

Chapter 4

Inference After Variable Selection

This chapter considers inference after variable selection including prediction intervals and bootstrap hypothesis testing. 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. Several of the sections of this chapter are much more technical than the rest of the book.

4.1 Variable Selection

Review Section 2.4 for 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.

When there is a sequence of M submodels, the final submodel I_d needs to be selected with a_d terms. 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 . 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 the Cox (1972) proportional hazards regression, to the predictors that had nonzero lasso or elastic net coefficients. See Tibshirani (1997) and Simon et al. (2011) for lasso and elastic net.

Forward selection and backward elimination both form a sequence of submodels I_1, \dots, I_p where I_j uses j predictors. Heuristically, backward elimination tries to delete the variable that will increase AIC the least, while forward selection tries to add the variable that will decrease AIC the most. Let I_{min} minimize the criterion such as AIC, BIC, or lasso. Often I_{min} from forward selection will differ from I_{min} from backward elimination, especially if the predictors are correlated.

Now suppose $p = 6$ and S in Equation (2.4) corresponds to x_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 . Suppose $I_{min} = I_d$. Compared to selecting a model I_d before examining the data, the model I_{min} fits the data a bit too well. 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**.

If $\hat{\beta}_{I_{min}}$ is $a \times 1$, form the $p \times 1$ vector $\hat{\beta}_{I_{min},0}$ from $\hat{\beta}_{I_{min}}$ by adding 0s corresponding to the omitted variables. For example, if $p = 4$ and $\hat{\beta}_{I_{min}} = (\hat{\beta}_1, \hat{\beta}_3)^T$, then $\hat{\beta}_{I_{min},0} = (\hat{\beta}_1, 0, \hat{\beta}_3, 0)^T$.

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 Some Tools for Large Sample Theory

This section gives some tools that are useful for inference after variable selection. The multivariate normal distribution is important. The last four subsections are more technical than most of this book. They can be omitted on first reading and refer to relevant theorems as needed.

4.2.1 The Multivariate Normal Distribution

For much of this book, \mathbf{X} is an $n \times p$ design matrix, but this subsection will usually use the notation $\mathbf{X} = (X_1, \dots, X_p)^T$ and \mathbf{Y} for the random vectors, and $\mathbf{x} = (x_1, \dots, x_p)^T$ for the observed value of the random vector. This notation will be useful to avoid confusion when studying conditional distributions such as $\mathbf{Y}|\mathbf{X} = \mathbf{x}$. It can be shown that Σ is positive semidefinite and symmetric.

Definition 4.1: Rao (1965, p. 437). A $p \times 1$ random vector \mathbf{X} has a p -dimensional *multivariate normal distribution* $N_p(\boldsymbol{\mu}, \Sigma)$ iff $\mathbf{t}^T \mathbf{X}$ has a univariate normal distribution for any $p \times 1$ vector \mathbf{t} .

If Σ is positive definite, then \mathbf{X} has a pdf

$$f(\mathbf{z}) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} e^{-(1/2)(\mathbf{z}-\boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{z}-\boldsymbol{\mu})} \quad (4.1)$$

where $|\boldsymbol{\Sigma}|^{1/2}$ is the square root of the determinant of $\boldsymbol{\Sigma}$. Note that if $p = 1$, then the quadratic form in the exponent is $(z - \mu)(\sigma^2)^{-1}(z - \mu)$ and X has the univariate $N(\mu, \sigma^2)$ pdf. If $\boldsymbol{\Sigma}$ is positive semidefinite but not positive definite, then \mathbf{X} has a degenerate distribution. For example, the univariate $N(0, 0^2)$ distribution is degenerate (the point mass at 0).

Definition 4.2. The *population mean* of a random $p \times 1$ vector $\mathbf{X} = (X_1, \dots, X_p)^T$ is

$$E(\mathbf{X}) = (E(X_1), \dots, E(X_p))^T$$

and the $p \times p$ *population covariance matrix*

$$\text{Cov}(\mathbf{X}) = E(\mathbf{X} - E(\mathbf{X}))(\mathbf{X} - E(\mathbf{X}))^T = (\sigma_{ij}).$$

That is, the ij entry of $\text{Cov}(\mathbf{X})$ is $\text{Cov}(X_i, X_j) = \sigma_{ij}$.

The covariance matrix is also called the variance-covariance matrix and variance matrix. Sometimes the notation $\text{Var}(\mathbf{X})$ is used. Note that $\text{Cov}(\mathbf{X})$ is a symmetric positive semidefinite matrix. If \mathbf{X} and \mathbf{Y} are $p \times 1$ random vectors, \mathbf{a} a conformable constant vector, and \mathbf{A} and \mathbf{B} are conformable constant matrices, then

$$E(\mathbf{a} + \mathbf{X}) = \mathbf{a} + E(\mathbf{X}) \quad \text{and} \quad E(\mathbf{X} + \mathbf{Y}) = E(\mathbf{X}) + E(\mathbf{Y}) \quad (4.2)$$

and

$$E(\mathbf{A}\mathbf{X}) = \mathbf{A}E(\mathbf{X}) \quad \text{and} \quad E(\mathbf{A}\mathbf{X}\mathbf{B}) = \mathbf{A}E(\mathbf{X})\mathbf{B}. \quad (4.3)$$

Thus

$$\text{Cov}(\mathbf{a} + \mathbf{A}\mathbf{X}) = \text{Cov}(\mathbf{A}\mathbf{X}) = \mathbf{A}\text{Cov}(\mathbf{X})\mathbf{A}^T. \quad (4.4)$$

Some important properties of multivariate normal (MVN) distributions are given in the following three theorems. These theorems can be proved using results from Johnson and Wichern (1988, pp. 127-132) or Severini (2005, ch. 8).

Theorem 4.1. a) If $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, then $E(\mathbf{X}) = \boldsymbol{\mu}$ and

$$\text{Cov}(\mathbf{X}) = \boldsymbol{\Sigma}.$$

b) If $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, then any linear combination $\mathbf{t}^T \mathbf{X} = t_1 X_1 + \dots + t_p X_p \sim N_1(\mathbf{t}^T \boldsymbol{\mu}, \mathbf{t}^T \boldsymbol{\Sigma} \mathbf{t})$. Conversely, if $\mathbf{t}^T \mathbf{X} \sim N_1(\mathbf{t}^T \boldsymbol{\mu}, \mathbf{t}^T \boldsymbol{\Sigma} \mathbf{t})$ for every $p \times 1$ vector \mathbf{t} , then $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$.

c) **The joint distribution of independent normal random variables is MVN.** If X_1, \dots, X_p are independent univariate normal $N(\mu_i, \sigma_i^2)$ random vectors, then $\mathbf{X} = (X_1, \dots, X_p)^T$ is $N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ where $\boldsymbol{\mu} = (\mu_1, \dots, \mu_p)^T$ and $\boldsymbol{\Sigma} = \text{diag}(\sigma_1^2, \dots, \sigma_p^2)$ (so the off diagonal entries $\sigma_{ij} = 0$ while the diagonal entries of $\boldsymbol{\Sigma}$ are $\sigma_{ii} = \sigma_i^2$).

d) If $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ and if \mathbf{A} is a $q \times p$ matrix, then $\mathbf{AX} \sim N_q(\mathbf{A}\boldsymbol{\mu}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)$. If \mathbf{a} is a $p \times 1$ vector of constants and b is a constant, then $\mathbf{a} + b\mathbf{X} \sim N_p(\mathbf{a} + b\boldsymbol{\mu}, b^2\boldsymbol{\Sigma})$. (Note that $b\mathbf{X} = b\mathbf{I}_p\mathbf{X}$ with $\mathbf{A} = b\mathbf{I}_p$.)

It will be useful to partition \mathbf{X} , $\boldsymbol{\mu}$, and $\boldsymbol{\Sigma}$. Let \mathbf{X}_1 and $\boldsymbol{\mu}_1$ be $q \times 1$ vectors, let \mathbf{X}_2 and $\boldsymbol{\mu}_2$ be $(p - q) \times 1$ vectors, let $\boldsymbol{\Sigma}_{11}$ be a $q \times q$ matrix, let $\boldsymbol{\Sigma}_{12}$ be a $q \times (p - q)$ matrix, let $\boldsymbol{\Sigma}_{21}$ be a $(p - q) \times q$ matrix, and let $\boldsymbol{\Sigma}_{22}$ be a $(p - q) \times (p - q)$ matrix. Then

$$\mathbf{X} = \begin{pmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{pmatrix}, \quad \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix}, \quad \text{and} \quad \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix}.$$

Theorem 4.2. a) **All subsets of a MVN are MVN:** $(X_{k_1}, \dots, X_{k_q})^T \sim N_q(\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}})$ where $\tilde{\boldsymbol{\mu}}_i = E(X_{k_i})$ and $\tilde{\boldsymbol{\Sigma}}_{ij} = \text{Cov}(X_{k_i}, X_{k_j})$. In particular, $\mathbf{X}_1 \sim N_q(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11})$ and $\mathbf{X}_2 \sim N_{p-q}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_{22})$.

b) If \mathbf{X}_1 and \mathbf{X}_2 are independent, then $\text{Cov}(\mathbf{X}_1, \mathbf{X}_2) = \boldsymbol{\Sigma}_{12} = E[(\mathbf{X}_1 - E(\mathbf{X}_1))(\mathbf{X}_2 - E(\mathbf{X}_2))^T] = \mathbf{0}$, a $q \times (p - q)$ matrix of zeroes.

c) If $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, then \mathbf{X}_1 and \mathbf{X}_2 are independent iff $\boldsymbol{\Sigma}_{12} = \mathbf{0}$.

d) If $\mathbf{X}_1 \sim N_q(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11})$ and $\mathbf{X}_2 \sim N_{p-q}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_{22})$ are independent, then

$$\begin{pmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{pmatrix} \sim N_p \left(\begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix}, \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma}_{22} \end{pmatrix} \right).$$

Theorem 4.3. The conditional distribution of a MVN is MVN. If $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, then the conditional distribution of \mathbf{X}_1 given that $\mathbf{X}_2 = \mathbf{x}_2$ is multivariate normal with mean $\boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{x}_2 - \boldsymbol{\mu}_2)$ and covariance matrix $\boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}$. That is,

$$\mathbf{X}_1 | \mathbf{X}_2 = \mathbf{x}_2 \sim N_q(\boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{x}_2 - \boldsymbol{\mu}_2), \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}).$$

Example 4.1. Let $p = 2$ and let $(Y, X)^T$ have a bivariate normal distribution. That is,

$$\begin{pmatrix} Y \\ X \end{pmatrix} \sim N_2 \left(\begin{pmatrix} \mu_Y \\ \mu_X \end{pmatrix}, \begin{pmatrix} \sigma_Y^2 & \text{Cov}(Y, X) \\ \text{Cov}(X, Y) & \sigma_X^2 \end{pmatrix} \right).$$

Also, recall that the population correlation between X and Y is given by

$$\rho(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{VAR}(X)}\sqrt{\text{VAR}(Y)}} = \frac{\sigma_{X,Y}}{\sigma_X\sigma_Y}$$

if $\sigma_X > 0$ and $\sigma_Y > 0$. Then $Y|X = x \sim N(E(Y|X = x), \text{VAR}(Y|X = x))$ where the conditional mean

$$E(Y|X = x) = \mu_Y + \text{Cov}(Y, X) \frac{1}{\sigma_X^2} (x - \mu_X) = \mu_Y + \rho(X, Y) \sqrt{\frac{\sigma_Y^2}{\sigma_X^2}} (x - \mu_X)$$

and the conditional variance

$$\begin{aligned} \text{VAR}(Y|X = x) &= \sigma_Y^2 - \text{Cov}(X, Y) \frac{1}{\sigma_X^2} \text{Cov}(X, Y) \\ &= \sigma_Y^2 - \rho(X, Y) \sqrt{\frac{\sigma_Y^2}{\sigma_X^2}} \rho(X, Y) \sqrt{\sigma_X^2} \sqrt{\sigma_Y^2} \\ &= \sigma_Y^2 - \rho^2(X, Y) \sigma_Y^2 = \sigma_Y^2 [1 - \rho^2(X, Y)]. \end{aligned}$$

Also $aX + bY$ is univariate normal with mean $a\mu_X + b\mu_Y$ and variance

$$a^2 \sigma_X^2 + b^2 \sigma_Y^2 + 2ab \text{Cov}(X, Y).$$

Remark 4.1. There are several common misconceptions. First, **it is not true that every linear combination $t^T \mathbf{X}$ of normal random variables is a normal random variable**, and **it is not true that all uncorrelated normal random variables are independent**. The key condition in Theorem 4.1b and Theorem 4.2c is that the joint distribution of \mathbf{X} is MVN. It is possible that X_1, X_2, \dots, X_p each has a marginal distribution that is univariate normal, but the joint distribution of \mathbf{X} is not MVN. Examine the following example from Rohatgi (1976, p. 229). Suppose that the joint pdf of X and Y is a mixture of two bivariate normal distributions both with $EX = EY = 0$ and $\text{VAR}(X) = \text{VAR}(Y) = 1$, but $\text{Cov}(X, Y) = \pm\rho$. Hence $f(x, y) =$

$$\begin{aligned} &\frac{1}{2} \frac{1}{2\pi\sqrt{1-\rho^2}} \exp\left(\frac{-1}{2(1-\rho^2)}(x^2 - 2\rho xy + y^2)\right) + \\ &\frac{1}{2} \frac{1}{2\pi\sqrt{1-\rho^2}} \exp\left(\frac{-1}{2(1-\rho^2)}(x^2 + 2\rho xy + y^2)\right) \equiv \frac{1}{2} f_1(x, y) + \frac{1}{2} f_2(x, y) \end{aligned}$$

where x and y are real and $0 < \rho < 1$. Since both marginal distributions of $f_i(x, y)$ are $N(0,1)$ for $i = 1$ and 2 by Theorem 4.2 a), the marginal distributions of X and Y are $N(0,1)$. Since $\int \int xy f_i(x, y) dx dy = \rho$ for $i = 1$ and $-\rho$ for $i = 2$, X and Y are uncorrelated, but X and Y are not independent since $f(x, y) \neq f_X(x) f_Y(y)$.

Remark 4.2. In Theorem 4.3, suppose that $\mathbf{X} = (Y, X_2, \dots, X_p)^T$. Let $X_1 = Y$ and $\mathbf{X}_2 = (X_2, \dots, X_p)^T$. Then $E[Y|\mathbf{X}_2] = \beta_1 + \beta_2 X_2 + \dots + \beta_p X_p$ and $\text{VAR}[Y|\mathbf{X}_2]$ is a constant that does not depend on \mathbf{X}_2 . Hence $Y|\mathbf{X}_2 = \beta_1 + \beta_2 X_2 + \dots + \beta_p X_p + e$ follows the multiple linear regression model.

4.2.2 The CLT and the Delta Method

The next three subsections will review large sample theory for the univariate case, then multivariate theory will be given.

Large sample theory, also called asymptotic theory, is used to approximate the distribution of an estimator when the sample size n is large. This theory is extremely useful if the exact sampling distribution of the estimator is complicated or unknown. To use this theory, one must determine what the estimator is estimating, the rate of convergence, the asymptotic distribution, and how large n must be for the approximation to be useful. Moreover, the (asymptotic) standard error (SE), an estimator of the asymptotic standard deviation, must be computable if the estimator is to be useful for inference. Often the bootstrap can be used to compute the SE.

Theorem 4.4: the Central Limit Theorem (CLT). Let Y_1, \dots, Y_n be iid with $E(Y) = \mu$ and $\text{VAR}(Y) = \sigma^2$. Let the sample mean $\bar{Y}_n = \frac{1}{n} \sum_{i=1}^n Y_i$. Then

$$\sqrt{n}(\bar{Y}_n - \mu) \xrightarrow{D} N(0, \sigma^2).$$

Hence

$$\sqrt{n} \left(\frac{\bar{Y}_n - \mu}{\sigma} \right) = \sqrt{n} \left(\frac{\sum_{i=1}^n Y_i - n\mu}{n\sigma} \right) \xrightarrow{D} N(0, 1).$$

Note that the sample mean is estimating the *population mean* μ with a \sqrt{n} convergence rate, the asymptotic distribution is normal, and the $\text{SE} = S/\sqrt{n}$ where S is the *sample standard deviation*. For distributions “close” to the normal distribution, the central limit theorem provides a good approximation if the sample size $n \geq 30$. Hesterberg (2014, pp. 41, 66) suggests $n \geq 5000$ is needed for moderately skewed distributions. A special case of the CLT is proven after Theorem 4.17.

Notation. The notation $X \sim Y$ and $X \stackrel{D}{=} Y$ both mean that the random variables X and Y have the same distribution. Hence $F_X(x) = F_Y(y)$ for all real y . The notation $Y_n \stackrel{D}{\rightarrow} X$ means that for large n we can approximate the cdf of Y_n by the cdf of X . The distribution of X is the limiting distribution or asymptotic distribution of Y_n . For the CLT, notice that

$$Z_n = \sqrt{n} \left(\frac{\bar{Y}_n - \mu}{\sigma} \right) = \left(\frac{\bar{Y}_n - \mu}{\sigma/\sqrt{n}} \right)$$

is the z-score of \bar{Y} . If $Z_n \stackrel{D}{\rightarrow} N(0, 1)$, then the notation $Z_n \approx N(0, 1)$, also written as $Z_n \sim AN(0, 1)$, means approximate the cdf of Z_n by the standard normal cdf. See Definition 4.3. Similarly, the notation

$$\bar{Y}_n \approx N(\mu, \sigma^2/n),$$

also written as $\bar{Y}_n \sim AN(\mu, \sigma^2/n)$, means approximate the cdf of \bar{Y}_n as if $\bar{Y}_n \sim N(\mu, \sigma^2/n)$. The distribution of X does not depend on n , but the approximate distribution $\bar{Y}_n \approx N(\mu, \sigma^2/n)$ does depend on n .

The two main applications of the CLT are to give the limiting distribution of $\sqrt{n}(\bar{Y}_n - \mu)$ and the limiting distribution of $\sqrt{n}(Y_n/n - \mu_X)$ for a random variable Y_n such that $Y_n = \sum_{i=1}^n X_i$ where the X_i are iid with $E(X) = \mu_X$ and $\text{VAR}(X) = \sigma_X^2$.

Example 4.2. a) Let Y_1, \dots, Y_n be iid $\text{Ber}(\rho)$. Then $E(Y) = \rho$ and $\text{VAR}(Y) = \rho(1 - \rho)$. (The Bernoulli (ρ) distribution is the binomial ($1, \rho$) distribution.) Hence

$$\sqrt{n}(\bar{Y}_n - \rho) \xrightarrow{D} N(0, \rho(1 - \rho))$$

by the CLT.

b) Now suppose that $Y_n \sim \text{BIN}(n, \rho)$. Then $Y_n \stackrel{D}{=} \sum_{i=1}^n X_i$ where X_1, \dots, X_n are iid $\text{Ber}(\rho)$. Hence

$$\sqrt{n} \left(\frac{Y_n}{n} - \rho \right) \xrightarrow{D} N(0, \rho(1 - \rho))$$

since

$$\sqrt{n} \left(\frac{Y_n}{n} - \rho \right) \stackrel{D}{=} \sqrt{n}(\bar{X}_n - \rho) \xrightarrow{D} N(0, \rho(1 - \rho))$$

by a).

c) Now suppose that $Y_n \sim \text{BIN}(k_n, \rho)$ where $k_n \rightarrow \infty$ as $n \rightarrow \infty$. Then

$$\sqrt{k_n} \left(\frac{Y_n}{k_n} - \rho \right) \approx N(0, \rho(1 - \rho))$$

or

$$\frac{Y_n}{k_n} \approx N \left(\rho, \frac{\rho(1 - \rho)}{k_n} \right) \quad \text{or} \quad Y_n \approx N(k_n \rho, k_n \rho(1 - \rho)).$$

Theorem 4.5: the Delta Method. If g does not depend on n , $g'(\theta) \neq 0$, and

$$\sqrt{n}(T_n - \theta) \xrightarrow{D} N(0, \sigma^2),$$

then

$$\sqrt{n}(g(T_n) - g(\theta)) \xrightarrow{D} N(0, \sigma^2 [g'(\theta)]^2).$$

Example 4.3. Let Y_1, \dots, Y_n be iid with $E(Y) = \mu$ and $\text{VAR}(Y) = \sigma^2$. Then by the CLT,

$$\sqrt{n}(\bar{Y}_n - \mu) \xrightarrow{D} N(0, \sigma^2).$$

Let $g(\mu) = \mu^2$. Then $g'(\mu) = 2\mu \neq 0$ for $\mu \neq 0$. Hence

$$\sqrt{n}((\bar{Y}_n)^2 - \mu^2) \xrightarrow{D} N(0, 4\sigma^2\mu^2)$$

for $\mu \neq 0$ by the delta method.

Example 4.4. Let $X \sim \text{Binomial}(n, p)$ where the positive integer n is large and $0 < p < 1$. Find the limiting distribution of $\sqrt{n} \left[\left(\frac{X}{n} \right)^2 - p^2 \right]$.

Solution. Example 4.2b gives the limiting distribution of $\sqrt{n}(\frac{X}{n} - p)$. Let $g(p) = p^2$. Then $g'(p) = 2p$ and by the delta method,

$$\sqrt{n} \left[\left(\frac{X}{n} \right)^2 - p^2 \right] = \sqrt{n} \left(g \left(\frac{X}{n} \right) - g(p) \right) \xrightarrow{D}$$

$$N(0, p(1-p)(g'(p))^2) = N(0, p(1-p)4p^2) = N(0, 4p^3(1-p)).$$

Example 4.5. Let $X_n \sim \text{Poisson}(n\lambda)$ where the positive integer n is large and $\lambda > 0$.

a) Find the limiting distribution of $\sqrt{n} \left(\frac{X_n}{n} - \lambda \right)$.

b) Find the limiting distribution of $\sqrt{n} \left[\sqrt{\frac{X_n}{n}} - \sqrt{\lambda} \right]$.

Solution. a) $X_n \stackrel{D}{=} \sum_{i=1}^n Y_i$ where the Y_i are iid Poisson(λ). Hence $E(Y) = \lambda = \text{Var}(Y)$. Thus by the CLT,

$$\sqrt{n} \left(\frac{X_n}{n} - \lambda \right) \stackrel{D}{=} \sqrt{n} \left(\frac{\sum_{i=1}^n Y_i}{n} - \lambda \right) \xrightarrow{D} N(0, \lambda).$$

b) Let $g(\lambda) = \sqrt{\lambda}$. Then $g'(\lambda) = \frac{1}{2\sqrt{\lambda}}$ and by the delta method,

$$\sqrt{n} \left[\sqrt{\frac{X_n}{n}} - \sqrt{\lambda} \right] = \sqrt{n} \left(g \left(\frac{X_n}{n} \right) - g(\lambda) \right) \xrightarrow{D}$$

$$N(0, \lambda (g'(\lambda))^2) = N \left(0, \lambda \frac{1}{4\lambda} \right) = N \left(0, \frac{1}{4} \right).$$

Example 4.6. Let Y_1, \dots, Y_n be independent and identically distributed (iid) from a Gamma(α, β) distribution.

a) Find the limiting distribution of $\sqrt{n} (\bar{Y} - \alpha\beta)$.

b) Find the limiting distribution of $\sqrt{n} ((\bar{Y})^2 - c)$ for appropriate constant c .

Solution: a) Since $E(Y) = \alpha\beta$ and $V(Y) = \alpha\beta^2$, by the CLT $\sqrt{n} (\bar{Y} - \alpha\beta) \xrightarrow{D} N(0, \alpha\beta^2)$.
 b) Let $\mu = \alpha\beta$ and $\sigma^2 = \alpha\beta^2$. Let $g(\mu) = \mu^2$ so $g'(\mu) = 2\mu$ and $[g'(\mu)]^2 = 4\mu^2 = 4\alpha^2\beta^2$. Then by the delta method, $\sqrt{n} ((\bar{Y})^2 - c) \xrightarrow{D} N(0, \sigma^2[g'(\mu)]^2) = N(0, 4\alpha^3\beta^4)$ where $c = \mu^2 = \alpha^2\beta^2$.

4.2.3 Modes of Convergence and Consistency

Definition 4.3. Let $\{Z_n, n = 1, 2, \dots\}$ be a sequence of random variables with cdfs F_n , and let X be a random variable with cdf F . Then Z_n **converges in distribution to X** , written

$$Z_n \xrightarrow{D} X,$$

or Z_n *converges in law to X* , written $Z_n \xrightarrow{L} X$, if

$$\lim_{n \rightarrow \infty} F_n(t) = F(t)$$

at each continuity point t of F . The distribution of X is called the **limiting distribution** or the **asymptotic distribution** of Z_n .

An important fact is that **the limiting distribution does not depend on the sample size n** . Notice that the CLT and delta method give the limiting distributions of $Z_n = \sqrt{n}(\bar{Y}_n - \mu)$ and $Z_n = \sqrt{n}(g(\bar{T}_n) - g(\theta))$, respectively.

Convergence in distribution is useful if the distribution of X_n is unknown or complicated and the distribution of X is easy to use. Then for large n we can approximate the probability that X_n is in an interval by the probability that X is in the interval. To see this, notice that if $X_n \xrightarrow{D} X$, then $P(a < X_n \leq b) = F_n(b) - F_n(a) \rightarrow F(b) - F(a) = P(a < X \leq b)$ if F is continuous at a and b .

Warning: convergence in distribution says that the cdf $F_n(t)$ of X_n gets close to the cdf of $F(t)$ of X as $n \rightarrow \infty$ provided that t is a continuity point of F . Hence for any $\epsilon > 0$ there exists N_ϵ such that if $n > N_\epsilon$, then $|F_n(t) - F(t)| < \epsilon$. Notice that N_ϵ depends on the value of t . Convergence in distribution does not imply that the random variables $X_n \equiv X_n(\omega)$ converge to the random variable $X \equiv X(\omega)$ for all ω .

Example 4.7. Suppose that $X_n \sim U(-1/n, 1/n)$. Then the cdf $F_n(x)$ of X_n is

$$F_n(x) = \begin{cases} 0, & x \leq -\frac{1}{n} \\ \frac{nx}{2} + \frac{1}{2}, & -\frac{1}{n} \leq x \leq \frac{1}{n} \\ 1, & x \geq \frac{1}{n}. \end{cases}$$

Sketching $F_n(x)$ shows that it has a line segment rising from 0 at $x = -1/n$ to 1 at $x = 1/n$ and that $F_n(0) = 0.5$ for all $n \geq 1$. Examining the cases $x < 0$, $x = 0$, and $x > 0$ shows that as $n \rightarrow \infty$,

$$F_n(x) \rightarrow \begin{cases} 0, & x < 0 \\ \frac{1}{2}, & x = 0 \\ 1, & x > 0. \end{cases}$$

Notice that the right hand side is not a cdf since right continuity does not hold at $x = 0$. Notice that if X is a random variable such that $P(X = 0) = 1$, then X has cdf

$$F_X(x) = \begin{cases} 0, & x < 0 \\ 1, & x \geq 0. \end{cases}$$

Since $x = 0$ is the only discontinuity point of $F_X(x)$ and since $F_n(x) \rightarrow F_X(x)$ for all continuity points of $F_X(x)$ (i.e. for $x \neq 0$),

$$X_n \xrightarrow{D} X.$$

Example 4.8. Suppose $Y_n \sim U(0, n)$. Then $F_n(t) = t/n$ for $0 < t \leq n$ and $F_n(t) = 0$ for $t \leq 0$. Hence $\lim_{n \rightarrow \infty} F_n(t) = 0$ for $t \leq 0$. If $t > 0$ and $n > t$, then $F_n(t) = t/n \rightarrow 0$ as $n \rightarrow \infty$. Thus $\lim_{n \rightarrow \infty} F_n(t) = 0$ for all t , and Y_n does not converge in distribution to any random variable Y since $H(t) \equiv 0$ is not a cdf.

Definition 4.4. A sequence of random variables X_n converges in distribution to a constant $\tau(\theta)$, written

$$X_n \xrightarrow{D} \tau(\theta), \quad \text{if } X_n \xrightarrow{D} X$$

where $P(X = \tau(\theta)) = 1$. The distribution of the random variable X is said to be *degenerate at $\tau(\theta)$* or to be a *point mass at $\tau(\theta)$* .

Definition 4.5. A sequence of random variables X_n converges in probability to a constant $\tau(\theta)$, written

$$X_n \xrightarrow{P} \tau(\theta),$$

if for every $\epsilon > 0$,

$$\lim_{n \rightarrow \infty} P(|X_n - \tau(\theta)| < \epsilon) = 1 \quad \text{or, equivalently,} \quad \lim_{n \rightarrow \infty} P(|X_n - \tau(\theta)| \geq \epsilon) = 0.$$

The sequence X_n **converges in probability to X** , written

$$X_n \xrightarrow{P} X,$$

if $X_n - X \xrightarrow{P} 0$.

Notice that $X_n \xrightarrow{P} X$ if for every $\epsilon > 0$,

$$\lim_{n \rightarrow \infty} P(|X_n - X| < \epsilon) = 1, \quad \text{or, equivalently,} \quad \lim_{n \rightarrow \infty} P(|X_n - X| \geq \epsilon) = 0.$$

Definition 4.6. Let the *parameter space* Θ be the set of possible values of θ . A sequence of estimators T_n of $\tau(\theta)$ is **consistent** for $\tau(\theta)$ if

$$T_n \xrightarrow{P} \tau(\theta)$$

for every $\theta \in \Theta$. If T_n is consistent for $\tau(\theta)$, then T_n is a **consistent estimator** of $\tau(\theta)$.

Consistency is a weak property that is usually satisfied by good estimators. T_n is a consistent estimator for $\tau(\theta)$ if the probability that T_n falls in any neighborhood of $\tau(\theta)$ goes to one, regardless of the value of $\theta \in \Theta$.

Definition 4.7. For a real number $r > 0$, Y_n *converges in r th mean* to a random variable Y , written

$$Y_n \xrightarrow{r} Y,$$

if

$$E(|Y_n - Y|^r) \rightarrow 0$$

as $n \rightarrow \infty$. In particular, if $r = 2$, Y_n **converges in quadratic mean** to Y , written

$$Y_n \xrightarrow{2} Y \quad \text{or} \quad Y_n \xrightarrow{\text{qm}} Y,$$

if

$$E[(Y_n - Y)^2] \rightarrow 0$$

as $n \rightarrow \infty$.

Theorem 4.6: Generalized Chebyshev's Inequality. Let $u : \mathbb{R} \rightarrow [0, \infty)$ be a nonnegative function. If $E[u(Y)]$ exists then for any $c > 0$,

$$P[u(Y) \geq c] \leq \frac{E[u(Y)]}{c}.$$

If $\mu = E(Y)$ exists, then taking $u(y) = |y - \mu|^r$ and $\tilde{c} = c^r$ gives **Markov's Inequality**: for $r > 0$ and any $c > 0$,

$$P[|Y - \mu| \geq c] = P[|Y - \mu|^r \geq c^r] \leq \frac{E[|Y - \mu|^r]}{c^r}.$$

If $r = 2$ and $\sigma^2 = \text{VAR}(Y)$ exists, then we obtain **Chebyshev's Inequality**:

$$P[|Y - \mu| \geq c] \leq \frac{\text{VAR}(Y)}{c^2}.$$

Proof. The proof is given for pdfs. For pmfs, replace the integrals by sums. Now

$$\begin{aligned} E[u(Y)] &= \int_{\mathbb{R}} u(y)f(y)dy = \int_{\{y:u(y)\geq c\}} u(y)f(y)dy + \int_{\{y:u(y)<c\}} u(y)f(y)dy \\ &\geq \int_{\{y:u(y)\geq c\}} u(y)f(y)dy \end{aligned}$$

since the integrand $u(y)f(y) \geq 0$. Hence

$$E[u(Y)] \geq c \int_{\{y:u(y)\geq c\}} f(y)dy = cP[u(Y) \geq c]. \quad \square$$

The following theorem gives sufficient conditions for T_n to be a consistent estimator of $\tau(\theta)$. Notice that $E_{\theta}[(T_n - \tau(\theta))^2] = MSE_{\tau(\theta)}(T_n) \rightarrow 0$ for all $\theta \in \Theta$ is equivalent to $T_n \xrightarrow{qm} \tau(\theta)$ for all $\theta \in \Theta$.

Theorem 4.7. a) If

$$\lim_{n \rightarrow \infty} MSE_{\tau(\theta)}(T_n) = 0$$

for all $\theta \in \Theta$, then T_n is a consistent estimator of $\tau(\theta)$.

b) If

$$\lim_{n \rightarrow \infty} \text{VAR}_{\theta}(T_n) = 0 \quad \text{and} \quad \lim_{n \rightarrow \infty} E_{\theta}(T_n) = \tau(\theta)$$

for all $\theta \in \Theta$, then T_n is a consistent estimator of $\tau(\theta)$.

Proof. a) Using Theorem 4.6 with $Y = T_n$, $u(T_n) = (T_n - \tau(\theta))^2$ and $c = \epsilon^2$ shows that for any $\epsilon > 0$,

$$P_{\theta}(|T_n - \tau(\theta)| \geq \epsilon) = P_{\theta}[(T_n - \tau(\theta))^2 \geq \epsilon^2] \leq \frac{E_{\theta}[(T_n - \tau(\theta))^2]}{\epsilon^2}.$$

Hence

$$\lim_{n \rightarrow \infty} E_{\theta}[(T_n - \tau(\theta))^2] = \lim_{n \rightarrow \infty} MSE_{\tau(\theta)}(T_n) \rightarrow 0$$

is a sufficient condition for T_n to be a consistent estimator of $\tau(\theta)$.

b) Recall that

$$MSE_{\tau(\theta)}(T_n) = \text{VAR}_{\theta}(T_n) + [\text{Bias}_{\tau(\theta)}(T_n)]^2$$

where $\text{Bias}_{\tau(\theta)}(T_n) = E_{\theta}(T_n) - \tau(\theta)$. Since $MSE_{\tau(\theta)}(T_n) \rightarrow 0$ if both $\text{VAR}_{\theta}(T_n) \rightarrow 0$ and $\text{Bias}_{\tau(\theta)}(T_n) = E_{\theta}(T_n) - \tau(\theta) \rightarrow 0$, the result follows from a). \square

The following result shows estimators that converge at a \sqrt{n} rate are consistent. Use this result and the delta method to show that $g(T_n)$ is a consistent

estimator of $g(\theta)$. Note that b) follows from a) with $X_\theta \sim N(0, v(\theta))$. The WLLN shows that \bar{Y} is a consistent estimator of $E(Y) = \mu$ if $E(Y)$ exists.

Theorem 4.8. a) Let X_θ be a random variable with distribution depending on θ , and $0 < \delta \leq 1$. If

$$n^\delta(T_n - \tau(\theta)) \xrightarrow{D} X_\theta$$

then $T_n \xrightarrow{P} \tau(\theta)$.

b) If

$$\sqrt{n}(T_n - \tau(\theta)) \xrightarrow{D} N(0, v(\theta))$$

for all $\theta \in \Theta$, then T_n is a consistent estimator of $\tau(\theta)$.

Definition 4.8. A sequence of random variables X_n converges almost everywhere (or almost surely, or with probability 1) to X if

$$P(\lim_{n \rightarrow \infty} X_n = X) = 1.$$

This type of convergence will be denoted by

$$X_n \xrightarrow{ae} X.$$

Notation such as “ X_n converges to X ae” will also be used. Sometimes “ae” will be replaced with “as” or “wp1.” We say that X_n converges almost everywhere to $\tau(\theta)$, written

$$X_n \xrightarrow{ae} \tau(\theta),$$

if $P(\lim_{n \rightarrow \infty} X_n = \tau(\theta)) = 1$.

Theorem 4.9. Let Y_n be a sequence of iid random variables with $E(Y_i) = \mu$. Then

a) **Strong Law of Large Numbers (SLLN):** $\bar{Y}_n \xrightarrow{ae} \mu$, and

b) **Weak Law of Large Numbers (WLLN):** $\bar{Y}_n \xrightarrow{P} \mu$.

Proof of WLLN when $V(Y_i) = \sigma^2$: By Chebyshev’s inequality, for every $\epsilon > 0$,

$$P(|\bar{Y}_n - \mu| \geq \epsilon) \leq \frac{V(\bar{Y}_n)}{\epsilon^2} = \frac{\sigma^2}{n\epsilon^2} \rightarrow 0$$

as $n \rightarrow \infty$. \square

In proving consistency results, there is an infinite sequence of estimators that depend on the sample size n . Hence the subscript n will be added to the estimators.

Definition 4.9. Lehmann (1999, pp. 53-54): a) A sequence of random variables W_n is *tight* or *bounded in probability*, written $W_n = O_P(1)$, if for every $\epsilon > 0$ there exist positive constants D_ϵ and N_ϵ such that

$$P(|W_n| \leq D_\epsilon) \geq 1 - \epsilon$$

for all $n \geq N_\epsilon$. Also $W_n = O_P(X_n)$ if $|W_n/X_n| = O_P(1)$.

b) The sequence $W_n = o_P(n^{-\delta})$ if $n^\delta W_n = o_P(1)$ which means that

$$n^\delta W_n \xrightarrow{P} 0.$$

c) W_n has the same order as X_n in probability, written $W_n \asymp_P X_n$, if for every $\epsilon > 0$ there exist positive constants N_ϵ and $0 < d_\epsilon < D_\epsilon$ such that

$$P\left(d_\epsilon \leq \left|\frac{W_n}{X_n}\right| \leq D_\epsilon\right) = P\left(\frac{1}{D_\epsilon} \leq \left|\frac{X_n}{W_n}\right| \leq \frac{1}{d_\epsilon}\right) \geq 1 - \epsilon$$

for all $n \geq N_\epsilon$.

d) Similar notation is used for a $k \times r$ matrix $\mathbf{A}_n = \mathbf{A} = [a_{i,j}(n)]$ if each element $a_{i,j}(n)$ has the desired property. For example, $\mathbf{A} = O_P(n^{-1/2})$ if each $a_{i,j}(n) = O_P(n^{-1/2})$.

Definition 4.10. Let $W_n = \|\hat{\boldsymbol{\mu}}_n - \boldsymbol{\mu}\|$.

a) If $W_n \asymp_P n^{-\delta}$ for some $\delta > 0$, then both W_n and $\hat{\boldsymbol{\mu}}_n$ have (tightness) rate n^δ .

b) If there exists a constant κ such that

$$n^\delta(W_n - \kappa) \xrightarrow{D} X$$

for some nondegenerate random variable X , then both W_n and $\hat{\boldsymbol{\mu}}_n$ have convergence rate n^δ .

Theorem 4.10. Suppose there exists a constant κ such that

$$n^\delta(W_n - \kappa) \xrightarrow{D} X.$$

a) Then $W_n = O_P(n^{-\delta})$.

b) If X is not degenerate, then $W_n \asymp_P n^{-\delta}$.

The above result implies that if W_n has convergence rate n^δ , then W_n has tightness rate n^δ , and the term “tightness” will often be omitted. Part a) is proved, for example, in Lehmann (1999, p. 67).

The following result shows that if $W_n \asymp_P X_n$, then $X_n \asymp_P W_n$, $W_n = O_P(X_n)$, and $X_n = O_P(W_n)$. Notice that if $W_n = O_P(n^{-\delta})$, then n^δ is a lower bound on the rate of W_n . As an example, if the CLT holds then $\bar{Y}_n = O_P(n^{-1/3})$, but $\bar{Y}_n \asymp_P n^{-1/2}$.

Theorem 4.11. a) If $W_n \asymp_P X_n$, then $X_n \asymp_P W_n$.

b) If $W_n \asymp_P X_n$, then $W_n = O_P(X_n)$.

c) If $W_n \asymp_P X_n$, then $X_n = O_P(W_n)$.

d) $W_n \asymp_P X_n$ iff $W_n = O_P(X_n)$ and $X_n = O_P(W_n)$.

Proof. a) Since $W_n \asymp_P X_n$,

$$P\left(d_\epsilon \leq \left|\frac{W_n}{X_n}\right| \leq D_\epsilon\right) = P\left(\frac{1}{D_\epsilon} \leq \left|\frac{X_n}{W_n}\right| \leq \frac{1}{d_\epsilon}\right) \geq 1 - \epsilon$$

for all $n \geq N_\epsilon$. Hence $X_n \asymp_P W_n$.

b) Since $W_n \asymp_P X_n$,

$$P(|W_n| \leq |X_n D_\epsilon|) \geq P\left(d_\epsilon \leq \left|\frac{W_n}{X_n}\right| \leq D_\epsilon\right) \geq 1 - \epsilon$$

for all $n \geq N_\epsilon$. Hence $W_n = O_P(X_n)$.

c) Follows by a) and b).

d) If $W_n \asymp_P X_n$, then $W_n = O_P(X_n)$ and $X_n = O_P(W_n)$ by b) and c). Now suppose $W_n = O_P(X_n)$ and $X_n = O_P(W_n)$. Then

$$P(|W_n| \leq |X_n| D_{\epsilon/2}) \geq 1 - \epsilon/2$$

for all $n \geq N_1$, and

$$P(|X_n| \leq |W_n| 1/d_{\epsilon/2}) \geq 1 - \epsilon/2$$

for all $n \geq N_2$. Hence

$$P(A) \equiv P\left(\left|\frac{W_n}{X_n}\right| \leq D_{\epsilon/2}\right) \geq 1 - \epsilon/2$$

and

$$P(B) \equiv P\left(d_{\epsilon/2} \leq \left|\frac{W_n}{X_n}\right|\right) \geq 1 - \epsilon/2$$

for all $n \geq N = \max(N_1, N_2)$. Since $P(A \cap B) = P(A) + P(B) - P(A \cup B) \geq P(A) + P(B) - 1$,

$$P(A \cap B) = P(d_{\epsilon/2} \leq \left|\frac{W_n}{X_n}\right| \leq D_{\epsilon/2}) \geq 1 - \epsilon/2 + 1 - \epsilon/2 - 1 = 1 - \epsilon$$

for all $n \geq N$. Hence $W_n \asymp_P X_n$. \square

The following result is used to prove the following Theorem 4.13 which says that if there are K estimators $T_{j,n}$ of a parameter β , such that $\|T_{j,n} - \beta\| = O_P(n^{-\delta})$ where $0 < \delta \leq 1$, and if T_n^* picks one of these estimators, then $\|T_n^* - \beta\| = O_P(n^{-\delta})$.

Theorem 4.12: Pratt (1959). Let $X_{1,n}, \dots, X_{K,n}$ each be $O_P(1)$ where K is fixed. Suppose $W_n = X_{i_n,n}$ for some $i_n \in \{1, \dots, K\}$. Then

$$W_n = O_P(1). \tag{4.5}$$

Proof.

$$P(\max\{X_{1,n}, \dots, X_{K,n}\} \leq x) = P(X_{1,n} \leq x, \dots, X_{K,n} \leq x) \leq$$

$$F_{W_n}(x) \leq P(\min\{X_{1,n}, \dots, X_{K,n}\} \leq x) = 1 - P(X_{1,n} > x, \dots, X_{K,n} > x).$$

Since K is finite, there exists $B > 0$ and N such that $P(X_{i,n} \leq B) > 1 - \epsilon/2K$ and $P(X_{i,n} > -B) > 1 - \epsilon/2K$ for all $n > N$ and $i = 1, \dots, K$. Bonferroni's inequality states that $P(\cap_{i=1}^K A_i) \geq \sum_{i=1}^K P(A_i) - (K - 1)$. Thus

$$F_{W_n}(B) \geq P(X_{1,n} \leq B, \dots, X_{K,n} \leq B) \geq$$

$$K(1 - \epsilon/2K) - (K - 1) = K - \epsilon/2 - K + 1 = 1 - \epsilon/2$$

and

$$-F_{W_n}(-B) \geq -1 + P(X_{1,n} > -B, \dots, X_{K,n} > -B) \geq$$

$$-1 + K(1 - \epsilon/2K) - (K - 1) = -1 + K - \epsilon/2 - K + 1 = -\epsilon/2.$$

Hence

$$F_{W_n}(B) - F_{W_n}(-B) \geq 1 - \epsilon \text{ for } n > N. \quad \square$$

Theorem 4.13. Suppose $\|T_{j,n} - \beta\| = O_P(n^{-\delta})$ for $j = 1, \dots, K$ where $0 < \delta \leq 1$. Let $T_n^* = T_{i_n,n}$ for some $i_n \in \{1, \dots, K\}$ where, for example, $T_{i_n,n}$ is the $T_{j,n}$ that minimized some criterion function. Then

$$\|T_n^* - \beta\| = O_P(n^{-\delta}). \quad (4.6)$$

Proof. Let $X_{j,n} = n^\delta \|T_{j,n} - \beta\|$. Then $X_{j,n} = O_P(1)$ so by Theorem 4.12, $n^\delta \|T_n^* - \beta\| = O_P(1)$. Hence $\|T_n^* - \beta\| = O_P(n^{-\delta})$. \square

4.2.4 Slutsky's Theorem and Related Results

Theorem 4.14: Slutsky's Theorem. Suppose $Y_n \xrightarrow{D} Y$ and $W_n \xrightarrow{P} w$ for some constant w . Then

- a) $Y_n + W_n \xrightarrow{D} Y + w$,
- b) $Y_n W_n \xrightarrow{D} wY$, and
- c) $Y_n/W_n \xrightarrow{D} Y/w$ if $w \neq 0$.

Theorem 4.15. a) If $X_n \xrightarrow{P} X$, then $X_n \xrightarrow{D} X$.

b) If $X_n \xrightarrow{ae} X$, then $X_n \xrightarrow{P} X$ and $X_n \xrightarrow{D} X$.

c) If $X_n \xrightarrow{r} X$, then $X_n \xrightarrow{P} X$ and $X_n \xrightarrow{D} X$.

d) $X_n \xrightarrow{P} \tau(\theta)$ iff $X_n \xrightarrow{D} \tau(\theta)$.

e) If $X_n \xrightarrow{P} \theta$ and τ is continuous at θ , then $\tau(X_n) \xrightarrow{P} \tau(\theta)$.

f) If $X_n \xrightarrow{D} \theta$ and τ is continuous at θ , then $\tau(X_n) \xrightarrow{D} \tau(\theta)$.

Suppose that for all $\theta \in \Theta$, $T_n \xrightarrow{D} \tau(\theta)$, $T_n \xrightarrow{r} \tau(\theta)$, or $T_n \xrightarrow{ae} \tau(\theta)$. Then T_n is a consistent estimator of $\tau(\theta)$ by Theorem 4.15. We are assuming that the function τ does not depend on n .

Example 4.9. Let Y_1, \dots, Y_n be iid with mean $E(Y_i) = \mu$ and variance $V(Y_i) = \sigma^2$. Then the sample mean \bar{Y}_n is a consistent estimator of μ since i) the SLLN holds (use Theorems 4.9 and 4.15), ii) the WLLN holds, and iii) the CLT holds (use Theorem 4.8). Since

$$\lim_{n \rightarrow \infty} \text{VAR}_\mu(\bar{Y}_n) = \lim_{n \rightarrow \infty} \sigma^2/n = 0 \quad \text{and} \quad \lim_{n \rightarrow \infty} E_\mu(\bar{Y}_n) = \mu,$$

\bar{Y}_n is also a consistent estimator of μ by Theorem 4.7b. By the delta method and Theorem 4.8b, $T_n = g(\bar{Y}_n)$ is a consistent estimator of $g(\mu)$ if $g'(\mu) \neq 0$ for all $\mu \in \Theta$. By Theorem 4.15e, $g(\bar{Y}_n)$ is a consistent estimator of $g(\mu)$ if g is continuous at μ for all $\mu \in \Theta$.

Theorem 4.16. Assume that the function g does not depend on n .

a) **Generalized Continuous Mapping Theorem:** If $X_n \xrightarrow{D} X$ and the function g is such that $P[X \in C(g)] = 1$ where $C(g)$ is the set of points where g is continuous, then $g(X_n) \xrightarrow{D} g(X)$.

b) **Continuous Mapping Theorem:** If $X_n \xrightarrow{D} X$ and the function g is continuous, then $g(X_n) \xrightarrow{D} g(X)$.

Remark 4.3. For Theorem 4.15, a) follows from Slutsky's Theorem by taking $Y_n \equiv X = Y$ and $W_n = X_n - X$. Then $Y_n \xrightarrow{D} Y = X$ and $W_n \xrightarrow{P} 0$. Hence $X_n = Y_n + W_n \xrightarrow{D} Y + 0 = X$. The convergence in distribution parts of b) and c) follow from a). Part f) follows from d) and e). Part e) implies that if T_n is a consistent estimator of θ and τ is a continuous function, then $\tau(T_n)$ is a consistent estimator of $\tau(\theta)$. Theorem 4.16 says that convergence in distribution is preserved by continuous functions, and even some discontinuities are allowed as long as the set of continuity points is assigned probability 1 by the asymptotic distribution. Equivalently, the set of discontinuity points is assigned probability 0.

Example 4.10. (Ferguson 1996, p. 40): If $X_n \xrightarrow{D} X$, then $1/X_n \xrightarrow{D} 1/X$ if X is a continuous random variable since $P(X = 0) = 0$ and $x = 0$ is the only discontinuity point of $g(x) = 1/x$.

Example 4.11. Show that if $Y_n \sim t_n$, a t distribution with n degrees of freedom, then $Y_n \xrightarrow{D} Z$ where $Z \sim N(0, 1)$.

Solution: $Y_n \stackrel{D}{=} Z/\sqrt{V_n/n}$ where $Z \perp V_n \sim \chi_n^2$. If $W_n = \sqrt{V_n/n} \xrightarrow{P} 1$, then the result follows by Slutsky's Theorem. But $V_n \stackrel{D}{=} \sum_{i=1}^n X_i$ where the

iid $X_i \sim \chi_1^2$. Hence $V_n/n \xrightarrow{P} 1$ by the WLLN and $\sqrt{V_n/n} \xrightarrow{P} 1$ by Theorem 4.15e.

Theorem 4.17: Continuity Theorem. Let Y_n be sequence of random variables with characteristic functions $\phi_n(t)$. Let Y be a random variable with characteristic function (cf) $\phi(t)$.

a)

$$Y_n \xrightarrow{D} Y \text{ iff } \phi_n(t) \rightarrow \phi(t) \forall t \in \mathbb{R}.$$

b) Also assume that Y_n has moment generating function (mgf) m_n and Y has mgf m . Assume that all of the mgfs m_n and m are defined on $|t| \leq d$ for some $d > 0$. Then if $m_n(t) \rightarrow m(t)$ as $n \rightarrow \infty$ for all $|t| < c$ where $0 < c < d$, then $Y_n \xrightarrow{D} Y$.

Application: Proof of a Special Case of the CLT. Following Rohatgi (1984, pp. 569-9), let Y_1, \dots, Y_n be iid with mean μ , variance σ^2 , and mgf $m_Y(t)$ for $|t| < t_o$. Then

$$Z_i = \frac{Y_i - \mu}{\sigma}$$

has mean 0, variance 1, and mgf $m_Z(t) = \exp(-t\mu/\sigma)m_Y(t/\sigma)$ for $|t| < \sigma t_o$. We want to show that

$$W_n = \sqrt{n} \left(\frac{\bar{Y}_n - \mu}{\sigma} \right) \xrightarrow{D} N(0, 1).$$

Notice that $W_n =$

$$n^{-1/2} \sum_{i=1}^n Z_i = n^{-1/2} \sum_{i=1}^n \left(\frac{Y_i - \mu}{\sigma} \right) = n^{-1/2} \frac{\sum_{i=1}^n Y_i - n\mu}{\sigma} = \frac{n^{-1/2}}{\frac{1}{n}} \frac{\bar{Y}_n - \mu}{\sigma}.$$

Thus

$$\begin{aligned} m_{W_n}(t) &= E(e^{tW_n}) = E[\exp(tn^{-1/2} \sum_{i=1}^n Z_i)] = E[\exp(\sum_{i=1}^n tZ_i/\sqrt{n})] \\ &= \prod_{i=1}^n E[e^{tZ_i/\sqrt{n}}] = \prod_{i=1}^n m_Z(t/\sqrt{n}) = [m_Z(t/\sqrt{n})]^n. \end{aligned}$$

Set $\psi(x) = \log(m_Z(x))$. Then

$$\log[m_{W_n}(t)] = n \log[m_Z(t/\sqrt{n})] = n\psi(t/\sqrt{n}) = \frac{\psi(t/\sqrt{n})}{\frac{1}{n}}.$$

Now $\psi(0) = \log[m_Z(0)] = \log(1) = 0$. Thus by L'Hôpital's rule (where the derivative is with respect to n), $\lim_{n \rightarrow \infty} \log[m_{W_n}(t)] =$

$$\lim_{n \rightarrow \infty} \frac{\psi(t/\sqrt{n})}{\frac{1}{n}} = \lim_{n \rightarrow \infty} \frac{\psi'(t/\sqrt{n}) \left[\frac{-t/2}{n^{3/2}} \right]}{\left(\frac{-1}{n^2} \right)} = \frac{t}{2} \lim_{n \rightarrow \infty} \frac{\psi'(t/\sqrt{n})}{\frac{1}{\sqrt{n}}}.$$

Now

$$\psi'(0) = \frac{m'_Z(0)}{m_Z(0)} = E(Z_i)/1 = 0,$$

so L'Hôpital's rule can be applied again, giving $\lim_{n \rightarrow \infty} \log[m_{W_n}(t)] =$

$$\frac{t}{2} \lim_{n \rightarrow \infty} \frac{\psi''(t/\sqrt{n}) \left[\frac{-t}{2n^{3/2}} \right]}{\left(\frac{-1}{2n^{3/2}} \right)} = \frac{t^2}{2} \lim_{n \rightarrow \infty} \psi''(t/\sqrt{n}) = \frac{t^2}{2} \psi''(0).$$

Now

$$\psi''(t) = \frac{d m'_Z(t)}{dt m_Z(t)} = \frac{m''_Z(t)m_Z(t) - (m'_Z(t))^2}{[m_Z(t)]^2}.$$

So

$$\psi''(0) = m''_Z(0) - [m'_Z(0)]^2 = E(Z_i^2) - [E(Z_i)]^2 = 1.$$

Hence $\lim_{n \rightarrow \infty} \log[m_{W_n}(t)] = t^2/2$ and

$$\lim_{n \rightarrow \infty} m_{W_n}(t) = \exp(t^2/2)$$

which is the $N(0,1)$ mgf. Thus by the continuity theorem,

$$W_n = \sqrt{n} \left(\frac{\bar{Y}_n - \mu}{\sigma} \right) \xrightarrow{D} N(0, 1). \quad \square$$

4.2.5 Multivariate Limit Theorems

Many of the univariate results of the previous 3 subsections can be extended to random vectors. For the limit theorems, the vector \mathbf{X} is typically a $k \times 1$ column vector and \mathbf{X}^T is a row vector. Let $\|\mathbf{x}\| = \sqrt{x_1^2 + \cdots + x_k^2}$ be the Euclidean norm of \mathbf{x} .

Definition 4.11. Let \mathbf{X}_n be a sequence of random vectors with joint cdfs $F_n(\mathbf{x})$ and let \mathbf{X} be a random vector with joint cdf $F(\mathbf{x})$.

a) \mathbf{X}_n converges in distribution to \mathbf{X} , written $\mathbf{X}_n \xrightarrow{D} \mathbf{X}$, if $F_n(\mathbf{x}) \rightarrow F(\mathbf{x})$ as $n \rightarrow \infty$ for all points \mathbf{x} at which $F(\mathbf{x})$ is continuous. The distribution of \mathbf{X} is the **limiting distribution** or **asymptotic distribution** of \mathbf{X}_n .

b) \mathbf{X}_n converges in probability to \mathbf{X} , written $\mathbf{X}_n \xrightarrow{P} \mathbf{X}$, if for every $\epsilon > 0$, $P(\|\mathbf{X}_n - \mathbf{X}\| > \epsilon) \rightarrow 0$ as $n \rightarrow \infty$.

c) Let $r > 0$ be a real number. Then \mathbf{X}_n converges in r th mean to \mathbf{X} , written $\mathbf{X}_n \xrightarrow{r} \mathbf{X}$, if $E(\|\mathbf{X}_n - \mathbf{X}\|^r) \rightarrow 0$ as $n \rightarrow \infty$.

d) \mathbf{X}_n converges almost everywhere to \mathbf{X} , written $\mathbf{X}_n \xrightarrow{ae} \mathbf{X}$, if $P(\lim_{n \rightarrow \infty} \mathbf{X}_n = \mathbf{X}) = 1$.

Theorems 4.18 and 4.19 below are the multivariate extensions of the limit theorems in subsection 4.2.2. When the limiting distribution of $\mathbf{Z}_n = \sqrt{n}(\mathbf{g}(\mathbf{T}_n) - \mathbf{g}(\boldsymbol{\theta}))$ is multivariate normal $N_k(\mathbf{0}, \boldsymbol{\Sigma})$, approximate the joint cdf of \mathbf{Z}_n with the joint cdf of the $N_k(\mathbf{0}, \boldsymbol{\Sigma})$ distribution. Thus to find probabilities, manipulate \mathbf{Z}_n as if $\mathbf{Z}_n \approx N_k(\mathbf{0}, \boldsymbol{\Sigma})$. To see that the CLT is a special case of the MCLT below, let $k = 1$, $E(X) = \mu$, and $V(X) = \boldsymbol{\Sigma}_x = \sigma^2$.

Theorem 4.18: the Multivariate Central Limit Theorem (MCLT). If $\mathbf{X}_1, \dots, \mathbf{X}_n$ are iid $k \times 1$ random vectors with $E(\mathbf{X}) = \boldsymbol{\mu}$ and $\text{Cov}(\mathbf{X}) = \boldsymbol{\Sigma}_x$, then

$$\sqrt{n}(\bar{\mathbf{X}}_n - \boldsymbol{\mu}) \xrightarrow{D} N_k(\mathbf{0}, \boldsymbol{\Sigma}_x)$$

where the sample mean

$$\bar{\mathbf{X}}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i.$$

To see that the delta method is a special case of the multivariate delta method, note that if T_n and parameter θ are real valued, then $\mathbf{D}_g(\boldsymbol{\theta}) = g'(\theta)$.

Theorem 4.19: the Multivariate Delta Method. If \mathbf{g} does not depend on n and

$$\sqrt{n}(T_n - \theta) \xrightarrow{D} N_k(\mathbf{0}, \boldsymbol{\Sigma}),$$

then

$$\sqrt{n}(\mathbf{g}(T_n) - \mathbf{g}(\boldsymbol{\theta})) \xrightarrow{D} N_d(\mathbf{0}, \mathbf{D}_g(\boldsymbol{\theta}) \boldsymbol{\Sigma} \mathbf{D}_g^T(\boldsymbol{\theta}))$$

where the $d \times k$ Jacobian matrix of partial derivatives

$$\mathbf{D}_g(\boldsymbol{\theta}) = \begin{bmatrix} \frac{\partial}{\partial \theta_1} g_1(\boldsymbol{\theta}) & \dots & \frac{\partial}{\partial \theta_k} g_1(\boldsymbol{\theta}) \\ \vdots & & \vdots \\ \frac{\partial}{\partial \theta_1} g_d(\boldsymbol{\theta}) & \dots & \frac{\partial}{\partial \theta_k} g_d(\boldsymbol{\theta}) \end{bmatrix}.$$

Here the mapping $\mathbf{g} : \mathbb{R}^k \rightarrow \mathbb{R}^d$ needs to be differentiable in a neighborhood of $\boldsymbol{\theta} \in \mathbb{R}^k$.

Definition 4.12. If the estimator $\mathbf{g}(T_n) \xrightarrow{P} \mathbf{g}(\boldsymbol{\theta})$ for all $\boldsymbol{\theta} \in \Theta$, then $\mathbf{g}(T_n)$ is a **consistent estimator** of $\mathbf{g}(\boldsymbol{\theta})$.

Theorem 4.20. If $0 < \delta \leq 1$, \mathbf{X} is a random vector, and

$$n^\delta(\mathbf{g}(T_n) - \mathbf{g}(\boldsymbol{\theta})) \xrightarrow{D} \mathbf{X},$$

then $\mathbf{g}(\mathbf{T}_n) \xrightarrow{P} \mathbf{g}(\boldsymbol{\theta})$.

Theorem 4.21. If $\mathbf{X}_1, \dots, \mathbf{X}_n$ are iid, $E(\|\mathbf{X}\|) < \infty$, and $E(\mathbf{X}) = \boldsymbol{\mu}$, then

a) WLLN: $\overline{\mathbf{X}}_n \xrightarrow{P} \boldsymbol{\mu}$ and

b) SLLN: $\overline{\mathbf{X}}_n \xrightarrow{ae} \boldsymbol{\mu}$.

Theorem 4.22: Continuity Theorem. Let \mathbf{X}_n be a sequence of $k \times 1$ random vectors with characteristic functions $\phi_n(\mathbf{t})$, and let \mathbf{X} be a $k \times 1$ random vector with cf $\phi(\mathbf{t})$. Then

$$\mathbf{X}_n \xrightarrow{D} \mathbf{X} \text{ iff } \phi_n(\mathbf{t}) \rightarrow \phi(\mathbf{t})$$

for all $\mathbf{t} \in \mathbb{R}^k$.

Theorem 4.23: Cramér Wold Device. Let \mathbf{X}_n be a sequence of $k \times 1$ random vectors, and let \mathbf{X} be a $k \times 1$ random vector. Then

$$\mathbf{X}_n \xrightarrow{D} \mathbf{X} \text{ iff } \mathbf{t}^T \mathbf{X}_n \xrightarrow{D} \mathbf{t}^T \mathbf{X}$$

for all $\mathbf{t} \in \mathbb{R}^k$.

Application: Proof of the MCLT Theorem 4.18. Note that for fixed \mathbf{t} , the $\mathbf{t}^T \mathbf{X}_i$ are iid random variables with mean $\mathbf{t}^T \boldsymbol{\mu}$ and variance $\mathbf{t}^T \boldsymbol{\Sigma} \mathbf{t}$. Hence by the CLT, $\mathbf{t}^T \sqrt{n}(\overline{\mathbf{X}}_n - \boldsymbol{\mu}) \xrightarrow{D} N(0, \mathbf{t}^T \boldsymbol{\Sigma} \mathbf{t})$. The right hand side has distribution $\mathbf{t}^T \mathbf{X}$ where $\mathbf{X} \sim N_k(\mathbf{0}, \boldsymbol{\Sigma})$. Hence by the Cramér Wold Device, $\sqrt{n}(\overline{\mathbf{X}}_n - \boldsymbol{\mu}) \xrightarrow{D} N_k(\mathbf{0}, \boldsymbol{\Sigma})$. \square

Theorem 4.24. a) If $\mathbf{X}_n \xrightarrow{P} \mathbf{X}$, then $\mathbf{X}_n \xrightarrow{D} \mathbf{X}$.

b)

$$\mathbf{X}_n \xrightarrow{P} \mathbf{g}(\boldsymbol{\theta}) \text{ iff } \mathbf{X}_n \xrightarrow{D} \mathbf{g}(\boldsymbol{\theta}).$$

Let $g(n) \geq 1$ be an increasing function of the sample size n : $g(n) \uparrow \infty$, e.g. $g(n) = \sqrt{n}$. See White (1984, p. 15). If a $k \times 1$ random vector $\mathbf{T}_n - \boldsymbol{\mu}$ converges to a nondegenerate multivariate normal distribution with convergence rate \sqrt{n} , then \mathbf{T}_n has (tightness) rate \sqrt{n} .

Definition 4.13. Let $\mathbf{A}_n = [a_{i,j}(n)]$ be an $r \times c$ random matrix.

a) $\mathbf{A}_n = O_P(\mathbf{X}_n)$ if $a_{i,j}(n) = O_P(\mathbf{X}_n)$ for $1 \leq i \leq r$ and $1 \leq j \leq c$.

b) $\mathbf{A}_n = o_p(\mathbf{X}_n)$ if $a_{i,j}(n) = o_p(\mathbf{X}_n)$ for $1 \leq i \leq r$ and $1 \leq j \leq c$.

c) $\mathbf{A}_n \asymp_P (1/g(n))$ if $a_{i,j}(n) \asymp_P (1/g(n))$ for $1 \leq i \leq r$ and $1 \leq j \leq c$.

d) Let $\mathbf{A}_{1,n} = \mathbf{T}_n - \boldsymbol{\mu}$ and $\mathbf{A}_{2,n} = \mathbf{C}_n - c\boldsymbol{\Sigma}$ for some constant $c > 0$. If $\mathbf{A}_{1,n} \asymp_P (1/g(n))$ and $\mathbf{A}_{2,n} \asymp_P (1/g(n))$, then $(\mathbf{T}_n, \mathbf{C}_n)$ has (tightness) rate $g(n)$.

Theorem 4.25: Continuous Mapping Theorem. Let $\mathbf{X}_n \in \mathbb{R}^k$. If $\mathbf{X}_n \xrightarrow{D} \mathbf{X}$ and if the function $\mathbf{g} : \mathbb{R}^k \rightarrow \mathbb{R}^j$ is continuous, then $\mathbf{g}(\mathbf{X}_n) \xrightarrow{D} \mathbf{g}(\mathbf{X})$.

The following two theorems are taken from Severini (2005, pp. 345-349, 354).

Theorem 4.26. Let $\mathbf{X}_n = (X_{1n}, \dots, X_{kn})^T$ be a sequence of $k \times 1$ random vectors, let \mathbf{Y}_n be a sequence of $k \times 1$ random vectors, and let $\mathbf{X} = (X_1, \dots, X_k)^T$ be a $k \times 1$ random vector. Let \mathbf{W}_n be a sequence of $k \times k$ nonsingular random matrices, and let \mathbf{C} be a $k \times k$ constant nonsingular matrix.

- a) $\mathbf{X}_n \xrightarrow{P} \mathbf{X}$ iff $X_{in} \xrightarrow{P} X_i$ for $i = 1, \dots, k$.
- b) **Slutsky's Theorem:** If $\mathbf{X}_n \xrightarrow{D} \mathbf{X}$ and $\mathbf{Y}_n \xrightarrow{P} \mathbf{c}$ for some constant $k \times 1$ vector \mathbf{c} , then i) $\mathbf{X}_n + \mathbf{Y}_n \xrightarrow{D} \mathbf{X} + \mathbf{c}$ and ii) $\mathbf{Y}_n^T \mathbf{X}_n \xrightarrow{D} \mathbf{c}^T \mathbf{X}$.
- c) If $\mathbf{X}_n \xrightarrow{D} \mathbf{X}$ and $\mathbf{W}_n \xrightarrow{P} \mathbf{C}$, then $\mathbf{W}_n \mathbf{X}_n \xrightarrow{D} \mathbf{C} \mathbf{X}$, $\mathbf{X}_n^T \mathbf{W}_n \xrightarrow{D} \mathbf{X}^T \mathbf{C}$, $\mathbf{W}_n^{-1} \mathbf{X}_n \xrightarrow{D} \mathbf{C}^{-1} \mathbf{X}$, and $\mathbf{X}_n^T \mathbf{W}_n^{-1} \xrightarrow{D} \mathbf{X}^T \mathbf{C}^{-1}$.

Theorem 4.27. Let W_n, X_n, Y_n , and Z_n be sequences of random variables such that $Y_n > 0$ and $Z_n > 0$. (Often Y_n and Z_n are deterministic, e.g. $Y_n = n^{-1/2}$.)

- a) If $W_n = O_P(1)$ and $X_n = O_P(1)$, then $W_n + X_n = O_P(1)$ and $W_n X_n = O_P(1)$, thus $O_P(1) + O_P(1) = O_P(1)$ and $O_P(1)O_P(1) = O_P(1)$.
- b) If $W_n = O_P(1)$ and $X_n = o_P(1)$, then $W_n + X_n = O_P(1)$ and $W_n X_n = o_P(1)$, thus $O_P(1) + o_P(1) = O_P(1)$ and $O_P(1)o_P(1) = o_P(1)$.
- c) If $W_n = O_P(Y_n)$ and $X_n = O_P(Z_n)$, then $W_n + X_n = O_P(\max(Y_n, Z_n))$ and $W_n X_n = O_P(Y_n Z_n)$, thus $O_P(Y_n) + O_P(Z_n) = O_P(\max(Y_n, Z_n))$ and $O_P(Y_n)O_P(Z_n) = O_P(Y_n Z_n)$.

Theorem 4.28. i) Suppose $\sqrt{n}(T_n - \boldsymbol{\mu}) \xrightarrow{D} N_p(\boldsymbol{\theta}, \boldsymbol{\Sigma})$. Let \mathbf{A} be a $q \times p$ constant matrix. Then $\mathbf{A}\sqrt{n}(T_n - \boldsymbol{\mu}) = \sqrt{n}(\mathbf{A}T_n - \mathbf{A}\boldsymbol{\mu}) \xrightarrow{D} N_q(\mathbf{A}\boldsymbol{\theta}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)$.

ii) Let $\boldsymbol{\Sigma} > 0$. If (T, \mathbf{C}) is a consistent estimator of $(\boldsymbol{\mu}, s \boldsymbol{\Sigma})$ where $s > 0$ is some constant, then $D_{\mathbf{x}}^2(T, \mathbf{C}) = (\mathbf{x} - T)^T \mathbf{C}^{-1} (\mathbf{x} - T) = s^{-1} D_{\mathbf{x}}^2(\boldsymbol{\mu}, \boldsymbol{\Sigma}) + o_P(1)$, so $D_{\mathbf{x}}^2(T, \mathbf{C})$ is a consistent estimator of $s^{-1} D_{\mathbf{x}}^2(\boldsymbol{\mu}, \boldsymbol{\Sigma})$.

iii) Let $\boldsymbol{\Sigma} > 0$. If $\sqrt{n}(T - \boldsymbol{\mu}) \xrightarrow{D} N_p(\mathbf{0}, \boldsymbol{\Sigma})$ and if \mathbf{C} is a consistent estimator of $\boldsymbol{\Sigma}$, then $n(T - \boldsymbol{\mu})^T \mathbf{C}^{-1} (T - \boldsymbol{\mu}) \xrightarrow{D} \chi_p^2$. In particular, $n(\bar{\mathbf{x}} - \boldsymbol{\mu})^T \mathbf{S}^{-1} (\bar{\mathbf{x}} - \boldsymbol{\mu}) \xrightarrow{D} \chi_p^2$.

Proof: ii) $D_{\mathbf{x}}^2(T, \mathbf{C}) = (\mathbf{x} - T)^T \mathbf{C}^{-1} (\mathbf{x} - T) = (\mathbf{x} - \boldsymbol{\mu} + \boldsymbol{\mu} - T)^T [\mathbf{C}^{-1} - s^{-1} \boldsymbol{\Sigma}^{-1} + s^{-1} \boldsymbol{\Sigma}^{-1}] (\mathbf{x} - \boldsymbol{\mu} + \boldsymbol{\mu} - T) = (\mathbf{x} - \boldsymbol{\mu})^T [s^{-1} \boldsymbol{\Sigma}^{-1}] (\mathbf{x} - \boldsymbol{\mu}) + (\mathbf{x} - T)^T [\mathbf{C}^{-1} - s^{-1} \boldsymbol{\Sigma}^{-1}] (\mathbf{x} - T) + (\mathbf{x} - \boldsymbol{\mu})^T [s^{-1} \boldsymbol{\Sigma}^{-1}] (\boldsymbol{\mu} - T) + (\boldsymbol{\mu} - T)^T [s^{-1} \boldsymbol{\Sigma}^{-1}] (\mathbf{x} - \boldsymbol{\mu}) + (\boldsymbol{\mu} - T)^T [s^{-1} \boldsymbol{\Sigma}^{-1}] (\boldsymbol{\mu} - T) = s^{-1} D_{\mathbf{x}}^2(\boldsymbol{\mu}, \boldsymbol{\Sigma}) + O_P(1)$.

(Note that $D_{\mathbf{x}}^2(T, \mathbf{C}) = s^{-1}D_{\mathbf{x}}^2(\boldsymbol{\mu}, \boldsymbol{\Sigma}) + O_P(n^{-\delta})$ if (T, \mathbf{C}) is a consistent estimator of $(\boldsymbol{\mu}, s \boldsymbol{\Sigma})$ with rate n^δ where $0 < \delta \leq 0.5$ if $[\mathbf{C}^{-1} - s^{-1}\boldsymbol{\Sigma}^{-1}] = O_P(n^{-\delta})$.)

Alternatively, $D_{\mathbf{x}}^2(T, \mathbf{C})$ is a continuous function of (T, \mathbf{C}) if $\mathbf{C} > 0$ for $n > 10p$. Hence $D_{\mathbf{x}}^2(T, \mathbf{C}) \xrightarrow{P} D_{\mathbf{x}}^2(\boldsymbol{\mu}, s\boldsymbol{\Sigma})$.

iii) Note that $\mathbf{Z}_n = \sqrt{n} \boldsymbol{\Sigma}^{-1/2}(T - \boldsymbol{\mu}) \xrightarrow{D} N_p(\mathbf{0}, \mathbf{I}_p)$. Thus $\mathbf{Z}_n^T \mathbf{Z}_n = n(T - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(T - \boldsymbol{\mu}) \xrightarrow{D} \chi_p^2$. Now $n(T - \boldsymbol{\mu})^T \mathbf{C}^{-1}(T - \boldsymbol{\mu}) = n(T - \boldsymbol{\mu})^T [\mathbf{C}^{-1} - \boldsymbol{\Sigma}^{-1} + \boldsymbol{\Sigma}^{-1}](T - \boldsymbol{\mu}) = n(T - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(T - \boldsymbol{\mu}) + n(T - \boldsymbol{\mu})^T [\mathbf{C}^{-1} - \boldsymbol{\Sigma}^{-1}](T - \boldsymbol{\mu}) = n(T - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(T - \boldsymbol{\mu}) + o_P(1) \xrightarrow{D} \chi_p^2$ since $\sqrt{n}(T - \boldsymbol{\mu})^T [\mathbf{C}^{-1} - \boldsymbol{\Sigma}^{-1}] \sqrt{n}(T - \boldsymbol{\mu}) = O_P(1)O_P(1)O_P(1) = o_P(1)$. \square

Example 4.12. Suppose that $\mathbf{x}_n \perp \mathbf{y}_n$ for $n = 1, 2, \dots$. Suppose $\mathbf{x}_n \xrightarrow{D} \mathbf{x}$, and $\mathbf{y}_n \xrightarrow{D} \mathbf{y}$ where $\mathbf{x} \perp \mathbf{y}$. Then

$$\begin{bmatrix} \mathbf{x}_n \\ \mathbf{y}_n \end{bmatrix} \xrightarrow{D} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}$$

by Theorem 4.22. To see this, let $\mathbf{t} = (t_1^T, t_2^T)^T$, $\mathbf{z}_n = (\mathbf{x}_n^T, \mathbf{y}_n^T)^T$, and $\mathbf{z} = (\mathbf{x}^T, \mathbf{y}^T)^T$. Since $\mathbf{x}_n \perp \mathbf{y}_n$ and $\mathbf{x} \perp \mathbf{y}$, the characteristic function

$$\phi_{\mathbf{z}_n}(\mathbf{t}) = \phi_{\mathbf{x}_n}(t_1)\phi_{\mathbf{y}_n}(t_2) \rightarrow \phi_{\mathbf{x}}(t_1)\phi_{\mathbf{y}}(t_2) = \phi_{\mathbf{z}}(\mathbf{t}).$$

Hence $\mathbf{g}(\mathbf{z}_n) \xrightarrow{D} \mathbf{g}(\mathbf{z})$ by Theorem 4.25.

4.3 Mixture Distributions

Mixture distributions are useful for model and variable selection since $\hat{\boldsymbol{\beta}}_{I_{min},0}$ is a mixture distribution of $\hat{\boldsymbol{\beta}}_{I_j,0}$, and the lasso estimator $\hat{\boldsymbol{\beta}}_L$ is a mixture distribution of $\hat{\boldsymbol{\beta}}_{L,\lambda_i}$ for $i = 1, \dots, M$. See the second to last paragraph of Section 4.1 for $\hat{\boldsymbol{\beta}}_{I_{min},0}$. A random vector \mathbf{u} has a mixture distribution if \mathbf{u} equals a random vector \mathbf{u}_j with probability π_j for $j = 1, \dots, J$. See Definition 4.2 for the population mean and population covariance matrix of a random vector.

Definition 4.14. The distribution of a $g \times 1$ random vector \mathbf{u} is a mixture distribution if the cumulative distribution function (cdf) of \mathbf{u} is

$$F_{\mathbf{u}}(\mathbf{t}) = \sum_{j=1}^J \pi_j F_{\mathbf{u}_j}(\mathbf{t}) \quad (4.7)$$

where the probabilities π_j satisfy $0 \leq \pi_j \leq 1$ and $\sum_{j=1}^J \pi_j = 1$, $J \geq 2$, and $F_{\mathbf{u}_j}(\mathbf{t})$ is the cdf of a $g \times 1$ random vector \mathbf{u}_j . Then \mathbf{u} has a mixture distribution of the \mathbf{u}_j with probabilities π_j .

Theorem 4.29. Suppose $E(h(\mathbf{u}))$ and the $E(h(\mathbf{u}_j))$ exist. Then

$$E(h(\mathbf{u})) = \sum_{j=1}^J \pi_j E[h(\mathbf{u}_j)]. \quad (4.8)$$

Hence

$$E(\mathbf{u}) = \sum_{j=1}^J \pi_j E[\mathbf{u}_j], \quad (4.9)$$

and $Cov(\mathbf{u}) = E(\mathbf{u}\mathbf{u}^T) - E(\mathbf{u})E(\mathbf{u}^T) = E(\mathbf{u}\mathbf{u}^T) - E(\mathbf{u})[E(\mathbf{u})]^T = \sum_{j=1}^J \pi_j E[\mathbf{u}_j\mathbf{u}_j^T] - E(\mathbf{u})[E(\mathbf{u})]^T =$

$$\sum_{j=1}^J \pi_j Cov(\mathbf{u}_j) + \sum_{j=1}^J \pi_j E(\mathbf{u}_j)[E(\mathbf{u}_j)]^T - E(\mathbf{u})[E(\mathbf{u})]^T. \quad (4.10)$$

If $E(\mathbf{u}_j) = \boldsymbol{\theta}$ for $j = 1, \dots, J$, then $E(\mathbf{u}) = \boldsymbol{\theta}$ and

$$Cov(\mathbf{u}) = \sum_{j=1}^J \pi_j Cov(\mathbf{u}_j).$$

This theorem is easy to prove if the \mathbf{u}_j are continuous random vectors with (joint) probability density functions (pdfs) $f_{\mathbf{u}_j}(\mathbf{t})$. Then \mathbf{u} is a continuous random vector with pdf

$$\begin{aligned} f_{\mathbf{u}}(\mathbf{t}) &= \sum_{j=1}^J \pi_j f_{\mathbf{u}_j}(\mathbf{t}), \quad \text{and} \quad E(h(\mathbf{u})) = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} h(\mathbf{t}) f_{\mathbf{u}}(\mathbf{t}) d\mathbf{t} \\ &= \sum_{j=1}^J \pi_j \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} h(\mathbf{t}) f_{\mathbf{u}_j}(\mathbf{t}) d\mathbf{t} = \sum_{j=1}^J \pi_j E[h(\mathbf{u}_j)] \end{aligned}$$

where $E[h(\mathbf{u}_j)]$ is the expectation with respect to the random vector \mathbf{u}_j . Note that

$$E(\mathbf{u})[E(\mathbf{u})]^T = \sum_{j=1}^J \sum_{k=1}^J \pi_j \pi_k E(\mathbf{u}_j)[E(\mathbf{u}_k)]^T. \quad (4.11)$$

4.4 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 $\mathbf{Z}_n = \sqrt{n}\mathbf{A}(\hat{\boldsymbol{\beta}}_{I_{min}} - \boldsymbol{\beta}_I)$ after variable selection. One reason is that with positive probability, $\hat{\boldsymbol{\beta}}_{I_{min}}$ does not have the same dimension as $\boldsymbol{\beta}_I$ if AIC is used. Hence \mathbf{Z}_n is not defined with positive probability. Also, the dimension of a random vector is $k \times 1$, say, while the dimension of $\hat{\boldsymbol{\beta}}_{I_{min}}$ is $K \times 1$ where K is a random variable. Hence the random quantity $\hat{\boldsymbol{\beta}}_{I_{min}}$ is not a random vector and not a statistic.

We will show that large sample theory becomes simple by using zero padding. 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 $\hat{\boldsymbol{\beta}}_{I_{min},0} = (\hat{\beta}_1, 0, \hat{\beta}_3, 0)^T$. Assume p is fixed, and $n \rightarrow \infty$.

The Rathnayake and Olive (2019) theory in this section applies to many regression models including many generalized linear models, some time series models, some survival regression models such as the Cox (1972) proportional hazards survival regression model and AFTs, and the multiple linear regression model where the error distribution is unknown.

Suppose the regression model satisfies $Y \perp\!\!\!\perp \mathbf{x} | \mathbf{x}^T \boldsymbol{\beta}$, that model (2.4) holds, and that if $S \subseteq I_j$, then $\sqrt{n}(\hat{\boldsymbol{\beta}}_{I_j} - \boldsymbol{\beta}_{I_j}) \xrightarrow{D} N_{a_j}(\mathbf{0}, \mathbf{V}_j)$. Also assume that a variable selection criterion, such as AIC or relaxed lasso, is used such that $P(S \subseteq I_{min}) \rightarrow 1$ as $n \rightarrow \infty$. Hence

$$\sqrt{n}(\hat{\boldsymbol{\beta}}_{I_{j,0}} - \boldsymbol{\beta}) \xrightarrow{D} N_p(\mathbf{0}, \mathbf{V}_{j,0}) \quad (4.12)$$

where $\mathbf{V}_{j,0}$ adds columns and rows of zeros corresponding to the x_i not in I_j . Hence $\mathbf{V}_{j,0}$ is singular unless I_j corresponds to the full model. Since fewer than 2^p regression models I contain the true model S , and each such model gives a \sqrt{n} consistent estimator $\hat{\boldsymbol{\beta}}_{I,0}$ of $\boldsymbol{\beta}$, the probability that I_{min} picks one of these models goes to one as $n \rightarrow \infty$. Then $\hat{\boldsymbol{\beta}}_{I_{min},0}$ is a \sqrt{n} consistent estimator of $\boldsymbol{\beta}$ under model (2.4) if the variable selection criterion is used with forward selection, backward elimination, or all subsets. This result holds since picking from a fixed number of \sqrt{n} consistent estimators results in a \sqrt{n} consistent estimator by Pratt (1959). See Theorem 4.12 and Theorem 4.13. This section will use mixture distributions to find the limiting distribution of $\sqrt{n}(\hat{\boldsymbol{\beta}}_{I_{min},0} - \boldsymbol{\beta})$.

Under regularity conditions, $P(S \subseteq I_{min}) \rightarrow 1$ as $n \rightarrow \infty$ if BIC or AIC is used with forward selection, backward elimination, or all subsets. See Charkhi and Claeskens (2018), Claeskens and Hjort (2008, pp. 70, 85-86, 101, 102, 114, 232), and Hjort and Claeskens (2006).

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 4.3. 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 4.30, Mixture Distribution CLT. Suppose the $g \times 1$ statistic T_n is equal to the estimator T_{jn} with probability π_{jn} for $j = 1, \dots, J$ where $\sum_j \pi_{jn} = 1$, $\pi_{jn} \rightarrow \pi_j$ as $n \rightarrow \infty$, and $\mathbf{u}_{jn} = \sqrt{n}(T_{jn} - \theta) \xrightarrow{D} \mathbf{u}_j$ with $E(\mathbf{u}_j) = \mathbf{0}$ and $\text{Cov}(\mathbf{u}_j) = \Sigma_j$. Then

$$\sqrt{n}(T_n - \theta) \xrightarrow{D} \mathbf{u} \quad (4.13)$$

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 π_j , $E(\mathbf{u}) = \mathbf{0}$, and $\text{Cov}(\mathbf{u}) = \Sigma \mathbf{u} = \sum_j \pi_j \Sigma_j$.

Proof: Note that T_n has a mixture distribution of the T_{jn} with probabilities π_{jn} . Hence $\sqrt{n}(T_n - \theta)$ has a mixture distribution of the $\mathbf{u}_{jn} = \sqrt{n}(T_{jn} - \theta)$, and the cdf of $\sqrt{n}(T_n - \theta)$ 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

Applying the above results makes large sample theory for $\hat{\beta}_{I_{min},0}$ simple. The following theorem is due to Rathnayake and Olive (2019), generalizing the Pelawa Watagoda and Olive (2019a) result for multiple linear regression.

Theorem 4.31, Variable Selection CLT. Assume $P(S \subseteq I_{min}) \rightarrow 1$ as $n \rightarrow \infty$, and let $\hat{\beta}_{I_{min},0} = \hat{\beta}_{I_k,0}$ with probabilities π_{kn} where $\pi_{kn} \rightarrow \pi_k$ as $n \rightarrow \infty$. Denote the positive π_k by π_j . Assume $\mathbf{u}_{jn} = \sqrt{n}(\hat{\beta}_{I_j,0} - \beta) \xrightarrow{D} \mathbf{u}_j \sim N_p(\mathbf{0}, \mathbf{V}_{j,0})$. a) Then

$$\mathbf{u}_n = \sqrt{n}(\hat{\beta}_{I_{min},0} - \beta) \xrightarrow{D} \mathbf{u} \quad (4.14)$$

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 π_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

$$\mathbf{v}_n = \mathbf{A} \mathbf{u}_n = \sqrt{n}(\mathbf{A} \hat{\beta}_{I_{min},0} - \mathbf{A} \beta) \xrightarrow{D} \mathbf{A} \mathbf{u} = \mathbf{v} \quad (4.15)$$

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 π_j .

Proof. a) Since \mathbf{u}_n has a mixture distribution of the \mathbf{u}_{kn} with probabilities π_{kn} , the cdf of \mathbf{u}_n is $F_{\mathbf{u}_n}(\mathbf{z}) = \sum_k \pi_{kn} F_{\mathbf{u}_{kn}}(\mathbf{z}) \rightarrow F_{\mathbf{u}}(\mathbf{z}) = \sum_j \pi_j F_{\mathbf{u}_j}(\mathbf{z})$ at

continuity points of the $F_{\mathbf{u}_j}(\mathbf{z})$ as $n \rightarrow \infty$.

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

Remark 4.4. If $P(S \subseteq I_{min}) \rightarrow 1$ as $n \rightarrow \infty$ (e.g. for AIC, BIC, or relaxed lasso), the values of π_j depend on the regression variable selection method such as backward elimination, forward selection, all subsets, and lasso. Typically the mixture distribution is not asymptotically normal. There are two exceptions. First, suppose $\pi_d = 1$ with $\mathbf{u} \sim \mathbf{u}_d \sim N_p(\mathbf{0}, \mathbf{V}_{d,0})$. This exception occurs if $a_S = p$ so S is the full model, and for methods like BIC that choose I_S with probability going to one under strong regularity conditions.

The second exception occurs for each $\pi_j > 0$, $\mathbf{A}\mathbf{u}_j \sim N_g(\mathbf{0}, \mathbf{A}\Sigma_j\mathbf{A}^T) = N_g(\mathbf{0}, \mathbf{A}\Sigma\mathbf{A}^T)$. Then $\sqrt{n}(\mathbf{A}\hat{\beta}_{I_{min},0} - \mathbf{A}\beta) \xrightarrow{D} \mathbf{A}\mathbf{u} \sim N_g(\mathbf{0}, \mathbf{A}\Sigma\mathbf{A}^T)$. This exception occurs for $\hat{\beta}_S$ if $\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{D} N_p(\mathbf{0}, \mathbf{V})$ where the asymptotic covariance matrix \mathbf{V} is diagonal and nonsingular. Then $\hat{\beta}_S$ has the same multivariate normal limiting distribution for I_{min} and for the full model.

Remark 4.5. This theory has several applications. First, the theory gives the asymptotic distribution for many variable selection estimators, which are some of the most used estimators in Statistics. Second, the theory is useful for explaining why $\hat{\beta}_{I_{min}}$ is not a good estimator, but $\hat{\beta}_{I_{min},0}$ is a good estimator. Suppose $I_{min} = I_j$ is observed. Due to selection bias, the model using predictors I_j underestimates the variability of the responses Y_1, \dots, Y_n , and $Cov(\mathbf{A}\hat{\beta}_{I_j})$ is not the correct covariance matrix for $\mathbf{A}\hat{\beta}_{I_{min}}$. Typically $\hat{\beta}_{I_{min}}$ is not a consistent estimator for any parameter vector β_{I_j} , since in general $P(I_{min} = I_j)$ does not go to one as $n \rightarrow \infty$, and the dimension of I_{min} is a random variable. Selection bias occurs from acting as if $\hat{\beta}_{I_{min}}$ is the “full model” (using large sample theory as if the “full model” was selected before gathering the data), when $\hat{\beta}_{I_{min},0}$ has large sample theory given by Theorem 4.31.

A third application will be bootstrap inference for hypothesis testing. See Section 4.8. Fourth, the theory can be used to justify prediction intervals after variable selection. See Section 4.5, and Olive et al. (2020). Fifth, recall p is fixed. Suppose a shrinkage method, such as lasso or elastic net, does variable selection. Let $\hat{\beta}_{I_{min}}$ be the regression estimator, such as a Cox regression, applied to a constant and the variables with nonzero shrinkage estimator coefficients. If the shrinkage estimator is consistent, then $P(S \subseteq I_{min}) \rightarrow 1$ as $n \rightarrow \infty$, and thus the relaxed shrinkage estimator $\hat{\beta}_{I_{min},0}$ is \sqrt{n} consistent. In particular, relaxed lasso and relaxed elastic net are \sqrt{n} consistent if lasso and elastic net are consistent.

Remark 4.6. If $\pi_d = 1$ corresponds to β_d , then $\hat{\beta}_{I_{min}}$ can give useful information about β_d , but information is lost about the parameters estimated to be zero if S is not the full model. There is a large literature on *variable selection consistency* and the *oracle property* where $P(I_{min} = S) \rightarrow 1$ as $n \rightarrow$

∞ . See Claeskens and Hjort (2008, pp. 99-114) 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 for all subsets regression, but only under very strong regularity conditions for fast methods such as forward selection, backward elimination, and lasso.

4.5 Prediction Intervals

Prediction intervals for regression and prediction regions for multivariate data 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 4.15. 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 $[\hat{L}_n, \hat{U}_n] \rightarrow [L_s, U_s]$ 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 $Y_{fi} \in PI_i$ 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 $Y_{fi} \in PI_i$ 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 ∞ 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 PIs for parametric 1D regression models, which include many parametric survival regression models.

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

$$[Z_{(k_1)}, Z_{(k_2)}] \quad (4.16)$$

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

The $\text{shorth}(c)$ 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

$$\text{shorth}(c) = [Z_{(s)}, Z_{(s+c-1)}]. \quad (4.17)$$

Let

$$k_n = \lceil n(1 - \delta) \rceil. \quad (4.18)$$

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

$$c = \min(n, \lceil n[1 - \delta + 1.12\sqrt{\delta/n}] \rceil). \quad (4.19)$$

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 4.13. 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

$$13 = 89 - 76$$

$$33 = 111 - 78$$

$$689 = 778 - 89$$

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

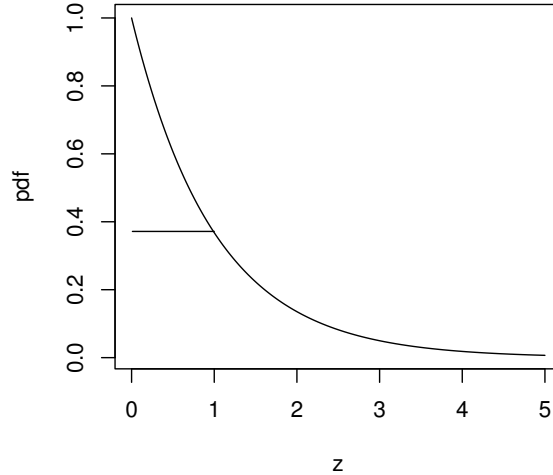


Fig. 4.1 The 36.8% Highest Density Region is $[0,1]$

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 $(-\infty, \infty)$.

A parametric 1D regression model is $Y|\mathbf{x} \sim D(h(\mathbf{x}), \boldsymbol{\gamma})$ for some real valued function, such as $h(\mathbf{x}) = \mathbf{x}^T \boldsymbol{\beta}$, where D is a parametric distribution that depends on the $p \times 1$ vector of predictors \mathbf{x} only through $h(\mathbf{x})$, and $\boldsymbol{\gamma}$ is a $q \times 1$ vector of parameters.

The first new large sample $100(1-\delta)\%$ prediction interval for Y_f applies the shorth(c) prediction interval to the parametric bootstrap sample Y_1^*, \dots, Y_B^* where the Y_i^* are iid from the distribution $D(\hat{h}(\mathbf{x}_f), \hat{\gamma})$ with

$$c = \min(B, \lceil B[1 - \delta + 1.12\sqrt{\delta/B}] \rceil). \quad (4.20)$$

If $Y|\mathbf{x}_f \sim D(h(\mathbf{x}_f), \gamma)$ and the regression method produces a consistent estimator $(\hat{h}(\mathbf{x}_f), \hat{\gamma})$ of $(h(\mathbf{x}_f), \gamma)$, then this new prediction interval is a large sample $100(1 - \delta)\%$ PI.

For models with a linear predictor, we will want prediction intervals after variable selection or model selection. The prediction interval (4.20) can have undercoverage if n is small compared to the number of estimated parameters. The modified shorth PI (4.21) inflates PI (4.20) to compensate for parameter estimation and model selection. Let d be the number of variables x_1^*, \dots, x_d^* used by the full model, forward selection, lasso, or relaxed lasso. We want $n \geq 10d$, and the prediction interval length will be increased (penalized) if n/d is not large. For the second new prediction interval, let $q_n = \min(1 - \delta + 0.05, 1 - \delta + d/n)$ for $\delta > 0.1$ and

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

If $1 - \delta < 0.999$ and $q_n < 1 - \delta + 0.001$, set $q_n = 1 - \delta$. Then compute the shorth PI with

$$c_{mod} = \min(B, \lceil B[q_n + 1.12\sqrt{\delta/B}] \rceil). \quad (4.21)$$

Olive (2007, 2018) and Pelawa Watagoda and Olive (2019b) used similar correction factors for regression models with an additive error since the maximum simulated undercoverage was about 0.05 when $n = 20d$. If a $q \times 1$ vector of parameters γ is also estimated, we may need to replace d by $d_q = d + q$.

Hong et al. (2018) explain why classical PIs after AIC variable selection may not work. Fix p and let I_{min} correspond to the predictors used after variable selection. To show that (4.20) and (4.21) are large sample prediction intervals, we need to show that $(\hat{\beta}_{I_{min},0}, \hat{\gamma}_{I_{min}})$ is a consistent estimator of (β, γ) . Theorem 4.31 shows that $\hat{\beta}_{I_{min},0}$ is a consistent estimator of β . Suppose $P(S \subseteq I_{min}) \rightarrow 1$ as $n \rightarrow \infty$. Suppose model (2.4) holds with $S \subseteq I_j$. Then under regularity conditions that are often mild, $(\hat{\beta}_{I_j}, \hat{\gamma}_{I_j})$ is a consistent estimator of (β_{I_j}, γ) . Then $\hat{\gamma}_{I_{min}}$ is a consistent estimator of γ . Hence if $P(S \subseteq I_{min}) \rightarrow 1$ as $n \rightarrow \infty$ (AIC, BIC, or relaxed lasso), then (4.20) and (4.21) are large sample PIs.

As an example, consider the *Weibull proportional hazards regression model*

$$Y|SP \sim W(\gamma = 1/\sigma, \lambda_0 \exp(SP)),$$

where $\lambda_0 = \exp(-\alpha/\sigma)$, and Y has a Weibull $W(\gamma, \lambda)$ distribution if the pdf of Y is

$$f(y) = \lambda \gamma y^{\gamma-1} \exp[-\lambda y^\gamma]$$

for $y > 0$. Note that the PI is for survival times Y , not censored survival times.

The *survpack* function `wpsim` simulates PI (4.21) for the WPH full model with $d = p$. WPH data $(\mathbf{x}_1, T_1), \dots, (\mathbf{x}_n, T_n)$ is generated as described for variable selection in Section 4.8. The T_i are right censored survival times corresponding to Y_i . Hence for the output below, $\boldsymbol{\beta} = (1, 1, 1, 1, 0, \dots, 0)^T$ with $p = 10$, and `psi = 0.9` means the ten predictor variables are highly correlated. The Weibull AFT is fit and used to get $\hat{\gamma}$ and $\hat{\boldsymbol{\beta}}_W$ for the WPH. Then $B = 1000$ values Y_1^*, \dots, Y_B^* are generated for $Y|\mathbf{x}_f \sim W(\hat{\gamma}, \hat{\lambda}_0 \exp(\hat{\boldsymbol{\beta}}_W^T \mathbf{x}_f))$. The large sample 95% PI (4.21) is used for Y_f with $d = p = 10$. 5000 WPH data sets are generated with 5000 values of (\mathbf{x}_f, Y_f) . The values of Y_f and Y_i^* are not censored. Then 94.76% of the 5000 PIs contained Y_f , with an average length of 1.0554.

```
wpsim(n=1000,p=10,k=4,nruns=5000,psi=0.9,gam=4,B=1000)
$int
(Intercept)
  0.0169485
$beta
 [1] 1 1 1 1 0 0 0 0 0 0
$fullpicov
 [1] 0.9476
$fullpimenlen
 [1] 1.0554
```

4.6 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 $\hat{\boldsymbol{\theta}}$ of the parameter $\boldsymbol{\theta}$, 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 4.16. 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 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 (4.29). 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 4.17. 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 4.18. 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 4.2 for the population mean and population covariance matrix. The Mahalanobis distance in Definition 4.8 is a random variable that estimates the population Mahalanobis distance defined after Definition 4.8.

Definition 4.19. The *i*th Mahalanobis distance $D_i = \sqrt{D_i^2}$ where the *i*th squared Mahalanobis distance is

$$D_i^2 = D_i^2(T(\mathbf{W}), \mathbf{C}(\mathbf{W})) = (\mathbf{x}_i - T(\mathbf{W}))^T \mathbf{C}^{-1}(\mathbf{W})(\mathbf{x}_i - T(\mathbf{W})) \quad (4.22)$$

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 is

$$D_{\mathbf{x}}^2(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \quad (4.23)$$

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

$$\mathcal{A}_n = \{\mathbf{x} : D_{\mathbf{x}}^2(\bar{\mathbf{x}}, \mathbf{S}) \leq D_{(c)}^2\} = \{\mathbf{x} : D_{\mathbf{x}}(\bar{\mathbf{x}}, \mathbf{S}) \leq D_{(c)}\}. \quad (4.24)$$

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 (4.28). Olive (2013b) showed that (4.24) 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 (4.24) 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. Empirically for many distributions, for $n \approx 20p$, the prediction region (4.24) 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

$$q_n = \min(1 - \delta/2, 1 - \delta + 10\delta p/n), \quad \text{otherwise.} \quad (4.25)$$

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

$$c = \lceil nq_n \rceil \quad (4.26)$$

in (4.24) decreased the undercoverage.

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 $(T, \mathbf{C}) = (\bar{\mathbf{x}}, \mathbf{S})$ satisfies this assumption. For $h > 0$, the hyperellipsoid

$$\{\mathbf{z} : (\mathbf{z} - T)^T \mathbf{C}^{-1} (\mathbf{z} - T) \leq h^2\} = \{\mathbf{z} : D_{\mathbf{z}}^2 \leq h^2\} = \{\mathbf{z} : D_{\mathbf{z}} \leq h\} \quad (4.27)$$

has volume equal to

$$\frac{2\pi^{p/2}}{p\Gamma(p/2)} h^p \sqrt{\det(\mathbf{C})} = \frac{2\pi^{p/2}}{p\Gamma(p/2)} h^p b^{p/2} \sqrt{\det(\mathbf{S}_M)}. \quad (4.28)$$

A future observation (random vector) \mathbf{x}_f is in the region (4.27) 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 (4.27) 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 (4.26). If $\mathbf{x}_1, \dots, \mathbf{x}_n$ and \mathbf{x}_f are iid, then prediction region (4.29) is asymptotically optimal for a large class of elliptically contoured

distributions since the volume of (4.29) 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.)

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

Definition 4.20. 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

$$\{\mathbf{z} : D_{\mathbf{z}}^2(\bar{\mathbf{x}}, \mathbf{S}) \leq D_{(U_n)}^2\}, \quad (4.29)$$

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

$$\{\mathbf{z} : D_{\mathbf{z}}^2(\bar{\mathbf{x}}, \mathbf{S}) \leq \chi_{p, 1-\delta}^2\}. \quad (4.30)$$

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 (4.19) and (4.26) do not differ by much more than 3% for $0.01 \leq \delta \leq 0.5$. See Figure 4.2 where $ol = (\text{Eq. 4.26})/n$ is plotted versus $fr = (\text{Eq. 4.19})/n$ 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 4.7. The nonparametric prediction region (4.29) 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 (4.29) contains U_n of the training data cases \mathbf{x}_i provided

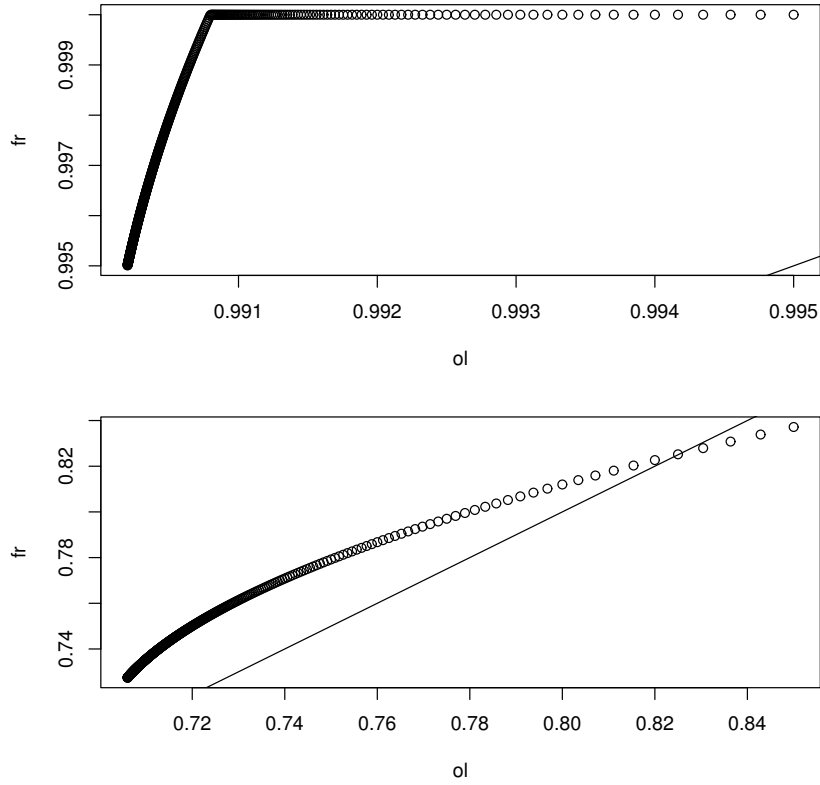


Fig. 4.2 Correction Factor Comparison when $\delta = 0.01$ (Top Plot) and $\delta = 0.3$ (Bottom Plot)

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 .

Remark 4.8. The most used prediction regions assume that the error vectors are iid from a multivariate normal distribution. Using (4.28), the ratio of the volumes of regions (4.30) and (4.29) 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 $\boldsymbol{\Sigma}_{\mathbf{x}} = \text{Cov}(\mathbf{x})$. See Budny (2014), Chen (2011), and Navarro (2014, 2016). Using $\gamma = h^2 = p/\delta$ in (4.27) usually results in prediction regions with volume and coverage that is too large.

Remark 4.9. The nonparametric prediction region (4.29) starts to have good coverage for $n \geq 10p$ for a large class of distributions. Olive (2013b) 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 (4.29) 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 $(\boldsymbol{\mu}, d\boldsymbol{\Sigma})$ where $d > 0$ and $\boldsymbol{\Sigma}$ is nonsingular, then (4.27) with $h = D_{(U_n)}$ is a large sample prediction region, but taking q_n given by (4.25) 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 (4.25) helped reduce undercoverage for small $n \geq 20p$ due to the unknown variability of (T, \mathbf{C}) .

4.7 Bootstrapping Hypothesis Tests and Confidence Regions

This section shows that, under regularity conditions, applying the nonparametric prediction region of Section 4.6 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 θ 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 θ 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 θ . Note that the lengths of the two intervals are the same. Where the interval is centered, at the parameter θ or the statistic \bar{Y}_n , determines whether the interval is a prediction or a confidence interval. See Theorem 4.32 for a similar relationship between confidence regions and prediction regions.

Definition 4.21. 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 (4.29) to the bootstrap sample will result in a confidence region for $\boldsymbol{\theta}$. When $g = 1$, applying the shorth PI (4.19) or PI (4.16) to the bootstrap sample results in a confidence interval for θ . Section 4.7.2 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 θ . 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 $U_n < 0.05$. 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 4.22. The interval $[L_n, U_n]$ is a large sample $100(1 - \delta)\%$ *confidence interval* for θ 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 θ 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 θ 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 (4.16) and (4.19) to the bootstrap sample. Some additional bootstrap CIs are obtained from bootstrap confidence regions from Section 4.7.2 when $g = 1$. See Efron (1982) and Chen (2016) for the percentile method CI. Let T_n be an estimator of a parameter θ 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 (4.31) is obtained by applying PI (4.16) to the bootstrap sample with B used instead of n . Hence (4.31) 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 4.23. The bootstrap percentile method large sample $100(1 - \delta)\%$ confidence interval for θ 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

$$[T_{(k_1)}^*, T_{(k_2)}^*]. \quad (4.31)$$

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

Definition 4.24. The large sample $100(1 - \delta)\%$ *lower shorth* CI for θ is $(-\infty, T_{(c)}^*]$, while the large sample $100(1 - \delta)\%$ *upper shorth* CI for θ 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

$$c = \min(B, \lceil B[1 - \delta + 1.12\sqrt{\delta/B}] \rceil). \quad (4.32)$$

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 4.13 for some theory for bootstrap CIs such as (4.31) and (4.32).

4.7.1 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 4.25. 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 4.14. Let \mathbf{w} be a random variable having the empirical distribution given by Definition 4.25. 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 4.15. 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 4.16. 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 4.7.3 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 $\sqrt{n}(T_n^* - T_n) \xrightarrow{D} \mathbf{u}$.

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

$$\bar{T}^* = \frac{1}{B} \sum_{i=1}^B T_i^* \quad \text{and} \quad \mathbf{S}_T^* = \frac{1}{B-1} \sum_{i=1}^B (T_i^* - \bar{T}^*)(T_i^* - \bar{T}^*)^T \quad (4.33)$$

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 $n\hat{\boldsymbol{\Sigma}}_T \xrightarrow{P} \boldsymbol{\Sigma}_A$. 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 4.10. For Example 4.14, the bootstrap works but is expensive compared to large sample theory. Fix n , then $\bar{T}^* \xrightarrow{P} \boldsymbol{\theta}_n = \bar{\boldsymbol{x}}$ and $\mathbf{S}_T^* \xrightarrow{P} (n-1)\mathbf{S}/n$ as $B \rightarrow \infty$, but using $(\bar{\boldsymbol{x}}, \mathbf{S})$ makes more sense. For Example 4.14, 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.

4.7.2 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 4.4. 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 (4.33). A useful result is $d_n F_{g,d_n,1-\delta} \rightarrow \chi_{g,1-\delta}^2$ 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 4.26. a) The standard bootstrap large sample $100(1 - \delta)\%$ confidence region for $\boldsymbol{\theta}$ is $\{\mathbf{w} : (\mathbf{w} - T_n)^T [\mathbf{S}_T^*]^{-1} (\mathbf{w} - T_n) \leq D_{1-\delta}^2\} =$

$$\{\mathbf{w} : D_{\mathbf{w}}^2(T_n, \mathbf{S}_T^*) \leq D_{1-\delta}^2\} \quad (4.34)$$

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 $\{\mathbf{w} : (\mathbf{w} - T_n)^T [\hat{\boldsymbol{\Sigma}}_A/n]^{-1} (\mathbf{w} - T_n) \leq D_{(k_B, T)}^2\} =$

$$\{\mathbf{w} : D_{\mathbf{w}}^2(T_n, \hat{\boldsymbol{\Sigma}}_A/n) \leq D_{(k_B, T)}^2\} \quad (4.35)$$

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 (4.34) 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 (4.35) 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 (4.29) to the bootstrap sample. Olive (2017ab, 2018) also gave the modified Bickel and Ren confidence region that uses $\hat{\boldsymbol{\Sigma}}_A = n\mathbf{S}_T^*$. The hybrid confidence region is due to Pelawa Watagoda and Olive (2019). Let $q_B = \min(1 - \delta + 0.05, 1 - \delta + g/B)$ for $\delta > 0.1$ and

$$q_B = \min(1 - \delta/2, 1 - \delta + 10\delta g/B), \quad \text{otherwise.} \quad (4.36)$$

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 (4.36) as a correction factor for finite $B \geq 50p$.

Definition 4.27. a) The prediction region method large sample $100(1 - \delta)\%$ confidence region for $\boldsymbol{\theta}$ is $\{\boldsymbol{w} : (\boldsymbol{w} - \bar{\boldsymbol{T}}^*)^T [\boldsymbol{S}_T^*]^{-1} (\boldsymbol{w} - \bar{\boldsymbol{T}}^*) \leq D_{(U_B)}^2\} =$

$$\{\boldsymbol{w} : D_{\boldsymbol{w}}^2(\bar{\boldsymbol{T}}^*, \boldsymbol{S}_T^*) \leq D_{(U_B)}^2\} \quad (4.37)$$

where $D_{(U_B)}^2$ is computed from $D_i^2 = (T_i^* - \bar{T}^*)^T [\boldsymbol{S}_T^*]^{-1} (T_i^* - \bar{T}^*)$ for $i = 1, \dots, B$. Note that the corresponding test for $H_0 : \boldsymbol{\theta} = \boldsymbol{\theta}_0$ rejects H_0 if $(\bar{\boldsymbol{T}}^* - \boldsymbol{\theta}_0)^T [\boldsymbol{S}_T^*]^{-1} (\bar{\boldsymbol{T}}^* - \boldsymbol{\theta}_0) > D_{(U_B)}^2$. (This procedure is basically the one sample Hotelling's T^2 test applied to the T_i^* using \boldsymbol{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 $\{\boldsymbol{w} : (\boldsymbol{w} - T_n)^T [\boldsymbol{S}_T^*]^{-1} (\boldsymbol{w} - T_n) \leq D_{(U_B, T)}^2\} =$

$$\{\boldsymbol{w} : D_{\boldsymbol{w}}^2(T_n, \boldsymbol{S}_T^*) \leq D_{(U_B, T)}^2\} \quad (4.38)$$

where the cutoff $D_{(U_B, T)}^2$ is the $100q_B$ th sample quantile of the $D_i^2 = (T_i^* - T_n)^T [\boldsymbol{S}_T^*]^{-1} (T_i^* - T_n)$. Note that the corresponding test for $H_0 : \boldsymbol{\theta} = \boldsymbol{\theta}_0$ rejects H_0 if $(T_n - \boldsymbol{\theta}_0)^T [\boldsymbol{S}_T^*]^{-1} (T_n - \boldsymbol{\theta}_0) > D_{(U_B, T)}^2$. c) Shift region (4.37) to have center T_n , or equivalently, change the cutoff of region (4.38) to $D_{(U_B)}^2$ to get the hybrid large sample $100(1 - \delta)\%$ confidence region: $\{\boldsymbol{w} : (\boldsymbol{w} - T_n)^T [\boldsymbol{S}_T^*]^{-1} (\boldsymbol{w} - T_n) \leq D_{(U_B)}^2\} =$

$$\{\boldsymbol{w} : D_{\boldsymbol{w}}^2(T_n, \boldsymbol{S}_T^*) \leq D_{(U_B)}^2\}. \quad (4.39)$$

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

Hyperellipsoids (4.37) and (4.39) have the same volume since they are the same region shifted to have a different center. The ratio of the volumes of regions (4.37) and (4.38) is

$$\frac{|\boldsymbol{S}_T^*|^{1/2}}{|\boldsymbol{S}_T^*|^{1/2}} \left(\frac{D_{(U_B)}}{D_{(U_B, T)}} \right)^g = \left(\frac{D_{(U_B)}}{D_{(U_B, T)}} \right)^g. \quad (4.40)$$

The volume of confidence region (4.38) tends to be greater than that of (4.37) 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 θ , and there is a bootstrap sample T_1^*, \dots, T_B^* where the statistic T_n is an estimator of θ 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 4.23. Efron (2014) used a similar large sample $100(1-\delta)\%$ confidence interval assuming that \bar{T}^* is asymptotically normal. The CI corresponding to (4.38) 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 (4.37)–(4.39) 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 4.13 for some theory for bootstrap CIs.

Remark 4.11. We may need $n \gg p$ before the S_T^* is a good estimator of $\text{Cov}(T) = \Sigma_T$. 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 θ , and the bootstrap sample covariance matrix S_T^* estimates $c_n \widehat{\text{Cov}}(T_n) \approx c_n \Sigma_T$ where c_n increases to 1 as $n \rightarrow \infty$. For multiple linear regression, this result happens for the residual bootstrap for least squares (OLS) with $c_n = (n-p)/n$. Then 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 $D_{\mathbf{w}}^{2*}(\bar{T}^*, S_T^*) \approx D_{\mathbf{w}}^2(\theta, \Sigma_T)/c_n$ 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 4.12. 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 S_T^* is a good estimator of $\text{Cov}(T_n^*)$. Prediction region theory uses correction factors like (4.26) and (4.19) to compensate for finite n . The bootstrap confidence regions (4.37)–(4.39) and the shorth CI use the correction factors (4.36) and (4.32) 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.

4.7.3 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 (4.38) is a large sample confidence region by Bickel and Ren (2001), (4.37) and (4.39) 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) $n\mathbf{S}_T^* \xrightarrow{P} \text{Cov}(\mathbf{u})$.

Suppose i) and ii) hold with $E(\mathbf{u}) = \mathbf{0}$ and $\text{Cov}(\mathbf{u}) = \boldsymbol{\Sigma}\mathbf{u}$. 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 4.13. 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 (4.37)–(4.39), the shorth confidence interval (4.32), and the “usual” percentile method confi-

dence interval (4.31) 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

$$\sqrt{n} (T_n - \boldsymbol{\theta}) \xrightarrow{D} \mathbf{u} \quad (4.41)$$

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

$$(T_n - \boldsymbol{\theta})^T [\mathbf{S}_T^*]^{-1} (T_n - \boldsymbol{\theta}) \xrightarrow{D} D^2 \quad (4.42)$$

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 (4.34) $\{\mathbf{w} : (\mathbf{w} - T_n)^T [\mathbf{S}_T^*]^{-1} (\mathbf{w} - T_n) \leq D_{1-\delta}^2\}$. 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 4.14. 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) $\sqrt{n}(T_i^* - \bar{T}^*) \xrightarrow{D} \mathbf{u}$. Suppose $(n\mathbf{S}_T^*)^{-1}$ is “not too ill conditioned.” Then

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

$$D_2^2 = D_{\boldsymbol{\theta}}^2(T_n, \mathbf{S}_T^*) = \sqrt{n}(T_n - \boldsymbol{\theta})^T (n\mathbf{S}_T^*)^{-1} \sqrt{n}(T_n - \boldsymbol{\theta}),$$

$$D_3^2 = D_{\boldsymbol{\theta}}^2(\bar{T}^*, \mathbf{S}_T^*) = \sqrt{n}(\bar{T}^* - \boldsymbol{\theta})^T (n\mathbf{S}_T^*)^{-1} \sqrt{n}(\bar{T}^* - \boldsymbol{\theta}), \quad \text{and}$$

$$D_4^2 = D_{T_i^*}^2(T_n, \mathbf{S}_T^*) = \sqrt{n}(T_i^* - T_n)^T (n\mathbf{S}_T^*)^{-1} \sqrt{n}(T_i^* - T_n),$$

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 (4.37), (4.39), and (4.39) will have coverage near $1 - \delta$. The regularity conditions for (4.37)–(4.39) 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 (4.34) com-

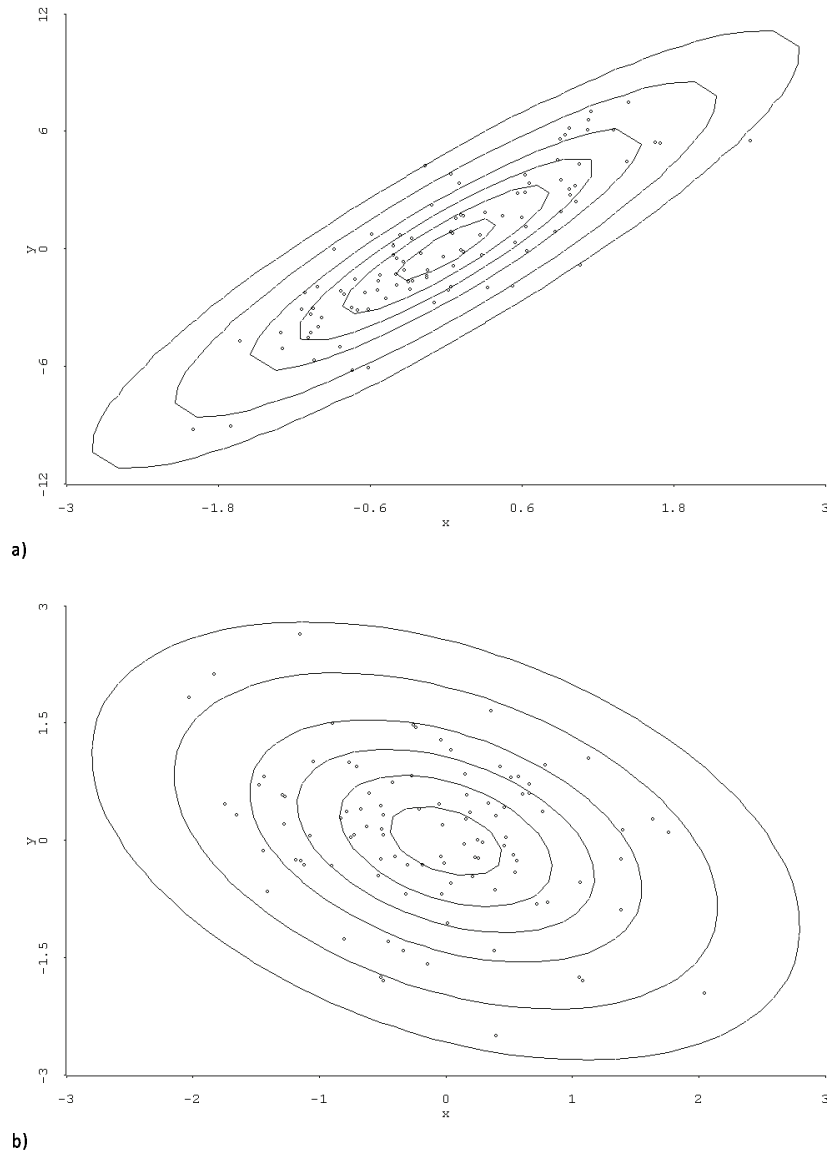


Fig. 4.3 Confidence Regions for 2 Statistics with MVN Distributions

puted 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 4.32: 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 4.22 and 4.23, and see Example 4.12.) 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 4.25. 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

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 $\sqrt{n}(T_n - \boldsymbol{\theta})$ 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 (4.37) applies the prediction region to the bootstrap. Then the hybrid region (4.39) centers that region at T_n . Hence (4.39) is a confidence region by the geometric argument, and (4.37) 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 (4.38) tend to be at least as large as the coverage and volume of (4.39).

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 4.32, 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 4.15. Remark 4.11 suggests that even if the statistic T_n is asymptotically normal so the Mahalanobis distances are asymptotically χ_g^2 , the prediction 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 4.32 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 method confidence region compensates for undercoverage when B is moderate, say $B \geq Jg$ where $J = 20$ or 50 . See Remark 4.12. 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 method 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 .

4.8 Bootstrapping Variable Selection

This section considers bootstrapping some survival regression models after variable selection, with emphasis on Cox PH regression. This section will explain why the bootstrap confidence regions (4.37), (4.38), and (4.39) give useful results. Much of the theory in Section 4.7.3 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_{j_n} with probability π_{j_n} for $j = 1, \dots, J$. Here \mathbf{A} is a known full rank $g \times p$ matrix with $1 \leq g \leq p$. We have $\sqrt{n}(T_n - \boldsymbol{\theta}) \xrightarrow{D} \mathbf{v}$ by (4.15) in Theorem 4.31 where $E(\mathbf{v}) = \mathbf{0}$, and $\boldsymbol{\Sigma}\mathbf{v} = \sum_j \sigma^2 \mathbf{A}\mathbf{V}_{j,0}\mathbf{A}^T$. Hence the geometric argument Theorem 4.32 holds: applying the prediction region (4.29) to an iid sample T_1, \dots, T_B and then centering the region at T_n gives a large sample confidence region for $\boldsymbol{\theta}$. Hence if $n\mathbf{S}_T^*$ is “not too ill conditioned,” there exists a cutoff $\hat{D}_{1-\delta}^2$ such that $\{\mathbf{w} : (\mathbf{w} - T_n)^T [\mathbf{S}_T^*]^{-1} (\mathbf{w} - T_n) \leq \hat{D}_{1-\delta}^2\}$ has coverage close to or higher than $1 - \delta$. See Remark 4.14.

We will denote the i th case by $(Z_i, \delta_i, \mathbf{x}_i)$ where $Z_i = Y_i$ if $\delta_i = 1$ so that the survival time is uncensored, and $Z_i = Y_i^*$ if $\delta_i = 0$ so that the survival time is right censored. In R , “time” is often used for the vector of Z_i and “status” for the vector of δ_i . Sometimes $T_i = Z_i$ is used for a possibly censored survival time, but in this chapter $T = \mathbf{A}\hat{\boldsymbol{\beta}}_{I_{min},0}$ is a test statistic.

Suppose the regression model satisfies $Y \perp \mathbf{x} | \mathbf{x}^T \boldsymbol{\beta}$, that Equation (2.4) holds, and that if $S \subseteq I_j$, then $\sqrt{n}(\hat{\boldsymbol{\beta}}_{I_j} - \boldsymbol{\beta}_{I_j}) \xrightarrow{D} N_{a_j}(\mathbf{0}, \mathbf{V}_j)$. Also assume that a variable selection criterion, such as AIC or relaxed lasso, is used such that $P(S \subseteq I_{min}) \rightarrow 1$ as $n \rightarrow \infty$. Hence

$$\sqrt{n}(\hat{\boldsymbol{\beta}}_{I_{j,0}} - \boldsymbol{\beta}) \xrightarrow{D} N_p(\mathbf{0}, \mathbf{V}_{j,0}) \quad (4.43)$$

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

For variable selection with $P(S \subseteq I_{min}) \rightarrow 1$ as $n \rightarrow \infty$, let $T_n = T_{k_n} = \hat{\boldsymbol{\beta}}_{I_{k,0}}$ with probabilities π_{k_n} where $\pi_{k_n} \rightarrow \pi_k$ as $n \rightarrow \infty$. Denote the π_k with $S \subseteq I_k$ by π_j . The other $\pi_k = 0$. Then Theorem 4.31 holds: $\sqrt{n}(\hat{\boldsymbol{\beta}}_{I_{min},0} - \boldsymbol{\beta}) \xrightarrow{D} \mathbf{u}$.

Note that $\mathbf{V}_{j,0}$ is singular unless I_j corresponds to the full model. For example, if $p = 3$ and model I_j uses a constant $x_1 \equiv 1$ and x_3 with

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

For variable selection, this section will show that the bootstrap sample data cloud T_1^*, \dots, T_B^* tends to be slightly more variable than the data cloud of iid T_1, \dots, T_B for large n . This result will hold for the parametric bootstrap and

nonparametric bootstrap, which are discussed in the next two subsections. Hence by the geometric argument, we expect $D_{(UB)}^2$ or $D_{(UB,T)}^2$ can be used as $\hat{D}_{1-\delta}^2$.

4.8.1 The Parametric Bootstrap

Suppose $Y_i|\mathbf{x}_i \sim D(\mathbf{x}_i^T\boldsymbol{\beta}, \gamma)$, $\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \xrightarrow{D} N_p(\mathbf{0}, \mathbf{V}(\boldsymbol{\beta}))$, and that $\mathbf{V}(\hat{\boldsymbol{\beta}}) \xrightarrow{P} \mathbf{V}(\boldsymbol{\beta})$ as $n \rightarrow \infty$. These assumptions tend to be mild for a parametric regression model where the maximum likelihood estimator (MLE) $\hat{\boldsymbol{\beta}}$ is used. Then $\mathbf{V}(\boldsymbol{\beta}) = \mathbf{I}^{-1}(\boldsymbol{\beta})$, the inverse Fisher information matrix. If $\mathbf{I}_n(\boldsymbol{\beta})$ is the Fisher information matrix based on a sample of size n , then $\mathbf{I}_n(\boldsymbol{\beta})/n \xrightarrow{P} \mathbf{I}(\boldsymbol{\beta})$. For the parametric regression model, we regress \mathbf{Y} on \mathbf{X} to obtain $(\hat{\boldsymbol{\beta}}, \hat{\gamma})$ where the $n \times 1$ vector $\mathbf{Y} = (Y_i)$ and the i th row of the $n \times p$ design matrix \mathbf{X} is \mathbf{x}_i^T .

The parametric bootstrap uses $\mathbf{Y}_j^* = (Y_i^*)$ where $Y_i^*|\mathbf{x}_i \sim D(\mathbf{x}_i^T\hat{\boldsymbol{\beta}}, \hat{\gamma})$ for $i = 1, \dots, n$. Regress \mathbf{Y}_j^* on \mathbf{X} to get $\hat{\boldsymbol{\beta}}_j^*$ for $j = 1, \dots, B$. The large sample theory for $\hat{\boldsymbol{\beta}}^*$ is simple. Note that if $Y_i^*|\mathbf{x}_i \sim D(\mathbf{x}_i^T\mathbf{b}, \hat{\gamma})$ where \mathbf{b} does not depend on n , then $(\mathbf{Y}^*, \mathbf{X})$ follows the parametric regression model with parameters $(\mathbf{b}, \hat{\gamma})$. Hence $\sqrt{n}(\hat{\boldsymbol{\beta}}^* - \mathbf{b}) \xrightarrow{D} N_p(\mathbf{0}, \mathbf{V}(\mathbf{b}))$. Now fix large integer n_0 , and let $\mathbf{b} = \hat{\boldsymbol{\beta}}_{n_0}$. Then $\sqrt{n}(\hat{\boldsymbol{\beta}}^* - \hat{\boldsymbol{\beta}}_{n_0}) \xrightarrow{D} N_p(\mathbf{0}, \mathbf{V}(\hat{\boldsymbol{\beta}}_{n_0}))$. Since $N_p(\mathbf{0}, \mathbf{V}(\hat{\boldsymbol{\beta}})) \xrightarrow{D} N_p(\mathbf{0}, \mathbf{V}(\boldsymbol{\beta}))$, we have

$$\sqrt{n}(\hat{\boldsymbol{\beta}}^* - \hat{\boldsymbol{\beta}}) \xrightarrow{D} N_p(\mathbf{0}, \mathbf{V}(\boldsymbol{\beta})) \quad (4.44)$$

as $n \rightarrow \infty$.

Now suppose $S \subseteq I$. Without loss of generality, let $\boldsymbol{\beta} = (\boldsymbol{\beta}_I^T, \boldsymbol{\beta}_O^T)^T$ and $\hat{\boldsymbol{\beta}} = (\hat{\boldsymbol{\beta}}(I)^T, \hat{\boldsymbol{\beta}}(O)^T)^T$. Then $(\mathbf{Y}, \mathbf{X}_I)$ follows the parametric regression model with parameters $(\boldsymbol{\beta}_I, \gamma)$. Hence $\sqrt{n}(\hat{\boldsymbol{\beta}}_I - \boldsymbol{\beta}_I) \xrightarrow{D} N_{a_I}(\mathbf{0}, \mathbf{V}(\boldsymbol{\beta}_I))$. Now $(\mathbf{Y}^*, \mathbf{X}_I)$ only follows the parametric regression model asymptotically, since $\hat{\boldsymbol{\beta}}(O) \neq \mathbf{0}$. However, under regularity conditions, $E(\hat{\boldsymbol{\beta}}_I^*) \approx \hat{\boldsymbol{\beta}}_I$ and $\text{Cov}(\hat{\boldsymbol{\beta}}_I^*) - \text{Cov}(\hat{\boldsymbol{\beta}}_I) \rightarrow \mathbf{0}$ as $n, B \rightarrow \infty$.

The parametric bootstrap should be useful for bootstrapping parametric survival regression models such as the Weibull PH regression model or the Weibull AFT.

4.8.2 The Nonparametric Bootstrap

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^* . This procedure is also called the *empirical bootstrap* or *naive bootstrap*.

Under regularity conditions, $\sqrt{n}(\hat{\beta}^* - \hat{\beta}) \xrightarrow{D} N_p(\mathbf{0}, \sigma^2 \mathbf{W}) \sim N_p(\mathbf{0}, \mathbf{V})$. 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 survival regression model holds for $Y_i | \mathbf{x}_i$, the \mathbf{x}_i are iid from some population with a nonsingular covariance matrix, the cases are independent, and the survival times are right censored. The cases $(Z_i, \delta_i, \mathbf{x}_i)$ are sampled with replacement. This method can be useful with proportional hazards regression models. See Burr (1994), Efron and Tibshirani (1986), and Shao, and Tu (1995).

4.8.3 Bootstrapping Variable Selection

Let the $g \times 1$ vector T_n be an estimator of the $g \times 1$ parameter vector θ . 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 $H_0 : \mathbf{A}\beta = \theta_0$ 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 explanation for why the bootstrap confidence regions (4.37), (4.38), and (4.39) give useful results after variable selection is due to Rathnayake and Olive (2019). Let the variable selection estimator $T_n = \mathbf{A}\hat{\beta}_{I_{min},0}$ with $\theta = \mathbf{A}\beta$. Then T_n is not smooth since T_n is equal to the estimator T_{jn} with probability π_{jn} for $j = 1, \dots, J$. Here \mathbf{A} is a known full rank $g \times p$ matrix with $1 \leq g \leq p$. We have $\sqrt{n}(T_n - \theta) \xrightarrow{D} \mathbf{v}$ by (4.15) where $E(\mathbf{v}) = \mathbf{0}$, and $\Sigma\mathbf{v} = \sum_j \pi_j \mathbf{A}\mathbf{V}_{j,0}\mathbf{A}^T$. Hence the geometric argument Theorem 4.32 holds: if we had iid data T_1, \dots, T_B , then R_c would be a large sample confidence region for θ . For variable selection, this section will show that the bootstrap sample data cloud T_1^*, \dots, T_B^* tends to be slightly more variable than the data cloud of iid T_1, \dots, T_B for large n . Empirically, for a mixture distribution, the bagging estimator \overline{T}^* tends to estimate θ at least as well as T_n . See Breiman (1996) and Yang (2003).

The full model should be checked before doing variable selection inference. Assume p is fixed and $n \geq 20p$. See Chapter 3 and 4. For the bootstrap, suppose that T_i^* is equal to T_{ij}^* with probability ρ_{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 $\mathbf{S}_{T,j}$.

Suppose the j th component of an iid sample T_1, \dots, T_B and the j th component of the bootstrap sample T_1^*, \dots, T_B^* have the same variability asymptotically. Since $E(T_{jn}) \approx \boldsymbol{\theta}$, each component of the iid sample is approximately centered at $\boldsymbol{\theta}$. The bootstrap components are centered at $E(T_{jn}^*)$, and often $E(T_{jn}^*) = T_{jn}$. Geometrically, separating the component clouds so that they are no longer centered at one value makes the overall data cloud larger. Thus the variability of T_n^* is larger than that of T_n for a mixture distribution, asymptotically. Hence the prediction region applied to the bootstrap sample is slightly larger than the prediction region applied to the iid sample, asymptotically (we want $n \geq 20p$). Hence cutoff $\hat{D}_{1,1-\delta}^2 = D_{(U_B)}^2$ gives coverage close to or higher than the nominal coverage for confidence regions (4.37) and (4.38), using the geometric argument. The deviation $T_i^* - T_n$ tends to be larger in magnitude than the deviation and $T_i^* - \bar{T}^*$. Hence the cutoff $\hat{D}_{2,1-\delta}^2 = D_{(U_{B,T})}^2$ tends to be larger than $D_{(U_B)}^2$, and region (4.38) tends to have higher coverage than region (4.39) for a mixture distribution.

Assume $P(S \subseteq I_{min}) \rightarrow 1$ as $n \rightarrow \infty$, and that $S \subseteq I_j$. The components of the iid sample are centered at $\mathbf{A}\boldsymbol{\beta}$ while the components of the bootstrap sample are centered at $\mathbf{A}\hat{\boldsymbol{\beta}}_{I_j,0}$. Consider regression models with $Y \perp\!\!\!\perp \mathbf{x} | \mathbf{x}^T \boldsymbol{\beta}$. Assume $\sqrt{n}\mathbf{A}(\hat{\boldsymbol{\beta}}_{I_j,0} - \boldsymbol{\beta}) \xrightarrow{D} N_{a_j}(\mathbf{0}, \boldsymbol{\Sigma}_j)$ where $\boldsymbol{\Sigma}_j = \mathbf{A}\mathbf{V}_{j,0}\mathbf{A}^T$. For the nonparametric bootstrap, assume $\sqrt{n}(\mathbf{A}\hat{\boldsymbol{\beta}}_{I_j,0}^* - \mathbf{A}\hat{\boldsymbol{\beta}}_{I_j,0}) \xrightarrow{D} N_{a_j}(\mathbf{0}, \boldsymbol{\Sigma}_j)$. Then the components of the iid sample and bootstrap sample have the same variability asymptotically. The components of iid sample are centered at $\mathbf{A}\boldsymbol{\beta}$ while the components of the bootstrap sample are centered at $\mathbf{A}\hat{\boldsymbol{\beta}}_{I_j,0}$. For the nonparametric bootstrap, the above results tend to hold if $\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \xrightarrow{D} N_p(\mathbf{0}, \mathbf{V})$ and if $\sqrt{n}(\hat{\boldsymbol{\beta}}^* - \hat{\boldsymbol{\beta}}) \xrightarrow{D} N_p(\mathbf{0}, \mathbf{V})$. Assumptions for the nonparametric bootstrap tend to be rather strong: often one assumption is that the n cases are iid from some population.

For the parametric bootstrap, Section 4.8.1 noted that under regularity conditions, $\text{Cov}(\hat{\beta}_I^*) - \text{Cov}(\hat{\beta}_I) \rightarrow \mathbf{0}$ as $n, B \rightarrow \infty$ if $S \subseteq I$. Hence $\text{Cov}(T_{jn}) - \text{Cov}(T_{jn}^*) \rightarrow \mathbf{0}$ as $n, B \rightarrow \infty$ if $S \subseteq I$. Here $T_n = \mathbf{A}\hat{\beta}_{I_{min},0}$, $T_{jn} = \mathbf{A}\hat{\beta}_{I_j,0}$, $T_n^* = \mathbf{A}\hat{\beta}_{I_{min},0}^*$, and $T_{jn}^* = \mathbf{A}\hat{\beta}_{I_j,0}^*$. Then $E(T_{jn}) \approx \mathbf{A}\beta = \theta$ while the $E(T_{jn}^*)$ are more variable than the $E(T_{jn})$ with $E(T_{jn}^*) \approx \mathbf{A}\hat{\beta}(I_j, 0)$, roughly, where $\hat{\beta}(I_j, 0)$ is formed from $\hat{\beta}(I_j)$ by adding zeros corresponding to variables not in I_j . Hence the j th component of an iid sample T_1, \dots, T_B and the j th component of the bootstrap sample T_1^*, \dots, T_B^* have the same variability asymptotically.

In simulations for $n \geq 20p$ for $H_0 : \mathbf{A}\beta_S = \theta_0$, 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 the submodel I_{min} was often more precise than inference with the full model if $n \geq 20p$ and $B \geq 50p$. It is possible that \mathbf{S}_T^* is singular if a column of the bootstrap sample is equal to $\mathbf{0}$. If the regression model has a $q \times 1$ vector of parameters γ , we may need to replace p by $p + q$.

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.

To see the effect of zero padding, consider $H_0 : \mathbf{A}\beta = \beta_O = \mathbf{0}$ where $\beta_O = (\beta_{i_1}, \dots, \beta_{i_q})^T$ and $O \subseteq E$ in Equation (2.4) so that H_0 is true. Suppose a nominal 95% confidence region is used and U_B is the 96th percentile. Hence the confidence region (4.37) or (4.38) covers at least 96% of the bootstrap sample. If $\hat{\beta}_{O,j}^* = \mathbf{0}$ for more than 4% of the $\hat{\beta}_{O,1}^*, \dots, \hat{\beta}_{O,B}^*$, 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 $\hat{\beta}_{O,j}^* = \mathbf{0}$ 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 (4.37), (4.38), or (4.39) for β_O and $\delta \in (0, 1)$, and the pvalue for $H_0 : \beta_O = \mathbf{0}$ is one. (This result holds since $\{\mathbf{0}\}$ contains 100% of the $\hat{\beta}_{O,j}^*$ 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 $\beta = (\beta_I^T, \beta_O^T)^T$. For the I_{min} model from variable 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 $\{\mathbf{0}\}$, $n \geq 20p$, and $B \geq 50p$. (Since the pvalue is one, this technique may be useful for data snooping: applying MLE theory to submodel I may have negligible selection bias.)

Remark 4.16. Note that there are several important variable selection models, including the model given by Equation (2.4) where $\mathbf{x}^T\beta = \mathbf{x}_S^T\beta_S$. Another model is $\mathbf{x}^T\beta = \mathbf{x}_{S_i}^T\beta_{S_i}$ for $i = 1, \dots, K$. Then there are $K \geq 2$ competing “true” nonnested submodels where β_{S_i} is $a_{S_i} \times 1$. For example, suppose the $K = 2$ models have predictors x_1, x_2, x_3 for S_1 and x_1, x_2, x_4 for

S_2 . Then x_3 and x_4 are likely to be selected and omitted often by forward selection for the B bootstrap samples. Hence omitting all predictors x_i that have a $\beta_{ij}^* = 0$ for at least one of the bootstrap samples $j = 1, \dots, B$ 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 β_O .

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

Another way to look at the bootstrap confidence region for 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 $100(1 - \delta)\%$.

4.8.4 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 β_i . Here $T_n = \hat{\beta}_i$ and $\theta = \beta_i$.

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 ψ 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 $\theta = \mathbf{A}\beta = \beta_E = \mathbf{0}$.

In the simulations, for $i = 1, \dots, n$, we generated $\mathbf{w}_i \sim N_p(\mathbf{0}, \mathbf{I})$ where the p elements of the vector \mathbf{w}_i are iid $N(0,1)$. Let the $p \times p$ 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 $\mathbf{z}_i = \mathbf{A}\mathbf{w}_i$ so that $Cov(\mathbf{z}_i) = \Sigma\mathbf{z} = \mathbf{A}\mathbf{A}^T = (\sigma_{ij})$ where the diagonal entries $\sigma_{ii} = [1 + (p - 1)\psi^2]$ and the off diagonal entries $\sigma_{ij} = [2\psi + (p - 2)\psi^2]$. Then $\sum_{j=1}^k z_j \sim N(0, k\sigma_{ii} + k(k - 1)\sigma_{ij}) = N(0, v^2)$. Let $\mathbf{x} = \mathbf{a}\mathbf{z}/v$. Hence the correlations are $Cor(x_i, x_j) = \rho = (2\psi + (p - 2)\psi^2)/(1 + (p - 1)\psi^2)$ for $i \neq j$. If $\psi = 1/\sqrt{cp}$, then $\rho \rightarrow 1/(c + 1)$ as $p \rightarrow \infty$ where $c > 0$. As ψ gets close to 1, the predictor vectors cluster about the line in the direction of $(1, \dots, 1)^T$. Let $SP = \mathbf{x}_i^T \beta = 1x_{i,1} + \dots + 1x_{i,k} \sim N(0, a^2)$ for $i = 1, \dots, n$. The simulations use $a = 1$ where $\beta = (1, \dots, 1, 0, \dots, 0)^T$ with k ones and $p - k$ zeros.


```

$beta
[1] 1 1 0 0
$k
[1] 2
PHbootsim(nruns=100,B=200,k=2) #fairly fast
$scicov
[1] 0.96 0.95 0.92 0.92 0.91 0.94 0.94 0.95 0.99 0.99
$avelen
[1] 0.8571470 0.8582906 0.7541797 0.7416362 2.5247451
      2.5247451 2.5558537 2.5021201 2.5021201 2.6243971
$beta
[1] 1 1 0 0
$k
[1] 2

```

The simulation computed the Frey shorth(c) interval for each β_i and used bootstrap confidence regions to test $H_0 : \beta_S = \mathbf{1}$ (whether first k $\beta_i = 1$) and $H_0 : \beta_E = \mathbf{0}$ (whether the last $p - k$ $\beta_i = 0$). 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 number of runs = 100 is tiny since the relaxed lasso simulation is slow. Using 5000 runs would be much better.

The regression models used the nonparametric bootstrap on the relaxed lasso estimator $\hat{\beta}_{I_{min},0}$. Table 4.1 gives results with $n = 100$, $p = 4$, and $k = 1$. Table 4.1 shows two rows for each model giving the observed confidence interval coverages and average lengths of the confidence intervals. The term “reg” is for the full model regression, and the term “vs” is for variable selection with relaxed lasso. The last six columns give results for the tests. The terms pr, hyb, and br are for the prediction region method (4.37), hybrid region (4.38), and Bickel and Ren region (4.39). 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 (4.40), but there is not enough information in Table 4.1 to compare the volume of the confidence region for the full model regression versus that for the relaxed lasso since the two methods have different determinants $|\mathbf{S}_T^*|$. Table 4.1 corresponds to the above R output with $k = 2$.

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 $\psi = 0$.

Table 4.1 Bootstrapping Cox PH Regression With Relaxed Lasso

ψ	β_1	β_2	β_{p-1}	β_p	pr0	hyb0	br0	pr1	hyb1	br1
reg,0	0.96	0.95	0.92	0.92	0.91	0.94	0.94	0.95	0.99	0.99
len	0.857	0.858	0.754	0.742	2.525	2.525	2.556	2.502	2.502	2.624
vs,0	0.94	0.96	0.97	0.99	0.95	0.97	0.97	0.93	0.95	0.95
len	0.864	0.847	0.733	0.722	2.556	2.556	2.662	2.512	2.512	2.625

Then it may be true that $\hat{\beta}_S$ 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 β_1 , β_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 (4.38) cutoffs and coverages were at least as high as those of the hybrid region (4.39).

4.9 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 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. For example, if $n = 500000$ and $p = 90$, using $n_1 = 900$ would result in a much smaller loss of efficiency than $n_1 = 250000$.

4.10 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, \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, $\boldsymbol{\beta}_E = \mathbf{0}$. 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 $\hat{\boldsymbol{\beta}}_{I_{min},0} = (\hat{\beta}_1, 0, \hat{\beta}_3, 0)^T$. If $S \subseteq I$, then $\sqrt{n}(\hat{\boldsymbol{\beta}}_I - \boldsymbol{\beta}_I) \xrightarrow{D} N_{a_I}(\mathbf{0}, \mathbf{V}_I)$.

3) **Theorem 4.31, Variable Selection CLT.** Assume $P(S \subseteq I_{min}) \rightarrow 1$ as $n \rightarrow \infty$, and let $T_n = \hat{\boldsymbol{\beta}}_{I_{min},0}$ and $T_{jn} = \hat{\boldsymbol{\beta}}_{I_j,0}$. Let $T_n = T_{kn} = \hat{\boldsymbol{\beta}}_{I_k,0}$ with probabilities π_{kn} where $\pi_{kn} \rightarrow \pi_k$ as $n \rightarrow \infty$. Denote the π_k with $S \subseteq I_k$ by π_j . The other $\pi_k = 0$ since $P(S \subseteq I_{min}) \rightarrow 1$ as $n \rightarrow \infty$. Assume $\sqrt{n}(\hat{\boldsymbol{\beta}}_{I_j} - \boldsymbol{\beta}_{I_j}) \xrightarrow{D} N_{a_j}(\mathbf{0}, \mathbf{V}_j)$ and $\mathbf{u}_{jn} = \sqrt{n}(\hat{\boldsymbol{\beta}}_{I_j,0} - \boldsymbol{\beta}) \xrightarrow{D} \mathbf{u}_j \sim N_p(\mathbf{0}, \mathbf{V}_{j,0})$.
a) Then

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

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 π_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

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

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 π_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 $\boldsymbol{\theta}$ is a set \mathcal{A}_n such that $P(\boldsymbol{\theta} \in \mathcal{A}_n)$ is eventually bounded below by $1 - \delta$ as $n \rightarrow \infty$.

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 $(\boldsymbol{\mu}, d\boldsymbol{\Sigma})$, 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 $H_0 : \boldsymbol{\theta} = \boldsymbol{\theta}_0$ versus $H_1 : \boldsymbol{\theta} \neq \boldsymbol{\theta}_0$ where $\boldsymbol{\theta}_0$ is a known $g \times 1$ vector. Make a confidence region and reject H_0 if $\boldsymbol{\theta}_0$ 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 $\boldsymbol{\theta}$ 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 $H_0 : \boldsymbol{\theta} = \boldsymbol{\theta}_0$ rejects H_0 if $(\bar{T}^* - \boldsymbol{\theta}_0)^T [\mathbf{S}_T^*]^{-1} (\bar{T}^* - \boldsymbol{\theta}_0) > D_{(U_B)}^2$. 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 β_1, \dots, β_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 $\boldsymbol{\beta}_O = (\beta_{i_1}, \dots, \beta_{i_g})^T$, and consider testing $H_0 : \mathbf{A}\boldsymbol{\beta}_O = \mathbf{0}$. If $\mathbf{A}\hat{\boldsymbol{\beta}}_{O,i}^* = \mathbf{0}$ 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, the $100(1 - \delta)\%$ prediction region method confidence region contains $\mathbf{0}$.)

7) **Theorem 4.32: 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 $\boldsymbol{\theta}$.

8) Applying the nonparametric prediction region (4.29) 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 $\boldsymbol{\theta}$. The j th bootstrap component data cloud is centered at $E(T_{ij}^*)$ and often $E(T_{jn}^*) = T_{jn}$. Confidence region (4.37) is the prediction region (4.29) applied to the bootstrap sample, and (4.37) is slightly larger in volume than (4.29) applied to the iid sample, asymptotically. The hybrid region (4.39) shifts (4.37) to be centered at T_n . Shifting the component clouds slightly and computing (4.29) does not change the axes of the prediction region (4.29) much compared to not shifting the component clouds. Hence by the geometric argument, we expect (4.39) 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 (4.38) tends to have higher coverage and volume than (4.39). Since \bar{T}^* tends to be closer to $\boldsymbol{\theta}$ than T_n , (4.37) 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.

4.11 Complements

Some variable selection methods for the Cox PH regression model include Fan and Li (2002), Huang et al. (2013) who give KKT conditions, Simon et al. (2013), and Tibshirani (1997). Also see Claeskens and Hjort (2008). For bootstrapping the Cox PH regression model, see Burr (1994), Efron and Tibshirani (1986), Rathnayake (2019), Rathnayake and Olive (2019), and Shao and Tu (1995). For bootstrapping some other survival analysis models, see Efron (1981), Gross and Lai (1996), and Li and Datta (2001).

This chapter followed Olive (2017b, ch. 5), Pelawa Watagoda and Olive (2019ab) and Rathnayake and Olive (2020) closely. Also see Olive (2013a, 2018), and Rathnayake (2019). 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 (1982), Efron and Hastie (2016, ch. 10–11), and Efron and Tibshirani (1993). Also see Chen (2016) and Hesterberg (2014).

There is a massive literature on variable selection and a fairly large literature for inference after variable selection. See, for example, Leeb and Pötscher (2006, 2008). Inference techniques for the variable selection model, other than data splitting, have not had much success. The methods are often inferior to data splitting, or are asymptotically equivalent to using the full model, or find a quantity to test that is not $\mathbf{A}\beta$. See Ewald and Schneider (2018).

4.12 Problems

4.1. Consider the Cushny and Peebles data set (see Staudte and Sheather 1990, p. 97) listed below. Find $\text{shorth}(7)$. 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 $\text{shorth}(5)$ for the following data set. Show work.

6 76 90 90 94 94 95 97 97 1008

4.3. Find $\text{shorth}(5)$ 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 θ of losses with the maximum likelihood estimator (MLE) \bar{X} assuming an exponential (θ) 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 $\text{shorth}(97)$ 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 μ 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.

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*.

4.8. 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%?