}\,\boldsymbol{I}_{p-1} \end{pmatrix}, \quad \text{and} \quad \boldsymbol{Z}_A = \begin{pmatrix} \boldsymbol{Z} \\ \mathbf{0} \end{pmatrix},$$

where $\mathbf{0}$ is the $(p-1) \times 1$ zero vector. Let $RSS_A(\boldsymbol{\eta}) = \|\boldsymbol{Z}_A - \boldsymbol{W}_A\boldsymbol{\eta}\|_2^2$. Then $\hat{\boldsymbol{\eta}}_{EN}$ can be obtained from the lasso of $\boldsymbol{Z}_A$ on $\boldsymbol{W}_A$: that is, $\hat{\boldsymbol{\eta}}_{EN}$ minimizes

$$Q_L(\boldsymbol{\eta}) = RSS_A(\boldsymbol{\eta}) + \lambda_2 \|\boldsymbol{\eta}\|_1 = Q_{EN}(\boldsymbol{\eta}). \tag{7.29}$$

Proof: We need to show that $Q_L(\boldsymbol{\eta}) = Q_{EN}(\boldsymbol{\eta})$. Note that $\boldsymbol{Z}_A^T \boldsymbol{Z}_A = \boldsymbol{Z}^T \boldsymbol{Z}$,

$$\boldsymbol{W}_A\,\boldsymbol{\eta} = \begin{pmatrix} \boldsymbol{W}\boldsymbol{\eta} \\ \sqrt{\lambda_1}\,\boldsymbol{\eta} \end{pmatrix},$$

and $\boldsymbol{Z}_A^T \boldsymbol{W}_A\,\boldsymbol{\eta} = \boldsymbol{Z}^T \boldsymbol{W}\boldsymbol{\eta}$. Then

$$RSS_A(\boldsymbol{\eta}) = \|\boldsymbol{Z}_A - \boldsymbol{W}_A\boldsymbol{\eta}\|_2^2 = (\boldsymbol{Z}_A - \boldsymbol{W}_A\boldsymbol{\eta})^T (\boldsymbol{Z}_A - \boldsymbol{W}_A\boldsymbol{\eta}) =$$

$$\boldsymbol{Z}_A^T \boldsymbol{Z}_A - \boldsymbol{Z}_A^T \boldsymbol{W}_A\boldsymbol{\eta} - \boldsymbol{\eta}^T \boldsymbol{W}_A^T \boldsymbol{Z}_A + \boldsymbol{\eta}^T \boldsymbol{W}_A^T \boldsymbol{W}_A\boldsymbol{\eta} =$$

$$\boldsymbol{Z}^T \boldsymbol{Z} - \boldsymbol{Z}^T \boldsymbol{W}\boldsymbol{\eta} - \boldsymbol{\eta}^T \boldsymbol{W}^T \boldsymbol{Z} + \begin{pmatrix} \boldsymbol{\eta}^T \boldsymbol{W}^T & \sqrt{\lambda_1}\ \boldsymbol{\eta}^T \end{pmatrix} \begin{pmatrix} \boldsymbol{W}\boldsymbol{\eta} \\ \sqrt{\lambda_1}\,\boldsymbol{\eta} \end{pmatrix}.$$

Thus

$$Q_L(\boldsymbol{\eta}) = \boldsymbol{Z}^T \boldsymbol{Z} - \boldsymbol{Z}^T \boldsymbol{W}\boldsymbol{\eta} - \boldsymbol{\eta}^T \boldsymbol{W}^T \boldsymbol{Z} + \boldsymbol{\eta}^T \boldsymbol{W}^T \boldsymbol{W}\boldsymbol{\eta} + \lambda_1 \boldsymbol{\eta}^T \boldsymbol{\eta} + \lambda_2 \|\boldsymbol{\eta}\|_1 =$$

$$RSS(\boldsymbol{\eta}) + \lambda_1\|\boldsymbol{\eta}\|_2^2 + \lambda_2\|\boldsymbol{\eta}\|_1 = Q_{EN}(\boldsymbol{\eta}). \quad \square$$

**Remark 7.20.** i) You could compute the elastic net estimator using a grid of 100 $\lambda_{1,n}$ values and a grid of $J \geq 10$ $\alpha$ values, which would take about $J \geq 10$ times as long to compute as lasso. The above equivalent lasso problem (7.29) still needs a grid of $\lambda_1 = (1-\alpha)\lambda_{1,n}$ and $\lambda_2 = 2\alpha\lambda_{1,n}$ values. Often $J = 11$, 21, 51, or 101. The elastic net estimator tends to be computed with fast methods for optimizing convex problems, such as coordinate descent. ii) Like lasso and ridge regression, the elastic net estimator is asymptotically equivalent to the OLS full model if $p$ is fixed and $\hat{\lambda}_{1,n} = o_P(\sqrt{n})$, but behaves worse than the OLS full model otherwise. See Theorem 7.8. iii) For prediction intervals, let $d$ be the number of nonzero coefficients from the equivalent augmented lasso problem (7.29). Alternatively, use $d_2$ with $d \approx d_2 = tr[\boldsymbol{W}_{AS}(\boldsymbol{W}_{AS}^T\boldsymbol{W}_{AS} + \lambda_{2,n}\boldsymbol{I}_{p-1})^{-1}\boldsymbol{W}_{AS}^T]$ where $\boldsymbol{W}_{AS}$ corresponds to the active set (not the augmented matrix). See Tibshirani and Taylor (2012, p. 1214). Again $\lambda_{2,n}$ may not be the $\lambda_2$ given by the software. iv) The number of nonzero lasso components (not including the constant) is at most $\min(n, p-1)$. Elastic net tends to do variable selection, but the number of nonzero components can equal $p-1$ (make the elastic net equal to ridge regression). Note that the number of nonzero components in the augmented lasso problem (7.29) is at most $\min(n+p-1, p-1) = p-1$. vi) The elastic net can be computed with `glmnet`, and there is an $R$ package `elasticnet`. vii) For fixed $\alpha > 0$, we could get $\lambda_M$ for elastic net from the equivalent lasso problem. For ridge regression, we could use the $\lambda_M$ for an $\alpha$ near 0.

Since lasso uses at most $\min(n, p-1)$ nontrivial predictors, elastic net and ridge regression can perform better than lasso if the true number of active nontrivial predictors $a_S > \min(n, p-1)$. For example, suppose $n = 1000$, $p = 5000$, and $a_S = 1500$.

Following Jia and Yu (2010), by standard Karush-Kuhn-Tucker (KKT) conditions for convex optimality for Equation (7.28), $\hat{\boldsymbol{\eta}}_{EN}$ is optimal if

$$2\boldsymbol{W}^T\boldsymbol{W}\hat{\boldsymbol{\eta}}_{EN} - 2\boldsymbol{W}^T\boldsymbol{Z} + 2\lambda_1\hat{\boldsymbol{\eta}}_{EN} + \lambda_2\boldsymbol{s}_n = 0, \quad \text{or}$$

$$(\boldsymbol{W}^T\boldsymbol{W} + \lambda_1\boldsymbol{I}_{p-1})\hat{\boldsymbol{\eta}}_{EN} = \boldsymbol{W}^T\boldsymbol{Z} - \frac{\lambda_2}{2}\boldsymbol{s}_n, \quad \text{or}$$

$$\hat{\boldsymbol{\eta}}_{EN} = \hat{\boldsymbol{\eta}}_R - n(\boldsymbol{W}^T\boldsymbol{W} + \lambda_1\boldsymbol{I}_{p-1})^{-1}\frac{\lambda_2}{2n}\boldsymbol{s}_n. \quad (7.30)$$

Hence

$$\hat{\boldsymbol{\eta}}_{EN} = \hat{\boldsymbol{\eta}}_{OLS} - \frac{\lambda_1}{n}\, n(\boldsymbol{W}^T\boldsymbol{W} + \lambda_1\boldsymbol{I}_{p-1})^{-1}\,\hat{\boldsymbol{\eta}}_{OLS} - \frac{\lambda_2}{2n}\, n(\boldsymbol{W}^T\boldsymbol{W} + \lambda_1\boldsymbol{I}_{p-1})^{-1}\,\boldsymbol{s}_n$$

$$= \hat{\boldsymbol{\eta}}_{OLS} - n(\boldsymbol{W}^T\boldsymbol{W} + \lambda_1\boldsymbol{I}_{p-1})^{-1}\,[\frac{\lambda_1}{n}\hat{\boldsymbol{\eta}}_{OLS} + \frac{\lambda_2}{2n}\boldsymbol{s}_n].$$

Note that if $\hat{\lambda}_{1,n}/\sqrt{n} \xrightarrow{P} \tau$ and $\hat{\alpha} \xrightarrow{P} \psi$, then $\hat{\lambda}_1/\sqrt{n} \xrightarrow{P} (1-\psi)\tau$ and $\hat{\lambda}_2/\sqrt{n} \xrightarrow{P}$ $2\psi\tau$. The following theorem shows elastic net is asymptotically equivalent to the OLS full model if $\hat{\lambda}_{1,n}/\sqrt{n} \xrightarrow{P} 0$. Note that we get the RR CLT if $\psi = 0$ and the lasso CLT (using $2\hat{\lambda}_{1,n}/\sqrt{n} \xrightarrow{P} 2\tau$) if $\psi = 1$. Under these conditions,

$$\sqrt{n}(\hat{\boldsymbol{\eta}}_{EN} - \boldsymbol{\eta}) = \sqrt{n}(\hat{\boldsymbol{\eta}}_{OLS} - \boldsymbol{\eta}) - n(\boldsymbol{W}^T\boldsymbol{W} + \hat{\lambda}_1\boldsymbol{I}_{p-1})^{-1}[\frac{\hat{\lambda}_1}{\sqrt{n}}\hat{\boldsymbol{\eta}}_{OLS} + \frac{\hat{\lambda}_2}{2\sqrt{n}}\boldsymbol{s}_n].$$

The following theorem is due to Slawski et al. (2010), and summarized in Pelawa Watagoda and Olive (2020).

**Theorem 7.8, Elastic Net CLT.** Assume $p$ is fixed and that the conditions of the OLS CLT Equation (7.14) hold for the model $\boldsymbol{Z} = \boldsymbol{W}\boldsymbol{\eta} + \boldsymbol{e}$.
a) If $\hat{\lambda}_{1,n}/\sqrt{n} \xrightarrow{P} 0$, then

$$\sqrt{n}(\hat{\boldsymbol{\eta}}_{EN} - \boldsymbol{\eta}) \xrightarrow{D} N_{p-1}(\boldsymbol{0}, \sigma^2\boldsymbol{V}).$$

b) If $\hat{\lambda}_{1,n}/\sqrt{n} \xrightarrow{P} \tau \geq 0$, $\hat{\alpha} \xrightarrow{P} \psi \in [0,1]$, and $\boldsymbol{s}_n \xrightarrow{P} \boldsymbol{s} = \boldsymbol{s}_{\boldsymbol{\eta}}$, then

$$\sqrt{n}(\hat{\boldsymbol{\eta}}_{EN} - \boldsymbol{\eta}) \xrightarrow{D} N_{p-1}\left(-\boldsymbol{V}[(1-\psi)\tau\boldsymbol{\eta} + \psi\tau\boldsymbol{s}], \sigma^2\boldsymbol{V}\right).$$

**Proof.** By the above remarks and the RR CLT Theorem 7.6,

$$\sqrt{n}(\hat{\boldsymbol{\eta}}_{EN} - \boldsymbol{\eta}) = \sqrt{n}(\hat{\boldsymbol{\eta}}_{EN} - \hat{\boldsymbol{\eta}}_R + \hat{\boldsymbol{\eta}}_R - \boldsymbol{\eta}) = \sqrt{n}(\hat{\boldsymbol{\eta}}_R - \boldsymbol{\eta}) + \sqrt{n}(\hat{\boldsymbol{\eta}}_{EN} - \hat{\boldsymbol{\eta}}_R)$$

$$\xrightarrow{D} N_{p-1}\left(-(1-\psi)\tau\boldsymbol{V}\boldsymbol{\eta}, \sigma^2\boldsymbol{V}\right) \quad - \quad \frac{2\psi\tau}{2}\boldsymbol{V}\boldsymbol{s}$$

$$\sim N_{p-1}\left(-\boldsymbol{V}[(1-\psi)\tau\boldsymbol{\eta} + \psi\tau\boldsymbol{s}], \sigma^2\boldsymbol{V}\right).$$

The mean of the normal distribution is $\boldsymbol{0}$ under a) since $\hat{\alpha}$ and $\boldsymbol{s}_n$ are bounded. $\square$

**Example 7.5**, continued. The rpack function enet does elastic net using 10-fold CV and a grid of $\alpha$ values $\{0, 1/am, 2/am, ..., am/am = 1\}$. The default uses $am = 10$. The default chose lasso with $alph = 1$. The function also makes a response plot, but does not add the lines for the pointwise prediction intervals since the false degrees of freedom $d$ is not computed.

```
library(glmnet); y <- marry[,3]; x <- marry[,-3]
tem <- enet(x,y)
tem$alph
[1] 1  #elastic net was lasso
tem<-enet(x,y,am=100)
tem$alph
[1] 0.97 #elastic net was not lasso with a finer grid
```

The *elastic net variable selection* estimator applies OLS to a constant and the active predictors that have nonzero elastic net $\hat{\eta}_i$. Hence elastic net is used as a variable selection method. Let $\boldsymbol{X}_A$ denote the matrix with a column of ones and the unstandardized active nontrivial predictors. Hence the elastic net variable selection estimator is $\hat{\boldsymbol{\beta}}_{ENVS} = (\boldsymbol{X}_A^T \boldsymbol{X}_A)^{-1} \boldsymbol{X}_A^T \boldsymbol{Y}$, and relaxed elastic net is an alternative to forward selection. Let $k$ be the number of active (nontrivial) predictors so $\hat{\boldsymbol{\beta}}_{ENVS}$ is $(k+1) \times 1$. Let $I_{min}$ correspond to the elastic net variable selection estimator and $\hat{\boldsymbol{\beta}}_{VS} = \hat{\boldsymbol{\beta}}_{ENVS,0} = \hat{\boldsymbol{\beta}}_{I_{min},0}$ to the zero padded elastic net variable selection estimator. Then by Remark 7.5 where $p$ is fixed, $\hat{\boldsymbol{\beta}}_{ENVS,0}$ is $\sqrt{n}$ consistent when elastic net is consistent, with the limiting distribution for $\hat{\boldsymbol{\beta}}_{ENVS,0}$ given by Theorem 7.4. Hence elastic net variable selection can be bootstrapped with the same methods used for forward selection in Section 7.4. Elastic net variable selection will often be better than elastic net when the model is sparse or if $n \geq 10(k+1)$. The elastic net can be better than elastic net variable selection if $(\boldsymbol{X}_A^T \boldsymbol{X}_A)$ is ill conditioned or if $n/(k+1) < 10$. Also see Olive (2019) and Rathnayake and Olive (2020).

## 7.12 Prediction Intervals

This section will develop prediction intervals after variable selection. Prediction intervals were considered in Sections 2.4 and 5.4.

The additive error regression model is $Y = m(\boldsymbol{x}) + e$ where $m(\boldsymbol{x})$ is a real valued function and the $e_i$ are iid, often with zero mean and constant variance $V(e) = \sigma^2$. The large sample theory for prediction intervals is simple for this model, and variable selection models for the multiple linear regression model have this form with $m(\boldsymbol{x}) = \boldsymbol{x}^T \boldsymbol{\beta} = \boldsymbol{x}_I^T \boldsymbol{\beta}_I$ if $S \subseteq I$. Let the residuals $r_i = Y_i - \hat{m}(\boldsymbol{x}_i) = Y_i - \hat{Y}_i$ for $i = 1, ..., n$. Assume $\hat{m}(\boldsymbol{x})$ is a consistent estimator of $m(\boldsymbol{x})$ such that the sample percentiles $[\hat{L}_n(r), \hat{U}_n(r)]$ of the residuals are consistent estimators of the population percentiles $[L, U]$ of the error distribution where $P(e \in [L, U]) = 1 - \delta$. Let $\hat{Y}_f = \hat{m}(\boldsymbol{x}_f)$. Then $P(Y_f \in [\hat{Y}_f + \hat{L}_n(r), \hat{Y}_f + \hat{U}_n(r)]) \to P(Y_f \in [m(\boldsymbol{x}_f) + L, m(\boldsymbol{x}_f) + U]) = P(e \in [L, U]) = 1 - \delta$ as $n \to \infty$. Three common choices are a) $P(e \leq U) = 1 - \delta/2$ and $P(e \leq L) = \delta/2$, b) $P(e^2 \leq U^2) = P(|e| \leq U) = P(-U \leq e \leq U) = 1 - \delta$ with $L = -U$, and c) the population shorth is the shortest interval (with length $U - L$) such that $P[e \in [L, U]) = 1 - \delta$. The PI c) is asymptotically optimal while a) and b) are asymptotically optimal on the class of symmetric zero mean unimodal error distributions. The split conformal PI (7.36), described below, estimates $[-U, U]$ in b).

Prediction intervals based on the shorth of the residuals need a correction factor for good coverage since the residuals tend to underestimate the errors in magnitude. With the exception of ridge regression, let $d$ be the number

of "variables" used by the method. For MLR, forward selection, lasso, and relaxed lasso use variables $x_1^*, ..., x_d^*$ while PCR and PLS use variables that are linear combinations of the predictors $V_j = \gamma_j^T x$ for $j = 1, ..., d$. (We could let $d = j$ if $j$ is the degrees of freedom of the selected model if that model was chosen in advance without model or variable selection. Hence $d = j$ is not the model degrees of freedom if model selection was used.) See Hong et al. (2018) for why classical prediction intervals after variable selection fail to work.

For $n/p$ large and $d = p$, Olive (2013a) developed prediction intervals for models of the form $Y_i = m(x_i) + e_i$, and variable selection models for MLR have this form, as noted by Olive (2018). Pelawa Watagoda and Olive (2020) gave two prediction intervals that can be useful even if $n/p$ is not large. These PIs will be defined below. The first PI modifies the Olive (2013a) PI that can only be computed if $n > p$. Olive (2007, 2017a, 2017b, 2018) used similar correction factors for several prediction intervals and prediction regions with $d = p$. We want $n \geq 10d$ so that the model does not overfit.

If the OLS model $I$ has $d$ predictors, and $S \subseteq I$, then

$$E(MSE(I)) = E\left(\sum_{i=1}^{n} \frac{r_i^2}{n-d}\right) = \sigma^2 = E\left(\sum_{i=1}^{n} \frac{e_i^2}{n}\right)$$

and $MSE(I)$ is a $\sqrt{n}$ consistent estimator of $\sigma^2$ for many error distributions by Su and Cook (2012). Also see Freedman (1981). For a wide range of regression models, extrapolation occurs if the leverage $h_f = x_{I,f}^T (X_I^T X_I)^{-1} x_{I,f} > 2d/n$: if $x_{I,f}$ is too far from the data $x_{I,1}, ..., x_{I,n}$, then the model may not hold and prediction can be arbitrarily bad. These results suggests that

$$\sqrt{\frac{n}{n-d}}\sqrt{(1+h_f)} \ \ r_i \approx \sqrt{\frac{n+2d}{n-d}} \ \ r_i \approx e_i.$$

In simulations for prediction intervals and prediction regions with $n = 20d$, the maximum simulated undercoverage was near 5% if $q_n$ in (7.31) is changed to $q_n = 1 - \delta$.

Next we give the correction factor and the first prediction interval. Let $q_n = \min(1 - \delta + 0.05, 1 - \delta + d/n)$ for $\delta > 0.1$ and

$$q_n = \min(1 - \delta/2, 1 - \delta + 10\delta d/n), \quad \text{otherwise.} \tag{7.31}$$

If $1 - \delta < 0.999$ and $q_n < 1 - \delta + 0.001$, set $q_n = 1 - \delta$. Let

$$c = \lceil n q_n \rceil, \tag{7.32}$$

and let

$$b_n = \left(1 + \frac{15}{n}\right)\sqrt{\frac{n+2d}{n-d}} \tag{7.33}$$

if $d \leq 8n/9$, and

$$b_n = 5\left(1 + \frac{15}{n}\right),$$

otherwise. As $d$ gets close to $n$, the model overfits and the coverage will be less than the nominal. The piecewise formula for $b_n$ allows the prediction interval to be computed even if $d \geq n$. Compute the shorth($c$) of the residuals $= [r_{(s)}, r_{(s+c-1)}] = [\tilde{\xi}_{\delta_1}, \tilde{\xi}_{1-\delta_2}]$. Then the first 100 $(1 - \delta)\%$ large sample PI for $Y_f$ is

$$[\hat{m}(\boldsymbol{x}_f) + b_n\tilde{\xi}_{\delta_1}, \hat{m}(\boldsymbol{x}_f) + b_n\tilde{\xi}_{1-\delta_2}]. \tag{7.34}$$

The second PI randomly divides the data into two half sets $H$ and $V$ where $H$ has $n_H = \lceil n/2 \rceil$ of the cases and $V$ has the remaining $n_V = n - n_H$ cases $i_1, ..., i_{n_V}$. The estimator $\hat{m}_H(\boldsymbol{x})$ is computed using the training data set $H$. Then the validation residuals $v_j = Y_{i_j} - \hat{m}_H(\boldsymbol{x}_{i_j})$ are computed for the $j = 1, ..., n_V$ cases in the validation set $V$. Find the Frey PI $[v_{(s)}, v_{(s+c-1)}]$ of the validation residuals (replacing $n$ in (2.11) by $n_V = n - n_H$). Then the second new $100(1 - \delta)\%$ large sample PI for $Y_f$ is

$$[\hat{m}_H(\boldsymbol{x}_f) + v_{(s)}, \hat{m}_H(\boldsymbol{x}_f) + v_{(s+c-1)}]. \tag{7.35}$$

**Remark 7.21.** Note that correction factors $b_n \to 1$ are used in large sample confidence intervals and tests if the limiting distribution is N(0,1) or $\chi_p^2$, but a $t_{d_n}$ or $pF_{p,d_n}$ cutoff is used: $t_{d_n,1-\delta}/z_{1-\delta} \to 1$ and $pF_{p,d_n,1-\delta}/\chi_{p,1-\delta}^2 \to 1$ if $d_n \to \infty$ as $n \to 1$. Using correction factors for large sample confidence intervals, tests, prediction intervals, prediction regions, and bootstrap confidence regions improves the performance for moderate sample size $n$.

**Remark 7.22.** For a good fitting model, residuals $r_i$ tend to be smaller in magnitude than the errors $e_i$, while validation residuals $v_i$ tend to be larger in magnitude than the $e_i$. Thus the Frey correction factor can be used for PI (7.35) while PI (7.34) needs a stronger correction factor.

We can also motivate PI (7.35) by modifying the justification for the Lei et al. (2018) split conformal prediction interval

$$[\hat{m}_H(\boldsymbol{x}_f) - a_q, \hat{m}_H(\boldsymbol{x}_f) + a_q] \tag{7.36}$$

where $a_q$ is the $100(1 - \alpha)$th quantile of the absolute validation residuals. PI (7.35) is a modification of the split conformal PI that is asymptotically optimal. Suppose $(Y_i, \boldsymbol{x}_i)$ are iid for $i = 1, ..., n, n + 1$ where $(Y_f, \boldsymbol{x}_f) = (Y_{n+1}, \boldsymbol{x}_{n+1})$. Compute $\hat{m}_H(\boldsymbol{x})$ from the cases in $H$. For example, get $\hat{\boldsymbol{\beta}}_H$ from the cases in $H$. Consider the validation residuals $v_i$ for $i = 1, ..., n_V$ and the validation residual $v_{n_V+1}$ for case $(Y_f, \boldsymbol{x}_f)$. Since these $n_V + 1$ cases are iid, the probability that $v_t$ has rank $j$ for $j = 1, ..., n_V + 1$ is $1/(n_V + 1)$ for each $t$, i.e., the ranks follow the discrete uniform distribution. Let $t = n_V + 1$ and let the $v_{(j)}$ be the ordered residuals using $j = 1, ..., n_V$. That is, get the

order statistics without using the unknown validation residual $v_{n_V+1}$. Then $v_{(i)}$ has rank $i$ if $v_{(i)} < v_{n_V+1}$ but rank $i+1$ if $v_{(i)} > v_{n_V+1}$. Thus

$$P(Y_f \in [\hat{m}_H(\boldsymbol{x}_f)+v_{(k)}, \hat{m}_H(\boldsymbol{x}_f)+v_{(k+b-1)}]) = P(v_{(k)} \leq v_{n_V+1} \leq v_{(k+b-1)}) \geq$$

$P(v_{n_V+1}$ has rank between $k+1$ and $k+b-1$ and there are no tied ranks) $\geq (b-1)/(n_V+1) \approx 1-\delta$ if $b = \lceil (n_V+1)(1-\delta) \rceil + 1$ and $k+b-1 \leq n_V$. This probability statement holds for a fixed $k$ such as $k = \lceil n_V \, \delta/2 \rceil$. The statement is not true when the shorth($b$) estimator is used since the shortest interval using $k = s$ can have $s$ change with the data set. That is, $s$ is not fixed. Hence if PI's were made from $J$ independent data sets, the PI's with fixed $k$ would contain $Y_f$ about $J(1-\delta)$ times, but this value would be smaller for the shorth($b$) prediction intervals where $s$ can change with the data set. The above argument works if the estimator $\hat{m}(\boldsymbol{x})$ is "symmetric in the data," which is satisfied for multiple linear regression estimators.

The PIs (7.34) to (7.36) can be used with $\hat{m}(\boldsymbol{x}) = \hat{Y}_f = \boldsymbol{x}_{I_d}^T \hat{\boldsymbol{\beta}}_{I_d}$ where $I_d$ denotes the index of predictors selected from the model or variable selection method. If $\hat{\boldsymbol{\beta}}$ is a consistent estimator of $\boldsymbol{\beta}$, the PIs (7.34) and (7.35) are asymptotically optimal for a large class of error distributions while the split conformal PI (7.36) needs the error distribution to be unimodal and symmetric for asymptotic optimality. Since $\hat{m}_H$ uses $n/2$ cases, $\hat{m}_H$ has about half the efficiency of $\hat{m}$. When $p \geq n$, the regularity conditions for consistent estimators are strong. For example, EBIC and lasso can have $P(S \subseteq I_{min}) \to 1$ as $n \to \infty$. Then forward selection with EBIC and lasso variable selection can produce consistent estimators. PLS can be $\sqrt{n}$ consistent.

None of the three prediction intervals (7.34), (7.35), and (7.36) dominates the other two. Recall that $\boldsymbol{\beta}_S$ is an $a_S \times 1$ vector in (7.1). If a good fitting method, such as lasso or forward selection with EBIC, is used, and $1.5a_S \leq n \leq 5a_S$, then PI (7.34) can be much shorter than PIs (7.35) and (7.36). For $n/d$ large, PIs (7.34) and (7.35) can be shorter than PI (7.36) if the error distribution is not unimodal and symmetric; however, PI (7.36) is often shorter if $n/d$ is not large since the sample shorth converges to the population shorth rather slowly. Grübel (1982) shows that for iid data, the length and center the shorth($k_n$) interval are $\sqrt{n}$ consistent and $n^{1/3}$ consistent estimators of the length and center of the population shorth interval. For a unimodal and symmetric error distribution, the three PIs are asymptotically equivalent, but PI (4.16) can be the shortest PI due to different correction factors.

If the estimator is poor, the split conformal PI (7.36) and PI (7.35) can have coverage closer to the nominal coverage than PI (7.34). For example, if $\hat{m}$ interpolates the data and $\hat{m}_H$ interpolates the training data from $H$, then the validation residuals will be huge. Hence PI (7.35) will be long compared to PI (7.36).

Asymptotically optimal PIs estimate the population shorth of the zero mean error distribution. Hence PIs that use the shorth of the residuals, such

as PIs (7.34) and (7.35), are the only easily computed asymptotically optimal PIs for a wide range of consistent estimators $\hat{\boldsymbol{\beta}}$ of $\boldsymbol{\beta}$ for the multiple linear regression model. If the error distribution is $e \sim EXP(1) - 1$, then the asymptotic length of the 95% PI (7.34) or (7.35) is 2.966 while that of the split conformal PI is $2(1.966) = 3.992$. For more about these PIs applied to MLR models, see Section 5.4 and Pelawa Watagoda and Olive (2020).

## 7.13 Outlier Resistant MLR Methods

Several methods from Section 6.1 can be modified to give outlier resistant MLR methods. Replace OLS by the MLR method such as lasso, elastic net, ridge regression, or forward selection.

   The first outlier resistant regression method was given by Application 3.3. Call the estimator the *MLD set MLR estimator*. Let the $i$th case $\boldsymbol{w}_i = (Y_i, \boldsymbol{x}_i^T)^T$ where the continuous predictors from $\boldsymbol{x}_i$ are denoted by $\boldsymbol{u}_i$ for $i = 1, ..., n$. Now let $D$ be the RMVN set $U$, the RFCH set $V$, or the covmb2 set $B$. Find $D$ by applying the MLD estimator to the $\boldsymbol{u}_i$, and then run the MLR method on the $m$ cases $\boldsymbol{w}_i$ corresponding to the set $D$ indices $i_1, ..., i_m$, where $m \geq n/2$. The set $B$ can be used even if $p > n$. The theory of the MLR method applies to the cleaned data set since $Y$ was not used to pick the subset of the data. Efficiency can be much lower since $m$ cases are used where $n/2 \leq m \leq n$, and the trimmed cases tend to be the "farthest" from the center of $\boldsymbol{u}$. The *rpack* function getu gets the RMVN set $U$. See the following $R$ code for the Buxton (1920) data where we could use the covmb2 set $B$ instead of the RMVN set $U$ by replacing the command *getu(x)* by getB(x). See Example 3.9.

   Second, replace OLS by the MLR method for the trimmed views or tvreg estimator. For $p > n$ or $n/p$ not large, trimming could be use the Euclidean distance from the coordinatewise median with $\boldsymbol{C}^{-1} = \boldsymbol{I}$ or use a regularized version of $\boldsymbol{C}_{covmb2}$ from Definition 3.26.

   Third, the MLR estimator can be applied to the RMVN set when RMVN is computed from the vectors $\boldsymbol{u}_i = (x_{i2}, ..., x_{ip}, Y_i)^T$ for $i = 1, ..., n$. Hence $\boldsymbol{u}_i$ is the $i$th case with $x_{i1} = 1$ deleted. This estimator is similar to the rmreg2 estimator that used OLS.

## 7.14 Summary

1) A *model for variable selection* can be described by $\boldsymbol{x}^T\boldsymbol{\beta} = \boldsymbol{x}_S^T\boldsymbol{\beta}_S + \boldsymbol{x}_E^T\boldsymbol{\beta}_E = \boldsymbol{x}_S^T\boldsymbol{\beta}_S$ where $\boldsymbol{x} = (\boldsymbol{x}_S^T, \boldsymbol{x}_E^T)^T$ is a $p \times 1$ vector of predictors, $\boldsymbol{x}_S$ is an $a_S \times 1$ vector, and $\boldsymbol{x}_E$ is a $(p - a_S) \times 1$ vector. Given that $\boldsymbol{x}_S$ is in the model, $\boldsymbol{\beta}_E = \boldsymbol{0}$. Assume $p$ is fixed while $n \to \infty$.

2) If $\hat{\boldsymbol{\beta}}_I$ is $a \times 1$, form the $p \times 1$ vector $\hat{\boldsymbol{\beta}}_{I,0}$ from $\hat{\boldsymbol{\beta}}_I$ by adding 0s corresponding to the omitted variables. For example, if $p = 4$ and $\hat{\boldsymbol{\beta}}_{I_{min}} = (\hat{\beta}_1, \hat{\beta}_3)^T$, then $\hat{\boldsymbol{\beta}}_{I_{min},0} = (\hat{\beta}_1, 0, \hat{\beta}_3, 0)^T$. For the OLS model with $S \subseteq I$, $\sqrt{n}(\hat{\boldsymbol{\beta}}_I - \boldsymbol{\beta}_I) \xrightarrow{D} N_{a_I}(\mathbf{0}, \boldsymbol{V}_I)$ where $(\boldsymbol{X}_I^T \boldsymbol{X}_I)/(n\sigma^2) \xrightarrow{P} \boldsymbol{V}_I^{-1}$.

3) **Theorem 7.3.** Assume $P(S \subseteq I_{min}) \to 1$ as $n \to \infty$, and let $\hat{\boldsymbol{\beta}}_{MIX} = \hat{\boldsymbol{\beta}}_{I_k,0}$ with probabilities $\pi_{kn}$ where $\pi_{kn} \to \pi_k$ as $n \to \infty$. Denote the positive $\pi_k$ by $\pi_j$. Assume $\boldsymbol{u}_{jn} = \sqrt{n}(\hat{\boldsymbol{\beta}}_{I_j,0} - \boldsymbol{\beta}) \xrightarrow{D} \boldsymbol{u}_j \sim N_p(\mathbf{0}, \boldsymbol{V}_{j,0})$. a) Then

$$\boldsymbol{u}_n = \sqrt{n}(\hat{\boldsymbol{\beta}}_{MIX} - \boldsymbol{\beta}) \xrightarrow{D} \boldsymbol{u} \tag{7.37}$$

where the cdf of $\boldsymbol{u}$ is $F_{\boldsymbol{u}}(\boldsymbol{t}) = \sum_j \pi_j F_{\boldsymbol{u}_j}(\boldsymbol{t})$. Thus $\boldsymbol{u}$ has a mixture distribution of the $\boldsymbol{u}_j$ with probabilities $\pi_j$, $E(\boldsymbol{u}) = \mathbf{0}$, and $\text{Cov}(\boldsymbol{u}) = \boldsymbol{\Sigma}_{\boldsymbol{u}} = \sum_j \pi_j \boldsymbol{V}_{j,0}$.

b) Let $\boldsymbol{A}$ be a $g \times p$ full rank matrix with $1 \le g \le p$. Then

$$\boldsymbol{v}_n = \boldsymbol{A} \boldsymbol{u}_n = \sqrt{n}(\boldsymbol{A}\hat{\boldsymbol{\beta}}_{MIX} - \boldsymbol{A}\boldsymbol{\beta}) \xrightarrow{D} \boldsymbol{A}\boldsymbol{u} = \boldsymbol{v} \tag{7.38}$$

where $\boldsymbol{v}$ has a mixture distribution of the $\boldsymbol{v}_j = \boldsymbol{A}\boldsymbol{u}_j \sim N_g(\mathbf{0}, \boldsymbol{A}\boldsymbol{V}_{j,0}\boldsymbol{A}^T)$ with probabilities $\pi_j$.

c) The estimator $\hat{\boldsymbol{\beta}}_{VS}$ is a $\sqrt{n}$ consistent estimator of $\boldsymbol{\beta}$. Hence $\sqrt{n}(\hat{\boldsymbol{\beta}}_{VS} - \boldsymbol{\beta}) = O_P(1)$.

d) If $\pi_d = 1$, then $\sqrt{n}(\hat{\boldsymbol{\beta}}_{SEL} - \boldsymbol{\beta}) \xrightarrow{D} \boldsymbol{u} \sim N_p(\mathbf{0}, \boldsymbol{V}_{d,0})$ where $SEL$ is $VS$ or $MIX$.

4) **Theorem 7.4, Variable Selection CLT.** Assume $P(S \subseteq I_{min}) \to 1$ as $n \to \infty$, and let $\hat{\boldsymbol{\beta}}_{VS} = \hat{\boldsymbol{\beta}}_{I_k,0}$ with probabilities $\pi_{kn}$ where $\pi_{kn} \to \pi_k$ as $n \to \infty$. Denote the positive $\pi_k$ by $\pi_j$. Assume $\boldsymbol{w}_{jn} = \sqrt{n}(\hat{\boldsymbol{\beta}}_{I_j,0}^C - \boldsymbol{\beta}) \xrightarrow{D} \boldsymbol{w}_j$. Then

$$\boldsymbol{w}_n = \sqrt{n}(\hat{\boldsymbol{\beta}}_{VS} - \boldsymbol{\beta}) \xrightarrow{D} \boldsymbol{w} \tag{7.39}$$

where the cdf of $\boldsymbol{w}$ is $F_{\boldsymbol{w}}(\boldsymbol{t}) = \sum_j \pi_j F_{\boldsymbol{w}_j}(\boldsymbol{t})$. Thus $\boldsymbol{w}$ is a mixture distribution of the $\boldsymbol{w}_j$ with probabilities $\pi_j$.

| Label | coef | SE | shorth 95% CI for $\beta_i$ |
|---|---|---|---|
| 5)  Constant=intercept= $x_1$ | $\hat{\beta}_1$ | $SE(\hat{\beta}_1)$ | $[\hat{L}_1, \hat{U}_1]$ |
| $x_2$ | $\hat{\beta}_2$ | $SE(\hat{\beta}_2)$ | $[\hat{L}_2, \hat{U}_2]$ |
| $\vdots$ | | | |
| $x_p$ | $\hat{\beta}_p$ | $SE(\hat{\beta}_p)$ | $[\hat{L}_p, \hat{U}_p]$ |

The classical OLS large sample 95% CI for $\beta_i$ is $\hat{\beta}_i \pm 1.96 SE(\hat{\beta}_i)$. Consider testing $H_0 : \beta_i = 0$ versus $H_A : \beta_i \ne 0$. If $0 \in$ CI for $\beta_i$, then fail to reject $H_0$, and conclude $x_i$ is not needed in the MLR model given the other predictors are in the model. If $0 \notin$ CI for $\beta_i$, then reject $H_0$, and conclude $x_i$ is needed in the MLR model.

6) A *model for variable selection* is $\boldsymbol{x}^T\boldsymbol{\beta} = \boldsymbol{x}_S^T\boldsymbol{\beta}_S + \boldsymbol{x}_E^T\boldsymbol{\beta}_E = \boldsymbol{x}_S^T\boldsymbol{\beta}_S$ where $\boldsymbol{x} = (\boldsymbol{x}_S^T, \boldsymbol{x}_E^T)^T$, $\boldsymbol{x}_S$ is an $a_S \times 1$ vector, and $\boldsymbol{x}_E$ is a $(p - a_S) \times 1$ vector. Let $\boldsymbol{x}_I$ be the vector of $a$ terms from a candidate subset indexed by $I$, and let $\boldsymbol{x}_O$ be the vector of the remaining predictors (out of the candidate submodel). If $S \subseteq I$, then $\boldsymbol{x}^T\boldsymbol{\beta} = \boldsymbol{x}_S^T\boldsymbol{\beta}_S = \boldsymbol{x}_S^T\boldsymbol{\beta}_S + \boldsymbol{x}_{I/S}^T\boldsymbol{\beta}_{(I/S)} + \boldsymbol{x}_O^T\boldsymbol{0} = \boldsymbol{x}_I^T\boldsymbol{\beta}_I$ where $\boldsymbol{x}_{I/S}$ denotes the predictors in $I$ that are not in $S$. Since this is true regardless of the values of the predictors, $\boldsymbol{\beta}_O = \boldsymbol{0}$ if $S \subseteq I$. Note that $\boldsymbol{\beta}_E = \boldsymbol{0}$. Let $k_S = a_S - 1 = $ the number of population active nontrivial predictors. Then $k = a - 1$ is the number of active predictors in the candidate submodel $I$.

7)

| $I_j$ | model | $x_2$ | $x_3$ | $x_4$ | $x_5$ | $\hat{\boldsymbol{\beta}}_{I_j,0}$ if $\hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}}_{I_j}$ |
|---|---|---|---|---|---|---|
| $I_2$ | 1 | | $*$ | | | $(\hat{\beta}_1, 0, \hat{\beta}_3, 0, 0)^T$ |
| $I_3$ | 2 | | $*$ | $*$ | | $(\hat{\beta}_1, 0, \hat{\beta}_3, \hat{\beta}_4, 0)^T$ |
| $I_4$ | 3 | $*$ | $*$ | $*$ | | $(\hat{\beta}_1, \hat{\beta}_2, \hat{\beta}_3, \hat{\beta}_4, 0)^T$ |
| $I_5$ | 4 | $*$ | $*$ | $*$ | $*$ | $(\hat{\beta}_1, \hat{\beta}_2, \hat{\beta}_3, \hat{\beta}_4, \hat{\beta}_4)^T = \hat{\boldsymbol{\beta}}_{OLS}$ |

Model $I_{min}$ is the model, among $p$ candidates, that minimizes $C_p$ if $n \geq 10$, or EBIC if $n < 10p$. Model $I_j$ contains $j$ predictors, $x_1^*, x_2^*, ..., x_j^*$ where $x_1^* = x_1 \equiv 1$, the constant.

8) Variable selection is a search for a subset of predictors that can be deleted without important loss of information if $n \geq 10p$ and such that model $I$ (containing the remaining predictors that were not deleted) is good for prediction if $n < 10p$. Note that the "100%" shorth CI for a $\beta_i$ that is a component of $\boldsymbol{\beta}_O$ is $[0,0]$.

9) Underfitting occurs if $S \not\subseteq I$ so that $\boldsymbol{x}_I$ is missing important predictors. Underfitting will occur if $\boldsymbol{x}_I$ is $k \times 1$ with $d = k < a_S$. Overfitting occurs if $S \subset I$ with $S \neq I$ or if $n < 5k$.

10) In 7) sometimes TRUE = $*$ and FALSE = blank. The $x_i$ may be replaced by the variable name or letters like a  b  c  d.

| $I_j$ | model | $x_2$ | $x_3$ | $x_4$ | $x_5$ |
|---|---|---|---|---|---|
| $I_2$ | 1 | FALSE | TRUE | FALSE | FALSE |
| $I_3$ | 2 | FALSE | TRUE | TRUE | FALSE |
| $I_4$ | 3 | TRUE | TRUE | TRUE | FALSE |
| $I_5$ | 4 | TRUE | TRUE | TRUE | TRUE |

11) The out$cp line gives $C_p(I_2), C_p(I_3), ..., C_p(I_p) = p$ and $I_{min}$ is the $I_j$ with the smallest $C_p$.

12) Typical bootstrap output for forward selection, lasso, and elastic net is shown below. The SE column is usually omitted except possibly for forward selection. The term "coef" might be replaced by "Estimate." This column gives $\hat{\boldsymbol{\beta}}_{I,0}$ where $I = I_{min}$ for forward selection, $I = L$ for lasso, and $I = EN$ for elastic net. Note that the SE entry is omitted if $\hat{\beta}_i = 0$ so variable $x_i$ was omitted by the variable selection method. In the output below, $\hat{\beta}_2 = \hat{\beta}_3 = 0$. The SE column corresponds to the OLS SE obtained by acting as if the OLS full model contains a constant and the variables not omitted by the variable

selection method. The OLS SE is incorrect unless the variables were selected before looking at the data for forward selection.

| Label | Estimate or coef | SE | shorth 95% CI for $\beta_i$ |
|---|---|---|---|
| Constant=intercept= $x_1$ | $\hat{\beta}_1$ | $SE(\hat{\beta}_1)$ | $[\hat{L}_1, \hat{U}_1]$ |
| $x_2$ | $\hat{\beta}_2$ | $SE(\hat{\beta}_2)$ | $[\hat{L}_2, \hat{U}_2]$ |
| $x_3$ | 0 | | $[\hat{L}_3, \hat{U}_3]$ |
| $x_4$ | 0 | | $[\hat{L}_4, \hat{U}_4]$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| $x_p$ | $\hat{\beta}_p$ | $SE(\hat{\beta}_p)$ | $[\hat{L}_p, \hat{U}_p]$ |

13) The OLS SE is also accurate for forward selection with $C_p$ if $\boldsymbol{X}^T\boldsymbol{X}/n \to \boldsymbol{V}^{-1} = diag(d_1, ..., d_p)$ where all $d_i > 0$. The diagonal limit matrix will occur if the predictors are orthogonal or if the nontrivial predictors are independent with 0 mean and finite variance.

```
regbootsim3(nruns=500)
$cicov
0.942 0.954 0.950 0.948 0.944 0.946 0.946 0.940 0.938 0.940
$avelen
0.398 0.399 0.397 0.399 2.448 2.448 2.448 2.448 2.448 2.450
$beta
[1] 1 1 0 0
$k
[1] 1
```

14) Simulation output for regression is similar to that shown above. Usually want coverage near 0.95 since nominal 95% CIs are used and tests with nominal $\delta = 0.05$ are used. To suggest that the actual coverage is near the nominal coverage of 0.95, want cov in [0.94,0.96] with 5000 runs, want cov in [0.93,0.97], with 1000 runs, want cov in [0.92,0.98] with 500 runs, and want cov in [0.91,0.99] with 100 runs. Let $SP = \boldsymbol{x}^T\boldsymbol{\beta} = \beta_1 + 1x_{i,2} + \cdots + 1x_{i,k+1}$ for $i = 1, ..., n$. Hence $\boldsymbol{\beta} = (\beta_1, 1, ..., 1, 0, ..., 0)^T$ with $\beta_1$, $k$ ones, and $p - k - 1$ zeros. Then $S = \{1, ..., k+1\}$ and $E = \{k+2, ..., p\}$. Note that $S$ corresponds to the first $k + 1$ $\beta_i$ while $E$ corresponds to the last $p - k + 1$ $\beta_i$.

The first 4 numbers are the bootstrap shorth confidence intervals for $\beta_1, \beta_2, \beta_{p-1}$, and $\beta_p$. The average lengths of the CIs along with the proportion of times (coverage) the CI for $\beta_i$ contained $\beta_i$ are given. The next three numbers test $H_0 : \boldsymbol{\beta}_E = \boldsymbol{0}$. The prediction region method, hybrid method, and Bickel and Ren methods are used. Hence the fifth interval gives the length of the interval $[0, D_{(c)}]$ where $H_0$ is rejected if $D_0 > D_{(c)}$ and the fifth "coverage" is the proportion of times the prediction region method test fails to reject $H_0$. The last three numbers are similar but test $H_0 : \boldsymbol{\beta}_S = (\beta_1, 1, ..., 1)^T$. Hence the last length 2.450 corresponds to the Bickel and Ren method with cover-

age 0.940. For the output shown, lengths near 2.45 correspond to $\sqrt{\chi_2^2(0.95)}$ where $P(X \leq \chi_2^2(0.95)) = 0.95$ if $X \sim \chi_2^2$.

15) Let $\boldsymbol{x}_i^T = (1 \quad \boldsymbol{u}_i^T)$. It is often convenient to use the centered response $\boldsymbol{Z} = \boldsymbol{Y} - \overline{\boldsymbol{Y}}$ where $\overline{\boldsymbol{Y}} = \overline{Y}\boldsymbol{1}$, and the $n \times (p-1)$ matrix of standardized nontrivial predictors $\boldsymbol{W} = (W_{ij})$. For $j = 1, ..., p-1$, let $W_{ij}$ denote the $(j+1)$th variable standardized so that $\sum_{i=1}^n W_{ij} = 0$ and $\sum_{i=1}^n W_{ij}^2 = n$. Then the sample correlation matrix of the nontrivial predictors $\boldsymbol{u}_i$ is

$$\boldsymbol{R}\boldsymbol{u} = \frac{\boldsymbol{W}^T\boldsymbol{W}}{n}.$$

Then regression through the origin is used for the model $\boldsymbol{Z} = \boldsymbol{W}\boldsymbol{\eta} + \boldsymbol{e}$ where the vector of fitted values $\hat{\boldsymbol{Y}} = \overline{\boldsymbol{Y}} + \hat{\boldsymbol{Z}}$. Thus the centered response $Z_i = Y_i - \overline{Y}$ and $\hat{Y}_i = \hat{Z}_i + \overline{Y}$. Then $\hat{\boldsymbol{\eta}}$ does not depend on the units of measurement of the predictors. Linear combinations of the $\boldsymbol{u}_i$ can be written as linear combinations of the $\boldsymbol{x}_i$, hence $\hat{\boldsymbol{\beta}}$ can be found from $\hat{\boldsymbol{\eta}}$.

16) Consider choosing $\hat{\boldsymbol{\eta}}$ to minimize the criterion

$$Q(\boldsymbol{\eta}) = \frac{1}{a}(\boldsymbol{Z} - \boldsymbol{W}\boldsymbol{\eta})^T(\boldsymbol{Z} - \boldsymbol{W}\boldsymbol{\eta}) + \frac{\lambda_{1,n}}{a}\sum_{i=1}^{p-1}|\eta_i|^j \qquad (7.40)$$

where $\lambda_{1,n} \geq 0$, $a > 0$, and $j > 0$ are known constants. Then $j = 2$ corresponds to ridge regression $\hat{\boldsymbol{\eta}}_R$, $j = 1$ corresponds to lasso $\hat{\boldsymbol{\eta}}_L$, and $a = 1, 2, n$, and $2n$ are common. The residual sum of squares $RSS_W(\boldsymbol{\eta}) = (\boldsymbol{Z} - \boldsymbol{W}\boldsymbol{\eta})^T(\boldsymbol{Z} - \boldsymbol{W}\boldsymbol{\eta})$, and $\lambda_{1,n} = 0$ corresponds to the OLS estimator $\hat{\boldsymbol{\eta}}_{OLS} = (\boldsymbol{W}^T\boldsymbol{W})^{-1}\boldsymbol{W}^T\boldsymbol{Z}$. Note that for a $k \times 1$ vector $\boldsymbol{\eta}$, the squared (Euclidean) $L_2$ norm $\|\boldsymbol{\eta}\|_2^2 = \boldsymbol{\eta}^T\boldsymbol{\eta} = \sum_{i=1}^k \eta_i^2$ and the $L_1$ norm $\|\boldsymbol{\eta}\|_1 = \sum_{i=1}^k |\eta_i|$.

Lasso and ridge regression have a parameter $\lambda$. When $\lambda = 0$, the OLS full model is used. Otherwise, the centered response and scaled nontrivial predictors are used with $\boldsymbol{Z} = \boldsymbol{W}\boldsymbol{\eta} + \boldsymbol{e}$. See 5). These methods also use a maximum value $\lambda_M$ of $\lambda$ and a grid of $M$ $\lambda$ values $0 \leq \lambda_1 < \lambda_2 < \cdots < \lambda_{M-1} < \lambda_M$ where often $\lambda_1 = 0$. For lasso, $\lambda_M$ is the smallest value of $\lambda$ such that $\hat{\boldsymbol{\eta}}_{\lambda_M} = \boldsymbol{0}$. Hence $\hat{\boldsymbol{\eta}}_{\lambda_i} \neq \boldsymbol{0}$ for $i < M$.

17) The elastic net estimator $\hat{\boldsymbol{\eta}}_{EN}$ minimizes

$$Q_{EN}(\boldsymbol{\eta}) = RSS(\boldsymbol{\eta}) + \lambda_1\|\boldsymbol{\eta}\|_2^2 + \lambda_2\|\boldsymbol{\eta}\|_1 \qquad (7.41)$$

where $\lambda_1 = (1-\alpha)\lambda_{1,n}$ and $\lambda_2 = 2\alpha\lambda_{1,n}$ with $0 \leq \alpha \leq 1$.

18) Use $\boldsymbol{Z}_n \sim AN_g(\boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n)$ to indicate that a normal approximation is used: $\boldsymbol{Z}_n \approx N_g(\boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n)$. Let $a$ be a constant, let $\boldsymbol{A}$ be a $k \times g$ constant matrix, and let $\boldsymbol{c}$ be a $k \times 1$ constant vector. If $\sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}) \xrightarrow{D} N_g(\boldsymbol{0}, \boldsymbol{V})$, then $a\boldsymbol{Z}_n = a\boldsymbol{I}_g\boldsymbol{Z}_n$ with $\boldsymbol{A} = a\boldsymbol{I}_g$,

$$a\boldsymbol{Z}_n \sim AN_g\left(a\boldsymbol{\mu}_n, a^2\boldsymbol{\Sigma}_n\right), \quad \text{and} \quad \boldsymbol{A}\boldsymbol{Z}_n + \boldsymbol{c} \sim AN_k\left(\boldsymbol{A}\boldsymbol{\mu}_n + \boldsymbol{c}, \boldsymbol{A}\boldsymbol{\Sigma}_n\boldsymbol{A}^T\right),$$

$$\hat{\boldsymbol{\theta}}_n \sim AN_g\left(\boldsymbol{\theta}, \frac{\boldsymbol{V}}{n}\right), \quad \text{and} \quad \boldsymbol{A}\hat{\boldsymbol{\theta}}_n + \boldsymbol{c} \sim AN_k\left(\boldsymbol{A}\boldsymbol{\theta} + \boldsymbol{c}, \frac{\boldsymbol{A}\boldsymbol{V}\boldsymbol{A}^T}{n}\right).$$

19) Assume $\hat{\boldsymbol{\eta}}_{OLS} = (\boldsymbol{W}^T\boldsymbol{W})^{-1}\boldsymbol{W}^T\boldsymbol{Z}$. Let $\boldsymbol{s}_n = (s_{1n}, ..., s_{p-1,n})^T$ where $s_{in} \in [-1, 1]$ and $s_{in} = \text{sign}(\hat{\eta}_i)$ if $\hat{\eta}_i \neq 0$. Here $\text{sign}(\eta_i) = 1$ if $\eta_i > 1$ and $\text{sign}(\eta_i) = -1$ if $\eta_i < 1$. Then

i) $\hat{\boldsymbol{\eta}}_R = \hat{\boldsymbol{\eta}}_{OLS} - \dfrac{\lambda_{1n}}{n}n(\boldsymbol{W}^T\boldsymbol{W} + \lambda_{1,n}\boldsymbol{I}_{p-1})^{-1}\hat{\boldsymbol{\eta}}_{OLS}$.

ii) $\hat{\boldsymbol{\eta}}_L = \hat{\boldsymbol{\eta}}_{OLS} - \dfrac{\lambda_{1,n}}{2n} n(\boldsymbol{W}^T\boldsymbol{W})^{-1} \boldsymbol{s}_n$.

iii) $\hat{\boldsymbol{\eta}}_{EN} = \hat{\boldsymbol{\eta}}_{OLS} - n(\boldsymbol{W}^T\boldsymbol{W} + \lambda_1\boldsymbol{I}_{p-1})^{-1} \left[\dfrac{\lambda_1}{n}\hat{\boldsymbol{\eta}}_{OLS} + \dfrac{\lambda_2}{2n}\boldsymbol{s}_n\right]$.

20) Assume that the sample correlation matrix $\boldsymbol{R_u} = \dfrac{\boldsymbol{W}^T\boldsymbol{W}}{n} \xrightarrow{P} \boldsymbol{V}^{-1}$.

Let $\boldsymbol{H} = \boldsymbol{W}(\boldsymbol{W}^T\boldsymbol{W})^{-1}\boldsymbol{W}^T = (h_{ij})$, and assume that $\max_{i=1,...,n} h_{ii} \xrightarrow{P} 0$ as $n \to \infty$. Let $\hat{\boldsymbol{\eta}}_A$ be $\hat{\boldsymbol{\eta}}_{EN}, \hat{\boldsymbol{\eta}}_L$, or $\hat{\boldsymbol{\eta}}_R$. Let $p$ be fixed.

i) OLS CLT: $\sqrt{n}(\hat{\boldsymbol{\eta}}_{OLS} - \boldsymbol{\eta}) \xrightarrow{D} N_{p-1}(\boldsymbol{0}, \sigma^2\boldsymbol{V})$.

ii) If $\hat{\lambda}_{1,n}/\sqrt{n} \xrightarrow{P} 0$, then

$$\sqrt{n}(\hat{\boldsymbol{\eta}}_A - \boldsymbol{\eta}) \xrightarrow{D} N_{p-1}(\boldsymbol{0}, \sigma^2\boldsymbol{V}).$$

iii) If $\hat{\lambda}_{1,n}/\sqrt{n} \xrightarrow{P} \tau \geq 0$, $\hat{\alpha} \xrightarrow{P} \psi \in [0, 1]$, and $\boldsymbol{s}_n \xrightarrow{P} \boldsymbol{s} = \boldsymbol{s_\eta}$, then

$$\sqrt{n}(\hat{\boldsymbol{\eta}}_{EN} - \boldsymbol{\eta}) \xrightarrow{D} N_{p-1}\left(-\boldsymbol{V}[(1 - \psi)\tau\boldsymbol{\eta} + \psi\tau\boldsymbol{s}], \sigma^2\boldsymbol{V}\right).$$

iv) If $\hat{\lambda}_{1,n}/\sqrt{n} \xrightarrow{P} \tau \geq 0$, then

$$\sqrt{n}(\hat{\boldsymbol{\eta}}_R - \boldsymbol{\eta}) \xrightarrow{D} N_{p-1}(-\tau\boldsymbol{V}\boldsymbol{\eta}, \sigma^2\boldsymbol{V}).$$

v) If $\hat{\lambda}_{1,n}/\sqrt{n} \xrightarrow{P} \tau \geq 0$ and $\boldsymbol{s}_n \xrightarrow{P} \boldsymbol{s} = \boldsymbol{s_\eta}$, then

$$\sqrt{n}(\hat{\boldsymbol{\eta}}_L - \boldsymbol{\eta}) \xrightarrow{D} N_{p-1}\left(\dfrac{-\tau}{2}\boldsymbol{V}\boldsymbol{s}, \sigma^2\boldsymbol{V}\right).$$

ii) and v) are the Lasso CLT, ii) and iv) are the RR CLT, and ii) and iii) are the EN CLT.

## 7.15 Complements

This chapter followed Pelawa Watagoda and Olive (2019, 2020) closely. Also see Olive (2013a, 2018), and Rathnayake and Olive (2020). For MLR, Olive (2017a: p. 123, 2017b: p. 176) showed that $\hat{\boldsymbol{\beta}}_{VS} = \hat{\boldsymbol{\beta}}_{I_{min},0}$ is a consistent es-

timator. Olive (2014: p. 283, 2017ab, 2018) recommended using the shorth($c$) estimator as a confidence interval. Olive (2017a: p. 128, 2017b: p. 181, 2018) showed that the prediction region method can simulate well for the $p \times 1$ vector $\hat{\boldsymbol{\beta}}_{VS} = \hat{\boldsymbol{\beta}}_{I_{min},0}$. Hastie et al. (2009, p. 57) noted that variable selection is a shrinkage estimator: the coefficients are shrunk to 0 for the omitted variables.

There is a massive literature on variable selection and a fairly large literature for inference after variable selection. See, for example, Leeb and Pötscher (2006, 2008), and Tibshirani et al. (2018). Knight and Fu (2000) have some results on the residual bootstrap that uses residuals from one estimator, such as full model OLS, but fit another estimator, such as lasso.

Inference techniques for the variable selection model, other than data splitting, have not had much success. For multiple linear regression, the methods are often inferior to data splitting, often assume normality, or are asymptotically equivalent to using the full model, or find a quantity to test that is not $\boldsymbol{A\beta}$. See Ewald and Schneider (2018). Berk et al. (2013) assumes normality, needs $p$ no more than about 30, assumes $\sigma^2$ can be estimated independently of the data, and Leeb et al. (2015) say the method does not work. The bootstrap confidence region (4.32) is centered at $\overline{T}^* \approx \sum_j \rho_{jn} T_{jn}$, which is closely related to a model averaging estimator. Wang and Zhou (2013) show that the Hjort and Claeskens (2003) confidence intervals based on frequentist model averaging are asymptotically equivalent to those obtained from the full model. See Buckland et al. (1997) and Schomaker and Heumann (2014) for standard errors when using the bootstrap or model averaging for linear model confidence intervals.

Efron (2014) used the confidence interval $\overline{T}^* \pm z_{1-\delta} SE(\overline{T}^*)$ assuming $\overline{T}^*$ is asymptotically normal and using delta method techniques, which require nonsingular covariance matrices. There is not yet rigorous theory for this method. Section 7.2 proved that $\overline{T}^*$ is asymptotically normal: under regularity conditions: if $\sqrt{n}(T_n - \boldsymbol{\theta}) \xrightarrow{D} N_g(\boldsymbol{0}, \boldsymbol{\Sigma}_A)$ and $\sqrt{n}(T_i^* - T_n) \xrightarrow{D} N_g(\boldsymbol{0}, \boldsymbol{\Sigma}_A)$, then under regularity conditions $\sqrt{n}(\overline{T}^* - \boldsymbol{\theta}) \xrightarrow{D} N_g(\boldsymbol{0}, \boldsymbol{\Sigma}_A)$. If $g = 1$, then the prediction region method large sample $100(1 - \delta)\%$ CI for $\theta$ has $P(\theta \in [\overline{T}^* - a_{(U_B)}, \overline{T}^* + a_{(U_B)}]) \to 1 - \delta$ as $n \to \infty$. If the Frey CI also has coverage converging to $1 - \delta$, than the two methods have the same asymptotic length (scaled by multiplying by $\sqrt{n}$), since otherwise the shorter interval will have lower asymptotic coverage.

We can get a prediction region by randomly dividing the data into two half sets $H$ and $V$ where $H$ has $n_H = \lceil n/2 \rceil$ of the cases and $V$ has the remaining $m = n_V = n - n_H$ cases. See Section 4.4.

**Robust Versions of OLS Alternatives:** Hastie et al. (2015, pp. 26-27) discuss some modifications of lasso that are robust to certain types of outliers. Robust methods for forward selection and LARS are given by Uraibi et al. (2017, 2019) that need $n >> p$. If $n$ is not much larger than $p$, then Hoffman

et al. (2015) have a robust Partial Least Squares–Lasso type estimator that uses a clever weighting scheme.

## 7.16 Problems

**7.1.** For the output below, an asterisk means the variable is in the model. All models have a constant, so model 1 contains a constant and mmen.

    a) List the variables, including a constant, that models 2, 3, and 4 contain.

    b) The term out$cp lists the $C_p$ criterion. Which model (1, 2, 3, or 4) is the minimum $C_p$ model $I_{min}$?

    c) Suppose $\hat{\boldsymbol{\beta}}_{I_{min}} = (241.5445, 1.001)^T$. What is $\hat{\boldsymbol{\beta}}_{VS} = \hat{\boldsymbol{\beta}}_{I_{min},0}$?

```
Selection Algorithm: forward #output for Problem 7.1
         pop mmen mmilmen milwmn
1  ( 1 ) " " "*"   " "      " "
2  ( 1 ) " " "*"   "*"      " "
3  ( 1 ) "*" "*"   "*"      " "
4  ( 1 ) "*" "*"   "*"      "*"
out$cp
[1] -0.8268967  1.0151462  3.0029429  5.0000000

        large sample full model inference
        Est.    SE   t   Pr(>|t|)   nparboot        resboot
int -1.249 0.838 -1.49 0.14 [-2.93,-0.093][-3.045,0.473]
L   -0.001 0.002 -0.28 0.78 [-0.005,0.003][-0.005,0.004]
logW 0.130 0.374  0.35 0.73 [-0.457,0.829][-0.703,0.890]
H    0.008 0.005  1.50 0.14 [-0.002,0.018][-0.003,0.016]
logS 0.640 0.169  3.80 0.00 [ 0.244,1.040][ 0.336,1.012]
```

    **7.2** Consider the above output for the OLS full model. The column *resboot* gives the large sample 95% CI for $\beta_i$ using the shorth applied to the $\hat{\beta}_{ij}^*$ for $j = 1, ..., B$ using the residual bootstrap. The standard large sample 95% CI for $\beta_i$ is $\hat{\beta}_i \pm 1.96 SE(\hat{\beta}_i)$. Hence for $\beta_2$ corresponding to $L$, the standard large sample 95% CI is $-0.001 \pm 1.96(0.002) = -0.001 \pm 0.00392 = [-0.00492, 0.00292]$ while the shorth 95% CI is $[-0.005, 0.004]$.

    a) Compute the standard 95% CIs for $\beta_i$ corresponding to W, H, and S. Also write down the shorth 95% CI. Are the standard and shorth 95% CIs fairly close?

    b) Consider testing $H_0 : \beta_i = 0$ versus $H_A : \beta_i \neq 0$. If the corresponding 95% CI for $\beta_i$ does not contain 0, then reject $H_0$ and conclude that the predictor variable $X_i$ is needed in the MLR model. If 0 is in the CI then fail to reject $H_0$ and conclude that the predictor variable $X_i$ is not needed in the MLR model given that the other predictors are in the MLR model.

Which variables, if any, are needed in the MLR model? Use the standard CI if the shorth CI gives a different result. The nontrivial predictor variables are L, W, H, and S.

**7.3.** Tremearne (1911) presents a data set of about 17 measurements on 112 people of Hausa nationality. We used $Y = height$. Along with a constant $x_{i,1} \equiv 1$, the five additional predictor variables used were $x_{i,2} = height\ when\ sitting$, $x_{i,3} = height\ when\ kneeling$, $x_{i,4} = head\ length$, $x_{i,5} = nasal\ breadth$, and $x_{i,6} = span$ (perhaps from left hand to right hand). The output below is for the OLS full model.

```
              Estimate Std.Err 95% shorth CI
 Intercept -77.0042 65.2956 [-208.864,55.051]
 X2          0.0156  0.0992 [-0.177,   0.217]
 X3          1.1553  0.0832 [ 0.983,   1.312]
 X4          0.2186  0.3180 [-0.378,   0.805]
 X5          0.2660  0.6615 [-1.038,   1.637]
 X6          0.1396  0.0385 [0.0575,   0.217]
```

a) Give the shorth 95% CI for $\beta_2$.

b) Compute the standard 95% CI for $\beta_2$.

c) Which variables, if any, are needed in the MLR model given that the other variables are in the model?

Now we use forward selection and $I_{min}$ is the minimum $C_p$ model.

```
              Estimate Std.Err 95% shorth CI
 Intercept -42.4846 51.2863 [-192.281, 52.492]
 X2          0               [   0.000,  0.268]
 X3          1.1707  0.0598 [   0.992,  1.289]
 X4          0               [   0.000,  0.840]
 X5          0               [   0.000,  1.916]
 X6          0.1467  0.0368 [  0.0747,  0.215]
    (Intercept)     a     b     c     d     e
 1         TRUE FALSE  TRUE FALSE FALSE FALSE
 2         TRUE FALSE  TRUE FALSE FALSE  TRUE
 3         TRUE FALSE  TRUE  TRUE FALSE  TRUE
 4         TRUE FALSE  TRUE  TRUE  TRUE  TRUE
 5         TRUE  TRUE  TRUE  TRUE  TRUE  TRUE
 > tem2$cp
 [1] 14.389492  0.792566  2.189839  4.024738  6.000000
```

d) What is the value of $C_p(I_{min})$ and what is $\hat{\boldsymbol{\beta}}_{I_{min},0}$?

e) Which variables, if any, are needed in the MLR model given that the other variables are in the model?

f) List the variables, including a constant, that model 3 contains.

**7.4.** Suppose the full model has $p$ predictors including a constant. Let submodel $I$ have $k$ predictors. Then

$$C_p(I) = \frac{SSE(I)}{MSE} + 2k - n = (p - k)(F_I - 1) + k$$

where MSE is for the full model. Since $F_I \geq 0$, $C_p(I_{min}) \geq -p$ and $C_p(I) \geq -p$. Assume the full model is one of the submodels considered. Then $-p \leq C_p(I_{min}) \leq p$. Let $r$ be the residual vector for the full model and $r_I$ that for the submodel. Then the correlation

$$corr(r, r_I) = \sqrt{\frac{n - p}{C_p(I) + n - 2k}}.$$

a) Show $corr(r, r_{I_{min}}) \to 1$ as $n \to \infty$.

b) Suppose $S$ is not a subset of $I$. Under the model $x^T\beta = x_S^T\beta_S$, $corr(r, r_I)$ will not converge to 1 as $n \to \infty$. Suppose that for large enough $n$, $[corr(r, r_I)]^2 \leq \gamma < 1$. Show that $C_p(I) \to \infty$ as $n \to \infty$.

**7.5.** The table below shows simulation results for bootstrapping OLS (reg) and forward selection (vs) with $C_p$ when $\beta = (1, 1, 0, 0)^T$. The $\beta_i$ columns give coverage = the proportion of CIs that contained $\beta_i$ and the average length of the CI. The test is for $H_0 : (\beta_3, \beta_4)^T = 0$ and $H_0$ is true. The "coverage" is the proportion of times the prediction region method bootstrap test failed to reject $H_0$. Since 1000 runs were used, a cov in [0.93,0.97] is reasonable for a nominal value of 0.95. Output is given for three different error distributions. If the coverage for both methods $\geq 0.93$, the method with the shorter average CI length was more precise. (If one method had coverage $\geq 0.93$ and the other had coverage $< 0.93$, we will say the method with coverage $\geq 0.93$ was more precise.)

a) For $\beta_2$, $\beta_3$, and $\beta_4$, which method, forward selection or the OLS full model, was more precise?

**Table 7.3** Bootstrapping Forward Selection, $n = 100, p = 4, \psi = 0.9, B = 1000$

|         | $\beta_1$ | $\beta_2$ | $\beta_3$ | $\beta_4$ | test  |
|---------|-------|--------|--------|--------|-------|
| reg cov | 0.93  | 0.95   | 0.95   | 0.94   | 0.95  |
| len     | 1.266 | 10.703 | 10.666 | 10.650 | 2.547 |
| vs cov  | 0.95  | 0.93   | 0.997  | 0.995  | 0.989 |
| len     | 1.260 | 8.901  | 8.986  | 8.977  | 2.759 |
| reg cov | 0.94  | 0.93   | 0.95   | 0.94   | 0.95  |
| len     | 0.393 | 3.285  | 3.266  | 3.279  | 2.475 |
| vs cov  | 0.94  | 0.97   | 0.998  | 0.997  | 0.995 |
| len     | 0.394 | 2.773  | 2.721  | 2.733  | 2.703 |
| reg cov | 0.95  | 0.94   | 0.95   | 0.95   | 0.95  |
| len     | 0.656 | 5.493  | 5.465  | 5.427  | 2.493 |
| vs cov  | 0.93  | 0.95   | 0.998  | 0.998  | 0.977 |
| len     | 0.657 | 4.599  | 4.655  | 4.642  | 2.783 |

b) The test "length" is the average length of the interval $[0, D_{(U_B)}] = D_{(U_B)}$ where the test fails to reject $H_0$ if $D_0 \leq D_{(U_B)}$. The OLS full model is

asymptotically normal, and hence for large enough $n$ and $B$ the reg len row for the test column should be near $\sqrt{\chi^2_{2,0.95}} = 2.477$.

Were the three values in the test column for reg within 0.11 of 2.477?

**7.6.** The table below shows simulation results for bootstrapping OLS (reg), lasso, and ridge regression (RR) with 10-fold CV when $\boldsymbol{\beta} = (1,1,0,0)^T$. The $\beta_i$ columns give coverage = the proportion of CIs that contained $\beta_i$ and the average length of the CI. The test is for $H_0 : (\beta_3, \beta_4)^T = \mathbf{0}$ and $H_0$ is true. The "coverage" is the proportion of times the prediction region method bootstrap test failed to reject $H_0$. OLS used 1000 runs while 100 runs were used for lasso and ridge regression. Since 100 runs were used, a cov in [0.89, 1] is reasonable for a nominal value of 0.95. If the coverage for both methods $\geq 0.89$, the method with the shorter average CI length was more precise. (If one method had coverage $\geq 0.89$ and the other had coverage $< 0.89$, we will say the method with coverage $\geq 0.89$ was more precise.) (Lengths for the test column are not comparable unless the statistics have the same asymptotic distribution.)

**Table 7.4** Bootstrapping lasso and RR, $n = 100, \psi = 0, p = 4, B = 250$

|       |     | $\beta_1$ | $\beta_2$ | $\beta_3$ | $\beta_4$ | test  |
|-------|-----|-------|-------|-------|-------|-------|
| reg   | cov | 0.945 | 0.947 | 0.941 | 0.941 | 0.937 |
|       | len | 0.397 | 0.399 | 0.400 | 0.398 | 2.451 |
| RR    | cov | 0.95  | 0.89  | 0.95  | 0.95  | 0.94  |
|       | len | 0.401 | 0.366 | 0.377 | 0.382 | 2.451 |
| reg   | cov | 0.928 | 0.948 | 0.953 | 0.952 | 0.943 |
|       | len | 0.661 | 0.673 | 0.675 | 0.676 | 2.490 |
| lasso | cov | 0.97  | 0.90  | 0.99  | 0.98  | 0.97  |
|       | len | 0.684 | 0.741 | 0.612 | 0.610 | 2.650 |

a) For $\beta_3$ and $\beta_4$ which method, ridge regression or the OLS full model, was more precise?

b) For $\beta_3$ and $\beta_4$ which method, lasso or the OLS full model, was more precise?

**7.7.** For ridge regression, let $\boldsymbol{A}_n = (\boldsymbol{W}^T\boldsymbol{W} + \lambda_{1,n}\boldsymbol{I}_{p-1})^{-1}\boldsymbol{W}^T\boldsymbol{W}$ and $\boldsymbol{B}_n = [\boldsymbol{I}_{p-1} - \lambda_{1,n}(\boldsymbol{W}^T\boldsymbol{W} + \lambda_{1,n}\boldsymbol{I}_{p-1})^{-1}]$. Show $\boldsymbol{A}_n - \boldsymbol{B}_n = \mathbf{0}$.

**7.8.** Table 7.5 below shows simulation results for bootstrapping OLS (reg) and forward selection (vs) with $C_p$ when $\boldsymbol{\beta} = (1,1,0,0,0)^T$. The $\beta_i$ columns give coverage = the proportion of CIs that contained $\beta_i$ and the average length of the CI. The test is for $H_0 : (\beta_3, \beta_4, \beta_5)^T = \mathbf{0}$ and $H_0$ is true. The "coverage" is the proportion of times the prediction region method bootstrap test failed to reject $H_0$. Since 1000 runs were used, a cov in [0.93,0.97] is reasonable for a nominal value of 0.95. Output is given for three different error distributions. If the coverage for both methods $\geq 0.93$, the method

with the shorter average CI length was more precise. (If one method had coverage $\geq 0.93$ and the other had coverage $< 0.93$, we will say the method with coverage $\geq 0.93$ was more precise.)

a) For $\beta_3$, $\beta_4$, and $\beta_5$, which method, forward selection or the OLS full model, was more precise?

**Table 7.5** Bootstrapping Forward Selection, $n = 100, p = 5, \psi = 0, B = 1000$

|        |     | $\beta_1$ | $\beta_2$ | $\beta_3$ | $\beta_4$ | $\beta_5$ | test |
|--------|-----|-----------|-----------|-----------|-----------|-----------|-------|
| reg cov |    | 0.95 | 0.93 | 0.93 | 0.93 | 0.94 | 0.93 |
|        | len | 0.658 | 0.672 | 0.673 | 0.674 | 0.674 | 2.861 |
| vs cov |    | 0.95 | 0.94 | 0.998 | 0.998 | 0.999 | 0.993 |
|        | len | 0.661 | 0.679 | 0.546 | 0.548 | 0.544 | 3.11 |
| reg cov |    | 0.96 | 0.93 | 0.94 | 0.96 | 0.93 | 0.94 |
|        | len | 0.229 | 0.230 | 0.229 | 0.231 | 0.230 | 2.787 |
| vs cov |    | 0.95 | 0.94 | 0.999 | 0.997 | 0.999 | 0.995 |
|        | len | 0.228 | 0.229 | 0.185 | 0.187 | 0.186 | 3.056 |
| reg cov |    | 0.94 | 0.94 | 0.95 | 0.94 | 0.94 | 0.93 |
|        | len | 0.393 | 0.398 | 0.399 | 0.399 | 0.398 | 2.839 |
| vs cov |    | 0.94 | 0.95 | 0.997 | 0.997 | 0.996 | 0.990 |
|        | len | 0.392 | 0.400 | 0.320 | 0.322 | 0.321 | 3.077 |

b) The test "length" is the average length of the interval $[0, D_{(U_B)}] = D_{(U_B)}$ where the test fails to reject $H_0$ if $D_{\mathbf{0}} \leq D_{(U_B)}$. The OLS full model is asymptotically normal, and hence for large enough $n$ and $B$ the reg len row for the test column should be near $\sqrt{\chi^2_{3,0.95}} = 2.795$.

Were the three values in the test column for reg within 0.1 of 2.795?

**7.9.** Suppose the MLR model $\boldsymbol{Y} = \boldsymbol{X\beta} + \boldsymbol{e}$, and the regression method fits $\boldsymbol{Z} = \boldsymbol{W\eta} + \boldsymbol{e}$. Suppose $\hat{Z} = 245.63$ and $\overline{Y} = 105.37$. What is $\hat{Y}$?

**7.10.** To get a large sample 90% PI for a future value $Y_f$ of the response variable, find a large sample 90% PI for a future residual and add $\hat{Y}_f$ to the endpoints of the of that PI. Suppose forward selection is used and the large sample 90% PI for a future residual is $[-778.28, 1336.44]$. What is the large sample 90% PI for $Y_f$ if $\hat{\boldsymbol{\beta}}_{I_{min}} = (241.545, 1.001)^T$ used a constant and the predictor *mmen* with corresponding $\boldsymbol{x}_{I_{min},f} = (1, 75000)^T$?

**7.11.** For ridge regression, let $\boldsymbol{A}_n = (\boldsymbol{W}^T\boldsymbol{W} + \lambda_{1,n}\boldsymbol{I}_{p-1})^{-1}\boldsymbol{W}^T\boldsymbol{W}$ and $\boldsymbol{B}_n = [\boldsymbol{I}_{p-1} - \lambda_{1,n}(\boldsymbol{W}^T\boldsymbol{W} + \lambda_{1,n}\boldsymbol{I}_{p-1})^{-1}]$. Show $\boldsymbol{A}_n - \boldsymbol{B}_n = \boldsymbol{0}$.

**7.12.** Consider choosing $\hat{\boldsymbol{\eta}}$ to minimize the elastic net criterion

$$Q(\boldsymbol{\eta}) = RSS(\boldsymbol{\eta}) + \lambda_1\|\boldsymbol{\eta}\|_2^2 + \lambda_2\|\boldsymbol{\eta}\|_1$$

where $\lambda_i \geq 0$ for $i = 1, 2$.

a) Which values of $\lambda_1$ and $\lambda_2$ correspond to ridge regression? (For example, both are zero, $\lambda_1$ is zero, or $\lambda_2$ is zero.)

b) Which values of $\lambda_1$ and $\lambda_2$ correspond to the OLS full model?

**7.13.** Consider choosing $\hat{\boldsymbol{\eta}}$ to minimize the criterion

$$Q(\boldsymbol{\eta}) = \frac{1}{a}(\boldsymbol{Z} - \boldsymbol{W}\boldsymbol{\eta})^T(\boldsymbol{Z} - \boldsymbol{W}\boldsymbol{\eta}) + \frac{\lambda_{1,n}}{a}\sum_{i=1}^{p-1}|\eta_i|^j$$

where $\lambda_{1,n} \geq 0$, $a > 0$, and $j > 0$ are known constants. Consider the regression methods OLS, forward selection, lasso, ridge regression, and lasso variable selection.
a) Which method corresponds to $j = 1$?
b) Which method corresponds to $j = 2$?
c) Which method corresponds to $\lambda_{1,n} = 0$?

**7.14.**
**R Problems** Some $R$ code for homework problems is at (http://parker.ad.siu.edu/Olive/robRhw.txt).
**Warning: Use a command like** *source("G:/rpack.txt")* **to download the programs. See Preface or Section 14.2.** Typing the name of the rpack function, e.g. *regbootsim3*, will display the code for the function. Use the args command, e.g. *args(regbootsim3)*, to display the needed arguments for the function.

```
regbootsim3(nruns=500)
#output similar to that for Problem 7.15
$cicov
0.942 0.954 0.950 0.948 0.944 0.946 0.946 0.940 0.938 0.940
$avelen
0.398 0.399 0.397 0.399 2.448 2.448 2.448 2.448 2.448 2.450
$beta
[1] 1 1 0 0
$k
[1] 1
```

**7.15.** Use the $R$ command for this problem, and put the output in *Word*. The output should be similar to that shown above. Consider the multiple linear regression model $Y_i = \beta_1 + \beta_2 x_{i,2} + \beta_3 x_{i,3} + \beta_4 x_{i,4} + e_i$ where $\boldsymbol{\beta} = (1,1,0,0)^T$. The function regbootsim3 bootstraps the regression model with the residual bootstrap. Note that $S = \{1,2\}$ and $E = \{3,4\}$. The first 4 numbers are the bootstrap shorth confidence intervals for $\beta_i$. The lengths of the CIs along with the proportion of times (coverage) the CI for $\beta_i$ contained $\beta_i$ are given. The CI lengths for the first 4 intervals should be near 0.392. With 500 runs, coverage in [0.92,0.98] suggests that the actual coverage is near the nominal coverage of 0.95. The next three numbers test $H_0 : \boldsymbol{\beta}_E = \boldsymbol{0}$ where $E$ corresponds to the last $p - k + 1$ $\beta_i$. The prediction region method, hybrid method, and Bickel and Ren methods are used. Hence the fifth interval

gives the length of the interval $[0, D_{(c)}]$ where $H_0$ is rejected if $D_0 > D_{(c)}$ and the fifth "coverage" is the proportion of times the prediction region method test fails to reject $H_0$. The last three numbers are similar but test $H_0$ : $\boldsymbol{\beta}_S = \mathbf{1}$ where $S$ corresponds to the first $k+1$ $\beta_i$. Hence the last length 2.450 corresponds to the Bickel and Ren method with coverage 0.940. Want lengths near 2.45 which correspond to $\sqrt{\chi_2^2(0.95)}$ where $P(X \le \chi_2^2(0.95)) = 0.95$ if $X \sim \chi_2^2$.

**7.16.** The $R$ program generates data satisfying the MLR model

$$Y = \beta_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 + e$$

where $\boldsymbol{\beta} = (\beta_1, \beta_2, \beta_3, \beta_4)^T = (1, 1, 0, 0)$.

a) Copy and paste the commands for this part into $R$. The output gives $\hat{\boldsymbol{\beta}}_{OLS}$ for the OLS full model. Give $\hat{\boldsymbol{\beta}}_{OLS}$. Is $\hat{\boldsymbol{\beta}}_{OLS}$ close to $\boldsymbol{\beta} = 1, 1, 0, 0)^T$?

b) The commands for this part bootstrap the OLS full model using the residual bootstrap. Copy and paste the output into *Word*. The output shows $T_j^* = \hat{\boldsymbol{\beta}}_j^*$ for $j = 1, ..., 5$.

c) $B = 1000$ $T_j^*$ were generated. The commands for this part compute the sample mean $\overline{T}^*$ of the $T_j^*$. Copy and paste the output into *Word*. Is $\overline{T}^*$ close to $\hat{\boldsymbol{\beta}}_{OLS}$ found in a)?

d) The commands for this part bootstrap the forward selection using the residual bootstrap. Copy and paste the output into *Word*. The output shows $T_j^* = \hat{\boldsymbol{\beta}}_{VS,j} = \hat{\boldsymbol{\beta}}_{I_{min},0,j}^*$ for $j = 1, ..., 5$. The last two variables may have a few 0s.

e) $B = 1000$ $T_j^*$ were generated. The commands for this part compute the sample mean $\overline{T}^*$ of the $T_j^*$ where $T_j^*$ is as in d). Copy and paste the output into *Word*. Is $\overline{T}^*$ close to $\boldsymbol{\beta} = (1, 1, 0, 0)$?

**7.17.**

**7.18.**

**7.19.** For the Buxton (1920) data with multiple linear regression, *height* was the response variable while an intercept, *head length, nasal height, bigonal breadth,* and *cephalic index* were used as predictors in the multiple linear regression model. Observation 9 was deleted since it had missing values. Five individuals, cases 61–65, were reported to be about 0.75 inches tall with head lengths well over five feet!

a) Copy and paste the commands for this problem into $R$. Include the lasso response plot in *Word*. The identity line passes right through the outliers which are obvious because of the large gap. Prediction interval (PI) bands are also included in the plot.

b) Copy and paste the commands for this problem into $R$. Include the lasso response plot in *Word*. This did lasso for the cases in the `covmb2` set $B$ applied to the predictors which included all of the clean cases and omitted the 5 outliers. The response plot was made for all of the data, including the outliers.

c) Copy and paste the commands for this problem into $R$. Include the DD plot in *Word*. The outliers are in the upper right corner of the plot.

**7.20.** This problem is like Problem 7.19, except elastic net is used instead of lasso.

a) Copy and paste the commands for this problem into $R$. Include the elastic net response plot in *Word*. The identity line passes right through the outliers which are obvious because of the large gap. Prediction interval (PI) bands are also included in the plot.

b) Copy and paste the commands for this problem into $R$. Include the elastic net response plot in *Word*. This did elastic net for the cases in the `covmb2` set $B$ applied to the predictors which included all of the clean cases and omitted the 5 outliers. The response plot was made for all of the data, including the outliers. (Problem 7.19 c) shows the DD plot for the data.)

**7.21.** Consider the Gladstone (1905) data set that has 12 variables on 267 persons after death. There are 5 infants in the data set. The response variable was *brain weight*. Head measurements were *breadth, circumference, head height, length,* and *size* as well as *cephalic index* and *brain weight*. *Age*, *height*, and three categorical variables *cause, ageclass* (0: under 20, 1: 20-45, 2: over 45) and *sex* were also given. The constant $x_1$ was the first variable. The variables *cause* and *ageclass* were not coded as factors. Coding as factors might improve the fit.

a) Copy and paste the commands for this problem into $R$. Include the lasso response plot in *Word*. The identity line passes right through the infants which are obvious because of the large gap. Prediction interval (PI) bands are also included in the plot.

b) Copy and paste the commands for this problem into $R$. Include the lasso response plot in *Word*. This did lasso for the cases in the `covmb2` set $B$ applied to the nontrivial predictors which are not categorical (omit the *constant, cause, ageclass* and *sex*) which omitted 8 cases, including the 5 infants. The response plot was made for all of the data.

c) Copy and paste the commands for this problem into $R$. Include the DD plot in *Word*. The infants are in the upper right corner of the plot.

**7.22.** This simulation is similar to that used to form Table 7.5. Since 1000 runs are used, coverage in [0.93,0.97] suggests that the actual coverage is close to the nominal coverage of 0.95.

The model is $Y = \boldsymbol{x}^T\boldsymbol{\beta} + e = \boldsymbol{x}_S^T\boldsymbol{\beta}_S + e$ where $\boldsymbol{\beta}_S = (\beta_1, \beta_2, ..., \beta_{k+1})^T = (\beta_1, \beta_2)^T$ and $k = 1$ is the number of active nontrivial predictors in the population model. The output for *test* tests $H_0 : (\beta_{k+2}, ..., \beta_p)^T = (\beta_3, ..., \beta_p)^T = \boldsymbol{0}$ and $H_0$ is true. The output gives the proportion of times the prediction region method bootstrap test fails to reject $H_0$. The nominal proportion is 0.95.

After getting your output, make a table similar to Table 7.5 with 4 lines. Two lines are for reg (the OLS full model) and two lines are for vs (forward selection with $I_{min}$). The $\beta_i$ columns give the coverage and lengths of the 95% CIs for $\beta_i$. If the coverage $\geq 0.93$, then the shorter CI length is more

precise. Were the CIs for forward selection more precise than the CIs for the OLS full model for $\beta_3$ and $\beta_4$?