

## Chapter 4

# Prediction and Variable Selection When $n \gg p$

This chapter considers variable selection when  $n \gg p$  and prediction intervals that can work if  $n > p$  or  $p > n$ . 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.

### 4.1 Variable Selection

*Variable selection*, also called subset or model selection, is the search for a subset of predictor variables that can be deleted with little loss of information if  $n/p$  is large. Consider the 1D regression model where  $Y \perp\!\!\!\perp \mathbf{x} | SP$  where  $SP = \mathbf{x}^T \boldsymbol{\beta}$ . See Chapters 1 and 10. A *model for variable selection* can be described by

$$\mathbf{x}^T \boldsymbol{\beta} = \mathbf{x}_S^T \boldsymbol{\beta}_S + \mathbf{x}_E^T \boldsymbol{\beta}_E = \mathbf{x}_S^T \boldsymbol{\beta}_S \quad (4.1)$$

where  $\mathbf{x} = (\mathbf{x}_S^T, \mathbf{x}_E^T)^T$  is a  $p \times 1$  vector of predictors,  $\mathbf{x}_S$  is an  $a_S \times 1$  vector, and  $\mathbf{x}_E$  is a  $(p - a_S) \times 1$  vector. Given that  $\mathbf{x}_S$  is in the model,  $\boldsymbol{\beta}_E = \mathbf{0}$  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  $\mathbf{x}_I$  be the vector of  $a$  terms from a candidate subset indexed by  $I$ , and let  $\mathbf{x}_O$  be the vector of the remaining predictors (out of the candidate submodel). Then

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

Suppose that  $S$  is a subset of  $I$  and that model (4.1) holds. Then

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

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

**Definition 4.1.** The model  $Y \perp\!\!\!\perp \mathbf{x} | \mathbf{x}^T \beta$  that uses all of the predictors is called the *full model*. A model  $Y \perp\!\!\!\perp \mathbf{x}_I | \mathbf{x}_I^T \beta_I$  that uses a subset  $\mathbf{x}_I$  of the predictors is called a *submodel*. The **full model is always a submodel**. The full model has *sufficient predictor*  $SP = \mathbf{x}^T \beta$  and the submodel has  $SP = \mathbf{x}_I^T \beta_I$ .

Forward selection or backward elimination with the Akaike (1973) AIC criterion or Schwarz (1978) BIC criterion are often used for variable selection. The relaxed lasso or relaxed elastic net estimator fits the regression method, such as a GLM or Cox (1972) proportional hazards regression, to the predictors that had nonzero lasso or elastic net coefficients. See Chapters 5 and 10.

To clarify notation, suppose  $p = 4$ , a constant  $x_1 = 1$  corresponding to  $\beta_1$  is always in the model, and  $\beta = (\beta_1, \beta_2, 0, 0)^T$ . Then the  $J = 2^{p-1} = 8$  possible subsets of  $\{1, 2, \dots, p\}$  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  $\hat{\beta}_{I_7} = (\hat{\beta}_1, \hat{\beta}_3, \hat{\beta}_4)^T$  and  $\mathbf{x}_{I_7} = (x_1, x_3, x_4)^T$ .

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

$$\beta_I + (\mathbf{X}_I\mathbf{X}_I)^{-1}\mathbf{X}_I^T\mathbf{X}_O\beta_O = \beta_I + \mathbf{A}\mathbf{X}_O\beta_O.$$

If  $S \subseteq I$ , then  $\beta_O = \mathbf{0}$ , but if underfitting occurs then the bias vector  $\mathbf{A}\mathbf{X}_O\beta_O$  can be large.

#### 4.1.1 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 (4.1) holds,  $S \subseteq I$ ,  $\beta_S$  is an  $a_S \times 1$  vector, and  $\beta_I$  is a  $j \times 1$  vector with  $j > a_S$ , then

$$\frac{1}{n} \sum_{i=1}^n V(\hat{Y}_{Ii}) = \frac{\sigma^2 j}{n} > \frac{\sigma^2 a_S}{n} = \frac{1}{n} \sum_{i=1}^n V(\hat{Y}_{Si}). \quad (4.2)$$

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  $\mathbf{H} = \mathbf{I}_n$ , the  $n \times n$  identity matrix. Then  $\hat{Y} = Y$  so that  $\text{VAR}(\hat{Y}|\mathbf{x}) = \text{VAR}(Y)$ . 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 (4.2) holds, assume that the full model includes all  $p$  possible terms so the full model may overfit but does not underfit. Then  $\hat{Y} = \mathbf{H}\mathbf{Y}$  and  $\text{Cov}(\hat{Y}) = \sigma^2 \mathbf{H}\mathbf{H}^T = \sigma^2 \mathbf{H}$ . Thus

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

where  $\text{tr}(\mathbf{A})$  is the trace operation. Replacing  $p$  by  $j$  and  $a_S$  and replacing  $\mathbf{H}$  by  $\mathbf{H}_I$  and  $\mathbf{H}_S$  implies Equation (4.2). Hence if only  $a_S$  parameters are needed and  $p \gg a_S$ , 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 1.17 and 1.18. 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  $\mathbf{x}_I$  has  $k$  terms (including the constant), then the partial  $F$  statistic for testing whether the  $p - k$  predictor variables in  $\mathbf{x}_O$  can be deleted is

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

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 4.2.**

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

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

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

**Theorem 4.1.** Suppose that a numerical variable selection method suggests several submodels with  $k$  predictors, including a constant, where  $2 \leq k \leq p$ .

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

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

c) As  $\text{corr}(r, r_I) \rightarrow 1$ ,

$$\text{corr}(\mathbf{x}^T \hat{\beta}, \mathbf{x}_I^T \hat{\beta}_I) = \text{corr}(\text{ESP}, \text{ESP}(I)) = \text{corr}(\hat{Y}, \hat{Y}_I) \rightarrow 1.$$

**Proof.** These results are a corollary of Theorem 4.2 below.  $\square$

**Remark 4.1.** 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 4.2 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  $k \rightarrow 0$ , the partial  $F$  test will reject  $H_0: \beta_O = \mathbf{0}$  (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 4.3.** 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 (4.1) holds and that a good estimator (such as OLS) for  $\hat{\beta}$  and  $\hat{\beta}_I$  is used.

**Application 4.1.** To visualize whether a candidate submodel using predictors  $\mathbf{x}_I$  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  $\mathbf{X}$  is the full rank  $n \times p$  design matrix for the full model. Let the corresponding vectors of OLS fitted values and residuals be  $\hat{\mathbf{Y}} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} = \mathbf{H} \mathbf{Y}$  and  $\mathbf{r} = (\mathbf{I} - \mathbf{H}) \mathbf{Y}$ , respectively. Suppose that  $\mathbf{X}_I$  is the  $n \times k$  design matrix for the candidate submodel and that the corresponding vectors of OLS fitted values and residuals are  $\hat{\mathbf{Y}}_I = \mathbf{X}_I(\mathbf{X}_I^T \mathbf{X}_I)^{-1} \mathbf{X}_I^T \mathbf{Y} = \mathbf{H}_I \mathbf{Y}$  and  $\mathbf{r}_I = (\mathbf{I} - \mathbf{H}_I) \mathbf{Y}$ , 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  $\text{corr}(x, y)$  denote the correlation between  $x$  and  $y$ .

**Theorem 4.2.** Suppose that every submodel contains a constant and that  $\mathbf{X}$  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  $b = [\text{corr}(Y, \hat{Y}_I)]^2 = R^2(I)$  and intercept  $a = \bar{Y}(1 - R^2(I))$  where  $\bar{Y} = \sum_{i=1}^n Y_i/n$  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  $b = [\text{corr}(\hat{Y}, \hat{Y}_I)]^2 = SSR(I)/SSR$  and intercept  $a = \bar{Y}[1 - (SSR(I)/SSR)]$  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  $b = [\text{corr}(r, r_I)]^2$  and

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

**Proof:** Recall that  $\mathbf{H}$  and  $\mathbf{H}_I$  are symmetric idempotent matrices and that  $\mathbf{H}\mathbf{H}_I = \mathbf{H}_I$ . The mean of OLS fitted values is equal to  $\bar{Y}$  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  $a = \bar{z} - b\bar{w}$  and

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

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

(\*) 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  $[\text{corr}(z, w)]^2$ .

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

ii) By (\*), the slope

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

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

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

iv) From iii),

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

Hence

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

and the slope

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

Also the slope

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

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

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

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

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

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 4.2.** Let  $I_{min}$  be the model than minimizes  $C_p(I)$  among the models  $I$  generated from the variable selection method such as forward se-

lection. 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  $\text{corr}(r, r_{I_{min}}) \rightarrow 1$  as  $n \rightarrow \infty$  by Theorem 4.2 vi). Referring to Equation (4.1), if  $P(S \subseteq I_{min})$  does not go to 1 as  $n \rightarrow \infty$ , then the above correlation would not go to one. Hence  $P(S \subseteq I_{min}) \rightarrow 1$  as  $n \rightarrow \infty$ .

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 4.3.** 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 4.2 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  $n \rightarrow \infty$ . 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 4.2 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 4.1.** 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  $I_1, \dots, I_p$  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  $SSE(I) = RSS(I) = \mathbf{r}^T(I)\mathbf{r}(I) = \sum_{i=1}^n r_i^2(I) = \sum_{i=1}^n (Y_i - \hat{Y}_i(I))^2$ . 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  $x_1^*, \dots, x_{j-1}^*$ . Compute  $SSE(I)$  and let  $I_j$  minimize  $SSE(I)$  for the  $p - j + 1$  models  $I$  that contain  $x_1^*, \dots, x_{j-1}^*$  and one other predictor not already selected. Denote the predictors in  $I_j$  by  $x_1^*, \dots, x_j^*$ . Continue in this manner for  $j = 2, \dots, M = p$ .

Backward elimination also forms a sequence of submodels  $I_1, \dots, I_p$  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  $Q_{p-1}(I)$  for the  $p - 1$  models  $I$  that exclude one of the predictors  $x_2, \dots, x_p$ . Denote the predictors in  $I_{p-1}$  by  $x_1^*, x_2^*, \dots, x_{p-1}^*$ . In general, to form  $I_j$  consider all models  $I$  with  $j$  predictors including variables  $x_1^*, \dots, x_{j+1}^*$ . Compute  $SSE(I)$ , and let  $I_j$  minimize  $SSE(I)$  for the  $p - j + 1$  models  $I$  that exclude one of the predictors  $x_2^*, \dots, x_{j+1}^*$ . Denote the predictors in  $I_j$  by  $x_1^*, \dots, x_j^*$ . Continue in this manner for  $j = p = M, p - 1, \dots, 2, 1$  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  $\mathbf{x}_I$  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 4.2. 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  $\hat{\beta}_I$  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 (4.1) 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  $x_1^*, \dots, x_{k-1}^*$  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 2.27, 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  $\hat{\beta}_{I_{min},0}$  (defined two paragraphs below) and the bootstrap for inference after variable selection, and ii) use data splitting for inference after variable selection.

## 4.2 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  $\mathbf{Z}_n = \sqrt{n}\mathbf{A}(\hat{\beta}_{I_{min}} - \beta_I)$  after variable selection. One reason is that with positive probability,  $\hat{\beta}_{I_{min}}$  does not have the same dimension as  $\beta_I$  if AIC or  $C_p$  is used. Hence  $\mathbf{Z}_n$  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  $n \rightarrow \infty$ . Suppose that model (4.1) holds. Assume the maximum leverage

$$\max_{i=1, \dots, n} \mathbf{x}_{iI_j}^T (\mathbf{X}_{I_j}^T \mathbf{X}_{I_j})^{-1} \mathbf{x}_{iI_j} \rightarrow 0$$

in probability as  $n \rightarrow \infty$  for each  $I_j$  with  $S \subseteq I_j$  where the dimension of  $I_j$  is  $a_j$ . For the OLS model with  $S \subseteq I_j$ ,  $\sqrt{n}(\hat{\beta}_{I_j} - \beta_{I_j}) \xrightarrow{D} N_{a_j}(\mathbf{0}, \mathbf{V}_j)$  where  $\mathbf{V}_j = \sigma^2 \mathbf{W}_j$  and  $(\mathbf{X}_{I_j}^T \mathbf{X}_{I_j})/n \xrightarrow{P} \mathbf{W}_j^{-1}$  by the LS CLT Theorem 2.26. Then

$$\mathbf{u}_{jn} = \sqrt{n}(\hat{\beta}_{I_j,0} - \beta) \xrightarrow{D} \mathbf{u}_j \sim N_p(\mathbf{0}, \mathbf{V}_{j,0}) \tag{4.3}$$

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

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

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

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 4.4.** 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  $k = 1, \dots, J$  where there are  $J$  subsets.

**Definition 4.5.** 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 1.6.

The first assumption in Theorem 4.3 is  $P(S \subseteq I_{min}) \rightarrow 1$  as  $n \rightarrow \infty$ . 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  $P(S \subseteq I_{min}) \rightarrow 1$  as  $n \rightarrow \infty$ . See Zhao and Yu (2006, p. 2554). Hence Theorem 4.3c) 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 4.4 and Remark 4.5. The assumption on  $\mathbf{u}_{jn}$  in Theorem 4.3 is reasonable by (4.3) since  $S \subseteq I_j$  for each  $\pi_j$ , and since  $\hat{\boldsymbol{\beta}}_{MIX}$  uses random selection.

**Theorem 4.3.** Assume  $P(S \subseteq I_{min}) \rightarrow 1$  as  $n \rightarrow \infty$ , and let  $\hat{\boldsymbol{\beta}}_{MIX} = \hat{\boldsymbol{\beta}}_{I_k,0}$  with probabilities  $\pi_{kn}$  where  $\pi_{kn} \rightarrow \pi_k$  as  $n \rightarrow \infty$ . Denote the positive

$\pi_k$  by  $\pi_j$ . Assume  $\mathbf{u}_{jn} = \sqrt{n}(\hat{\boldsymbol{\beta}}_{I_j,0} - \boldsymbol{\beta}) \xrightarrow{D} \mathbf{u}_j \sim N_p(\mathbf{0}, \mathbf{V}_{j,0})$ . a) Then

$$\mathbf{u}_n = \sqrt{n}(\hat{\boldsymbol{\beta}}_{MIX} - \boldsymbol{\beta}) \xrightarrow{D} \mathbf{u} \tag{4.4}$$

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

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

$$\mathbf{v}_n = \mathbf{A}\mathbf{u}_n = \sqrt{n}(\mathbf{A}\hat{\boldsymbol{\beta}}_{MIX} - \mathbf{A}\boldsymbol{\beta}) \xrightarrow{D} \mathbf{A}\mathbf{u} = \mathbf{v} \tag{4.5}$$

where  $\mathbf{v}$  has a mixture distribution of the  $\mathbf{v}_j = \mathbf{A}\mathbf{u}_j \sim N_g(\mathbf{0}, \mathbf{A}\mathbf{V}_{j,0}\mathbf{A}^T)$  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  $\sqrt{n}(\hat{\boldsymbol{\beta}}_{SEL} - \boldsymbol{\beta}) \xrightarrow{D} \mathbf{u} \sim N_p(\mathbf{0}, \mathbf{V}_{d,0})$  where  $SEL$  is  $VS$  or  $MIX$ .

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

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

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  $\hat{\boldsymbol{\beta}}_{MIX} = (\hat{\beta}_1, \dots, \hat{\beta}_p)^T$ . Let  $\hat{\boldsymbol{\beta}}_{i,MIX} = \hat{\beta}_i$ . Similarly, if  $I = \{i_1, \dots, i_a\}$ , then  $\hat{\boldsymbol{\beta}}_{I,MIX} = (\hat{\beta}_{i_1}, \dots, \hat{\beta}_{i_a})^T$ . Subscripts after  $MIX$  denote the  $i$ th vector from a sample  $\hat{\boldsymbol{\beta}}_{MIX,1}, \dots, \hat{\boldsymbol{\beta}}_{MIX,B}$ . 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  $\pi_j$ ,  $\mathbf{A}\mathbf{u}_j \sim N_g(\mathbf{0}, \mathbf{A}\mathbf{V}_{j,0}\mathbf{A}^T) = N_g(\mathbf{0}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)$ . Then  $\sqrt{n}(\mathbf{A}\hat{\boldsymbol{\beta}}_{MIX} - \mathbf{A}\boldsymbol{\beta}) \xrightarrow{D} \mathbf{A}\mathbf{u} \sim N_g(\mathbf{0}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)$ . This special case occurs for  $\hat{\boldsymbol{\beta}}_{S,MIX}$  if  $\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \xrightarrow{D} N_p(\mathbf{0}, \mathbf{V})$  where the asymptotic covariance matrix  $\mathbf{V}$  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,  $\sqrt{n}(\mathbf{A}\hat{\boldsymbol{\beta}}_{VS} - \mathbf{A}\boldsymbol{\beta})$  is selecting from the  $\mathbf{A}\mathbf{u}_j$  which have the same distribution. Then the confidence regions applied to  $\mathbf{A}\hat{\boldsymbol{\beta}}_{SEL}^* = \mathbf{B}\hat{\boldsymbol{\beta}}_{S,SEL}^*$  should have similar volume and cutoffs where  $SEL$  is  $MIX$ ,  $VS$ , or  $FULL$ .

Theorem 4.3 can be used to justify prediction intervals after variable selection. See Pelawa Watagoda and Olive (2020). Theorem 4.3d) is useful for *variable selection consistency* and the *oracle property* where  $\pi_d = \pi_S = 1$  if  $P(I_{min} = S) \rightarrow 1$  as  $n \rightarrow \infty$ . See Claeskens and Hjort (2008, pp. 101-114) and Fan and Li (2001) for references. A necessary condition for  $P(I_{min} = S) \rightarrow 1$  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 4.4.** If  $A_1, A_2, \dots, A_k$  are pairwise disjoint and if  $\cup_{i=1}^k A_i = S$ , then the collection of sets  $A_1, A_2, \dots, A_k$  is a *partition* of  $S$ . Then the *Law of Total Probability* states that if  $A_1, A_2, \dots, A_k$  form a partition of  $S$  such that  $P(A_i) > 0$  for  $i = 1, \dots, k$ , 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  $A_{k+1}, \dots, A_m$  satisfy  $P(A_i) = 0$  for  $i = k+1, \dots, m$ . 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  $\hat{\beta}_{VS} | (\hat{\beta}_{VS} = \hat{\beta}_{I_k,0})$  to find the distribution of  $\mathbf{w}_n = \sqrt{n}(\hat{\beta}_{VS} - \beta)$ . Let  $W = W_{VS} = k$  if  $\hat{\beta}_{VS} = \hat{\beta}_{I_k,0}$  where  $P(W_{VS} = k) = \pi_{kn}$  for  $k = 1, \dots, J$ . Then  $(\hat{\beta}_{VS:n}, W_{VS:n}) = (\hat{\beta}_{VS}, W_{VS})$  has a joint distribution where the sample size  $n$  is usually suppressed. Note that  $\hat{\beta}_{VS} = \hat{\beta}_{I_w,0}$ . Define  $P(B|A_k)P(A_k) = 0$  if  $P(A_k) = 0$ . Let  $\hat{\beta}_{I_k,0}^C$  be a random vector from the conditional distribution  $\hat{\beta}_{I_k,0} | (W_{VS} = k)$ . Let  $\mathbf{w}_{kn} = \sqrt{n}(\hat{\beta}_{I_k,0} - \beta) | (W_{VS} = k) \sim \sqrt{n}(\hat{\beta}_{I_k,0}^C - \beta)$ . Denote  $F_{\mathbf{z}}(\mathbf{t}) = P(z_1 \leq t_1, \dots, z_p \leq t_p)$  by  $P(\mathbf{z} \leq \mathbf{t})$ . Then

$$\begin{aligned} F_{\mathbf{w}_n}(\mathbf{t}) &= P[n^{1/2}(\hat{\beta}_{VS} - \beta) \leq \mathbf{t}] = \\ &= \sum_{k=1}^J P[n^{1/2}(\hat{\beta}_{VS} - \beta) \leq \mathbf{t} | (\hat{\beta}_{VS} = \hat{\beta}_{I_k,0})] P(\hat{\beta}_{VS} = \hat{\beta}_{I_k,0}) = \\ &= \sum_{k=1}^J P[n^{1/2}(\hat{\beta}_{I_k,0} - \beta) \leq \mathbf{t} | (\hat{\beta}_{VS} = \hat{\beta}_{I_k,0})] \pi_{kn} \end{aligned}$$

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

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

Charkhi and Claeskens (2018) showed that  $\mathbf{w}_{jn} = \sqrt{n}(\hat{\beta}_{I_j,0}^C - \beta) \xrightarrow{D} \mathbf{w}_j$  if  $S \subseteq I_j$  for the MLE with AIC. Here  $\mathbf{w}_j$  is a multivariate truncated normal distribution (where no truncation is possible) that is symmetric about  $\mathbf{0}$ . Hence  $E(\mathbf{w}_j) = \mathbf{0}$ , and  $\text{Cov}(\mathbf{w}_j) = \Sigma_j$  exists. Referring to Definitions 4.4 and 4.5, note that both  $\sqrt{n}(\hat{\beta}_{MIX} - \beta)$  and  $\sqrt{n}(\hat{\beta}_{VS} - \beta)$  are selecting from the  $\mathbf{u}_{kn} = \sqrt{n}(\hat{\beta}_{I_k,0} - \beta)$  and asymptotically from the  $\mathbf{u}_j$  of Equation (4.3). The random selection for  $\hat{\beta}_{MIX}$  does not change the distribution of  $\mathbf{u}_{jn}$ , but selection bias does change the distribution of the selected  $\mathbf{u}_{jn}$  to that of  $\mathbf{w}_{jn}$ . Similarly, selection bias does change the distribution of the selected  $\mathbf{u}_j$  to that of  $\mathbf{w}_j$ . The reasonable Theorem 4.4 assumption that  $\mathbf{w}_{jn} \xrightarrow{D} \mathbf{w}_j$  may not be mild.

**Theorem 4.4, Variable Selection CLT.** Assume  $P(S \subseteq I_{min}) \rightarrow 1$  as  $n \rightarrow \infty$ , and let  $\hat{\beta}_{VS} = \hat{\beta}_{I_k,0}$  with probabilities  $\pi_{kn}$  where  $\pi_{kn} \rightarrow \pi_k$  as  $n \rightarrow \infty$ . Denote the positive  $\pi_k$  by  $\pi_j$ . Assume  $\mathbf{w}_{jn} = \sqrt{n}(\hat{\beta}_{I_j,0}^C - \beta) \xrightarrow{D} \mathbf{w}_j$ . Then

$$\mathbf{w}_n = \sqrt{n}(\hat{\beta}_{VS} - \beta) \xrightarrow{D} \mathbf{w} \quad (4.6)$$

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

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

**Remark 4.5.** If  $P(S \subseteq I_{min}) \rightarrow 1$  as  $n \rightarrow \infty$ , then  $\hat{\beta}_{VS}$  is a  $\sqrt{n}$  consistent estimator of  $\beta$  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). By both this result and Theorems 4.3 and 4.4, the lasso variable selection and elastic net variable selection estimators are  $\sqrt{n}$  consistent if lasso and elastic net are consistent.

Mixture distributions are useful for variable selection since  $\hat{\beta}_{I_{min},0}$  has a mixture distribution of the  $\hat{\beta}_{I_j,0}$ . Review mixture distributions from Section 1.6. The following theorem is due to Pelawa Watagoda and Olive (2019a). Note that the cdf of  $T_n$  is  $F_{T_n}(\mathbf{z}) = \sum_j \pi_{jn} F_{T_{jn}}(\mathbf{z})$  where  $F_{T_{jn}}(\mathbf{z})$  is the cdf of  $T_{jn}$ .

**Theorem 4.5, Mixture Distribution CLT.** Suppose the  $g \times 1$  statistic  $T_n$  is equal to the estimator  $T_{jn}$  with probability  $\pi_{jn}$  for  $j = 1, \dots, J$  where

$\sum_j \pi_{jn} = 1$ ,  $\pi_{jn} \rightarrow \pi_j$  as  $n \rightarrow \infty$ , and  $\mathbf{u}_{jn} = \sqrt{n}(T_{jn} - \boldsymbol{\theta}) \xrightarrow{D} \mathbf{u}_j$  with  $E(\mathbf{u}_j) = \mathbf{0}$  and  $\text{Cov}(\mathbf{u}_j) = \boldsymbol{\Sigma}_j$ . Then

$$\sqrt{n}(T_n - \boldsymbol{\theta}) \xrightarrow{D} \mathbf{u} \quad (4.7)$$

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

**Proof:** Note that  $T_n$  has a mixture distribution of the  $T_{jn}$  with probabilities  $\pi_{jn}$ . Hence  $\sqrt{n}(T_n - \boldsymbol{\theta})$  has a mixture distribution of the  $\mathbf{u}_{jn} = \sqrt{n}(T_{jn} - \boldsymbol{\theta})$ , and the cdf of  $\sqrt{n}(T_n - \boldsymbol{\theta})$  is  $\sum_j \pi_{jn} F_{\mathbf{u}_{jn}}(\mathbf{z}) \rightarrow \sum_j \pi_j F_{\mathbf{u}_j}(\mathbf{z})$  at continuity points  $\mathbf{z}$  of the  $F_{\mathbf{u}_j}$ .  $\square$

**Remark 4.6.** Another variable selection model is  $\mathbf{x}^T \boldsymbol{\beta} = \mathbf{x}_{S_i}^T \boldsymbol{\beta}_{S_i}$  for  $i = 1, \dots, K$ . Then submodel  $I$  underfits if no  $S_i \subseteq I$ . A necessary condition for an estimator to be consistent is  $P(\text{no } S_i \subseteq I_{min}) \rightarrow 0$  as  $n \rightarrow \infty$ . Then in Theorem 4.4, we can replace  $P(S \subseteq I_{min}) \rightarrow 1$  by  $P(\text{no } S_i \subseteq I_{min}) \rightarrow 0$  as  $n \rightarrow \infty$ .

### 4.3 Prediction Intervals

Prediction intervals for regression and prediction regions for multivariate regression are important topics. Inference after variable selection will consider bootstrap hypothesis testing. Applying certain prediction intervals or prediction regions to the bootstrap sample will result in confidence intervals or confidence regions. The prediction intervals and regions are based on samples of size  $n$ , while the bootstrap sample size is  $B = B_n$ . Hence this section and the following section are important.

**Definition 4.6.** Consider predicting a future test value  $Y_f$  given a  $p \times 1$  vector of predictors  $\mathbf{x}_f$  and training data  $(Y_1, \mathbf{x}_1), \dots, (Y_n, \mathbf{x}_n)$ . A large sample  $100(1 - \delta)\%$  prediction interval (PI) for  $Y_f$  has the form  $[\hat{L}_n, \hat{U}_n]$  where  $P(\hat{L}_n \leq Y_f \leq \hat{U}_n)$  is eventually bounded below by  $1 - \delta$  as the sample size  $n \rightarrow \infty$ . A large sample  $100(1 - \delta)\%$  PI is *asymptotically optimal* if it has the shortest asymptotic length: the length of  $[\hat{L}_n, \hat{U}_n]$  converges to  $U_s - L_s$  as  $n \rightarrow \infty$  where  $[L_s, U_s]$  is the population shorth: the shortest interval covering at least  $100(1 - \delta)\%$  of the mass.

If  $Y_f | \mathbf{x}_f$  has a pdf, we often want  $P(\hat{L}_n \leq Y_f \leq \hat{U}_n) \rightarrow 1 - \delta$  as  $n \rightarrow \infty$ . The interpretation of a  $100(1 - \delta)\%$  PI for a random variable  $Y_f$  is similar to that of a confidence interval (CI). Collect data, then form the PI, and repeat for a total of  $k$  times where the  $k$  trials are independent from the same population. If  $Y_{fi}$  is the  $i$ th random variable and  $PI_i$  is the  $i$ th PI,



then the probability that  $Y_{fj} \in PI_j$  for  $j$  of the PIs approximately follows a binomial( $k, \rho = 1 - \delta$ ) distribution. Hence if 100 95% PIs are made,  $\rho = 0.95$  and  $Y_{fj} \in PI_j$  happens about 95 times.

There are two big differences between CIs and PIs. First, the length of the CI goes to 0 as the sample size  $n$  goes to  $\infty$  while the length of the PI converges to some nonzero number  $J$ , say. Secondly, many confidence intervals work well for large classes of distributions while many prediction intervals assume that the distribution of the data is known up to some unknown parameters. Usually the  $N(\mu, \sigma^2)$  distribution is assumed, and the parametric PI may not perform well if the normality assumption is violated. This section will describe three nonparametric PIs for the additive error regression model,  $Y = m(\mathbf{x}) + e$ , that work well for a large class of unknown zero mean error distributions.

First we will consider the location model,  $Y_i = \mu + e_i$ , where  $Y_1, \dots, Y_n, Y_f$  are iid and there are no vectors of predictors  $\mathbf{x}_i$  and  $\mathbf{x}_f$ . Let  $Z_{(1)} \leq Z_{(2)} \leq \dots \leq Z_{(n)}$  be the order statistics of  $n$  iid random variables  $Z_1, \dots, Z_n$ . Let a future random variable  $Z_f$  be such that  $Z_1, \dots, Z_n, Z_f$  are iid. Let  $k_1 = \lceil n\delta/2 \rceil$  and  $k_2 = \lceil n(1 - \delta/2) \rceil$  where  $\lceil x \rceil$  is the smallest integer  $\geq x$ . For example,  $\lceil 7.7 \rceil = 8$ . Then a common nonparametric large sample  $100(1 - \delta)\%$  prediction interval for  $Z_f$  is

$$[Z_{(k_1)}, Z_{(k_2)}] \quad (4.8)$$

where  $0 < \delta < 1$ . See Frey (2013) for references.

The shorth( $c$ ) estimator of the population shorth is useful for making asymptotically optimal prediction intervals. With the  $Z_i$  and  $Z_{(i)}$  as in the above paragraph, let the shortest closed interval containing at least  $c$  of the  $Z_i$  be

$$\text{shorth}(c) = [Z_{(s)}, Z_{(s+c-1)}]. \quad (4.9)$$

Let

$$k_n = \lceil n(1 - \delta) \rceil. \quad (4.10)$$

Frey (2013) showed that for large  $n\delta$  and iid data, the shorth( $k_n$ ) prediction interval has maximum undercoverage  $\approx 1.12\sqrt{\delta/n}$ , and used the shorth( $c$ ) estimator as the large sample  $100(1 - \delta)\%$  PI where

$$c = \min(n, \lceil n[1 - \delta + 1.12\sqrt{\delta/n}] \rceil). \quad (4.11)$$

An interesting fact is that the maximum undercoverage occurs for the family of uniform  $U(\theta_1, \theta_2)$  distributions where such a distribution has pdf  $f(y) = 1/(\theta_2 - \theta_1)$  for  $\theta_1 \leq y \leq \theta_2$  where  $f(y) = 0$ , otherwise, and  $\theta_1 < \theta_2$ .

A problem with the prediction intervals that cover  $\approx 100(1 - \delta)\%$  of the training data cases  $Y_i$  (such as (4.8) using  $c = k_n$  given by (4.9)), is that they have coverage lower than the nominal coverage of  $1 - \delta$  for moderate  $n$ . This result is not surprising since empirically statistical methods perform worse on test data. For iid data, Frey (2013) used (4.10) to correct for undercoverage.

**Example 4.1.** Given below were votes for preseason 1A basketball poll from Nov. 22, 2011 WSIL News where the 778 was a typo: the actual value was 78. As shown below, finding  $\text{shorth}(3)$  from the ordered data is simple. If the outlier was corrected,  $\text{shorth}(3) = [76, 78]$ .

111    89    778    78    76

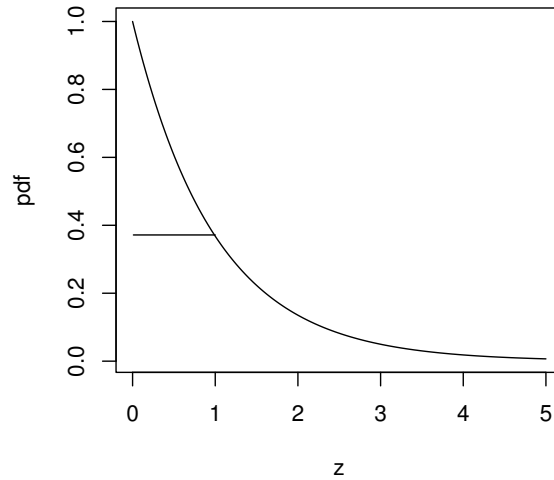
order data: 76 78 89 111 778

$$13 = 89 - 76$$

$$33 = 111 - 78$$

$$689 = 778 - 89$$

$\text{shorth}(3) = [76, 89]$



**Fig. 4.1** The 36.8% Highest Density Region is  $[0, 1]$

**Remark. 4.7.** The large sample  $100(1 - \delta)\%$  shorth PI (4.10) may or may not be asymptotically optimal if the  $100(1 - \delta)\%$  population shorth is  $[L_s, U_s]$  and  $F(x)$  is not strictly increasing in intervals  $(L_s - \delta, L_s + \delta)$  and  $(U_s - \delta, U_s + \delta)$  for some  $\delta > 0$ . To see the issue, suppose  $Y$  has probability mass function (pmf)  $p(0) = 0.4$ ,  $p(1) = 0.3$ ,  $p(2) = 0.2$ ,  $p(3) = 0.06$ , and  $p(4) = 0.04$ . Then the 90% population shorth is  $[0, 2]$  and the  $100(1 - \delta)\%$

population shorth is  $[0,3]$  for  $(1 - \delta) \in (0.9, 0.96]$ . Let  $W_i = I(Y_i \leq x) = 1$  if  $Y_i \leq x$  and 0, otherwise. The empirical cdf

$$\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n I(Y_i \leq x) = \frac{1}{n} \sum_{i=1}^n I(Y_{(i)} \leq x)$$

is the sample proportion of  $Y_i \leq x$ . If  $Y_1, \dots, Y_n$  are iid, then for fixed  $x$ ,  $n\hat{F}_n(x) \sim \text{binomial}(n, F(x))$ . Thus  $\hat{F}_n(x) \sim AN(F(x), F(x)(1 - F(x))/n)$ . For the  $Y$  with the above pmf,  $\hat{F}_n(2) \xrightarrow{P} 0.9$  as  $n \rightarrow \infty$  with  $P(\hat{F}_n(2) < 0.9) \rightarrow 0.5$  and  $P(\hat{F}_n(2) \geq 0.9) \rightarrow 0.5$  as  $n \rightarrow \infty$ . Hence the large sample 90% PI (4.10) will be  $[0,2]$  or  $[0,3]$  with probabilities  $\rightarrow 0.5$  as  $n \rightarrow \infty$  with expected asymptotic length of 2.5 and expected asymptotic coverage converging to 0.93. However, the large sample  $100(1 - \delta)\%$  PI (4.10) converges to  $[0,3]$  and is asymptotically optimal with asymptotic coverage 0.96 for  $(1 - \delta) \in (0.9, 0.96)$ .

For a random variable  $Y$ , the  $100(1 - \delta)\%$  highest density region is a union of  $k \geq 1$  disjoint intervals such that the mass within the intervals  $\geq 1 - \delta$  and the sum of the  $k$  interval lengths is as small as possible. Suppose that  $f(z)$  is a unimodal pdf that has interval support, and that the pdf  $f(z)$  of  $Y$  decreases rapidly as  $z$  moves away from the mode. Let  $[a, b]$  be the shortest interval such that  $F_Y(b) - F_Y(a) = 1 - \delta$  where the cdf  $F_Y(z) = P(Y \leq z)$ . Then the interval  $[a, b]$  is the  $100(1 - \delta)\%$  highest density region. To find the  $100(1 - \delta)\%$  highest density region of a pdf, move a horizontal line down from the top of the pdf. The line will intersect the pdf or the boundaries of the support of the pdf at  $[a_1, b_1], \dots, [a_k, b_k]$  for some  $k \geq 1$ . Stop moving the line when the areas under the pdf corresponding to the intervals is equal to  $1 - \delta$ . As an example, let  $f(z) = e^{-z}$  for  $z > 0$ . See Figure 4.1 where the area under the pdf from 0 to 1 is 0.368. Hence  $[0,1]$  is the 36.8% highest density region. The shorth PI estimates the highest density interval which is the highest density region for a distribution with a unimodal pdf. Often the highest density region is an interval  $[a, b]$  where  $f(a) = f(b)$ , especially if the support where  $f(z) > 0$  is  $(-\infty, \infty)$ .

The additive error regression model is  $Y = m(\mathbf{x}) + e$  where  $m(\mathbf{x})$  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  $m(\mathbf{x}) = \mathbf{x}^T \boldsymbol{\beta} = \mathbf{x}_I^T \boldsymbol{\beta}_I$  if  $S \subseteq I$ . Let the residuals  $r_i = Y_i - \hat{m}(\mathbf{x}_i) = Y_i - \hat{Y}_i$  for  $i = 1, \dots, n$ . Assume  $\hat{m}(\mathbf{x})$  is a consistent estimator of  $m(\mathbf{x})$  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  $\hat{Y}_f = \hat{m}(\mathbf{x}_f)$ . Then  $P(Y_f \in [\hat{Y}_f + \hat{L}_n(r), \hat{Y}_f + \hat{U}_n(r)]) \rightarrow P(Y_f \in [m(\mathbf{x}_f) + L, m(\mathbf{x}_f) + U]) = P(e \in [L, U]) = 1 - \delta$  as  $n \rightarrow \infty$ . 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 (4.16), 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  $x_1^*, \dots, x_d^*$  while PCR and PLS use variables that are linear combinations of the predictors  $V_j = \gamma_j^T \mathbf{x}$  for  $j = 1, \dots, d$ . (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 Chapter 5 for more about these estimators. 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  $Y_i = m(\mathbf{x}_i) + e_i$ , and variable selection models for MLR have this form, as noted by Olive (2018). Pelawa Watagoda and Olive (2019b) 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  $h_f = \mathbf{x}_{I,f}^T (\mathbf{X}_I^T \mathbf{X}_I)^{-1} \mathbf{x}_{I,f} > 2d/n$ : if  $\mathbf{x}_{I,f}$  is too far from the data  $\mathbf{x}_{I,1}, \dots, \mathbf{x}_{I,n}$ , 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 (4.11) 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

$$q_n = \min(1 - \delta/2, 1 - \delta + 10\delta d/n), \text{ otherwise.} \quad (4.12)$$

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

$$c = \lceil nq_n \rceil, \quad (4.13)$$

and let

$$b_n = \left(1 + \frac{15}{n}\right) \sqrt{\frac{n+2d}{n-d}} \quad (4.14)$$

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

$$[\hat{m}(\mathbf{x}_f) + b_n \tilde{\xi}_{\delta_1}, \hat{m}(\mathbf{x}_f) + b_n \tilde{\xi}_{1-\delta_2}]. \quad (4.15)$$

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  $i_1, \dots, i_{n_V}$ . The estimator  $\hat{m}_H(\mathbf{x})$  is computed using the training data set  $H$ . Then the validation residuals  $v_j = Y_{i_j} - \hat{m}_H(\mathbf{x}_{i_j})$  are computed for the  $j = 1, \dots, n_V$  cases in the validation set  $V$ . Find the Frey PI  $[v_{(s)}, v_{(s+c-1)}]$  of the validation residuals (replacing  $n$  in (4.10) by  $n_V = n - n_H$ ). Then the second new 100 $(1 - \delta)\%$  large sample PI for  $Y_f$  is

$$[\hat{m}_H(\mathbf{x}_f) + v_{(s)}, \hat{m}_H(\mathbf{x}_f) + v_{(s+c-1)}]. \quad (4.16)$$

**Remark 4.8.** Note that correction factors  $b_n \rightarrow 1$  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:  $t_{d_n,1-\delta}/z_{1-\delta} \rightarrow 1$  and  $pF_{p,d_n,1-\delta}/\chi_{p,1-\delta}^2 \rightarrow 1$  if  $d_n \rightarrow \infty$  as  $n \rightarrow 1$ . 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 4.9.** 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 (4.15) while PI (4.14) needs a stronger correction factor.

We can also motivate PI (4.15) by modifying the justification for the Lei et al. (2018) split conformal prediction interval

$$[\hat{m}_H(\mathbf{x}_f) - a_q, \hat{m}_H(\mathbf{x}_f) + a_q] \quad (4.17)$$

where  $a_q$  is the  $100(1 - \alpha)$ th quantile of the absolute validation residuals. PI (4.15) is a modification of the split conformal PI that is asymptotically optimal. Suppose  $(Y_i, \mathbf{x}_i)$  are iid for  $i = 1, \dots, n, n + 1$  where  $(Y_f, \mathbf{x}_f) = (Y_{n+1}, \mathbf{x}_{n+1})$ . Compute  $\hat{m}_H(\mathbf{x})$  from the cases in  $H$ . For example, get  $\hat{\beta}_H$  from the cases in  $H$ . Consider the validation residuals  $v_i$  for  $i = 1, \dots, n_V$  and the validation residual  $v_{n_V+1}$  for case  $(Y_f, \mathbf{x}_f)$ . Since these  $n_V + 1$  cases are iid, the probability that  $v_t$  has rank  $j$  for  $j = 1, \dots, n_V + 1$  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  $j = 1, \dots, n_V$ . 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

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

$P(v_{n_V+1}$  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 PIs 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  $\hat{m}(\mathbf{x})$  is "symmetric in the data," which is satisfied for multiple linear regression estimators.

The PIs (4.14) to (4.16) can be used with  $\hat{m}(\mathbf{x}) = \hat{Y}_f = \mathbf{x}_{I_d}^T \hat{\beta}_{I_d}$  where  $I_d$  denotes the index of predictors selected from the model or variable selection method. If  $\hat{\beta}$  is a consistent estimator of  $\beta$ , the PIs (4.14) and (4.15) are asymptotically optimal for a large class of error distributions while the split conformal PI (4.16) 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  $P(S \subseteq I_{min}) \rightarrow 1$  as  $n \rightarrow \infty$ . Then forward selection with EBIC and relaxed lasso can produce consistent estimators. PLS can be  $\sqrt{n}$  consistent. See Chapter 5 for the large sample for many MLR estimators.

None of the three prediction intervals (4.14), (4.15), and (4.16) dominates the other two. Recall that  $\beta_S$  is an  $a_S \times 1$  vector in (4.1). 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 (4.14) can be much shorter than PIs (4.15) and (4.16). For  $n/d$  large, PIs (4.14) and (4.15) can be shorter than PI (4.16) if the error distribution is not unimodal and symmetric; however, PI (4.16) 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 (4.16) and PI (4.15) can have coverage closer to the nominal coverage than PI (4.14). 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 (4.15) will be long compared to PI (4.16).

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 (4.14) and (4.15), are the only easily computed asymptotically optimal PIs for a wide range of consistent estimators  $\hat{\beta}$  of  $\beta$  for the multiple linear regression model. If the error distribution is  $e \sim EXP(1) - 1$ , then the asymptotic length of the 95% PI (4.14) or (4.15) 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 5.10 and Pelawa Watagoda and Olive (2019b).

## 4.4 Prediction Regions

Consider predicting a  $p \times 1$  future test value  $\mathbf{x}_f$ , given past training data  $\mathbf{x}_1, \dots, \mathbf{x}_n$  where  $\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{x}_f$  are iid. Much as confidence regions and intervals give a measure of precision for the point estimator  $\hat{\theta}$  of the parameter  $\theta$ , prediction regions and intervals give a measure of precision of the point estimator  $T = \hat{\mathbf{x}}_f$  of the future random vector  $\mathbf{x}_f$ .

**Definition 4.7.** A large sample  $100(1 - \delta)\%$  prediction region is a set  $\mathcal{A}_n$  such that  $P(\mathbf{x}_f \in \mathcal{A}_n)$  is eventually bounded below by  $1 - \delta$  as  $n \rightarrow \infty$ . A prediction region is *asymptotically optimal* if its volume converges in probability to the volume of the minimum volume covering region or the highest density region of the distribution of  $\mathbf{x}_f$ .

If  $\mathbf{x}_f$  has a pdf, we often want  $P(\mathbf{x}_f \in \mathcal{A}_n) \rightarrow 1 - \delta$  as  $n \rightarrow \infty$ . A PI is a prediction region where  $p = 1$ . Highest density regions are usually hard to estimate for  $p$  not much larger than four, but many elliptically contoured distributions with a nonsingular population covariance matrix, including the multivariate normal distribution, have highest density regions that can be estimated by the nonparametric prediction region (4.24). For more about highest density regions, see Olive (2017b, pp. 148-155) and Hyndman (1996).

For multivariate data, sample Mahalanobis distances play a role similar to that of residuals in multiple linear regression. Let the observed training data be collected in an  $n \times p$  matrix  $\mathbf{W}$ . Let the  $p \times 1$  column vector  $T = T(\mathbf{W})$  be a multivariate location estimator, and let the  $p \times p$  symmetric positive definite matrix  $\mathbf{C} = \mathbf{C}(\mathbf{W})$  be a dispersion estimator.

**Definition 4.8.** Let  $x_{1j}, \dots, x_{nj}$  be measurements on the  $j$ th random variable  $X_j$  corresponding to the  $j$ th column of the data matrix  $\mathbf{W}$ . The  $j$ th *sample mean* is  $\bar{x}_j = \frac{1}{n} \sum_{k=1}^n x_{kj}$ . The *sample covariance*  $S_{ij}$  estimates  $\text{Cov}(X_i, X_j) = \sigma_{ij} = E[(X_i - E(X_i))(X_j - E(X_j))]$ , and

$$S_{ij} = \frac{1}{n-1} \sum_{k=1}^n (x_{ki} - \bar{x}_i)(x_{kj} - \bar{x}_j).$$

$S_{ii} = S_i^2$  is the *sample variance* that estimates the population variance  $\sigma_{ii} = \sigma_i^2$ . The *sample correlation*  $r_{ij}$  estimates the population correlation  $\text{Cor}(X_i, X_j) = \rho_{ij} = \frac{\sigma_{ij}}{\sigma_i \sigma_j}$ , and

$$r_{ij} = \frac{S_{ij}}{S_i S_j} = \frac{S_{ij}}{\sqrt{S_{ii} S_{jj}}} = \frac{\sum_{k=1}^n (x_{ki} - \bar{x}_i)(x_{kj} - \bar{x}_j)}{\sqrt{\sum_{k=1}^n (x_{ki} - \bar{x}_i)^2} \sqrt{\sum_{k=1}^n (x_{kj} - \bar{x}_j)^2}}.$$

**Definition 4.9.** Let  $\mathbf{x}_1, \dots, \mathbf{x}_n$  be the data where  $\mathbf{x}_i$  is a  $p \times 1$  vector. The **sample mean** or *sample mean vector*

$$\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i = (\bar{x}_1, \dots, \bar{x}_p)^T = \frac{1}{n} \mathbf{W}^T \mathbf{1}$$

where  $\mathbf{1}$  is the  $n \times 1$  vector of ones. The **sample covariance matrix**

$$\mathbf{S} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T = (S_{ij}).$$

That is, the  $ij$  entry of  $\mathbf{S}$  is the sample covariance  $S_{ij}$ . The *classical estimator of multivariate location and dispersion* is  $(T, \mathbf{C}) = (\bar{\mathbf{x}}, \mathbf{S})$ . The **sample correlation matrix**

$$\mathbf{R} = (r_{ij}).$$

That is, the  $ij$  entry of  $\mathbf{R}$  is the sample correlation  $r_{ij}$ .

It can be shown that  $(n-1)\mathbf{S} = \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T - \bar{\mathbf{x}} \bar{\mathbf{x}}^T =$

$$\mathbf{W}^T \mathbf{W} - \frac{1}{n} \mathbf{W}^T \mathbf{1} \mathbf{1}^T \mathbf{W}.$$

Hence if the *centering matrix*  $\mathbf{G} = \mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^T$ , then  $(n-1)\mathbf{S} = \mathbf{W}^T \mathbf{G} \mathbf{W}$ .

See Definition 1.24 for the population mean and population covariance matrix. Definition 2.18 also defined a sample covariance matrix. The Ma-



halanobis distance in Definition 4.9 is a random variable that estimates the population Mahalanobis distance of Definition 1.38.

**Definition 4.9.** The  $i$ th Mahalanobis distance  $D_i = \sqrt{D_i^2}$  where the  $i$ th squared Mahalanobis distance is

$$D_i^2 = D_i^2(T(\mathbf{W}), \mathbf{C}(\mathbf{W})) = (\mathbf{x}_i - T(\mathbf{W}))^T \mathbf{C}^{-1}(\mathbf{W})(\mathbf{x}_i - T(\mathbf{W})) \quad (4.18)$$

for each point  $\mathbf{x}_i$ . Notice that  $D_i^2$  is a random variable (scalar valued). Let  $(T, \mathbf{C}) = (T(\mathbf{W}), \mathbf{C}(\mathbf{W}))$ . Then

$$D_{\mathbf{x}}^2(T, \mathbf{C}) = (\mathbf{x} - T)^T \mathbf{C}^{-1}(\mathbf{x} - T).$$

Hence  $D_i^2$  uses  $\mathbf{x} = \mathbf{x}_i$ .

Let the  $p \times 1$  location vector be  $\boldsymbol{\mu}$ , often the population mean, and let the  $p \times p$  dispersion matrix be  $\boldsymbol{\Sigma}$ , often the population covariance matrix. Notice that if  $\mathbf{x}$  is a random vector, then the population squared Mahalanobis distance from Definition 1.38 is

$$D_{\mathbf{x}}^2(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \quad (4.19)$$

and that the term  $\boldsymbol{\Sigma}^{-1/2}(\mathbf{x} - \boldsymbol{\mu})$  is the  $p$ -dimensional analog to the  $z$ -score used to transform a univariate  $N(\mu, \sigma^2)$  random variable into a  $N(0, 1)$  random variable. Hence the sample Mahalanobis distance  $D_i = \sqrt{D_i^2}$  is an analog of the absolute value  $|Z_i|$  of the sample  $Z$ -score  $Z_i = (X_i - \bar{X})/\hat{\sigma}$ . Also notice that the Euclidean distance of  $\mathbf{x}_i$  from the estimate of center  $T(\mathbf{W})$  is  $D_i(T(\mathbf{W}), \mathbf{I}_p)$  where  $\mathbf{I}_p$  is the  $p \times p$  identity matrix.

Consider the hyperellipsoid

$$\mathcal{A}_n = \{\mathbf{x} : D_{\mathbf{x}}^2(\bar{\mathbf{x}}, \mathbf{S}) \leq D_{(c)}^2\} = \{\mathbf{x} : D_{\mathbf{x}}(\bar{\mathbf{x}}, \mathbf{S}) \leq D_{(c)}\}. \quad (4.20)$$

If  $n$  is large, we can use  $c = k_n = \lceil n(1 - \delta) \rceil$ . If  $n$  is not large, using  $c = U_n$  where  $U_n$  decreases to  $k_n$ , can improve small sample performance.  $U_n$  will be defined in the paragraph below Equation (4.23). Olive (2013a) showed that (4.19) is a large sample  $100(1 - \delta)\%$  prediction region under mild conditions, although regions with smaller volumes may exist. Note that the result follows since if  $\boldsymbol{\Sigma}_{\mathbf{x}}$  and  $\mathbf{S}$  are nonsingular, then the Mahalanobis distance is a continuous function of  $(\bar{\mathbf{x}}, \mathbf{S})$ . Let  $\boldsymbol{\mu} = E(\mathbf{x})$  and  $D = D(\boldsymbol{\mu}, \boldsymbol{\Sigma}_{\mathbf{x}})$ . Then  $D_i \xrightarrow{D} D$  and  $D_i^2 \xrightarrow{D} D^2$ . Hence the sample percentiles of the  $D_i$  are consistent estimators of the population percentiles of  $D$  at continuity points of the cumulative distribution function of  $D$ .

A problem with the prediction regions that cover  $\approx 100(1 - \delta)\%$  of the training data cases  $\mathbf{x}_i$  (such as (4.19) for  $c = k_n$ ), is that they have coverage lower than the nominal coverage of  $1 - \delta$  for moderate  $n$ . This result is not surprising since empirically statistical methods perform worse on test data.

Increasing  $c$  will improve the coverage for moderate samples. Also see Remark 4.8. Empirically for many distributions, for  $n \approx 20p$ , the prediction region (4.19) applied to iid data using  $c = k_n = \lceil n(1 - \delta) \rceil$  tended to have undercoverage as high as 5%. The undercoverage decreases rapidly as  $n$  increases. Let  $q_n = \min(1 - \delta + 0.05, 1 - \delta + p/n)$  for  $\delta > 0.1$  and

$$q_n = \min(1 - \delta/2, 1 - \delta + 10\delta p/n), \text{ otherwise.} \quad (4.21)$$

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

$$c = \lceil nq_n \rceil \quad (4.22)$$

in (4.19) decreased the undercoverage. Note that Equations (4.11) and (4.12) are similar to Equations (4.20) and (4.21), but replace  $p$  by  $d$ .

If  $(T, \mathbf{C})$  is a  $\sqrt{n}$  consistent estimator of  $(\boldsymbol{\mu}, d\boldsymbol{\Sigma})$  for some constant  $d > 0$  where  $\boldsymbol{\Sigma}$  is nonsingular, then  $D^2(T, \mathbf{C}) = (\mathbf{x} - T)^T \mathbf{C}^{-1} (\mathbf{x} - T) =$

$$\begin{aligned} & (\mathbf{x} - \boldsymbol{\mu} + \boldsymbol{\mu} - T)^T [\mathbf{C}^{-1} - d^{-1}\boldsymbol{\Sigma}^{-1} + d^{-1}\boldsymbol{\Sigma}^{-1}] (\mathbf{x} - \boldsymbol{\mu} + \boldsymbol{\mu} - T) \\ & = d^{-1}D^2(\boldsymbol{\mu}, \boldsymbol{\Sigma}) + o_p(1). \end{aligned}$$

Thus the sample percentiles of  $D_i^2(T, \mathbf{C})$  are consistent estimators of the percentiles of  $d^{-1}D^2(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  (at continuity points  $D_{1-\delta}$  of the cdf of  $D^2(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ ). If  $\mathbf{x} \sim N_m(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , then  $D_{\mathbf{x}}^2(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = D^2(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \sim \chi_m^2$ .

Suppose  $(T, \mathbf{C}) = (\bar{\mathbf{x}}_M, b\mathbf{S}_M)$  is the sample mean and scaled sample covariance matrix applied to some subset of the data. The classical estimator and RMVN estimator from Section 7.1 satisfy this assumption. For  $h > 0$ , the hyperellipsoid

$$\{\mathbf{z} : (\mathbf{z} - T)^T \mathbf{C}^{-1} (\mathbf{z} - T) \leq h^2\} = \{\mathbf{z} : D_{\mathbf{z}}^2 \leq h^2\} = \{\mathbf{z} : D_{\mathbf{z}} \leq h\} \quad (4.23)$$

has volume equal to

$$\frac{2\pi^{p/2}}{p\Gamma(p/2)} h^p \sqrt{\det(\mathbf{C})} = \frac{2\pi^{p/2}}{p\Gamma(p/2)} h^p b^{p/2} \sqrt{\det(\mathbf{S}_M)}. \quad (4.24)$$

A future observation (random vector)  $\mathbf{x}_f$  is in the region (4.22) if  $D_{\mathbf{x}_f} \leq h$ .

If  $(T, \mathbf{C})$  is a consistent estimator of  $(\boldsymbol{\mu}, d\boldsymbol{\Sigma})$  for some constant  $d > 0$  where  $\boldsymbol{\Sigma}$  is nonsingular, then (4.22) is a large sample  $100(1 - \delta)\%$  prediction region if  $h = D_{(U_n)}$  where  $D_{(U_n)}$  is the  $100q_n$ th sample quantile of the  $D_i$  where  $q_n$  is defined above (4.21). If  $\mathbf{x}_1, \dots, \mathbf{x}_n$  and  $\mathbf{x}_f$  are iid, then prediction region (4.24) is asymptotically optimal for a large class of elliptically contoured distributions since the volume of (4.24) converges in probability to the volume of the highest density region. (These distributions have a highest density region which is a hyperellipsoid determined by a population Mahalanobis distance. See Section 1.7.)

The Olive (2013a) nonparametric prediction region uses  $(T, \mathbf{C}) = (\bar{\mathbf{x}}, \mathbf{S})$ . For the classical prediction region, see Chew (1966) and Johnson and Wichern (1988, pp. 134, 151). Refer to the above paragraph for  $D_{(U_n)}$ .

**Definition 4.10.** The large sample  $100(1 - \delta)\%$  nonparametric prediction region for a future value  $\mathbf{x}_f$  given iid data  $\mathbf{x}_1, \dots, \mathbf{x}_n$  is

$$\{\mathbf{z} : D_{\mathbf{z}}^2(\bar{\mathbf{x}}, \mathbf{S}) \leq D_{(U_n)}^2\}, \quad (4.25)$$

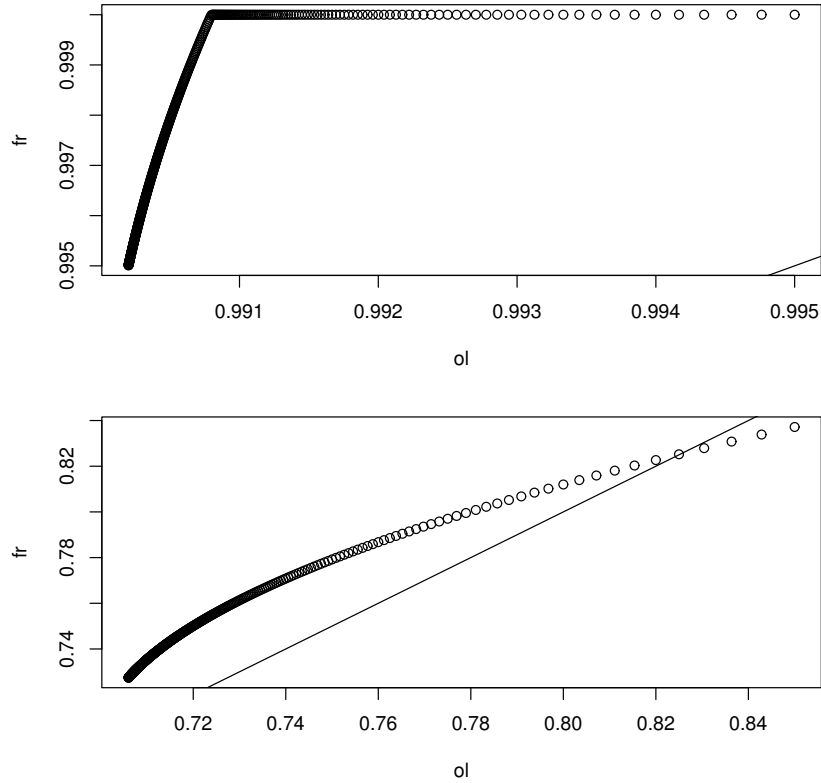
while the large sample  $100(1 - \delta)\%$  classical prediction region is

$$\{\mathbf{z} : D_{\mathbf{z}}^2(\bar{\mathbf{x}}, \mathbf{S}) \leq \chi_{p, 1-\delta}^2\}. \quad (4.26)$$

If  $p$  is small, Mahalanobis distances tend to be right skewed with a population shorth that discards the right tail. For  $p = 1$  and  $n \geq 20$ , the finite sample correction factors  $c/n$  for  $c$  given by (4.10) and (4.21) do not differ by much more than 3% for  $0.01 \leq \delta \leq 0.5$ . See Figure 4.2 where  $ol = (\text{Eq. 4.21})/n$  is plotted versus  $fr = (\text{Eq. 4.10})/n$  for  $n = 20, 21, \dots, 500$ . The top plot is for  $\delta = 0.01$ , while the bottom plot is for  $\delta = 0.3$ . The identity line is added to each plot as a visual aid. The value of  $n$  increases from 20 to 500 from the right of the plot to the left of the plot. Examining the axes of each plot shows that the correction factors do not differ greatly. *R* code to create Figure 4.2 is shown below.

```
cmar <- par("mar"); par(mfrow = c(2, 1))
par(mar=c(4.0, 4.0, 2.0, 0.5))
frey(0.01); frey(0.3)
par(mfrow = c(1, 1)); par(mar=cmar)
```

**Remark 4.10.** The nonparametric prediction region (4.24) is useful if  $\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{x}_f$  are iid from a distribution with a nonsingular covariance matrix, and the sample size  $n$  is large enough. The distribution could be continuous, discrete, or a mixture. The asymptotic coverage is  $1 - \delta$  if  $D$  has a pdf, although prediction regions with smaller volume may exist. If the  $100(1 - \delta)\%$ th percentile  $D_{1-\delta}$  of  $D$  is not a continuity point of the distribution of  $D$ , then the asymptotic coverage tends to be  $\geq 1 - \delta$  since a sample percentile with cutoff  $q_n$  that decreases to  $1 - \delta$  is used and a closed region is used. Often  $D$  has a continuous distribution and hence has no discontinuity points for  $0 < \delta < 1$ . (If there is a jump in the distribution from 0.9 to 0.96 at discontinuity point  $a$ , and the nominal coverage is 0.95, we want 0.96 coverage instead of 0.9. So we want the sample percentile to decrease to  $a$ .) The nonparametric prediction region (4.24) contains  $U_n$  of the training data cases  $\mathbf{x}_i$  provided that  $\mathbf{S}$  is nonsingular, even if the model is wrong. For many distributions, the coverage started to be close to  $1 - \delta$  for  $n \geq 10p$  where the coverage is the simulated percentage of times that the prediction region contained  $\mathbf{x}_f$ .



**Fig. 4.2** Correction Factor Comparison when  $\delta = 0.01$  (Top Plot) and  $\delta = 0.3$  (Bottom Plot)

**Remark 4.11.** The most used prediction regions assume that the error vectors are iid from a multivariate normal distribution. Using (4.23), the ratio of the volumes of regions (4.25) and (4.24) is

$$\left( \frac{\chi_{p,1-\delta}^2}{D_{(U_n)}^2} \right)^{p/2},$$

which can become close to zero rapidly as  $p$  gets large if the  $\mathbf{x}_i$  are not from the light tailed multivariate normal distribution. For example, suppose  $\chi_{4,0.5}^2 \approx 3.33$  and  $D_{(U_n)}^2 \approx D_{\mathbf{x},0.5}^2 = 6$ . Then the ratio is  $(3.33/6)^2 \approx 0.308$ . Hence if the data is not multivariate normal, severe undercoverage can occur if the classical prediction region is used, and the undercoverage tends to get worse as the dimension  $p$  increases. The coverage need not to go to 0, since by the multivariate Chebyshev's inequality,  $P(D_{\mathbf{x}}^2(\boldsymbol{\mu}, \boldsymbol{\Sigma}_{\mathbf{x}}) \leq \gamma) \geq 1 - p/\gamma > 0$

for  $\gamma > p$  where the population covariance matrix  $\Sigma_{\mathbf{x}} = \text{Cov}(\mathbf{x})$ . See Budny (2014), Chen (2011), and Navarro (2014, 2016). Using  $\gamma = h^2 = p/\delta$  in (4.22) usually results in prediction regions with volume and coverage that is too large.

**Remark 4.12.** The nonparametric prediction region (4.24) starts to have good coverage for  $n \geq 10p$  for a large class of distributions. Olive (2013a) suggests  $n \geq 50p$  may be needed for the prediction region to have a good volume. Of course for any  $n$  there are error distributions that will have severe undercoverage.

For the multivariate lognormal distribution with  $n = 20p$ , the large sample nonparametric 95% prediction region (4.24) had coverages 0.970, 0.959, and 0.964 for  $p = 100, 200$ , and 500. Some *R* code is below.

```
nruns=1000 #lognormal, p = 100, n = 20p = 2000
count<-0
for(i in 1:nruns){
x <- exp(matrix(rnorm(200000), ncol=100, nrow=2000))
xff <- exp(as.vector(rnorm(100)))
count <- count + predrgn(x,xf=xff)$inr}
count #970/1000, may take a few minutes
```

Notice that for the training data  $\mathbf{x}_1, \dots, \mathbf{x}_n$ , if  $\mathbf{C}^{-1}$  exists, then  $c \approx 100q_n\%$  of the  $n$  cases are in the prediction regions for  $\mathbf{x}_f = \mathbf{x}_i$ , and  $q_n \rightarrow 1 - \delta$  even if  $(T, \mathbf{C})$  is not a good estimator. Hence the coverage  $q_n$  of the training data is robust to model assumptions. Of course the volume of the prediction region could be large if a poor estimator  $(T, \mathbf{C})$  is used or if the  $\mathbf{x}_i$  do not come from an elliptically contoured distribution. Also notice that  $q_n = 1 - \delta/2$  or  $q_n = 1 - \delta + 0.05$  for  $n \leq 20p$  and  $q_n \rightarrow 1 - \delta$  as  $n \rightarrow \infty$ . If  $q_n \equiv 1 - \delta$  and  $(T, \mathbf{C})$  is a consistent estimator of  $(\boldsymbol{\mu}, d\Sigma)$  where  $d > 0$  and  $\Sigma$  is nonsingular, then (4.22) with  $h = D_{(U_n)}$  is a large sample prediction region, but taking  $q_n$  given by (4.20) improves the finite sample performance of the prediction region. Taking  $q_n \equiv 1 - \delta$  does not take into account variability of  $(T, \mathbf{C})$ , and for  $n = 20p$  the resulting prediction region tended to have undercoverage as high as  $\min(0.05, \delta/2)$ . Using (4.20) helped reduce undercoverage for small  $n \geq 20p$  due to the unknown variability of  $(T, \mathbf{C})$ .

## 4.5 Bootstrapping Hypothesis Tests and Confidence Regions

This section shows that, under regularity conditions, applying the nonparametric prediction region of Section 4.4 to a bootstrap sample results in a confidence region. The volume of a confidence region  $\rightarrow 0$  as  $n \rightarrow 0$ , while

the volume of a prediction region goes to that of a population region that would contain a new  $\mathbf{x}_f$  with probability  $1 - \delta$ . The nominal coverage is  $100(1 - \delta)$ . If the actual coverage  $100(1 - \delta_n) > 100(1 - \delta)$ , then the region is *conservative*. If  $100(1 - \delta_n) < 100(1 - \delta)$ , then the region is *liberal*. A region that is 5% conservative is considered “much better” than a region that is 5% liberal.

When teaching confidence intervals, it is often noted that by the central limit theorem, the probability that  $\bar{Y}_n$  is within two standard deviations ( $2SD(\bar{Y}_n) = 2\sigma/\sqrt{n}$ ) of  $\theta = \mu$  is about 95%. Hence the probability that  $\theta$  is within two standard deviations of  $\bar{Y}_n$  is about 95%. Thus the interval  $[\theta - 1.96S/\sqrt{n}, \theta + 1.96S/\sqrt{n}]$  is a large sample 95% prediction interval for a future value of the sample mean  $\bar{Y}_{n,f}$  if  $\theta$  is known, while  $[\bar{Y}_n - 1.96S/\sqrt{n}, \bar{Y}_n + 1.96S/\sqrt{n}]$  is a large sample 95% confidence interval for the population mean  $\theta$ . Note that the lengths of the two intervals are the same. Where the interval is centered, at the parameter  $\theta$  or the statistic  $\bar{Y}_n$ , determines whether the interval is a prediction or a confidence interval. See Theorem 4.7 for a similar relationship between confidence regions and prediction regions.

**Definition 4.11.** A large sample  $100(1 - \delta)\%$  confidence region for a vector of parameters  $\boldsymbol{\theta}$  is a set  $\mathcal{A}_n$  such that  $P(\boldsymbol{\theta} \in \mathcal{A}_n)$  is eventually bounded below by  $1 - \delta$  as  $n \rightarrow \infty$ .

If  $\mathcal{A}_n$  is based on a squared Mahalanobis distance  $D^2$  with a limiting distribution that has a pdf, we often want  $P(\boldsymbol{\theta} \in \mathcal{A}_n) \rightarrow 1 - \delta$  as  $n \rightarrow \infty$ .

There are several methods for obtaining a bootstrap sample  $T_1^*, \dots, T_B^*$  where the sample size  $n$  is suppressed:  $T_i^* = T_{in}^*$ . The parametric bootstrap, nonparametric bootstrap, and residual bootstrap will be used. Applying prediction region (4.24) to the bootstrap sample will result in a confidence region for  $\boldsymbol{\theta}$ . When  $g = 1$ , applying the shorth PI (4.10) or PI (4.7) to the bootstrap sample results in a confidence interval for  $\theta$ . Section 4.5.2 will help clarify ideas.

When  $g = 1$ , a confidence interval is a special case of a confidence region. One sided confidence intervals give a lower or upper confidence bound for  $\theta$ . A large sample  $100(1 - \delta)\%$  lower confidence interval  $(-\infty, U_n]$  uses an upper confidence bound  $U_n$  and is in the lower tail of the distribution of  $\hat{\theta}$ . A large sample  $100(1 - \delta)\%$  upper confidence interval  $[L_n, \infty)$  uses a lower confidence bound  $L_n$  and is in the upper tail of the distribution of  $\hat{\theta}$ . These CIs can be useful if  $\theta \in [a, b]$  and  $\theta = a$  or  $\theta = b$  is of interest for a hypothesis test. For example,  $[a, b] = [0, 1]$  if  $\theta = \rho^2$ , the squared population correlation. Then use  $[0, U_n]$  and  $[L_n, 1]$  as CIs, e.g. if we expect  $\theta = 0$  we might test  $H_0 : \theta \leq 0.05$  versus  $H_0 : \theta > 0.05$ , and fail to reject  $H_0$  if  $U_n < 0.05$ . See Section 4.5.4 for an illustration. Again we often want the probability to converge to  $1 - \delta$  if the confidence interval is based on a statistic with an asymptotic distribution that has a pdf.

**Definition 4.12.** The interval  $[L_n, U_n]$  is a large sample  $100(1 - \delta)\%$  *confidence interval* for  $\theta$  if  $P(L_n \leq \theta \leq U_n)$  is eventually bounded below by  $1 - \delta$  as  $n \rightarrow \infty$ . The interval  $(-\infty, U_n]$  is a large sample  $100(1 - \delta)\%$  *lower confidence interval* for  $\theta$  if  $P(\theta \leq U_n)$  is eventually bounded below by  $1 - \delta$  as  $n \rightarrow \infty$ . The interval  $[L_n, \infty)$  is large sample  $100(1 - \delta)\%$  *upper confidence interval* for  $\theta$  if  $P(\theta \geq L_n)$  is eventually bounded below by  $1 - \delta$  as  $n \rightarrow \infty$ .

Next we discuss bootstrap confidence intervals that are obtained by applying prediction intervals (4.7) and (4.10) to the bootstrap sample. Some additional bootstrap CIs are obtained from bootstrap confidence regions from Section 4.5.2 when  $g = 1$ . See Efron (1982) and Chen (2016) for the percentile method CI. Let  $T_n$  be an estimator of a parameter  $\theta$  such as  $T_n = \bar{Z} = \sum_{i=1}^n Z_i/n$  with  $\theta = E(Z_1)$ . Let  $T_1^*, \dots, T_B^*$  be a bootstrap sample for  $T_n$ . Let  $T_{(1)}^*, \dots, T_{(B)}^*$  be the order statistics of the the bootstrap sample. The CI (4.26) is obtained by applying PI (4.7) to the bootstrap sample with  $B$  used instead of  $n$ . Hence (4.26) is also a large sample prediction interval for a future value of  $T_f^*$  if the  $T_i^*$  are iid from the empirical distribution discussed in Section 4.5.1.

**Definition 4.13.** The bootstrap percentile method large sample  $100(1 - \delta)\%$  confidence interval for  $\theta$  is an interval  $[T_{(k_L)}^*, T_{(k_U)}^*]$  containing  $\approx [B(1 - \delta)]$  of the  $T_i^*$ . Let  $k_1 = \lceil B\delta/2 \rceil$  and  $k_2 = \lceil B(1 - \delta/2) \rceil$ . A common choice is

$$[T_{(k_1)}^*, T_{(k_2)}^*]. \quad (4.27)$$

The large sample  $100(1 - \delta)\%$  *lower percentile method* CI for  $\theta$  is  $(-\infty, T_{(\lceil B(1-\delta) \rceil)}^*]$ . The large sample  $100(1 - \delta)\%$  *upper percentile method* CI for  $\theta$  is  $[T_{(\lceil B\delta \rceil)}^*, \infty)$ .

**Definition 4.14.** The large sample  $100(1 - \delta)\%$  *lower shorth* CI for  $\theta$  is  $(-\infty, T_{(c)}^*]$ , while the large sample  $100(1 - \delta)\%$  *upper shorth* CI for  $\theta$  is  $[T_{(B-c+1)}^*, \infty)$ . The large sample  $100(1 - \delta)\%$  *shorth(c) CI* uses the interval  $[T_{(1)}^*, T_{(c)}^*], [T_{(2)}^*, T_{(c+1)}^*], \dots, [T_{(B-c+1)}^*, T_{(B)}^*]$  of shortest length. Here

$$c = \min(B, \lceil B[1 - \delta + 1.12\sqrt{\delta/B}] \rceil). \quad (4.28)$$

Applied to a bootstrap sample, the Frey shorth interval can be regarded as the shortest percentile method confidence interval, asymptotically. Hence the shorth confidence interval is a practical implementation of the Hall (1988) shortest bootstrap interval based on all possible bootstrap samples. See Remark 4.16 for some theory for bootstrap CIs such as (4.26) and (4.27).

### 4.5.1 The Bootstrap

This subsection illustrates the nonparametric bootstrap with some examples. Suppose a statistic  $T_n$  is computed from a data set of  $n$  cases. The nonparametric bootstrap draws  $n$  cases with replacement from that data set. Then  $T_1^*$  is the statistic  $T_n$  computed from the sample. This process is repeated  $B$  times to produce the bootstrap sample  $T_1^*, \dots, T_B^*$ . Sampling cases with replacement uses the empirical distribution.

**Definition 4.15.** Suppose that data  $\mathbf{x}_1, \dots, \mathbf{x}_n$  has been collected and observed. Often the data is a random sample (iid) from a distribution with cdf  $F$ . The *empirical distribution* is a discrete distribution where the  $\mathbf{x}_i$  are the possible values, and each value is equally likely. If  $\mathbf{w}$  is a random variable having the empirical distribution, then  $p_i = P(\mathbf{w} = \mathbf{x}_i) = 1/n$  for  $i = 1, \dots, n$ . The *cdf of the empirical distribution* is denoted by  $F_n$ .

**Example 4.2.** Let  $\mathbf{w}$  be a random variable having the empirical distribution given by Definition 4.15. Show that  $E(\mathbf{w}) = \bar{\mathbf{x}} \equiv \bar{\mathbf{x}}_n$  and  $\text{Cov}(\mathbf{w}) = \frac{n-1}{n} \mathbf{S} \equiv \frac{n-1}{n} \mathbf{S}_n$ .

Solution: Recall that for a discrete random vector, the population expected value  $E(\mathbf{w}) = \sum \mathbf{x}_i p_i$  where  $\mathbf{x}_i$  are the values that  $\mathbf{w}$  takes with positive probability  $p_i$ . Similarly, the population covariance matrix

$$\text{Cov}(\mathbf{w}) = E[(\mathbf{w} - E(\mathbf{w}))(\mathbf{w} - E(\mathbf{w}))^T] = \sum (\mathbf{x}_i - E(\mathbf{w}))(\mathbf{x}_i - E(\mathbf{w}))^T p_i.$$

Hence

$$E(\mathbf{w}) = \sum_{i=1}^n \mathbf{x}_i \frac{1}{n} = \bar{\mathbf{x}},$$

and

$$\text{Cov}(\mathbf{w}) = \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T \frac{1}{n} = \frac{n-1}{n} \mathbf{S}. \quad \square$$

**Example 4.3.** If  $W_1, \dots, W_n$  are iid from a distribution with cdf  $F_W$ , then the empirical cdf  $F_n$  corresponding to  $F_W$  is given by

$$F_n(y) = \frac{1}{n} \sum_{i=1}^n I(W_i \leq y)$$

where the indicator  $I(W_i \leq y) = 1$  if  $W_i \leq y$  and  $I(W_i \leq y) = 0$  if  $W_i > y$ . Fix  $n$  and  $y$ . Then  $nF_n(y) \sim \text{binomial}(n, F_W(y))$ . Thus  $E[F_n(y)] = F_W(y)$  and  $V[F_n(y)] = F_W(y)[1 - F_W(y)]/n$ . By the central limit theorem,

$$\sqrt{n}(F_n(y) - F_W(y)) \xrightarrow{D} N(0, F_W(y)[1 - F_W(y)]).$$



Thus  $F_n(y) - F_W(y) = O_P(n^{-1/2})$ , and  $F_n$  is a reasonable estimator of  $F_W$  if the sample size  $n$  is large.

Suppose there is data  $\mathbf{w}_1, \dots, \mathbf{w}_n$  collected into an  $n \times p$  matrix  $\mathbf{W}$ . Let the statistic  $T_n = t(\mathbf{W}) = T(F_n)$  be computed from the data. Suppose the statistic estimates  $\boldsymbol{\mu} = T(F)$ , and let  $t(\mathbf{W}^*) = t(F_n^*) = T_n^*$  indicate that  $t$  was computed from an iid sample from the empirical distribution  $F_n$ : a sample  $\mathbf{w}_1^*, \dots, \mathbf{w}_n^*$  of size  $n$  was drawn with replacement from the observed sample  $\mathbf{w}_1, \dots, \mathbf{w}_n$ . This notation is used for von Mises differentiable statistical functions in large sample theory. See Serfling (1980, ch. 6). The empirical distribution is also important for the influence function (widely used in robust statistics). The *nonparametric bootstrap* draws  $B$  samples of size  $n$  from the rows of  $\mathbf{W}$ , e.g. from the empirical distribution of  $\mathbf{w}_1, \dots, \mathbf{w}_n$ . Then  $T_{jn}^*$  is computed from the  $j$ th bootstrap sample for  $j = 1, \dots, B$ .

**Example 4.4.** Suppose the data is 1, 2, 3, 4, 5, 6, 7. Then  $n = 7$  and the sample median  $T_n$  is 4. Using  $R$ , we drew  $B = 2$  bootstrap samples (samples of size  $n$  drawn with replacement from the original data) and computed the sample median  $T_{1,n}^* = 3$  and  $T_{2,n}^* = 4$ .

```
b1 <- sample(1:7, replace=T)
b1
[1] 3 2 3 2 5 2 6
median(b1)
[1] 3
b2 <- sample(1:7, replace=T)
b2
[1] 3 5 3 4 3 5 7
median(b2)
[1] 4
```

The bootstrap has been widely used to estimate the population covariance matrix of the statistic  $\text{Cov}(T_n)$ , for testing hypotheses, and for obtaining confidence regions (often confidence intervals). An iid sample  $T_{1n}, \dots, T_{Bn}$  of size  $B$  of the statistic would be very useful for inference, but typically we only have one sample of data and one value  $T_n = T_{1n}$  of the statistic. Often  $T_n = t(\mathbf{w}_1, \dots, \mathbf{w}_n)$ , and the bootstrap sample  $T_{1n}^*, \dots, T_{Bn}^*$  is formed where  $T_{jn}^* = t(\mathbf{w}_{j1}^*, \dots, \mathbf{w}_{jn}^*)$ . Section 4.5.3 will show that  $T_{1n}^* - T_n, \dots, T_{Bn}^* - T_n$  is pseudodata for  $T_{1n} - \boldsymbol{\theta}, \dots, T_{Bn} - \boldsymbol{\theta}$  when  $n$  is large in that  $\sqrt{n}(T_n - \boldsymbol{\theta}) \xrightarrow{D} \mathbf{u}$  and  $\sqrt{n}(T_n^* - T_n) \xrightarrow{D} \mathbf{u}$ .

**Example 4.5.** Suppose there is training data  $(\mathbf{y}_i, \mathbf{x}_i)$  for the model  $\mathbf{y}_i = m(\mathbf{x}_i) + \boldsymbol{\epsilon}_i$  for  $i = 1, \dots, n$ , and it is desired to predict a future test value  $\mathbf{y}_f$  given  $\mathbf{x}_f$  and the training data. The model can be fit and the residual vectors formed. One method for obtaining a prediction region for  $\mathbf{y}_f$  is to form the pseudodata  $\hat{\mathbf{y}}_f + \hat{\boldsymbol{\epsilon}}_i$  for  $i = 1, \dots, n$ , and apply the nonparametric

prediction region (4.24) to the pseudodata. See Section 8.3 and Olive (2017b, 2018). The residual bootstrap could also be used to make a bootstrap sample  $\hat{\mathbf{y}}_f + \hat{\boldsymbol{\epsilon}}_1^*, \dots, \hat{\mathbf{y}}_f + \hat{\boldsymbol{\epsilon}}_B^*$  where the  $\hat{\boldsymbol{\epsilon}}_j^*$  are selected with replacement from the residual vectors for  $j = 1, \dots, B$ . As  $B \rightarrow \infty$ , the bootstrap sample will take on the  $n$  values  $\hat{\mathbf{y}}_f + \hat{\boldsymbol{\epsilon}}_i$  (the pseudodata) with probabilities converging to  $1/n$  for  $i = 1, \dots, n$ .

Suppose there is a statistic  $T_n$  that is a  $g \times 1$  vector. Let

$$\bar{T}^* = \frac{1}{B} \sum_{i=1}^B T_i^* \quad \text{and} \quad \mathbf{S}_T^* = \frac{1}{B-1} \sum_{i=1}^B (T_i^* - \bar{T}^*)(T_i^* - \bar{T}^*)^T \quad (4.29)$$

be the sample mean and sample covariance matrix of the bootstrap sample  $T_1^*, \dots, T_B^*$  where  $T_i^* = T_{i,n}^*$ . Fix  $n$ , and let  $E(T_{i,n}^*) = \boldsymbol{\theta}_n$  and  $\text{Cov}(T_{i,n}^*) = \boldsymbol{\Sigma}_n$ .

We will often assume that  $\text{Cov}(T_n) = \boldsymbol{\Sigma}_T$ , and  $\sqrt{n}(T_n - \boldsymbol{\theta}) \xrightarrow{D} N_g(\mathbf{0}, \boldsymbol{\Sigma}_A)$  where  $\boldsymbol{\Sigma}_A > 0$  is positive definite and nonsingular. Often  $n\hat{\boldsymbol{\Sigma}}_T \xrightarrow{P} \boldsymbol{\Sigma}_A$ . For example, using least squares and the residual bootstrap for the multiple linear regression model,  $\boldsymbol{\Sigma}_n = \frac{n-p}{n} \text{MSE}(\mathbf{X}^T \mathbf{X})^{-1}$ ,  $T_n = \boldsymbol{\theta}_n = \hat{\boldsymbol{\beta}}$ ,  $\boldsymbol{\theta} = \boldsymbol{\beta}$ ,  $\hat{\boldsymbol{\Sigma}}_T = \text{MSE}(\mathbf{X}^T \mathbf{X})^{-1}$  and  $\boldsymbol{\Sigma}_A = \sigma^2 \lim_{n \rightarrow \infty} (\mathbf{X}^T \mathbf{X}/n)^{-1}$ . See Example 4.6 in Section 4.6.

Suppose the  $T_i^* = T_{i,n}^*$  are iid from some distribution with cdf  $\tilde{F}_n$ . For example, if  $T_{i,n}^* = t(F_n^*)$  where iid samples from  $F_n$  are used, then  $\tilde{F}_n$  is the cdf of  $t(F_n^*)$ . With respect to  $\tilde{F}_n$ , both  $\boldsymbol{\theta}_n$  and  $\boldsymbol{\Sigma}_n$  are parameters, but with respect to  $F$ ,  $\boldsymbol{\theta}_n$  is a random vector and  $\boldsymbol{\Sigma}_n$  is a random matrix. For fixed  $n$ , by the multivariate central limit theorem,

$$\sqrt{B}(\bar{T}^* - \boldsymbol{\theta}_n) \xrightarrow{D} N_g(\mathbf{0}, \boldsymbol{\Sigma}_n) \quad \text{and} \quad \text{B}(\bar{T}^* - \boldsymbol{\theta}_n)^T [\mathbf{S}_T^*]^{-1} (\bar{T}^* - \boldsymbol{\theta}_n) \xrightarrow{D} \chi_r^2$$

as  $B \rightarrow \infty$ .

**Remark 4.13.** For Examples 4.2, 4.5, and 4.6, the bootstrap works but is expensive compared to alternative methods. For Example 4.2, fix  $n$ , then  $\bar{T}^* \xrightarrow{P} \boldsymbol{\theta}_n = \bar{\mathbf{x}}$  and  $\mathbf{S}_T^* \xrightarrow{P} (n-1)\mathbf{S}/n$  as  $B \rightarrow \infty$ , but using  $(\bar{\mathbf{x}}, \mathbf{S})$  makes more sense. For Example 4.5, use the pseudodata instead of the residual bootstrap. For Example 4.6, using  $\hat{\boldsymbol{\beta}}$  and the classical estimated covariance matrix  $\widehat{\text{Cov}}(\hat{\boldsymbol{\beta}}) = \text{MSE}(\mathbf{X}^T \mathbf{X})^{-1}$  makes more sense than using the bootstrap. For these three examples, it is known how the bootstrap sample behaves as  $B \rightarrow \infty$ . The bootstrap can be very useful when  $\sqrt{n}(T_n - \boldsymbol{\theta}) \xrightarrow{D} N_g(\mathbf{0}, \boldsymbol{\Sigma}_A)$ , but it not known how to estimate  $\boldsymbol{\Sigma}_A$  without using a resampling method like the bootstrap. The bootstrap may be useful when  $\sqrt{n}(T_n - \boldsymbol{\theta}) \xrightarrow{D} \mathbf{u}$ , but the limiting distribution (the distribution of  $\mathbf{u}$ ) is unknown.

### 4.5.2 Bootstrap Confidence Regions for Hypothesis Testing

When the bootstrap is used, a large sample  $100(1 - \delta)\%$  confidence region for a  $g \times 1$  parameter vector  $\boldsymbol{\theta}$  is a set  $\mathcal{A}_n = \mathcal{A}_{n,B}$  such that  $P(\boldsymbol{\theta} \in \mathcal{A}_{n,B})$  is eventually bounded below by  $1 - \delta$  as  $n, B \rightarrow \infty$ . The  $B$  is often suppressed. Consider testing  $H_0 : \boldsymbol{\theta} = \boldsymbol{\theta}_0$  versus  $H_1 : \boldsymbol{\theta} \neq \boldsymbol{\theta}_0$  where  $\boldsymbol{\theta}_0$  is a known  $g \times 1$  vector. Then reject  $H_0$  if  $\boldsymbol{\theta}_0$  is not in the confidence region  $\mathcal{A}_n$ . Let the  $g \times 1$  vector  $T_n$  be an estimator of  $\boldsymbol{\theta}$ . Let  $T_1^*, \dots, T_B^*$  be the bootstrap sample for  $T_n$ . Let  $\mathbf{A}$  be a full rank  $g \times p$  constant matrix. For variable selection, consider testing  $H_0 : \mathbf{A}\boldsymbol{\beta} = \boldsymbol{\theta}_0$  versus  $H_1 : \mathbf{A}\boldsymbol{\beta} \neq \boldsymbol{\theta}_0$  with  $\boldsymbol{\theta} = \mathbf{A}\boldsymbol{\beta}$  where often  $\boldsymbol{\theta}_0 = \mathbf{0}$ . Then let  $T_n = \mathbf{A}\hat{\boldsymbol{\beta}}_{I_{min},0}$  and let  $T_i^* = \mathbf{A}\hat{\boldsymbol{\beta}}_{I_{min},0,i}^*$  for  $i = 1, \dots, B$ . The statistic  $\hat{\boldsymbol{\beta}}_{I_{min},0}$  is the variable selection estimator padded with zeroes. See Section 4.2. Let  $\bar{T}^*$  and  $\mathbf{S}_T^*$  be the sample mean and sample covariance matrix of the bootstrap sample  $T_1^*, \dots, T_B^*$ . See Equation (4.28). See Theorem 2.25 for why  $d_n F_{g,d_n,1-\delta} \rightarrow \chi_{g,1-\delta}^2$  as  $d_n \rightarrow \infty$ . Here  $P(X \leq \chi_{g,1-\delta}^2) = 1 - \delta$  if  $X \sim \chi_g^2$ , and  $P(X \leq F_{g,d_n,1-\delta}) = 1 - \delta$  if  $X \sim F_{g,d_n}$ . Let  $k_B = \lceil B(1 - \delta) \rceil$ .

**Definition 4.16.** a) The standard bootstrap large sample  $100(1 - \delta)\%$  confidence region for  $\boldsymbol{\theta}$  is  $\{\mathbf{w} : (\mathbf{w} - T_n)^T [\mathbf{S}_T^*]^{-1} (\mathbf{w} - T_n) \leq D_{1-\delta}^2\} =$

$$\{\mathbf{w} : D_{\mathbf{w}}^2(T_n, \mathbf{S}_T^*) \leq D_{1-\delta}^2\} \quad (4.30)$$

where  $D_{1-\delta}^2 = \chi_{g,1-\delta}^2$  or  $D_{1-\delta}^2 = d_n F_{g,d_n,1-\delta}$  where  $d_n \rightarrow \infty$  as  $n \rightarrow \infty$ . b) The Bickel and Ren (2001) large sample  $100(1 - \delta)\%$  confidence region for  $\boldsymbol{\theta}$  is  $\{\mathbf{w} : (\mathbf{w} - T_n)^T [\hat{\boldsymbol{\Sigma}}_A/n]^{-1} (\mathbf{w} - T_n) \leq D_{(k_B, T)}^2\} =$

$$\{\mathbf{w} : D_{\mathbf{w}}^2(T_n, \hat{\boldsymbol{\Sigma}}_A/n) \leq D_{(k_B, T)}^2\} \quad (4.31)$$

where the cutoff  $D_{(k_B, T)}^2$  is the  $100k_B$ th sample quantile of the

$$D_i^2 = (T_i^* - T_n)^T [\hat{\boldsymbol{\Sigma}}_A/n]^{-1} (T_i^* - T_n) = n(T_i^* - T_n)^T [\hat{\boldsymbol{\Sigma}}_A]^{-1} (T_i^* - T_n).$$

Confidence region (4.29) needs  $\sqrt{n}(T_n - \boldsymbol{\theta}) \xrightarrow{D} N_g(\mathbf{0}, \boldsymbol{\Sigma}_A)$  and  $n\mathbf{S}_T^* \xrightarrow{P} \boldsymbol{\Sigma}_A > 0$  as  $n, B \rightarrow \infty$ . See Machado and Parente (2005) for regularity conditions for this assumption. Bickel and Ren (2001) have interesting sufficient conditions for (4.30) to be a confidence region when  $\hat{\boldsymbol{\Sigma}}_A$  is a consistent estimator of positive definite  $\boldsymbol{\Sigma}_A$ . Let the vector of parameters  $\boldsymbol{\theta} = T(F)$ , the statistic  $T_n = T(F_n)$ , and the bootstrapped statistic  $T^* = T(F_n^*)$  where  $F$  is the cdf of iid  $\mathbf{x}_1, \dots, \mathbf{x}_n$ ,  $F_n$  is the empirical cdf, and  $F_n^*$  is the empirical cdf of  $\mathbf{x}_1^*, \dots, \mathbf{x}_n^*$ , a sample from  $F_n$  using the nonparametric bootstrap. If  $\sqrt{n}(F_n - F) \xrightarrow{D} \mathbf{z}_F$ , a Gaussian random process, and if  $T$  is sufficiently smooth (has a Hadamard derivative  $\dot{T}(F)$ ), then  $\sqrt{n}(T_n - \boldsymbol{\theta}) \xrightarrow{D} \mathbf{u}$  and

$\sqrt{n}(T_i^* - T_n) \xrightarrow{D} \mathbf{u}$  with  $\mathbf{u} = \dot{T}(F)\mathbf{z}_F$ . Note that  $F_n$  is a perfectly good cdf “ $F$ ” and  $F_n^*$  is a perfectly good empirical cdf from  $F_n = “F.”$  Thus if  $n$  is fixed, and a sample of size  $m$  is drawn with replacement from the empirical distribution, then  $\sqrt{m}(T(F_m^*) - T_n) \xrightarrow{D} \dot{T}(F_n)\mathbf{z}_{F_n}$ . Now let  $n \rightarrow \infty$  with  $m = n$ . Then bootstrap theory gives  $\sqrt{n}(T_i^* - T_n) \xrightarrow{D} \lim_{n \rightarrow \infty} \dot{T}(F_n)\mathbf{z}_{F_n} = \dot{T}(F)\mathbf{z}_F \sim \mathbf{u}$ .

The following three confidence regions will be used for inference after variable selection. The Olive (2017ab, 2018) prediction region method applies prediction region (4.24) to the bootstrap sample. Olive (2017ab, 2018) also gave the modified Bickel and Ren confidence region that uses  $\hat{\Sigma}_A = n\mathbf{S}_T^*$ . The hybrid confidence region is due to Pelawa Watagoda and Olive (2019a). Let  $q_B = \min(1 - \delta + 0.05, 1 - \delta + g/B)$  for  $\delta > 0.1$  and

$$q_B = \min(1 - \delta/2, 1 - \delta + 10\delta g/B), \quad \text{otherwise.} \quad (4.32)$$

If  $1 - \delta < 0.999$  and  $q_B < 1 - \delta + 0.001$ , set  $q_B = 1 - \delta$ . Let  $D_{(U_B)}$  be the  $100q_B$ th sample quantile of the  $D_i$ . Use (4.31) as a correction factor for finite  $B \geq 50p$ .

**Definition 4.17.** a) The prediction region method large sample  $100(1 - \delta)\%$  confidence region for  $\boldsymbol{\theta}$  is  $\{\mathbf{w} : (\mathbf{w} - \bar{\mathbf{T}}^*)^T [\mathbf{S}_T^*]^{-1} (\mathbf{w} - \bar{\mathbf{T}}^*) \leq D_{(U_B)}^2\} =$

$$\{\mathbf{w} : D_{\mathbf{w}}^2(\bar{\mathbf{T}}^*, \mathbf{S}_T^*) \leq D_{(U_B)}^2\} \quad (4.33)$$

where  $D_{(U_B)}^2$  is computed from  $D_i^2 = (T_i^* - \bar{\mathbf{T}}^*)^T [\mathbf{S}_T^*]^{-1} (T_i^* - \bar{\mathbf{T}}^*)$  for  $i = 1, \dots, B$ . Note that the corresponding test for  $H_0 : \boldsymbol{\theta} = \boldsymbol{\theta}_0$  rejects  $H_0$  if  $(\bar{\mathbf{T}}^* - \boldsymbol{\theta}_0)^T [\mathbf{S}_T^*]^{-1} (\bar{\mathbf{T}}^* - \boldsymbol{\theta}_0) > D_{(U_B)}^2$ . (This procedure is basically the one sample Hotelling’s  $T^2$  test applied to the  $T_i^*$  using  $\mathbf{S}_T^*$  as the estimated covariance matrix and replacing the  $\chi_{g,1-\delta}^2$  cutoff by  $D_{(U_B)}^2$ .) b) The modified Bickel and Ren (2001) large sample  $100(1 - \delta)\%$  confidence region is  $\{\mathbf{w} : (\mathbf{w} - T_n)^T [\mathbf{S}_T^*]^{-1} (\mathbf{w} - T_n) \leq D_{(U_B, T)}^2\} =$

$$\{\mathbf{w} : D_{\mathbf{w}}^2(T_n, \mathbf{S}_T^*) \leq D_{(U_B, T)}^2\} \quad (4.34)$$

where the cutoff  $D_{(U_B, T)}^2$  is the  $100q_B$ th sample quantile of the  $D_i^2 = (T_i^* - T_n)^T [\mathbf{S}_T^*]^{-1} (T_i^* - T_n)$ . Note that the corresponding test for  $H_0 : \boldsymbol{\theta} = \boldsymbol{\theta}_0$  rejects  $H_0$  if  $(T_n - \boldsymbol{\theta}_0)^T [\mathbf{S}_T^*]^{-1} (T_n - \boldsymbol{\theta}_0) > D_{(U_B, T)}^2$ . c) Shift region (4.32) to have center  $T_n$ , or equivalently, change the cutoff of region (4.33) to  $D_{(U_B)}^2$  to get the hybrid large sample  $100(1 - \delta)\%$  confidence region:  $\{\mathbf{w} : (\mathbf{w} - T_n)^T [\mathbf{S}_T^*]^{-1} (\mathbf{w} - T_n) \leq D_{(U_B)}^2\} =$

$$\{\mathbf{w} : D_{\mathbf{w}}^2(T_n, \mathbf{S}_T^*) \leq D_{(U_B)}^2\}. \quad (4.35)$$

Note that the corresponding test for  $H_0 : \boldsymbol{\theta} = \boldsymbol{\theta}_0$  rejects  $H_0$  if  $(T_n - \boldsymbol{\theta}_0)^T [\mathbf{S}_T^*]^{-1} (T_n - \boldsymbol{\theta}_0) > D_{(U_B)}^2$ .

Hyperellipsoids (4.32) and (4.34) have the same volume since they are the same region shifted to have a different center. The ratio of the volumes of regions (4.32) and (4.33) is

$$\frac{|\mathbf{S}_T^*|^{1/2}}{|\mathbf{S}_T^*|^{1/2}} \left( \frac{D_{(U_B)}}{D_{(U_B, T)}} \right)^g = \left( \frac{D_{(U_B)}}{D_{(U_B, T)}} \right)^g. \quad (4.36)$$

The volume of confidence region (4.33) tends to be greater than that of (4.32) since the  $T_i^*$  are closer to  $\bar{T}^*$  than  $T_n$  on average.

If  $g = 1$ , then a hyperellipsoid is an interval, and confidence intervals are special cases of confidence regions. Suppose the parameter of interest is  $\theta$ , and there is a bootstrap sample  $T_1^*, \dots, T_B^*$  where the statistic  $T_n$  is an estimator of  $\theta$  based on a sample of size  $n$ . The percentile method uses an interval that contains  $U_B \approx k_B = \lceil B(1-\delta) \rceil$  of the  $T_i^*$ . Let  $a_i = |T_i^* - \bar{T}^*|$ . Let  $\bar{T}^*$  and  $S_T^{2*}$  be the sample mean and variance of the  $T_i^*$ . Then the squared Mahalanobis distance  $D_\theta^2 = (\theta - \bar{T}^*)^2 / S_T^{2*} \leq D_{(U_B)}^2$  is equivalent to  $\theta \in [\bar{T}^* - S_T^* D_{(U_B)}, \bar{T}^* + S_T^* D_{(U_B)}] = [\bar{T}^* - a_{(U_B)}, \bar{T}^* + a_{(U_B)}]$ , which is an interval centered at  $\bar{T}^*$  just long enough to cover  $U_B$  of the  $T_i^*$ . Hence the prediction region method is a special case of the percentile method if  $g = 1$ . See Definition 4.13. Efron (2014) used a similar large sample  $100(1-\delta)\%$  confidence interval assuming that  $\bar{T}^*$  is asymptotically normal. The CI corresponding to (4.33) is defined similarly, and  $[T_n - a_{(U_B)}, T_n + a_{(U_B)}]$  is the CI for (4.34). Note that the three CIs corresponding to (4.32)–(4.34) can be computed without finding  $S_T^*$  or  $D_{(U_B)}$  even if  $S_T^* = 0$ . The Frey (2013) shorth( $c$ ) CI (4.27) computed from the  $T_i^*$  can be much shorter than the Efron (2014) or prediction region method confidence intervals. See Remark 4.16 for some theory for bootstrap CIs.

**Remark 4.14.** From Example 4.6,  $\text{Cov}(\hat{\beta}^*) = \frac{n-p}{n} \text{MSE}(\mathbf{X}^T \mathbf{X})^{-1} = \frac{n-p}{n} \widehat{\text{Cov}}(\hat{\beta})$  where  $\widehat{\text{Cov}}(\hat{\beta}) = \text{MSE}(\mathbf{X}^T \mathbf{X})^{-1}$  starts to give good estimates of  $\text{Cov}(\hat{\beta}) = \boldsymbol{\Sigma}_T$  for many error distributions if  $n \geq 10p$  and  $T = \hat{\beta}$ . For the residual bootstrap with large  $B$ , note that  $\mathbf{S}_T^* \approx 0.95 \widehat{\text{Cov}}(\hat{\beta})$  for  $n = 20p$  and  $\mathbf{S}_T^* \approx 0.99 \widehat{\text{Cov}}(\hat{\beta})$  for  $n = 100p$ . Hence we may need  $n \gg p$  before the  $\mathbf{S}_T^*$  is a good estimator of  $\text{Cov}(T) = \boldsymbol{\Sigma}_T$ . The distribution of  $\sqrt{n}(T_n - \theta)$  is approximated by the distribution of  $\sqrt{n}(T^* - T_n)$  or by the distribution of  $\sqrt{n}(T^* - \bar{T}^*)$ , but  $n$  may need to be large before the approximation is good.

Suppose the bootstrap sample mean  $\bar{T}^*$  estimates  $\theta$ , and the bootstrap sample covariance matrix  $\mathbf{S}_T^*$  estimates  $c_n \widehat{\text{Cov}}(T_n) \approx c_n \boldsymbol{\Sigma}_T$  where  $c_n$  increases to 1 as  $n \rightarrow \infty$ . Then  $\mathbf{S}_T^*$  is not a good estimator of  $\widehat{\text{Cov}}(T_n)$  until  $c_n \approx 1$  ( $n \geq 100p$  for OLS  $\hat{\beta}$ ), but the squared Mahalanobis distance  $D_{\mathbf{w}}^{2*}(\bar{T}^*, \mathbf{S}_T^*) \approx D_{\mathbf{w}}^2(\theta, \boldsymbol{\Sigma}_T) / c_n$  and  $D_{(U_B)}^{2*} \approx D_{1-\delta}^2 / c_n$ . Hence the prediction region method has a cutoff  $D_{(U_B)}^{2*}$  that estimates the cutoff  $D_{1-\delta}^2 / c_n$ . Thus the prediction region method may give good results for much smaller  $n$  than

a bootstrap method that uses a  $\chi_{g,1-\delta}^2$  cutoff when a cutoff  $\chi_{g,1-\delta}^2/c_n$  should be used for moderate  $n$ .

**Remark 4.15.** For bootstrapping the  $p \times 1$  vector  $\hat{\beta}_{I_{min},0}$ , we will often want  $n \geq 20p$  and  $B \geq \max(100, n, 50p)$ . If  $T_n$  is  $g \times 1$ , we might replace  $p$  by  $g$  or replace  $p$  by  $d$  if  $d$  is the model degrees of freedom. Sometimes much larger  $n$  is needed to avoid undercoverage. We want  $B \geq 50g$  so that  $\mathbf{S}_T^*$  is a good estimator of  $Cov(T_n^*)$ . Prediction region theory uses correction factors like (4.21) and (4.10) to compensate for finite  $n$ . The bootstrap confidence regions (4.32)–(4.34) and the shorth CI use the correction factors (4.31) and (4.27) to compensate for finite  $B \geq 50g$ . Note that the correction factors make the volume of the confidence region larger as  $B$  decreases. Hence a test with larger  $B$  will have more power.

### 4.5.3 Theory for Bootstrap Confidence Regions

Consider testing  $H_0 : \boldsymbol{\theta} = \boldsymbol{\theta}_0$  versus  $H_1 : \boldsymbol{\theta} \neq \boldsymbol{\theta}_0$  where  $\boldsymbol{\theta}$  is  $g \times 1$ . This section gives some theory for bootstrap confidence regions and for the bagging estimator  $\bar{T}^*$ , also called the smoothed bootstrap estimator. Empirically, bootstrapping with the bagging estimator often outperforms bootstrapping with  $T_n$ . See Breiman (1996), Yang (2003), and Efron (2014). See Büchlmann and Yu (2002) and Friedman and Hall (2007) for theory and references for the bagging estimator. Since (4.33) is a large sample confidence region by Bickel and Ren (2001), (4.32) and (4.34) are too, provided  $\sqrt{n}(\bar{T}^* - T_n) \xrightarrow{P} \mathbf{0}$ .

If i)  $\sqrt{n}(T_n - \boldsymbol{\theta}) \xrightarrow{D} \mathbf{u}$ , then under regularity conditions, ii)  $\sqrt{n}(T_i^* - T_n) \xrightarrow{D} \mathbf{u}$ , iii)  $\sqrt{n}(\bar{T}^* - \boldsymbol{\theta}) \xrightarrow{D} \mathbf{u}$ , iv)  $\sqrt{n}(T_i^* - \bar{T}^*) \xrightarrow{D} \mathbf{u}$ , and v)  $n\mathbf{S}_T^* \xrightarrow{P} Cov(\mathbf{u})$ .

Suppose i) and ii) hold with  $E(\mathbf{u}) = \mathbf{0}$  and  $Cov(\mathbf{u}) = \boldsymbol{\Sigma}\mathbf{u}$ . With respect to the bootstrap sample,  $T_n$  is a constant and the  $\sqrt{n}(T_i^* - T_n)$  are iid for  $i = 1, \dots, B$ . Let  $\sqrt{n}(T_i^* - T_n) \xrightarrow{D} \mathbf{v}_i \sim \mathbf{u}$  where the  $\mathbf{v}_i$  are iid with the same distribution as  $\mathbf{u}$ . Fix  $B$ . Then the average of the  $\sqrt{n}(T_i^* - T_n)$  is

$$\sqrt{n}(\bar{T}^* - T_n) \xrightarrow{D} \frac{1}{B} \sum_{i=1}^B \mathbf{v}_i \sim AN_g \left( \mathbf{0}, \frac{\boldsymbol{\Sigma}\mathbf{u}}{B} \right)$$

where  $\mathbf{z} \sim AN_g(\mathbf{0}, \boldsymbol{\Sigma})$  is an asymptotic multivariate normal approximation. Hence as  $B \rightarrow \infty$ ,  $\sqrt{n}(\bar{T}^* - T_n) \xrightarrow{P} \mathbf{0}$ , and iii) and iv) hold. If  $B$  is fixed and  $\mathbf{u} \sim N_g(\mathbf{0}, \boldsymbol{\Sigma}\mathbf{u})$ , then

$$\frac{1}{B} \sum_{i=1}^B \mathbf{v}_i \sim N_g \left( \mathbf{0}, \frac{\boldsymbol{\Sigma}\mathbf{u}}{B} \right) \text{ and } \sqrt{B}\sqrt{n}(\bar{T}^* - T_n) \xrightarrow{D} N_g(\mathbf{0}, \boldsymbol{\Sigma}\mathbf{u}).$$

Hence the prediction region method gives a large sample confidence region for  $\boldsymbol{\theta}$  provided that the sample percentile  $\hat{D}_{1-\delta}^2$  of the  $D_{T_i^*}^2(\bar{T}^*, \mathbf{S}_T^*) = \sqrt{n}(T_i^* - \bar{T}^*)^T (n\mathbf{S}_T^*)^{-1} \sqrt{n}(T_i^* - \bar{T}^*)$  is a consistent estimator of the percentile  $D_{n,1-\delta}^2$  of the random variable  $D_{\boldsymbol{\theta}}^2(\bar{T}^*, \mathbf{S}_T^*) = \sqrt{n}(\boldsymbol{\theta} - \bar{T}^*)^T (n\mathbf{S}_T^*)^{-1} \sqrt{n}(\boldsymbol{\theta} - \bar{T}^*)$  in that  $\hat{D}_{1-\delta}^2 - D_{n,1-\delta}^2 \xrightarrow{P} 0$ . Since iii) and iv) hold, the sample percentile will be consistent under much weaker conditions than v) if  $\boldsymbol{\Sigma}_{\mathbf{u}}$  is nonsingular. Olive (2017b: § 5.3.3, 2018) proved that the prediction region method gives a large sample confidence region under the much stronger conditions of v) and  $\mathbf{u} \sim N_g(\mathbf{0}, \boldsymbol{\Sigma}_{\mathbf{u}})$ , but the above Pelawa Watagoda and Olive (2019a) proof is simpler.

**Remark 4.16.** Note that if  $\sqrt{n}(T_n - \boldsymbol{\theta}) \xrightarrow{D} U$  and  $\sqrt{n}(T_i^* - T_n) \xrightarrow{D} U$  where  $U$  has a unimodal probability density function symmetric about zero, then the confidence intervals from the three confidence regions (4.32)–(4.34), the shorth confidence interval (4.27), and the “usual” percentile method confidence interval (4.26) are asymptotically equivalent (use the central proportion of the bootstrap sample, asymptotically).

Assume  $n\mathbf{S}_T^* \xrightarrow{P} \boldsymbol{\Sigma}_A$  as  $n, B \rightarrow \infty$  where  $\boldsymbol{\Sigma}_A$  and  $\mathbf{S}_T^*$  are nonsingular  $g \times g$  matrices, and  $T_n$  is an estimator of  $\boldsymbol{\theta}$  such that

$$\sqrt{n} (T_n - \boldsymbol{\theta}) \xrightarrow{D} \mathbf{u} \quad (4.37)$$

as  $n \rightarrow \infty$ . Then

$$\sqrt{n} \boldsymbol{\Sigma}_A^{-1/2} (T_n - \boldsymbol{\theta}) \xrightarrow{D} \boldsymbol{\Sigma}_A^{-1/2} \mathbf{u} = \mathbf{z},$$

$$n (T_n - \boldsymbol{\theta})^T \hat{\boldsymbol{\Sigma}}_A^{-1} (T_n - \boldsymbol{\theta}) \xrightarrow{D} \mathbf{z}^T \mathbf{z} = D^2$$

as  $n \rightarrow \infty$  where  $\hat{\boldsymbol{\Sigma}}_A$  is a consistent estimator of  $\boldsymbol{\Sigma}_A$ , and

$$(T_n - \boldsymbol{\theta})^T [\mathbf{S}_T^*]^{-1} (T_n - \boldsymbol{\theta}) \xrightarrow{D} D^2 \quad (4.38)$$

as  $n, B \rightarrow \infty$ . Assume the cumulative distribution function of  $D^2$  is continuous and increasing in a neighborhood of  $D_{1-\delta}^2$  where  $P(D^2 \leq D_{1-\delta}^2) = 1 - \delta$ . If the distribution of  $D^2$  is known, then we could use the large sample confidence region (4.29)  $\{\mathbf{w} : (\mathbf{w} - T_n)^T [\mathbf{S}_T^*]^{-1} (\mathbf{w} - T_n) \leq D_{1-\delta}^2\}$ . Often by a central limit theorem or the multivariate delta method,  $\sqrt{n}(T_n - \boldsymbol{\theta}) \xrightarrow{D} N_g(\mathbf{0}, \boldsymbol{\Sigma}_A)$ , and  $D^2 \sim \chi_g^2$ . Note that  $[\mathbf{S}_T^*]^{-1}$  could be replaced by  $n\hat{\boldsymbol{\Sigma}}_A^{-1}$ .

**Remark 4.17.** Under reasonable conditions, i)  $\sqrt{n}(T_n - \boldsymbol{\theta}) \xrightarrow{D} \mathbf{u}$ , ii)  $\sqrt{n}(T_i^* - T_n) \xrightarrow{D} \mathbf{u}$ , iii)  $\sqrt{n}(\bar{T}^* - \boldsymbol{\theta}) \xrightarrow{D} \mathbf{u}$ , and iv)  $\sqrt{n}(T_i^* - \bar{T}^*) \xrightarrow{D} \mathbf{u}$ . Then

$$D_1^2 = D_{T_i^*}^2(\bar{T}^*, \mathbf{S}_T^*) = \sqrt{n}(T_i^* - \bar{T}^*)^T (n\mathbf{S}_T^*)^{-1} \sqrt{n}(T_i^* - \bar{T}^*),$$

$$\begin{aligned}
D_2^2 &= D_{\boldsymbol{\theta}}^2(T_n, \mathbf{S}_T^*) = \sqrt{n}(T_n - \boldsymbol{\theta})^T (n\mathbf{S}_T^*)^{-1} \sqrt{n}(T_n - \boldsymbol{\theta}), \\
D_3^2 &= D_{\boldsymbol{\theta}}^2(\bar{T}^*, \mathbf{S}_T^*) = \sqrt{n}(\bar{T}^* - \boldsymbol{\theta})^T (n\mathbf{S}_T^*)^{-1} \sqrt{n}(\bar{T}^* - \boldsymbol{\theta}), \quad \text{and} \\
D_4^2 &= D_{T_i^*}^2(T_n, \mathbf{S}_T^*) = \sqrt{n}(T_i^* - T_n)^T (n\mathbf{S}_T^*)^{-1} \sqrt{n}(T_i^* - T_n),
\end{aligned}$$

are well behaved. If  $(n\mathbf{S}_T^*)^{-1} \xrightarrow{P} \boldsymbol{\Sigma}_T^{-1}$ , then  $D_j^2 \xrightarrow{D} D^2 = \mathbf{u}^T \boldsymbol{\Sigma}_T^{-1} \mathbf{u}$ . If  $(n\mathbf{S}_T^*)^{-1}$  is “not too ill conditioned” then  $D_j^2 \approx \mathbf{u}^T (n\mathbf{S}_T^*)^{-1} \mathbf{u}$  for large  $n$ , and the confidence regions (4.32), (4.33), and (4.34) will have coverage near  $1 - \delta$ . The regularity conditions for (4.32)–(4.34) are weaker when  $g = 1$ , since  $\mathbf{S}_T^*$  does not need to be computed.

The following Pelawa Watagoda and Olive (2019a) theorem is very useful. Let  $D_{(U_B)}^2$  be the cutoff for the nonparametric prediction region (4.24) computed from the  $D_i^2(\bar{T}, \mathbf{S}_T)$  for  $i = 1, \dots, B$ . Hence  $n$  is replaced by  $B$ . Since  $T_n$  depends on the sample size  $n$ , we need  $(n\mathbf{S}_T)^{-1}$  to be fairly well behaved (“not too ill conditioned”) for each  $n \geq 20g$ , say. This condition is weaker than  $(n\mathbf{S}_T)^{-1} \xrightarrow{P} \boldsymbol{\Sigma}_A^{-1}$ . Note that  $T_i = T_{in}$ .

**Theorem 4.7: Geometric Argument.** Suppose  $\sqrt{n}(T_n - \boldsymbol{\theta}) \xrightarrow{D} \mathbf{u}$  with  $E(\mathbf{u}) = \mathbf{0}$  and  $Cov(\mathbf{u}) = \boldsymbol{\Sigma}\mathbf{u}$ . Assume  $T_1, \dots, T_B$  are iid with nonsingular covariance matrix  $\boldsymbol{\Sigma}_{T_n}$ . Then the large sample  $100(1 - \delta)\%$  prediction region  $R_p = \{\mathbf{w} : D_{\mathbf{w}}^2(\bar{T}, \mathbf{S}_T) \leq D_{(U_B)}^2\}$  centered at  $\bar{T}$  contains a future value of the statistic  $T_f$  with probability  $1 - \delta_B \rightarrow 1 - \delta$  as  $B \rightarrow \infty$ . Hence the region  $R_c = \{\mathbf{w} : D_{\mathbf{w}}^2(T_n, \mathbf{S}_T) \leq D_{(U_B)}^2\}$  is a large sample  $100(1 - \delta)\%$  confidence region for  $\boldsymbol{\theta}$  where  $T_n$  is a randomly selected  $T_i$ .

**Proof.** The region  $R_c$  centered at a randomly selected  $T_n$  contains  $\bar{T}$  with probability  $1 - \delta_B$  which is eventually bounded below by  $1 - \delta$  as  $B \rightarrow \infty$ . Since the  $\sqrt{n}(T_i - \boldsymbol{\theta})$  are iid,

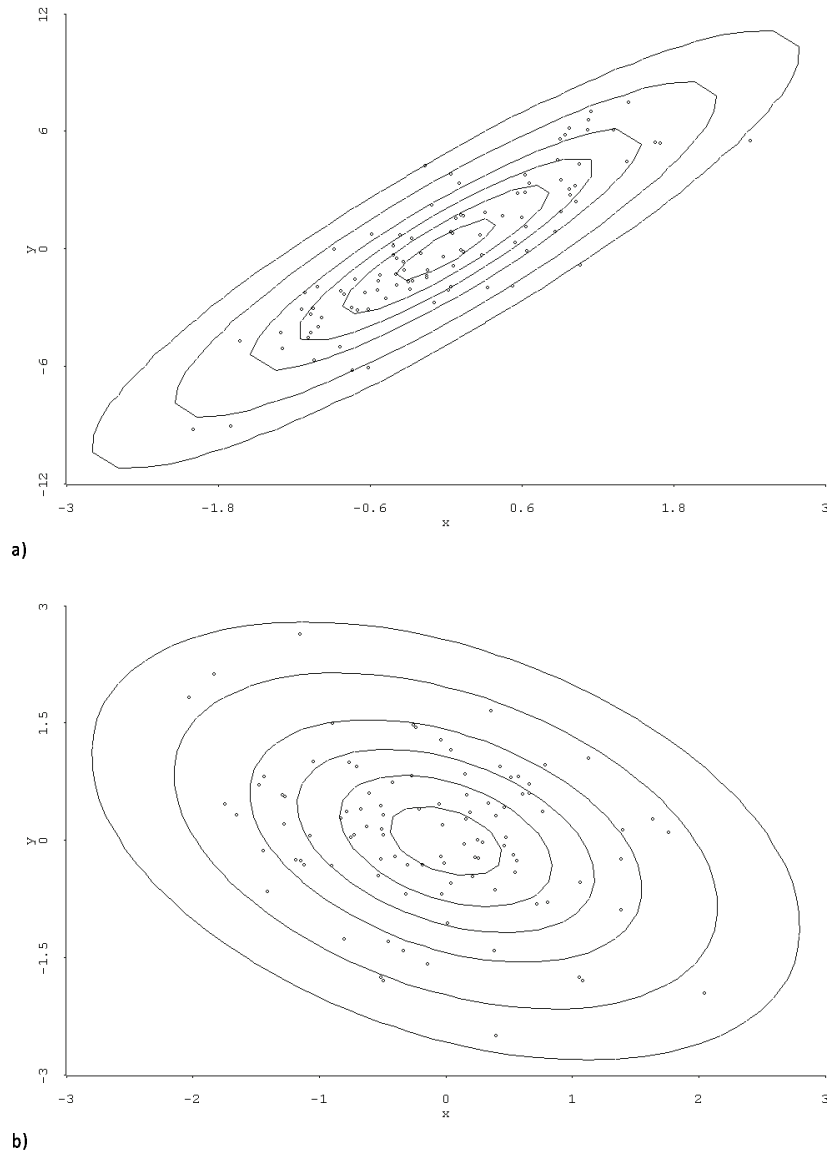
$$\begin{bmatrix} \sqrt{n}(T_1 - \boldsymbol{\theta}) \\ \vdots \\ \sqrt{n}(T_B - \boldsymbol{\theta}) \end{bmatrix} \xrightarrow{D} \begin{bmatrix} \mathbf{v}_1 \\ \vdots \\ \mathbf{v}_B \end{bmatrix}$$

where the  $\mathbf{v}_i$  are iid with the same distribution as  $\mathbf{u}$ . (Use Theorems 1.30 and 1.31, and see Example 1.16.) For fixed  $B$ , the average of these random vectors is

$$\sqrt{n}(\bar{T} - \boldsymbol{\theta}) \xrightarrow{D} \frac{1}{B} \sum_{i=1}^B \mathbf{v}_i \sim AN_g \left( \mathbf{0}, \frac{\boldsymbol{\Sigma}\mathbf{u}}{B} \right)$$

by Theorem 1.33. Hence  $(\bar{T} - \boldsymbol{\theta}) = O_P((nB)^{-1/2})$ , and  $\bar{T}$  gets arbitrarily close to  $\boldsymbol{\theta}$  compared to  $T_n$  as  $B \rightarrow \infty$ . Thus  $R_c$  is a large sample  $100(1 - \delta)\%$  confidence region for  $\boldsymbol{\theta}$  as  $n, B \rightarrow \infty$ .  $\square$





**Fig. 4.3** Confidence Regions for 2 Statistics with MVN Distributions

Examining the iid data cloud  $T_1, \dots, T_B$  and the bootstrap sample data cloud  $T_1^*, \dots, T_B^*$  is often useful for understanding the bootstrap. If  $\sqrt{n}(T_n - \theta)$  and  $\sqrt{n}(T_i^* - T_n)$  both converge in distribution to  $\mathbf{u}$ , then the bootstrap sample data cloud of  $T_1^*, \dots, T_B^*$  is like the data cloud of iid  $T_1, \dots, T_B$  shifted to be centered at  $T_n$ . The nonparametric confidence region (4.32) applies the prediction region to the bootstrap. Then the hybrid region (4.34) centers that region at  $T_n$ . Hence (4.34) is a confidence region by the geometric argument, and (4.32) is a confidence region if  $\sqrt{n}(\bar{T}^* - T_n) \xrightarrow{P} \mathbf{0}$ . Since the  $T_i^*$  are closer to  $\bar{T}^*$  than  $T_n$  on average,  $D_{(U_B, T)}^2$  tends to be greater than  $D_{(U_B)}^2$ . Hence the coverage and volume of (4.33) tend to be at least as large as the coverage and volume of (4.34).

The hyperellipsoid corresponding to the squared Mahalanobis distance  $D^2(T_n, \mathbf{C})$  is centered at  $T_n$ , while the hyperellipsoid corresponding to the squared Mahalanobis distance  $D^2(\bar{T}, \mathbf{C})$  is centered at  $\bar{T}$ . Note that  $D_{\bar{T}}^2(T_n, \mathbf{C}) = (\bar{T} - T_n)^T \mathbf{C}^{-1} (\bar{T} - T_n) = (T_n - \bar{T})^T \mathbf{C}^{-1} (T_n - \bar{T}) = D_{T_n}^2(\bar{T}, \mathbf{C})$ . Thus  $D_{\bar{T}}^2(T_n, \mathbf{C}) \leq D_{(U_B)}^2$  iff  $D_{T_n}^2(\bar{T}, \mathbf{C}) \leq D_{(U_B)}^2$ .

The prediction region method will often simulate well even if  $B$  is rather small. If the ellipses are centered at  $T_n$  or  $\bar{T}^*$ , Figure 4.3 shows confidence regions if the plotted points are  $T_1^*, \dots, T_B^*$  where the  $T_i^*$  are approximately multivariate normal. If the ellipses are centered at  $\bar{T}$ , Figure 4.3 shows 10%, 30%, 50%, 70%, 90%, and 98% prediction regions for a future value of  $T_f$  for two multivariate normal statistics. Then the plotted points are iid  $T_1, \dots, T_B$ . If  $n \text{Cov}(T) \xrightarrow{P} \Sigma_A$ , and the  $T_i^*$  are iid from the bootstrap distribution, then  $\text{Cov}(\bar{T}^*) \approx \text{Cov}(T)/B \approx \Sigma_A/(nB)$ . By Theorem 4.7, if  $\bar{T}^*$  is in the 90% prediction region with probability near 90%, then the confidence region should give simulated coverage near 90% and the volume of the confidence region should be near that of the 90% prediction region. If  $B = 100$ , then  $\bar{T}^*$  falls in a covering region of the same shape as the prediction region, but centered near  $T_n$  and the lengths of the axes are divided by  $\sqrt{B}$ . Hence if  $B = 100$ , then the axes lengths of this covering region are about one tenth of those in Figure 4.3. Hence when  $T_n$  falls within the 70% prediction region, the probability that  $\bar{T}^*$  falls in the 90% prediction region is near one. If  $T_n$  is just within or just without the boundary of the 90% prediction region,  $\bar{T}^*$  tends to be just within or just without of the 90% prediction region. Hence the coverage and volume of prediction region confidence region is near that of the nominal coverage 90% and near the volume of the 90% prediction region.

Hence  $B$  does not need to be large provided that  $n$  and  $B$  are large enough so that  $S_T^* \approx \text{Cov}(T^*) \approx \Sigma_A/n$ . If  $n$  is large, the sample covariance matrix starts to be a good estimator of the population covariance matrix when  $B \geq Jg$  where  $J = 20$  or  $50$ . For small  $g$ , using  $B = 1000$  often led to good simulations, but  $B = \max(50g, 100)$  may work well.

**Remark 4.18.** Remark 4.14 suggests that even if the statistic  $T_n$  is asymptotically normal so the Mahalanobis distances are asymptotically  $\chi_g^2$ , the pre-

diction region method can give better results for moderate  $n$  by using the cutoff  $D_{(U_B)}^2$  instead of the cutoff  $\chi_{g,1-\delta}^2$ . Theorem 4.7 says that the hyperellipsoidal prediction and confidence regions have exactly the same volume. We compensate for the prediction region undercoverage when  $n$  is moderate by using  $D_{(U_n)}^2$ . If  $n$  is large, by using  $D_{(U_B)}^2$ , the prediction region confidence region compensates for undercoverage when  $B$  is moderate, say  $B \geq Jg$  where  $J = 20$  or  $50$ . See Remark 4.15. This result can be useful if a simulation with  $B = 1000$  or  $B = 10000$  is much slower than a simulation with  $B = Jg$ . The price to pay is that the prediction region confidence region is inflated to have better coverage, so the power of the hypothesis test is decreased if moderate  $B$  is used instead of larger  $B$ .

#### 4.5.4 Bootstrapping the Population Coefficient of Multiple Determination

This subsection illustrates a case where the shorth( $c$ ) bootstrap CI fails, but the lower shorth CI can be useful. See Definition 4.14.

The multiple linear regression (MLR) model is

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

for  $i = 1, \dots, n$ . See Definition 1.17 for the *coefficient of multiple determination*

$$R^2 = [\text{corr}(Y_i, \hat{Y}_i)]^2 = \frac{\text{SSR}}{\text{SSTO}} = 1 - \frac{\text{SSE}}{\text{SSTO}}$$

where  $\text{corr}(Y_i, \hat{Y}_i)$  is the sample correlation of  $Y_i$  and  $\hat{Y}_i$ .

Assume that the variance of the errors is  $\sigma_e^2$  and that the variance of  $Y$  is  $\sigma_Y^2$ . Let the linear combination  $L = \sum_{i=2}^p x_i \beta_i$  where  $Y = \beta_1 + \sum_{i=2}^p x_i \beta_i + e = \beta_1 + L + e$ . Let the variance of  $L$  be  $\sigma_L^2$ . Then

$$R^2 = 1 - \frac{\sum_{i=1}^n r_i^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} \xrightarrow{P} \tau^2 = 1 - \frac{\sigma_e^2}{\sigma_Y^2} = 1 - \frac{\sigma_e^2}{\sigma_e^2 + \sigma_L^2}.$$

Here we assume that  $e$  is independent of the predictors  $x_2, \dots, x_p$ . Hence  $e$  is independent of  $L$  and the variance  $\sigma_Y^2 = V(L + e) = V(L) + V(e) = \sigma_L^2 + \sigma_e^2$ .

One of the sufficient conditions for the shorth( $c$ ) interval to be a large sample CI for  $\theta$  is  $\sqrt{n}(T - \theta) \xrightarrow{D} N(0, \sigma^2)$ . If the function  $t(\theta)$  has an inverse, and  $\sqrt{n}(t(T) - t(\theta)) \xrightarrow{D} N(0, v^2)$ , then the above condition typically holds by the delta method. See Remark 4.16.

For  $T = R^2$  and  $\theta = \tau^2$ , the test statistic  $F_0$  for testing  $H_0 : \beta_2 = \cdots = \beta_p = 0$  in the Anova  $F$  test has  $(p-1)F_0 \xrightarrow{D} \chi_{p-1}^2$  for a large class of error distributions when  $H_0$  is true, where

$$F_0 = \frac{R^2}{1 - R^2} \frac{n - p}{p - 1}$$

if the MLR model has a constant. If  $H_0$  is false, then  $F_0$  has an asymptotic scaled noncentral  $\chi^2$  distribution. These results suggest that the large sample distribution of  $\sqrt{n}(R^2 - \tau^2)$  may not be  $N(0, \sigma^2)$  if  $H_0$  is false so  $\tau^2 > 0$ . If  $\tau^2 = 0$ , we may have  $\sqrt{n}(R^2 - 0) \xrightarrow{D} N(0, 0)$ , the point mass at 0. Hence the shorth CI may not be a large sample CI for  $\tau^2$ . The lower shorth CI should be useful for testing  $H_0 : \tau^2 = 0$  versus  $H_A : \tau^2 > a$  where  $0 < a \leq 1$  since the coverage is 1 and the length of the CI converges to 0. So reject  $H_0$  if  $a$  is not in the CI.

The simulation simulated iid data  $\mathbf{w}$  with  $\mathbf{u} = \mathbf{A}\mathbf{w}$  and  $\mathbf{A}_{ij} = \psi$  for  $i \neq j$  and  $\mathbf{A}_{ii} = 1$  where  $0 \leq \psi < 1$  and  $\mathbf{u} = (x_2, \dots, x_p)^T$ . Hence  $\text{Cor}(x_i, x_j) = \rho = [2\psi + (p-3)\psi^2]/[1 + (p-2)\psi^2]$  for  $i \neq j$ . If  $\psi = 1/\sqrt{kp}$ , then  $\rho \rightarrow 1/(k+1)$  as  $p \rightarrow \infty$  where  $k > 0$ . We used  $\mathbf{w} \sim N_{p-1}(\mathbf{0}, \mathbf{I}_{p-1})$ . If  $\psi$  is high or if  $p$  is large with  $\psi \geq 0.5$ , then the data are clustered tightly about the line with direction  $\mathbf{1} = (1, \dots, 1)^T$ , and there is a dominant principal component with eigenvector  $\mathbf{1}$  and eigenvalue  $\lambda_1$ . We used  $\psi = 0, 1/\sqrt{p}$ , and 0.9. Then  $\rho = 0, \rho \rightarrow 0.5$ , or  $\rho \rightarrow 1$  as  $p \rightarrow \infty$ .

We also used  $V(x_2) = \dots = V(x_p) = \sigma_x^2$ . If  $p > 2$ , then  $\text{Cov}(x_i, x_j) = \rho\sigma_x^2$  for  $i \neq j$  and  $\text{Cov}(x_i, x_j) = V(x_i) = \sigma_x^2$  for  $i = j$ . Then  $V(Y) = \sigma_Y^2 = \sigma_L^2 + \sigma_e^2$  where

$$\begin{aligned} \sigma_L^2 = V(L) &= V\left(\sum_{i=2}^p \beta_i x_i\right) = \text{Cov}\left(\sum_{i=2}^p \beta_i x_i, \sum_{j=2}^p \beta_j x_j\right) = \sum_{i=2}^p \sum_{j=2}^p \beta_i \beta_j \text{Cov}(x_i, x_j) \\ &= \sum_{i=2}^p \beta_i^2 \sigma_x^2 + 2\rho\sigma_x^2 \sum_{i=2}^p \sum_{j=i+1}^p \beta_i \beta_j. \end{aligned}$$

The simulations took  $\beta_i \equiv 0$  or  $\beta_i \equiv 1$  for  $i = 2, \dots, p$ . For the latter case,

$$\sigma_L^2 = V(L) = (p-1)\sigma_x^2 + 2\rho\sigma_x^2 p(p-1)/2.$$

The zero mean errors  $e_i$  were from 5 distributions: i)  $N(0,1)$ , ii)  $t_3$ , iii)  $EXP(1) - 1$ , iv) uniform $(-1, 1)$ , and v)  $(1 - \epsilon)N(0, 1) + \epsilon N(0, (1 + s)^2)$  with  $\epsilon = 0.1$  and  $s = 9$  in the simulation. Then  $Y = 1 + bx_2 + bx_3 + \dots + bx_p + e$  with  $b = 0$  or  $b = 1$ .

**Remark 4.19.** Suppose the simulation uses  $K$  runs and  $W_i = 1$  if  $\mu$  is in the  $i$ th CI, and  $W_i = 0$  otherwise, for  $i = 1, \dots, K$ . Then the  $W_i$  are iid binomial $(1, 1 - \delta_n)$  where  $\rho_n = 1 - \delta_n$  is the true coverage of the CI when the sample size is  $n$ . Let  $\hat{\rho}_n = \overline{W}$ . Since  $\sum_{i=1}^K W_i \sim \text{binomial}(K, \rho_n)$ , the standard error  $SE(\overline{W}) = \sqrt{\rho_n(1 - \rho_n)/K}$ . For  $K = 5000$  and  $\rho_n$  near 0.9, we have  $3SE(\overline{W}) \approx 0.01$ . Hence an observed coverage of  $\hat{\rho}_n$  within 0.01 of the nominal coverage  $1 - \delta$  suggests that there is no reason to doubt that the nominal CI coverage is different from the observed coverage. So for a large sample

95% CI, we want the observed coverage to be between 0.94 and 0.96. Also a difference of 0.01 is not large. Coverage slightly higher than the nominal coverage is better than coverage slightly lower than the nominal coverage.

Bootstrapping confidence intervals for quantities like  $\rho^2$  and  $\tau^2$  is notoriously difficult. If  $\beta_2 = \dots = \beta_p = 0$ , then  $\sigma_L^2 = 0$  and  $\tau^2 = 0$ . However, the probability that  $R_i^{2*} > 0 = 1$ . Hence the usual two sided bootstrap percentile and shorth intervals for  $\tau^2$  will never contain 0. The one sided bootstrap CI  $[0, T_{(c)}^*]$  always contains 0, and is useful if the length of the CI goes to 0 as  $n \rightarrow \infty$ . In the table below,  $\beta_i = b$  for  $i = 2, \dots, p$ . If  $b = 0$ , then  $\tau^2 = 0$ .

The simulation for the table used 5000 runs with the bootstrap sample size  $B = 1000$ . When  $n = 400$ , the shorth(c) CI never contains  $\tau^2 = 0$  and the average length of the CI is 0.035. See *ccov* and *clen*. The lower shorth CI always contained  $\tau^2 = 0$  with *lcov* = 1, and the average CI length was *llen* = 0.036. The upper shorth CI never contains  $\tau^2 = 0$ , and the average length is near 1.

**Table 4.1** Bootstrapping  $\tau^2$  with  $R^2$  and  $B = 1000$

etype	n	p	b	$\psi$	$\tau^2$	ccov	clen	lcov	llen	ucov	ulen
1	100	4	0	0	0	0	0.135	1	0.137	0	0.990
1	200	4	0	0	0	0	0.0693	1	0.0702	0	0.995
1	400	4	0	0	0	0	0.0354	1	0.0358	0	0.988

Three *linmodpack* functions were used in the simulation. The function *shorthLU* gets the shorth(c) CI, the lower shorth CI, and the upper shorth CI. The function *Rsqrboot* bootstraps  $R^2$ , while the function *Rsqrbootsim* does the simulation. Some *R* code for the first line of Table 4.1 is below where  $b = cc$ .

```
Rsqrbootsim(n=100,p=4,BB=1000,nruns=5000,type=1,psi=0,
cc=0)
$rho
[1] 0
$sigesq
[1] 1
$sigLsq
[1] 0
$poprsq
[1] 0
$cccov
[1] 0
$avelen
[1] 0.1348881
$lcicov
```

```

[1] 1
$lavelen
[1] 0.13688
$ucicov
[1] 0
$uavelen
[1] 0.9896608

```

## 4.6 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 (4.32), (4.33), and (4.34) give useful results. Much of the theory in Section 4.5.3 does not apply to the variable selection estimator  $T_n = \mathbf{A}\hat{\boldsymbol{\beta}}_{I_{min},0}$  with  $\boldsymbol{\theta} = \mathbf{A}\boldsymbol{\beta}$ , because  $T_n$  is not smooth since  $T_n$  is equal to the estimator  $T_{jn}$  with probability  $\pi_{jn}$  for  $j = 1, \dots, J$ . Here  $\mathbf{A}$  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  $\mathbf{Y}^*$  and  $\mathbf{X}^*$  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  $\mathbf{Y}^*$  and  $\mathbf{X}^*$  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  $i = 1, \dots, B$ . 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 (4.33) 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  $j = 1, \dots, J$  where  $\sum_j \rho_{jn} = 1$ , and  $\rho_{jn} \rightarrow \pi_j$  as  $n \rightarrow \infty$ . Let  $B_{jn}$  count the number of times  $T_i^* = T_{ij}^*$  in the bootstrap sample. Then the bootstrap sample  $T_1^*, \dots, T_B^*$  can be written as

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

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

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

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

Let  $T_n = \hat{\beta}_{MIX}$  and  $T_{ij} = \hat{\beta}_{I_j,0}$ . If  $S \subseteq I_j$ , assume  $\sqrt{n}(\hat{\beta}_{I_j} - \beta_{I_j}) \xrightarrow{D} N_{a_j}(\mathbf{0}, \mathbf{V}_j)$  and  $\sqrt{n}(\hat{\beta}_{I_j}^* - \hat{\beta}_{I_j}) \xrightarrow{D} N_{a_j}(\mathbf{0}, \mathbf{V}_j)$ . Then by Equation (4.3),

$$\sqrt{n}(\hat{\beta}_{I_j,0} - \beta) \xrightarrow{D} N_p(\mathbf{0}, \mathbf{V}_{j,0}) \text{ and } \sqrt{n}(\hat{\beta}_{I_j,0}^* - \hat{\beta}_{I_j,0}) \xrightarrow{D} N_p(\mathbf{0}, \mathbf{V}_{j,0}). \quad (4.39)$$

This result means that the component clouds have the same variability asymptotically. The iid data component clouds are all centered at  $\beta$ . If the bootstrap data component clouds were all centered at the same value  $\tilde{\beta}$ , then the bootstrap cloud would be like an iid data cloud shifted to be centered at  $\tilde{\beta}$ , and (4.33) would be a confidence region for  $\theta = \beta$ . Instead, the bootstrap data component clouds are shifted slightly from a common center, and are each centered at a  $\hat{\beta}_{I_j,0}$ . 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 (4.33) having slightly higher asymptotic coverage than applying (4.33) to the iid data cloud. Also, (4.33) tends to have higher coverage than (4.34) since the cutoff for (4.33) tends to be larger than the cutoff for (4.34). Region (4.32) has the same volume as region (4.34), but tends to have higher coverage since empirically, the bagging estimator  $\bar{T}^*$  tends to estimate  $\theta$  at least as well as  $T_n$  for a mixture distribution. A similar argument holds if  $T_n = \mathbf{A}\hat{\beta}_{MIX}$ ,  $T_{ij} = \mathbf{A}\hat{\beta}_{I_j,0}$ , and  $\theta = \mathbf{A}\beta$ .

To see that  $T^*$  has more variability than  $T_n$ , asymptotically, look at Figure 4.3. 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  $j = 1, \dots, J$ . 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.

In the simulations for  $H_0 : \mathbf{A}\beta = \mathbf{B}\beta_S = \theta_0$  with  $n \geq 20p$ , the coverage tended to get close to  $1 - \delta$  for  $B \geq \max(200, 50p)$  so that  $\mathbf{S}_T^*$  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  $\mathbf{S}_T^*$  can be singular due to one or more columns of zeros in the bootstrap sample for  $\beta_1, \dots, \beta_p$ . 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  $\mathbf{S}_T^*$  for (4.32), (4.33), and (4.34). 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  $T_1, \dots, T_B$ , 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  $\hat{\boldsymbol{\beta}}_{MIX}^*$ .

#### 4.6.1 The Parametric Bootstrap

The parametric bootstrap generates  $\mathbf{Y}_j^* = (Y_i^*)$  from a parametric distribution. Then regress  $\mathbf{Y}_j^*$  on  $\mathbf{X}$  to get  $\hat{\boldsymbol{\beta}}_j^*$  for  $j = 1, \dots, B$ . Consider the parametric bootstrap for the MLR model with  $\mathbf{Y}^* \sim N_n(\mathbf{X}\boldsymbol{\beta}, \hat{\sigma}_n^2 \mathbf{I}) \sim N_n(\mathbf{H}\mathbf{Y}, \hat{\sigma}_n^2 \mathbf{I})$  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

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

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

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

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

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



### 4.6.2 The Residual Bootstrap

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

**Example 4.6.** For multiple linear regression,  $Y_i = \mathbf{x}_i^T \boldsymbol{\beta} + e_i$  is written in matrix form as  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$ . Regress  $\mathbf{Y}$  on  $\mathbf{X}$  to obtain  $\hat{\boldsymbol{\beta}}$ ,  $\mathbf{r}$ , and  $\hat{\mathbf{Y}}$  with  $i$ th element  $\hat{Y}_i = \hat{m}(\mathbf{x}_i) = \mathbf{x}_i^T \hat{\boldsymbol{\beta}}$ . For  $j = 1, \dots, B$ , regress  $\mathbf{Y}_j^*$  on  $\mathbf{X}$  to form  $\hat{\boldsymbol{\beta}}_{1,n}^*, \dots, \hat{\boldsymbol{\beta}}_{B,n}^*$  using the residual bootstrap.

Now examine the OLS model. Let  $\hat{\mathbf{Y}} = \hat{\mathbf{Y}}_{OLS} = \mathbf{X}\hat{\boldsymbol{\beta}}_{OLS} = \mathbf{H}\mathbf{Y}$  be the fitted values from the OLS full model. Let  $\mathbf{r}^W$  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),

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

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

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

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

Let  $\hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}}_{OLS}$ . Then  $\hat{\boldsymbol{\beta}}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}^*$  with  $\text{Cov}(\hat{\boldsymbol{\beta}}^*) = \sigma_n^2 (\mathbf{X}^T \mathbf{X})^{-1} = \frac{n-p}{n} MSE (\mathbf{X}^T \mathbf{X})^{-1}$ , and  $E(\hat{\boldsymbol{\beta}}^*) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T E(\mathbf{Y}^*) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{H}\mathbf{Y} = \hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}}_n$  since  $\mathbf{H}\mathbf{X} = \mathbf{X}$ . The expectations are with respect to the bootstrap distribution where  $\hat{\mathbf{Y}}$  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  $\widehat{\text{Cov}}(\hat{\boldsymbol{\beta}}_{OLS}) = MSE (\mathbf{X}^T \mathbf{X})^{-1}$ . The sample covariance matrix of the  $\hat{\boldsymbol{\beta}}^*$  is estimating  $\text{Cov}(\hat{\boldsymbol{\beta}}^*)$  as  $B \rightarrow \infty$ . Hence the residual bootstrap standard error  $SE(\hat{\beta}_i^*) \approx \sqrt{\frac{n-p}{n}} SE(\hat{\beta}_i)$  for  $i = 1, \dots, p$  where

$\hat{\boldsymbol{\beta}}_{OLS} = \hat{\boldsymbol{\beta}} = (\hat{\beta}_1, \dots, \hat{\beta}_p)^T$ . The LS CLT Theorem 2.26 says

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

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

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

See Freedman (1981).

For the above residual bootstrap,  $\hat{\boldsymbol{\beta}}_{I_j}^* = (\mathbf{X}_{I_j}^T \mathbf{X}_{I_j})^{-1} \mathbf{X}_{I_j}^T \mathbf{Y}^* = \mathbf{D}_j \mathbf{Y}^*$  with  $\text{Cov}(\hat{\boldsymbol{\beta}}_{I_j}^*) = \sigma_n^2 (\mathbf{X}_{I_j}^T \mathbf{X}_{I_j})^{-1}$  and  $E(\hat{\boldsymbol{\beta}}_{I_j}^*) = (\mathbf{X}_{I_j}^T \mathbf{X}_{I_j})^{-1} \mathbf{X}_{I_j}^T E(\mathbf{Y}^*) = (\mathbf{X}_{I_j}^T \mathbf{X}_{I_j})^{-1} \mathbf{X}_{I_j}^T \mathbf{H} \mathbf{Y} = \hat{\boldsymbol{\beta}}_{I_j}$  since  $\mathbf{H} \mathbf{X}_{I_j} = \mathbf{X}_{I_j}$ . The expectations are with respect to the bootstrap distribution where  $\hat{\mathbf{Y}}$  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  $n\text{Cov}(\hat{\boldsymbol{\beta}}_I^*) = n[(n-p)/n]\hat{\sigma}_n^2 (\mathbf{X}_I^T \mathbf{X}_I)^{-1} \xrightarrow{P} \mathbf{V}_I$  as  $n \rightarrow \infty$  with  $\hat{\sigma}_n^2 = \text{MSE}$ . Hence  $\hat{\boldsymbol{\beta}}_I^* - \hat{\boldsymbol{\beta}}_I \xrightarrow{P} \mathbf{0}$  as  $n \rightarrow \infty$  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 4.20.** The Cauchy Schwartz inequality says  $|\mathbf{a}^T \mathbf{b}| \leq \|\mathbf{a}\| \|\mathbf{b}\|$ . Suppose  $\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) = O_P(1)$  is bounded in probability. This will occur if  $\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \xrightarrow{D} N_p(\mathbf{0}, \boldsymbol{\Sigma})$ , e.g. if  $\hat{\boldsymbol{\beta}}$  is the OLS estimator. Then

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

Hence

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

since  $\max \|\mathbf{x}_i\| = O_P(1)$  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 4.21.** 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  $\sqrt{n}(\hat{\boldsymbol{\beta}}_I^* - \hat{\boldsymbol{\beta}}_I) \xrightarrow{D} N_{a_I}(\mathbf{0}, \mathbf{V}_I)$  as  $n, B \rightarrow \infty$ . For the model  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$ , the  $e_i$  are iid from a distribution that does not depend on  $n$ , and  $\boldsymbol{\beta}_E = \mathbf{0}$ . For  $\mathbf{Y}^* = \mathbf{X}\hat{\boldsymbol{\beta}} + \mathbf{r}^W$ , the distribution of the  $r_i^W$  depends on  $n$  and  $\hat{\boldsymbol{\beta}}_E \neq \mathbf{0}$  although  $\sqrt{n}\hat{\boldsymbol{\beta}}_E = O_P(1)$ .

### 4.6.3 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  $(Y_i^*, \mathbf{x}_i^*)$  with replacement from the  $n$  cases  $(Y_i, \mathbf{x}_i)$ , and regresses the  $Y_i^*$  on the  $\mathbf{x}_i^*$  to get  $\hat{\beta}_{VS,1}^*$ , and then draws another sample to get  $\hat{\beta}_{MIX,1}^*$ . This process is repeated  $B$  times to get the two bootstrap samples for  $i = 1, \dots, B$ .

Then for the full model,

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

and for a submodel  $I$ ,

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

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

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

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

One set of regularity conditions is that the MLR model holds, and if  $\mathbf{x}_i = (1 \ \mathbf{u}_i^T)^T$ , then the  $\mathbf{w}_i = (Y_i \ \mathbf{u}_i^T)^T$  are iid from some population with a nonsingular covariance matrix. Since cases are sampled with replacement, we have  $Y_i^* = \mathbf{x}_i^{*T} \beta + e_i^*$  for  $i = 1, \dots, n$ . In matrix form  $\mathbf{Y}^* = \mathbf{X}^* \beta + \mathbf{e}^*$ , but  $\mathbf{X}^*$  is a random matrix and the  $e_i^*$  are not iid from the distribution of the  $e_i$  since the  $e_i^*$  are “sampled with replacement” from the unknown  $e_1, \dots, e_n$ .

The nonparametric bootstrap uses  $\mathbf{w}_1^*, \dots, \mathbf{w}_n^*$  where the  $\mathbf{w}_i^*$  are sampled with replacement from  $\mathbf{w}_1, \dots, \mathbf{w}_n$ . By Example 4.2,  $E(\mathbf{w}^*) = \bar{\mathbf{w}}$ , and

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

Note that  $\hat{\beta}$  is a constant with respect to the bootstrap distribution. Assume all inverse matrices exist. Then by Theorem 2.20,

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

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

#### 4.6.4 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  $T_1, \dots, T_B$ , and ii) zero padding.

To see the effect of zero padding, consider  $H_0 : \mathbf{A}\boldsymbol{\beta} = \boldsymbol{\beta}_O = \mathbf{0}$  where  $\boldsymbol{\beta}_O = (\beta_{i_1}, \dots, \beta_{i_q})^T$  and  $O \subseteq E$  in (4.1) so that  $H_0$  is true. Suppose a nominal 95% confidence region is used and  $U_B = 0.96$ . Hence the confidence region (4.32) or (4.33) covers at least 96% of the bootstrap sample. If  $\hat{\boldsymbol{\beta}}_{O,j}^* = \mathbf{0}$  for more than 4% of the  $\hat{\boldsymbol{\beta}}_{O,1}^*, \dots, \hat{\boldsymbol{\beta}}_{O,B}^*$ , then  $\mathbf{0}$  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  $\hat{\boldsymbol{\beta}}_{O,j}^* = \mathbf{0}$  for  $j = 1, \dots, B$ . Then  $\mathbf{S}_T^*$  is singular, but the singleton set  $\{\mathbf{0}\}$  is the large sample  $100(1 - \delta)\%$  confidence region (4.32), (4.33), or (4.34) for  $\boldsymbol{\beta}_O$  and  $\delta \in (0, 1)$ , and the pvalue for  $H_0 : \boldsymbol{\beta}_O = \mathbf{0}$  is one. (This result holds since  $\{\mathbf{0}\}$  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  $\mathbf{x}_O$  is not needed in the model given  $\mathbf{x}_I$  is in the model if the “100%” confidence region is  $\{\mathbf{0}\}$ ,  $n \geq 20p$ ,  $B \geq 50p$ , 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 4.22.** The assumption  $\rho_{jn} \rightarrow \pi_j$  as  $n \rightarrow \infty$  seems to be the most reasonable for the residual bootstrap since  $|r_i - e_i| \rightarrow 0$  fast by Remark 4.20. The assumption may not hold for the parametric bootstrap of Section 4.6.1 if the  $e_i$  are not iid  $N(0, \sigma^2)$ . Another way to look at the bootstrap confidence region for OLS variable selection estimators is to consider the estimator  $T_{2,n}$  that chooses  $I_j$  with probability equal to the observed bootstrap proportion  $\hat{\rho}_{jn}$ . The bootstrap sample  $T_1^*, \dots, T_B^*$  tends to be slightly more variable than an iid sample  $T_{2,1}, \dots, T_{2,B}$ , and the geometric argument suggests that the large sample coverage of the nominal  $100(1 - \delta)\%$  confidence region will be at least as large as the nominal coverage  $100(1 - \delta)\%$ .

**Remark 4.23.** Note that there are several important variable selection models, including the model given by Equation (4.1) where  $\mathbf{x}^T \boldsymbol{\beta} = \mathbf{x}_S^T \boldsymbol{\beta}_S$ . Another model is  $\mathbf{x}^T \boldsymbol{\beta} = \mathbf{x}_{S_i}^T \boldsymbol{\beta}_{S_i}$  for  $i = 1, \dots, K$ . Then there are  $K \geq 2$  competing “true” nonnested submodels where  $\boldsymbol{\beta}_{S_i}$  is  $a_{S_i} \times \mathbf{1}$ . 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  $\beta_{i_j}^* = 0$  for at least one of the bootstrap samples  $j = 1, \dots, B$  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  $\beta_{\mathcal{O}}$ . See Remark 4.6.

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  $i = 1, \dots, p$ .

**Example 4.7.** Cook and Weisberg (1999, 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  $\approx \hat{\beta}_i \pm 1.96SE(\hat{\beta}_i)$ .

```

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  $H_0 : \mathbf{A}\beta = (\beta_2, \beta_3, \beta_4)^T = \mathbf{0}$ . Using the prediction region method with the

full model gave an interval  $[0, 2.930]$  with  $D_{\mathbf{0}} = 1.641$ . Note that  $\sqrt{\chi_{3,0.95}^2} = 2.795$ . 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  $D_{\mathbf{0}} = 1.134$ . 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  $D_{\mathbf{0}} = 1.577$ . 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  $D_{\mathbf{0}} = 1.245$ . So fail to reject  $H_0$ . The ratio of the volumes of the bootstrap confidence regions for this test was 0.392. (Use (4.35) with  $\mathbf{S}_T^*$  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$scov))*outfs$D0^3)/(sqrt(det(out$scov))*out$D0^3)
```

**Example 4.8.** 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 categor-

ical 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. 4.8
      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      0      0      0.000 [-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      0      0      0.000 [ -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      0      0      0.000 [-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      0      0      0.000 [ -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 <- modlboot(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

```

#### 4.6.5 Simulations

For variable selection with the  $p \times 1$  vector  $\hat{\beta}_{I_{min},0}$ , consider testing  $H_0 : \mathbf{A}\beta = \theta_0$  versus  $H_1 : \mathbf{A}\beta \neq \theta_0$  with  $\theta = \mathbf{A}\beta$  where often  $\theta_0 = \mathbf{0}$ . Then let  $T_n = \mathbf{A}\hat{\beta}_{I_{min},0}$  and let  $T_i^* = \mathbf{A}\hat{\beta}_{I_{min},0,i}^*$  for  $i = 1, \dots, B$ . The shorth estimator can be applied to a bootstrap sample  $\hat{\beta}_{i1}^*, \dots, \hat{\beta}_{iB}^*$  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 1.1. 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  $\theta = \mathbf{A}\beta = \beta_i$ ,  $\theta = \mathbf{A}\beta = \beta_S = \mathbf{1}$  and  $\theta = \mathbf{A}\beta = \beta_E = \mathbf{0}$ .

Let  $\mathbf{x} = (\mathbf{1} \mathbf{u}^T)^T$  where  $\mathbf{u}$  is the  $(p-1) \times 1$  vector of nontrivial predictors. In the simulations, for  $i = 1, \dots, n$ , we generated  $\mathbf{w}_i \sim N_{p-1}(\mathbf{0}, \mathbf{I})$  where the  $m = p-1$  elements of the vector  $\mathbf{w}_i$  are iid  $N(0,1)$ . Let the  $m \times m$  matrix  $\mathbf{A} = (a_{ij})$  with  $a_{ii} = 1$  and  $a_{ij} = \psi$  where  $0 \leq \psi < 1$  for  $i \neq j$ . Then the vector  $\mathbf{u}_i = \mathbf{A}\mathbf{w}_i$  so that  $Cov(\mathbf{u}_i) = \Sigma_{\mathbf{u}} = \mathbf{A}\mathbf{A}^T = (\sigma_{ij})$  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  $\rho \rightarrow 1/(c+1)$  as  $p \rightarrow \infty$  where  $c > 0$ . As  $\psi$  gets close to 1, the predictor vectors cluster about the line in the direction of  $(1, \dots, 1)^T$ . Let  $Y_i = 1 + 1x_{i,2} + \dots + 1x_{i,k+1} + e_i$  for  $i = 1, \dots, n$ . Hence  $\beta = (1, \dots, 1, 0, \dots, 0)^T$  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)  $\text{EXP}(1) - 1$ , iv)  $\text{uniform}(-1, 1)$ , 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  $\text{shorth}(c)$  interval for each  $\beta_i$  and used bootstrap confidence regions to test  $H_0 : \beta_S = \mathbf{1}$  (whether first  $k+1$   $\beta_i = 1$ ) and  $H_0 : \beta_E = \mathbf{0}$  (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  $\hat{\beta}_{I_{min},0}$ . Table 4.2 gives results for when the iid errors  $e_i \sim N(0,1)$  with  $n = 100$ ,  $p = 4$ , and  $k = 1$ . Table 4.2 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 (4.32), hybrid region (4.34), and Bickel and Ren region (4.33). The 0 indicates the test was  $H_0 : \beta_E = \mathbf{0}$ , while the 1 indicates that the test was  $H_0 : \beta_S = \mathbf{1}$ . The length and coverage =  $P(\text{fail to reject } H_0)$  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  $\sqrt{\chi_{g,0.95}^2}$  if the statistic  $T$  is asymptotically normal. Note that  $\sqrt{\chi_{2,0.95}^2} = 2.448$  is close to 2.45 for the full model regression bootstrap tests.

Volume ratios of the three confidence regions can be compared using (4.35), but there is not enough information in Table 4.2 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  $|\mathcal{S}_T^*|$ .

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 4.2,  $\hat{\beta}_S$  has 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  $\beta_S = (\beta_1, \beta_2)^T$ . Forward selection inference was more precise for  $\beta_E = (\beta_3, \beta_4)^T$ . The Bickel and Ren (4.33) cutoffs and coverages were at least as high as those of the hybrid region (4.34).

For  $\psi > 0$  and  $I_{min}$ , the coverages for the  $\beta_i$  corresponding to  $\beta_S$  were near 0.95, but the average length could be shorter since  $I_{min}$  tends to have

**Table 4.2** 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

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  $H_0 : \beta_E = \mathbf{0}$  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  $H_0 : \beta_S = \mathbf{1}$  was higher than that for the OLS full model.

**Table 4.3** 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 (4.34) 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  $\mathbf{S}_T^*$  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 (4.32) and (4.33) worked for much smaller  $n$ , but tended to be longer than the shorth CIs.

See Table 4.3 for one of the worst scenarios for the shorth, where shlen, prlen, and brlen are for the average CI lengths based on the shorth, (4.32), and (4.33), respectively. In Table 4.3,  $k = 8$  and the two nonzero  $\pi_j$  correspond to the full model  $\hat{\beta}$  and  $\hat{\beta}_{S,0}$ . Hence  $\beta_i = 1$  for  $i = 1, \dots, 9$  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 4.2 showed that the CI lengths are proportional to  $1/\sqrt{n}$ . When  $n = 25000$ , the shorth CI uses the 95.16th percentile while CI (4.32) uses the 95.00th percentile, allowing the average CI length of (4.32) 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 (4.32) and (4.33). 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

```

## 4.7 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  $p \gg n$  is possible. See Chapter 6, 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  $(H_1, V_1), (H_2, V_2), \dots, (H_J, V_J)$  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  $n_1 = 250000$ .

## 4.8 Summary

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

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  $S \subseteq I$ ,  $\sqrt{n}(\hat{\boldsymbol{\beta}}_I - \boldsymbol{\beta}_I) \xrightarrow{D} N_{a_I}(\mathbf{0}, \mathbf{V}_I)$  where  $(\mathbf{X}_I^T \mathbf{X}_I)/(n\sigma^2) \xrightarrow{P} \mathbf{V}_I^{-1}$ .

3) **Theorem 4.4, Variable Selection CLT.** Assume  $P(S \subseteq I_{min}) \rightarrow 1$  as  $n \rightarrow \infty$ , and let  $T_n = \hat{\beta}_{I_{min},0}$  and  $T_{jn} = \hat{\beta}_{I_j,0}$ . Let  $T_n = T_{kn} = \hat{\beta}_{I_k,0}$  with probabilities  $\pi_{kn}$  where  $\pi_{kn} \rightarrow \pi_k$  as  $n \rightarrow \infty$ . Denote the  $\pi_k$  with  $S \subseteq I_k$  by  $\pi_j$ . The other  $\pi_k = 0$  since  $P(S \subseteq I_{min}) \rightarrow 1$  as  $n \rightarrow \infty$ . Assume  $\sqrt{n}(\hat{\beta}_{I_j} - \beta_{I_j}) \xrightarrow{D} N_{a_j}(\mathbf{0}, \mathbf{V}_j)$  and  $\mathbf{u}_{jn} = \sqrt{n}(\hat{\beta}_{I_j,0} - \beta) \xrightarrow{D} \mathbf{u}_j \sim N_p(\mathbf{0}, \mathbf{V}_{j,0})$ . a) Then

$$\sqrt{n}(\hat{\beta}_{I_{min},0} - \beta) \xrightarrow{D} \mathbf{u}$$

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

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

$$\sqrt{n}(\mathbf{A}\hat{\beta}_{I_{min},0} - \mathbf{A}\beta) \xrightarrow{D} \mathbf{A}\mathbf{u} = \mathbf{v}$$

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

4) For  $h > 0$ , the hyperellipsoid  $\{\mathbf{z} : (\mathbf{z} - T)^T \mathbf{C}^{-1}(\mathbf{z} - T) \leq h^2\} = \{\mathbf{z} : D_{\mathbf{z}}^2 \leq h^2\} = \{\mathbf{z} : D_{\mathbf{z}} \leq h\}$ . A future observation (random vector)  $\mathbf{x}_f$  is in this region if  $D_{\mathbf{x}_f} \leq h$ . A large sample  $100(1 - \delta)\%$  prediction region is a set  $\mathcal{A}_n$  such that  $P(\mathbf{x}_f \in \mathcal{A}_n)$  is eventually bounded below by  $1 - \delta$  as  $n \rightarrow \infty$  where  $0 < \delta < 1$ . A *large sample*  $100(1 - \delta)\%$  *confidence region* for a vector of parameters  $\theta$  is a set  $\mathcal{A}_n$  such that  $P(\theta \in \mathcal{A}_n)$  is eventually bounded below by  $1 - \delta$  as  $n \rightarrow \infty$ .

5) Let  $q_n = \min(1 - \delta + 0.05, 1 - \delta + p/n)$  for  $\delta > 0.1$  and  $q_n = \min(1 - \delta/2, 1 - \delta + 10\delta p/n)$ , otherwise. If  $q_n < 1 - \delta + 0.001$ , set  $q_n = 1 - \delta$ . If  $(T, \mathbf{C})$  is a consistent estimator of  $(\mu, d\Sigma)$ , then  $\{\mathbf{z} : D_{\mathbf{z}}(T, \mathbf{C}) \leq h\}$  is a large sample  $100(1 - \delta)\%$  prediction regions if  $h = D_{(U_n)}$  where  $D_{(U_n)}$  is the  $100q_n$ th sample quantile of the  $D_i$ . The large sample  $100(1 - \delta)\%$  nonparametric prediction region  $\{\mathbf{z} : D_{\mathbf{z}}^2(\bar{\mathbf{x}}, \mathbf{S}) \leq D_{(U_n)}^2\}$  uses  $(T, \mathbf{C}) = (\bar{\mathbf{x}}, \mathbf{S})$ . We want  $n \geq 10p$  for good coverage and  $n \geq 50p$  for good volume.

6) Consider testing  $H_0 : \theta = \theta_0$  versus  $H_1 : \theta \neq \theta_0$  where  $\theta_0$  is a known  $g \times 1$  vector. Make a confidence region and reject  $H_0$  if  $\theta_0$  is not in the confidence region. Let  $q_B$  and  $U_B$  be as in 5) with  $n$  replaced by  $B$  and  $p$  replaced by  $g$ . Let  $\bar{T}^*$  and  $\mathbf{S}_T^*$  be the sample mean and sample covariance matrix of the bootstrap sample  $T_1^*, \dots, T_B^*$ . a) The prediction region method large sample  $100(1 - \delta)\%$  confidence region for  $\theta$  is  $\{\mathbf{w} : (\mathbf{w} - \bar{T}^*)^T [\mathbf{S}_T^*]^{-1} (\mathbf{w} - \bar{T}^*) \leq D_{(U_B)}^2\} = \{\mathbf{w} : D_{\mathbf{w}}^2(\bar{T}^*, \mathbf{S}_T^*) \leq D_{(U_B)}^2\}$  where  $D_{(U_B)}^2$  is computed from  $D_i^2 = (T_i^* - \bar{T}^*)^T [\mathbf{S}_T^*]^{-1} (T_i^* - \bar{T}^*)$  for  $i = 1, \dots, B$ . Note that the corresponding test for  $H_0 : \theta = \theta_0$  rejects  $H_0$  if  $(\bar{T}^* - \theta_0)^T [\mathbf{S}_T^*]^{-1} (\bar{T}^* - \theta_0) > D_{(U_B)}^2$ . This procedure applies the nonparametric prediction region to the bootstrap sample. b) The modified Bickel and Ren (2001) large sample  $100(1 - \delta)\%$  confidence region is  $\{\mathbf{w} : (\mathbf{w} - T_n)^T [\mathbf{S}_T^*]^{-1} (\mathbf{w} - T_n) \leq D_{(U_B, T)}^2\} = \{\mathbf{w} : D_{\mathbf{w}}^2(T_n, \mathbf{S}_T^*) \leq D_{(U_B, T)}^2\}$  where the cutoff  $D_{(U_B, T)}^2$  is the  $100q_B$ th sample

quantile of the  $D_i^2 = (T_i^* - T_n)^T [\mathbf{S}_T^*]^{-1} (T_i^* - T_n)$ . c) The hybrid large sample  $100(1 - \delta)\%$  confidence region:  $\{\mathbf{w} : (\mathbf{w} - T_n)^T [\mathbf{S}_T^*]^{-1} (\mathbf{w} - T_n) \leq D_{(U_B)}^2\} = \{\mathbf{w} : D_{\mathbf{w}}^2(T_n, \mathbf{S}_T^*) \leq D_{(U_B)}^2\}$ .

If  $g = 1$ , confidence intervals can be computed without  $\mathbf{S}_T^*$  or  $D^2$  for a), b), and c).

For some data sets,  $\mathbf{S}_T^*$  may be singular due to one or more columns of zeroes in the bootstrap sample for  $\beta_1, \dots, \beta_p$ . The variables corresponding to these columns are likely not needed in the model given that the other predictors are in the model if  $n$  and  $B$  are large enough. Let  $\beta_O = (\beta_{i_1}, \dots, \beta_{i_g})^T$ , and consider testing  $H_0 : \mathbf{A}\beta_O = \mathbf{0}$ . If  $\mathbf{A}\hat{\beta}_{O,i}^* = \mathbf{0}$  for greater than  $B\delta$  of the bootstrap samples  $i = 1, \dots, B$ , then fail to reject  $H_0$ . (If  $\mathbf{S}_T^*$  is nonsingular, then the  $100(1 - \delta)\%$  prediction region method confidence region contains  $\mathbf{0}$ .)

7) **Theorem 4.7: Geometric Argument.** Suppose  $\sqrt{n}(T_n - \theta) \xrightarrow{D} \mathbf{u}$  with  $E(\mathbf{u}) = \mathbf{0}$  and  $Cov(\mathbf{u}) = \Sigma \mathbf{u}$ . Assume  $T_1, \dots, T_B$  are iid with nonsingular covariance matrix  $\Sigma_{T_n}$ . Then the large sample  $100(1 - \delta)\%$  prediction region  $R_p = \{\mathbf{w} : D_{\mathbf{w}}^2(\bar{T}, \mathbf{S}_T) \leq D_{(U_B)}^2\}$  centered at  $\bar{T}$  contains a future value of the statistic  $T_f$  with probability  $1 - \delta_B \rightarrow 1 - \delta$  as  $B \rightarrow \infty$ . Hence the region  $R_c = \{\mathbf{w} : D_{\mathbf{w}}^2(T_n, \mathbf{S}_T) \leq D_{(U_B)}^2\}$  is a large sample  $100(1 - \delta)\%$  confidence region for  $\theta$ .

8) Applying the nonparametric prediction region (4.24) to the iid data  $T_1, \dots, T_B$  results in the  $100(1 - \delta)\%$  confidence region  $\{\mathbf{w} : (\mathbf{w} - T_n)^T \mathbf{S}_T^{-1} (\mathbf{w} - T_n) \leq D_{(U_B)}^2(T_n, \mathbf{S}_T)\}$  where  $D_{(U_B)}^2(T_n, \mathbf{S}_T)$  is computed from the  $(T_i - T_n)^T \mathbf{S}_T^{-1} (T_i - T_n)$  provided the  $\mathbf{S}_T = \mathbf{S}_{T_n}$  are “not too ill conditioned.” For OLS variable selection, assume there are two or more component clouds. The bootstrap component data clouds have the same asymptotic covariance matrix as the iid component data clouds, which are centered at  $\theta$ . The  $j$ th bootstrap component data cloud is centered at  $E(T_{ij}^*)$  and often  $E(T_{jn}^*) = T_{jn}$ . Confidence region (4.32) is the prediction region (4.24) applied to the bootstrap sample, and (4.32) is slightly larger in volume than (4.24) applied to the iid sample, asymptotically. The hybrid region (4.34) shifts (4.32) to be centered at  $T_n$ . Shifting the component clouds slightly and computing (4.24) does not change the axes of the prediction region (4.24) much compared to not shifting the component clouds. Hence by the geometric argument, we expect (4.34) to have coverage at least as high as the nominal, asymptotically, provided the  $\mathbf{S}_T^*$  are “not too ill conditioned.” The Bickel and Ren confidence region (4.33) tends to have higher coverage and volume than (4.34). Since  $\bar{T}^*$  tends to be closer to  $\theta$  than  $T_n$ , (4.32) tends to have good coverage.

9) Suppose  $m$  independent large sample  $100(1 - \delta)\%$  prediction regions are made where  $\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{x}_f$  are iid from the same distribution for each of the  $m$  runs. Let  $Y$  count the number of times  $\mathbf{x}_f$  is in the prediction region. Then  $Y \sim \text{binomial}(m, 1 - \delta_n)$  where  $1 - \delta_n$  is the true coverage. Simulation can be used to see if the true or actual coverage  $1 - \delta_n$  is close to the nominal coverage  $1 - \delta$ . A prediction region with  $1 - \delta_n < 1 - \delta$  is liberal and a region with  $1 - \delta_n > 1 - \delta$  is conservative. It is better to be conservative by 3% than

liberal by 3%. Parametric prediction regions tend to have large undercoverage and so are too liberal. Similar definitions are used for confidence regions.

10) For the bootstrap, perform variable selection on  $\mathbf{Y}_i^*$  and  $\mathbf{X}$  (or  $\mathbf{X}^*$  for the nonparametric bootstrap), fit the model that minimizes the criterion, and add 0s corresponding to the omitted variables, resulting in estimators  $\hat{\beta}_1^*, \dots, \hat{\beta}_B^*$  where  $\hat{\beta}_i^* = \hat{\beta}_{I_{min,0,i}^*}$ .

11) Let  $Z_1, \dots, Z_n$  be random variables, let  $Z_{(1)}, \dots, Z_{(n)}$  be the order statistics, and let  $c$  be a positive integer. Compute  $Z_{(c)} - Z_{(1)}, Z_{(c+1)} - Z_{(2)}, \dots, Z_{(n)} - Z_{(n-c+1)}$ . Let  $\text{shorth}(c) = [Z_{(d)}, Z_{(d+c-1)}]$  correspond to the interval with the shortest length.

The large sample  $100(1-\delta)\%$  *shorth*( $c$ ) CI uses the interval  $[T_{(1)}^*, T_{(c)}^*], [T_{(2)}^*, T_{(c+1)}^*], \dots, [T_{(B-c+1)}^*, T_{(B)}^*]$  of shortest length. Here  $c = \min(B, \lceil B[1 - \delta + 1.12\sqrt{\delta/B}] \rceil)$ . The shorth CI is computed by applying the shorth PI to the bootstrap sample.

## 4.9 Complements

This chapter followed Olive (2017b, ch. 5) and Pelawa Watagoda and Olive (2019ab) closely. Also see Olive (2013a, 2018), Pelawa Watagoda (2017), and Rathnayake and Olive (2019). For MLR, Olive (2017a: p. 123, 2017b: p. 176) showed that  $\hat{\beta}_{I_{min,0}}$  is a consistent estimator. Olive (2014: p. 283, 2017ab, 2018) recommended using the *shorth*( $c$ ) estimator for the percentile method. Olive (2017a: p. 128, 2017b: p. 181, 2018) showed that the prediction region method can simulate well for the  $p \times 1$  vector  $\hat{\beta}_{I_{min,0}}$ . Hastie et al. (2009, p. 57) noted that variable selection is a shrinkage estimator: the coefficients are shrunk to 0 for the omitted variables.

Good references for the bootstrap include Efron (1979, 1982), Efron and Hastie (2016, ch. 10–11), and Efron and Tibshirani (1993). Also see Chen (2016) and Hesterberg (2014). One of the sufficient conditions for the bootstrap confidence region is that  $T$  has a well behaved Hadamard derivative. Fréchet differentiability implies Hadamard differentiability, and many statistics are shown to be Hadamard differentiable in Bickel and Ren (2001), Clarke (1986, 2000), Fernholtz (1983), Gill (1989), Ren (1991), and Ren and Sen (1995). Bickel and Ren (2001) showed that their method can work when Hadamard differentiability fails.

There is a massive literature on variable selection and a fairly large literature for inference after variable selection. See, for example, Leeb and Pötscher (2005, 2006, 2008), Leeb et al. (2015), Tibshirani et al. (2016), 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  $\mathbf{A}\boldsymbol{\beta}$ . 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  $\bar{T}^* \approx \sum_j \rho_{jn} T_{jn}$ , 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  $\bar{T}^* \pm z_{1-\delta} SE(\bar{T}^*)$  assuming  $\bar{T}^*$  is asymptotically normal and using delta method techniques, which require nonsingular covariance matrices. There is not yet rigorous theory for this method. Section 4.2 proved that  $\bar{T}^*$  is asymptotically normal: under regularity conditions: if  $\sqrt{n}(T_n - \boldsymbol{\theta}) \xrightarrow{D} N_g(\mathbf{0}, \boldsymbol{\Sigma}_A)$  and  $\sqrt{n}(T_i^* - T_n) \xrightarrow{D} N_g(\mathbf{0}, \boldsymbol{\Sigma}_A)$ , then under regularity conditions  $\sqrt{n}(\bar{T}^* - \boldsymbol{\theta}) \xrightarrow{D} N_g(\mathbf{0}, \boldsymbol{\Sigma}_A)$ . If  $g = 1$ , then the prediction region method large sample  $100(1 - \delta)\%$  CI for  $\theta$  has  $P(\theta \in [\bar{T}^* - a_{(U_B)}, \bar{T}^* + a_{(U_B)}]) \rightarrow 1 - \delta$  as  $n \rightarrow \infty$ . If the Frey CI also has coverage converging to  $1 - \delta$ , then the two methods have the same asymptotic length (scaled by multiplying by  $\sqrt{n}$ ), since otherwise the shorter interval will have lower asymptotic coverage.

For the mixture distribution with two or more component groups,  $\sqrt{n}(T_n - \boldsymbol{\theta}) \xrightarrow{D} \mathbf{v}$  by Theorem 4.4 b). If  $\sqrt{n}(T_i^* - c_n) \xrightarrow{D} \mathbf{u}$  then  $c_n$  must be a value such as  $c_n = \bar{T}^*$ ,  $c_n = \sum_j \rho_{jn} T_{jn}$ , or  $c_n = \sum_j \pi_j T_{jn}$ . Next we will examine  $\bar{T}^*$ . If  $S \subseteq I_j$ , then  $\sqrt{n}(\hat{\boldsymbol{\beta}}_{I_j,0} - \boldsymbol{\beta}) \xrightarrow{D} N_p(\mathbf{0}, \mathbf{V}_{j,0})$ , and for the parametric and nonparametric bootstrap,  $\sqrt{n}(\hat{\boldsymbol{\beta}}_{I_j,0}^* - \hat{\boldsymbol{\beta}}_{I_j,0}) \xrightarrow{D} N_p(\mathbf{0}, \mathbf{V}_{j,0})$ . Let  $T_n = \mathbf{A}\hat{\boldsymbol{\beta}}_{I_{min},0}$  and  $T_{jn} = \mathbf{A}\hat{\boldsymbol{\beta}}_{I_j,0} = \mathbf{A}\mathbf{D}_{j0}\mathbf{Y}$  using notation from Section 4.6. Let  $\boldsymbol{\theta} = \mathbf{A}\boldsymbol{\beta}$ . Hence from Section 4.5.3,  $\sqrt{n}(\bar{T}_j^* - T_{jn}) \xrightarrow{P} \mathbf{0}$ . Assume  $\hat{\rho}_{in} \xrightarrow{P} \rho_i$  as  $n \rightarrow \infty$ . Then  $\sqrt{n}(\bar{T}^* - \boldsymbol{\theta}) =$

$$\sum_i \hat{\rho}_{in} \sqrt{n}(\bar{T}_i^* - \boldsymbol{\theta}) = \sum_j \hat{\rho}_{jn} \sqrt{n}(\bar{T}_j^* - \boldsymbol{\theta}) + \sum_k \hat{\rho}_{kn} \sqrt{n}(\bar{T}_k^* - \boldsymbol{\theta})$$

$$= d_n + a_n \text{ where } a_n \xrightarrow{P} \mathbf{0} \text{ since } \rho_k = 0. \text{ Now}$$

$$d_n = \sum_j \hat{\rho}_{jn} \sqrt{n}(\bar{T}_j^* - T_{jn} + T_{jn} - \boldsymbol{\theta}) = \sum_j \hat{\rho}_{jn} \sqrt{n}(T_{jn} - \boldsymbol{\theta}) + c_n$$



where  $c_n = o_P(1)$  since  $\sqrt{n}(\overline{T}_j^* - T_{jn}) = o_P(1)$ . Hence under regularity conditions, if  $\sqrt{n}(\overline{T}^* - \boldsymbol{\theta}) \xrightarrow{D} \mathbf{w}$  then  $\sum_j \rho_j \sqrt{n}(T_{jn} - \boldsymbol{\theta}) \xrightarrow{D} \mathbf{w}$ .

To examine the last term and  $\mathbf{w}$ , let the  $n \times 1$  vector  $\mathbf{Y}$  have characteristic function  $\phi_{\mathbf{Y}}$ ,  $E(\mathbf{Y}) = \mathbf{X}\boldsymbol{\beta}$ , and  $\text{Cov}(\mathbf{Y}) = \sigma^2 \mathbf{I}$ . Let  $\mathbf{Z} = (\mathbf{Y}^T, \dots, \mathbf{Y}^T)^T$  be a  $Jn \times 1$  vector with  $J$  copies of  $\mathbf{Y}$  stacked into a vector. Let  $\mathbf{t} = (\mathbf{t}_1^T, \dots, \mathbf{t}_J^T)^T$ . Then  $\mathbf{Z}$  has characteristic function  $\phi_{\mathbf{Z}}(\mathbf{t}) = \phi_{\mathbf{Y}}(\sum_{j=1}^J \mathbf{t}_j) = \phi_{\mathbf{Y}}(\mathbf{s})$ . Now assume  $\mathbf{Y} \sim N_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$ . Then  $\mathbf{t}^T \mathbf{Z} = \mathbf{s}^T \mathbf{Y} \sim N(\mathbf{s}^T \mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{s}^T \mathbf{s})$ . Hence  $\mathbf{Z}$  has a multivariate normal distribution by Definition 1.23 with  $E(\mathbf{Z}) = (\mathbf{X}\boldsymbol{\beta}^T, \dots, \mathbf{X}\boldsymbol{\beta}^T)^T$ , and  $\text{Cov}(\mathbf{Z})$  a block matrix with  $J \times J$  blocks each equal to  $\sigma^2 \mathbf{I}$ . Then

$$\begin{aligned} \sum_j \rho_j T_{jn} &= \sum_j \rho_j \mathbf{A} \mathbf{D}_{j0} \mathbf{Y} = \mathbf{B} \mathbf{Y} \sim N_g(\boldsymbol{\theta}, \sigma^2 \mathbf{B} \mathbf{B}^T) = \\ &N_g(\boldsymbol{\theta}, \sigma^2 \sum_j \sum_k \rho_j \rho_k \mathbf{A} \mathbf{D}_{j0} \mathbf{D}_{k0}^T \mathbf{A}) \end{aligned}$$

since  $E(T_{jn}) = E(\mathbf{A} \hat{\boldsymbol{\beta}}_{I_j, 0}) = \mathbf{A} \boldsymbol{\beta} = \boldsymbol{\theta}$  if  $S \subseteq I_j$ . Since  $(T_{1n}^T, \dots, T_{jn}^T)^T = \text{diag}(\mathbf{A} \mathbf{D}_{10}, \dots, \mathbf{A} \mathbf{D}_{j0}) \mathbf{Z}$ , then  $(T_{1n}^T, \dots, T_{jn}^T)^T$  is multivariate normal and

$$\sum_j \rho_j T_{jn} \sim N_g[\boldsymbol{\theta}, \sum_j \sum_k \pi_j \pi_k \text{Cov}(T_{jn}, T_{kn})].$$

Now assume  $n \mathbf{D}_{j0} \mathbf{D}_{k0}^T \xrightarrow{P} \mathbf{W}_{jk}$  as  $n \rightarrow \infty$ . Then

$$\sum_j \rho_j \sqrt{n}(T_{jn} - \boldsymbol{\theta}) \xrightarrow{D} \mathbf{w} \sim N_g(\mathbf{0}, \sigma^2 \sum_j \sum_k \rho_j \rho_k \mathbf{A} \mathbf{W}_{jk} \mathbf{A}).$$

We conjecture that this result may hold under milder conditions than  $\mathbf{Y} \sim N_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$ , but even the above results are not yet rigorous. If  $\sqrt{n}(T_{jn} - \boldsymbol{\theta}) \xrightarrow{D} \mathbf{w}_j \sim N_g(\mathbf{0}, \boldsymbol{\Sigma}_j)$ , then a possibly poor approximation is  $\overline{T}^* \approx \sum_j \rho_j T_{jn} \approx N_g[\boldsymbol{\theta}, \sum_j \sum_k \rho_j \rho_k \text{Cov}(T_{jn}, T_{kn})]$ , and estimating  $\sum_j \sum_k \rho_j \rho_k \text{Cov}(T_{jn}, T_{kn})$  with delta method techniques may not be possible.

The double bootstrap technique may be useful. See Hall (1986) and Chang and Hall (2015) for references. The double bootstrap for  $\overline{T}^* = \overline{T}_B^*$  says that  $T_n = \overline{T}^*$  is a statistic that can be bootstrapped. Let  $B_d \geq 50g_{max}$  where  $1 \leq g_{max} \leq p$  is the largest dimension of  $\boldsymbol{\theta}$  to be tested with the double bootstrap. Draw a bootstrap sample of size  $B$  and compute  $\overline{T}^* = T_1^*$ . Repeat for a total of  $B_d$  times. Apply the confidence region (4.32), (4.33), or (4.34) to the double bootstrap sample  $T_1^*, \dots, T_{B_d}^*$ . If  $D_{(U_{B_d})} \approx D_{(U_{B_d}, T)} \approx \sqrt{\chi_{g, 1-\delta}^2}$ , then  $\overline{T}^*$  may be approximately multivariate normal. The CI (4.32) applied to the double bootstrap sample could be regarded as a modified Frey CI

without delta method techniques. Of course the double bootstrap tends to be too computationally expensive to simulate.

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. Compute  $(\bar{\mathbf{x}}_H, \mathbf{S}_H)$  from the cases in  $H$ . Then compute the distances  $D_i^2 = (\mathbf{x}_i - \bar{\mathbf{x}}_H)^T \mathbf{S}_H^{-1} (\mathbf{x}_i - \bar{\mathbf{x}}_H)$  for the  $m$  vectors  $\mathbf{x}_i$  in  $V$ . Then a large sample  $100(1 - \delta)\%$  prediction region for  $\mathbf{x}_F$  is  $\{\mathbf{x} : D_{\mathbf{x}}^2(\bar{\mathbf{x}}_H, \mathbf{S}_H) \leq D_{(k_m)}^2\}$  where  $k_m = \lceil m(1 - \delta) \rceil$ . This prediction region may give better coverage than the nonparametric prediction region (4.24) if  $5p \leq n \leq 20p$ .

The iid sample  $T_1, \dots, T_B$  has sample mean  $\bar{T}$ . Let  $T_{in} = T_{ijn}$  if  $T_{jn}$  is chosen  $D_{jn}$  times where the random variables  $D_{jn}/B \xrightarrow{P} \pi_{jn}$ . The  $D_{jn}$  follow a multinomial distribution. Then the iid sample can be written as

$$T_{1,1}, \dots, T_{D_{1n},1}, \dots, T_{1,J}, \dots, T_{D_{Jn},J},$$

where the  $T_{ij}$  are not iid. Denote  $T_{1j}, \dots, T_{D_{jn},j}$  as the  $j$ th component of the iid sample with sample mean  $\bar{T}_j$  and sample covariance matrix  $\mathbf{S}_{T,j}$ . Thus

$$\bar{T} = \frac{1}{B} \sum_{i=1}^B T_{ijn} = \sum_j \frac{D_{jn}}{B} \frac{1}{D_{jn}} \sum_{i=1}^{D_{jn}} T_{ij} = \sum_j \hat{\pi}_{jn} \bar{T}_j.$$

Hence  $\bar{T}$  is a random linear combination of the  $\bar{T}_j$ . Conditionally on the  $D_{jn}$ , the  $T_{ij}$  are independent, and  $\bar{T}$  is a linear combination of the  $\bar{T}_j$ . Note that  $\text{Cov}(\bar{T}) = \text{Cov}(T_n)/B$ .

**Software.** The simulations were done in *R*. See R Core Team (2016). We used several *R* functions including forward selection as computed with the `regsubsets` function from the `leaps` library. Several `linmodpack` functions were used. The function `predrgn` makes the nonparametric prediction region and determines whether  $\mathbf{x}_f$  is in the region. The function `predreg` also makes the nonparametric prediction region, and determines if  $\mathbf{0}$  is in the region. For multiple linear regression, the function `regboot` does the residual bootstrap for multiple linear regression, `regbootsim` simulates the residual bootstrap for regression, and the function `rowboot` does the empirical nonparametric bootstrap. The function `vsbootsim` simulates the bootstrap for all subsets variable selection, so needs  $p$  small, while `vsbootsim2` simulates the prediction region method for forward selection. The functions `fselboot` and `vselboot` bootstrap the forward selection and all subsets variable selection estimators that minimize  $C_p$ . See Examples 4.7 and 4.8. The `shorth3` function computes the `shorth(c)` intervals with the Frey (2013) correction used when  $g = 1$ . Table 4.2 was made using `regbootsim3` for the OLS full model and `vsbootsim4` for forward selection. The functions `bicboot` and `bicbootsim` are useful if BIC is used instead of  $C_p$ . For forward selection

with  $C_p$ , the function `vscisim` was used to make Table 4.3, and can be used to compare the shorth, prediction region method, and Bickel and Ren CIs for  $\beta_i$ .

## 4.10 Problems

**4.1.** Consider the Cushny and Peebles data set (see Staudte and Sheather 1990, p. 97) listed below. Find `shorth(7)`. Show work.

0.0 0.8 1.0 1.2 1.3 1.3 1.4 1.8 2.4 4.6

**4.2.** Find `shorth(5)` for the following data set. Show work.

6 76 90 90 94 94 95 97 97 1008

**4.3.** Find `shorth(5)` for the following data set. Show work.

66 76 90 90 94 94 95 95 97 98

**4.4.** Suppose you are estimating the mean  $\theta$  of losses with the maximum likelihood estimator (MLE)  $\bar{X}$  assuming an exponential ( $\theta$ ) distribution. Compute the sample mean of the fourth bootstrap sample.

actual losses 1, 2, 5, 10, 50:  $\bar{X} = 13.6$

bootstrap samples:

2, 10, 1, 2, 2:  $\bar{X} = 3.4$

50, 10, 50, 2, 2:  $\bar{X} = 22.8$

10, 50, 2, 1, 1:  $\bar{X} = 12.8$

5, 2, 5, 1, 50:  $\bar{X} = ?$

**4.5.** The data below are a sorted residuals from a least squares regression where  $n = 100$  and  $p = 4$ . Find `shorth(97)` of the residuals.

number	1	2	3	4	...	97	98	99	100
residual	-2.39	-2.34	-2.03	-1.77	...	1.76	1.81	1.83	2.16

**4.6.** To find the sample median of a list of  $n$  numbers where  $n$  is odd, order the numbers from smallest to largest and the median is the middle ordered number. The sample median estimates the population median. Suppose the sample is  $\{14, 3, 5, 12, 20, 10, 9\}$ . Find the sample median for each of the three bootstrap samples listed below.

Sample 1: 9, 10, 9, 12, 5, 14, 3

Sample 2: 3, 9, 20, 10, 9, 5, 14

Sample 3: 14, 12, 10, 20, 3, 3, 5

**4.7.** Suppose you are estimating the mean  $\mu$  of losses with  $T = \bar{X}$ .

actual losses 1, 2, 5, 10, 50:  $\bar{X} = 13.6$ ,

a) Compute  $T_1^*, \dots, T_4^*$ , where  $T_i^*$  is the sample mean of the  $i$ th bootstrap sample. bootstrap samples:

2, 10, 1, 2, 2:

50, 10, 50, 2, 2:

10, 50, 2, 1, 1:

5, 2, 5, 1, 50:

b) Now compute the bagging estimator which is the sample mean of the

$T_i^*$ : the bagging estimator  $\bar{T}^* = \frac{1}{B} \sum_{i=1}^B T_i^*$  where  $B = 4$  is the number of

bootstrap samples.

**4.8.** Consider the output for Example 4.7 for the minimum  $C_p$  forward selection model.

a) What is  $\hat{\beta}_{I_{min}}$ ?

b) What is  $\hat{\beta}_{I_{min},0}$ ?

c) The large sample 95% shorth CI for  $H$  is  $[0, 0.016]$ . Is  $H$  needed in the minimum  $C_p$  model given that the other predictors are in the model?

d) The large sample 95% shorth CI for  $\log(S)$  is  $[0.324, 0.913]$  for all subsets. Is  $\log(S)$  needed in the minimum  $C_p$  model given that the other predictors are in the model?

e) Suppose  $x_1 = 1$ ,  $x_4 = H = 130$ , and  $x_5 = \log(S) = 5.075$ . Find  $\hat{Y} = (x_1 \ x_4 \ x_5) \hat{\beta}_{I_{min}}$ . Note that  $Y = \log(M)$ .

**4.9<sup>Q</sup>.** Suppose  $\mathbf{Y}^* = \mathbf{X}\hat{\beta} + \mathbf{r}^W$  where  $E(\mathbf{r}^W) = \mathbf{0}$  and  $Cov(\mathbf{r}^W) = Cov(\mathbf{Y}^*) = MSE \mathbf{I}_n$ . Then  $\hat{\beta}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}^*$ . Recall that  $\mathbf{X}$  is an  $n \times p$  constant matrix. Simplify quantities when possible.

a) What is  $E(\hat{\beta}^*)$ ?

b) What is  $Cov(\hat{\beta}^*)$ ?

c) Recall that  $\mathbf{X}\hat{\beta} = \mathbf{P}\mathbf{Y}$ . What is  $E(\hat{\beta}_I^*) = E[(\mathbf{X}_I^T \mathbf{X}_I)^{-1} \mathbf{X}_I^T \mathbf{Y}^*]$ ?

d) What is  $Cov(\hat{\beta}_I^*)$ ?

**4.10<sup>Q</sup>.** Suppose  $\mathbf{Y}^* \sim N_n(\mathbf{X}\hat{\beta}, \sigma_n^2 \mathbf{I}_n)$ . Hence  $Y_i^* = \mathbf{x}_i^T \hat{\beta} + \epsilon_i^P$  where  $E(\epsilon_i^P) = 0$  and  $V(\epsilon_i^P) = \sigma_n^2$ . Hence  $\mathbf{A}\mathbf{Y}^* \sim N_g(\mathbf{A}\mathbf{X}\hat{\beta}, \sigma_n^2 \mathbf{A}\mathbf{A}^T)$  if  $\mathbf{A}$  is a  $g \times n$  constant matrix. Recall that  $\mathbf{X}$  is an  $n \times p$  constant matrix. Simplify quantities when possible.

a) What is the distribution of  $\hat{\beta}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}^*$ ?

b) Using a), what is  $E(\hat{\beta}^*)$ ?

c) Recall that  $\mathbf{X}\hat{\beta} = \mathbf{P}\mathbf{Y}$ . What is the distribution of  $\hat{\beta}_I^* = (\mathbf{X}_I^T \mathbf{X}_I)^{-1} \mathbf{X}_I^T \mathbf{Y}^*$  if  $\hat{\beta}_I^*$  is  $k \times 1$ ?

**4.11<sup>Q</sup>.** Suppose  $\mathbf{Y}^* = \mathbf{X}\hat{\beta} + \mathbf{r}^W$  where  $E(\mathbf{r}^W) = \mathbf{0}$  and  $Cov(\mathbf{r}^W) = Cov(\mathbf{Y}^*) = \text{diag}(r_i^2) = \text{diag}(r_1^2, \dots, r_n^2)$ . Then  $\hat{\beta}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}^*$  is the least squares estimator from regressing  $\mathbf{Y}^*$  on  $\mathbf{X}$ , an  $n \times p$  constant matrix. This model is used for the wild bootstrap. Simplify quantities when possible. (Can simplify a) and c), but can't simplify b) and d) much.)

- a) What is  $E(\hat{\boldsymbol{\beta}}^*)$ ?  
 b) What is  $Cov(\hat{\boldsymbol{\beta}}^*)$ ?  
 c) Recall that  $\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{P}\mathbf{Y}$ . What is  $E(\hat{\boldsymbol{\beta}}_I^*) = E[(\mathbf{X}_I^T \mathbf{X}_I)^{-1} \mathbf{X}_I^T \mathbf{Y}^*]$ ?  
 d) What is  $Cov(\hat{\boldsymbol{\beta}}_I^*)$ ?

4.12.

4.13.

4.14.

4.15.

4.16.

4.17.

4.18.

4.19.

4.20.

### R Problems

Use the command `source("G:/linmodpack.txt")` to download the functions and the command `source("G:/linmoddata.txt")` to download the data. See Preface or Section 11.1. Typing the name of the `linmodpack` function, e.g. `regbootsim2`, will display the code for the function. Use the `args` command, e.g. `args(regbootsim2)`, to display the needed arguments for the function. For the following problem, the *R* command can be copied and pasted from (<http://parker.ad.siu.edu/Olive/linmodrhw.txt>) into *R*.

4.21. a) Type the *R* command `predsim()` and paste the output into *Word*.

This program computes  $\mathbf{x}_i \sim N_4(\mathbf{0}, \text{diag}(1, 2, 3, 4))$  for  $i = 1, \dots, 100$  and  $\mathbf{x}_f = \mathbf{x}_{101}$ . One hundred such data sets are made, and `ncvr`, `scvr`, and `mcvr` count the number of times  $\mathbf{x}_f$  was in the nonparametric, semiparametric, and parametric MVN 90% prediction regions. The volumes of the prediction regions are computed and `voln`, `vols`, and `volm` are the average ratio of the volume of the  $i$ th prediction region over that of the semiparametric region. Hence `vols` is always equal to 1. For multivariate normal data, these ratios should converge to 1 as  $n \rightarrow \infty$ .

b) Were the three coverages near 90%?

4.22. 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 `regbootsim2` bootstraps the regression model, finds bootstrap confidence intervals for  $\beta_i$  and a bootstrap confidence region for  $(\beta_3, \beta_4)^T$  corresponding to the test  $H_0 : \beta_3 = \beta_4 = 0$  versus  $H_A$ : not  $H_0$ . See the *R* code near Table 4.3. The lengths of the CIs along with the proportion of times the CI for  $\beta_i$  contained  $\beta_i$  are given. 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 test fails to reject  $H_0$ . Since nominal 95% CIs were used and the nominal

level of the test is 0.05 when  $H_0$  is true, we want the coverages near 0.95. The CI lengths for the first 4 intervals should be near 0.392. The residual bootstrap is used.

Copy and paste the commands for this problem into *R*, and include the output in *Word*.