

Chapter 4

Prediction Regions and Bootstrap Confidence Regions

In chapter two, it was shown that applying certain prediction intervals to the bootstrap sample results in confidence intervals. In this chapter, it will be shown that applying the nonparametric prediction region to the bootstrap sample results in a confidence region. Prediction intervals are a special case of prediction regions when $p = 1$ so the $p \times 1$ random vector is a random variable.

4.1 Prediction Regions

Consider predicting a $p \times 1$ future test value \mathbf{x}_f , given past training data $\mathbf{x}_1, \dots, \mathbf{x}_n$ where $\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{x}_f$ are iid. Much as confidence regions and intervals give a measure of precision for the point estimator $\hat{\theta}$ of the parameter θ , 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 4.1. 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 is from a distribution with a pdf, we often want $P(\mathbf{x}_f \in \mathcal{A}_n) \rightarrow 1 - \delta$ as $n \rightarrow \infty$. The following definition makes sense when the highest density region is unique. Section 2.4 discussed the highest density region for a random variable where $p = 1$. Then nonzero flat spots in the pdf can cause the region to have higher than nominal coverage. For example, the highest density region of a uniform(θ_1, θ_2) random variable is not unique. See Figure 2.1 where the area under the pdf from 0 to 1 gives the 36.8% highest density region. Figure 3.1 shows the highest density regions for two bivariate normal distributions.

Definition 4.2. When unique, the $100(1 - \delta)\%$ highest density region $R(f_{1-\delta}) = \{\mathbf{z} : f(\mathbf{z}) \geq f_\delta\}$ where f_δ is the largest constant such that $P[\mathbf{x} \in R(f_{1-\delta})] \geq 1 - \delta$ and $f(\mathbf{z})$ is the probability density function (pdf) of \mathbf{x} .

Highest density regions are usually hard to estimate for p not much larger than four, but for elliptically contoured distributions with continuous decreasing g , the highest density region is the hyperellipsoid

$$\{\mathbf{z} : (\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{z} - \boldsymbol{\mu}) \leq u_{1-\delta}\} = \{\mathbf{z} : D_{\mathbf{z}}^2(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \leq u_{1-\delta}\} \quad (4.1)$$

where $P(U \leq u_{1-\delta}) = 1 - \delta$, and U is given by (3.9). If $HDR_Y(1 - \delta)$ is the $100(1 - \delta)\%$ highest density region for a random variable Y , and $X \sim U(0, \theta) \perp\!\!\!\perp Y$ (meaning X is independent of Y), then the $100(1 - \delta)\%$ highest density region for (X_f, Y_f) is

$$\{(x, y) : x \in (0, \theta), y \in HDR_Y(1 - \delta)\}.$$

There is a moderate amount of literature for prediction regions that may perform well for small p . Let $\hat{f}_{(1)}, \dots, \hat{f}_{(n)}$ be the order statistics of $\hat{f}(\mathbf{x}_1), \dots, \hat{f}(\mathbf{x}_n)$. Hyndman (1996) used the estimated highest density region

$$\hat{R}(f_{1-\delta}) = \{\mathbf{z} : d\hat{f}(\mathbf{z}) \geq d\hat{f}_{(h)}\} \quad (4.2)$$

where $d > 0$ can be any constant, $h = \max(1, \lfloor n\delta \rfloor)$, and $\lfloor x \rfloor$ is the integer part of x . Here \hat{f} is a kernel density estimator. See Remark 4.3, and see Lei et al. (2013) for references.

Let $D_{(c)}^2$ be the c th order statistic of D_1^2, \dots, D_n^2 , and consider the hyperellipsoid

$$\mathcal{A}_n = \{\mathbf{x} : D_{\mathbf{x}}^2(\bar{\mathbf{x}}, \mathbf{S}) \leq D_{(c)}^2\} = \{\mathbf{x} : D_{\mathbf{x}}(\bar{\mathbf{x}}, \mathbf{S}) \leq D_{(c)}\}. \quad (4.3)$$

If n is large, we can use $c = k_n = \lceil n(1 - \delta) \rceil$. Olive (2013a) showed that (4.3) 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 $(\bar{\mathbf{x}}, \mathbf{S})$. See Theorem 11.29. 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 (cdf) of D .

A problem with the prediction regions that cover $\approx 100(1 - \delta)\%$ of the training data cases \mathbf{x}_i (such as (4.3) 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 than on training data*. Increasing c will improve the coverage for moderate samples. Empirically for many distributions, for $n \approx 20p$, the prediction

region (4.3) applied to iid data using $k_n = \lceil n(1 - \delta) \rceil$ tended to have undercoverage as high as 5%. The undercoverage decreases rapidly as n increases. Let $q_n = \min(1 - \delta + 0.05, 1 - \delta + p/n)$ for $\delta > 0.1$ and

$$q_n = \min(1 - \delta/2, 1 - \delta + 10\delta p/n), \text{ otherwise.} \quad (4.4)$$

If $1 - \delta < 0.999$ and $q_n < 1 - \delta + 0.001$, set $q_n = 1 - \delta$. Let $D_{(U_n)}$ be the $100q_n$ th sample quantile of the D_i where

$$D_i^2 = D_{\mathbf{x}_i}^2 = (\mathbf{x}_i - \bar{\mathbf{x}})^T \mathbf{S}^{-1} (\mathbf{x}_i - \bar{\mathbf{x}}).$$

Definition 4.3. 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

$$\{\mathbf{z} : D_{\mathbf{z}}^2(\bar{\mathbf{x}}, \mathbf{S}) \leq D_{(U_n)}^2\}, \quad (4.5)$$

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

$$\{\mathbf{z} : D_{\mathbf{z}}^2(\bar{\mathbf{x}}, \mathbf{S}) \leq \chi_{p, 1-\delta}^2\}. \quad (4.6)$$

Remark 4.1. The nonparametric prediction region (4.5) 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 nonparametric prediction region is asymptotically optimal on a large class of elliptically contoured distributions in that the prediction region's volume converges in probability to the volume of the highest density region (4.1). The asymptotic coverage is $1 - \delta$ if the $100(1 - \delta)$ th percentile $D_{1-\delta}$ of D is a continuity point of the distribution of D , although prediction regions with smaller volume may exist. If $D_{1-\delta}$ 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 (4.5) contains U_n of the training data cases \mathbf{x}_i provided that \mathbf{S} is nonsingular, even if the model is wrong. For many distributions, the coverage started to be close to $1 - \delta$ for $n \geq 10p$ where the coverage is the simulated percentage of times that the prediction region contained \mathbf{x}_f .

Remark 4.2. The most used prediction regions assume that the error vectors are iid from a multivariate normal distribution. The ratio of the volumes of regions (4.5) and (4.6) 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 (4.7) usually results in prediction regions with volume and coverage that is too large.

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 $D_{\mathbf{x}}^2(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = D^2(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \sim \chi_m^2$. The Olive (2013a) nonparametric prediction region uses $(T, \mathbf{C}) = (\bar{\mathbf{x}}, \mathbf{S})$. For the classical prediction region, see Chew (1966) and Johnson and Wichern (1988, pp. 134, 151).

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, RFCH, and RMVN estimators satisfy this assumption. For $h > 0$, the hyperellipsoid

$$\{\mathbf{z} : (\mathbf{z} - T)^T \mathbf{C}^{-1} (\mathbf{z} - T) \leq h^2\} = \{\mathbf{z} : D_{\mathbf{z}}^2 \leq h^2\} = \{\mathbf{z} : D\mathbf{z} \leq h\} \quad (4.7)$$

has volume equal to

$$\frac{2\pi^{p/2}}{p\Gamma(p/2)} h^p \sqrt{\det(\mathbf{C})} = \frac{2\pi^{p/2}}{p\Gamma(p/2)} h^p b^{p/2} \sqrt{\det(\mathbf{S}_M)}. \quad (4.8)$$

A future observation (random vector) \mathbf{x}_f is in the region (4.7) 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 (4.7) 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 near (4.4). For example, use $U_n = c = \lceil nq_n \rceil$.

Remark 4.3. There may not yet be any practical competing prediction regions that do not have the form of (4.7) if p is much larger than two and the distribution of the \mathbf{x}_i is unknown. The prediction region of Section 4.3 also has this form. Remark 4.1 suggests that the nonparametric prediction region (4.5) starts to have good coverage for $n \geq 10p$ for a large class of distributions. Of course for any n there are error distributions that will have severe undercoverage. Prediction regions that estimate the pdf $f(\mathbf{z})$ with a kernel density estimator quickly become impractical as p increases since large sample sizes are needed for good estimates. See Silverman (1986, p. 129).

For example, the Hyndman nominal 95% prediction region (4.2) was computed for iid $N_p(\mathbf{0}, \mathbf{I})$ data with 1000 runs. Let the coverage be the observed proportion of prediction regions that contained the future value \mathbf{x}_f . For $p = 1$, the coverage was 0.933 for $n = 40$. For $p = 2$, the coverage was 0.911 for $n = 50$, and 0.930 for $n = 150$. For $p = 4$, the coverage was 0.920 for $n = 250$. For $p = 5$ the coverage was 0.866 for $n = 200$ and 0.934 for $n = 2000$. For $p = 8$, the coverage was 0.735 for $n = 125$. For the multivariate lognormal distribution with $n = 20p$, the Olive (2013a) large sample nonparametric 95% prediction region (4.5) had coverages 0.970, 0.959, and 0.964 for $p = 100, 200$, and 500. Some *R* code is below.

```
nruns=1000 #p = 1
count<-0
for(i in 1:nruns){
x <- rnorm(40)
xff <- rnorm(1)
count <- count + hdr2(x,xf=xff)$inr}
count #933/1000

count<-0 #p = 5
for(i in 1:nruns){
x <- matrix(rnorm(1000),ncol=5,nrow=200)
xff <- as.vector(rnorm(5))
count <- count + hdr2(x,xf=xff)$inr}
count #886/1000

#lognormal, p = 100
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
```

Olive (2013a) used three prediction regions (4.7) that can be displayed with the DD plot. The nonparametric prediction region (4.5) uses the classical

estimator $(T, \mathbf{C}) = (\bar{\mathbf{x}}, \mathbf{S})$ and $h = D_{(U_n)}$. The other two prediction regions are defined below.

Definition 4.4. The *semiparametric prediction region* uses $(T, \mathbf{C}) = (T_{RMVN}, \mathbf{C}_{RMVN})$ and $h = D_{(U_n)}$. The *parametric MVN prediction region* uses $(T, \mathbf{C}) = (T_{RMVN}, \mathbf{C}_{RMVN})$ and $h^2 = \chi_{p, q_n}^2$ where $P(W \leq \chi_{p, \delta}^2) = \delta$ if $W \sim \chi_p^2$.

All three prediction regions are asymptotically optimal for MVN distributions with nonsingular Σ . The first two prediction regions are asymptotically optimal for a large class of $EC(\boldsymbol{\mu}, \Sigma, g)$ distributions given by Assumption (E1) used in Theorem 3.20, provided g is continuous and decreasing. For distributions with nonsingular covariance matrix $c_X \Sigma$, the nonparametric region is a large sample $(1 - \delta)100\%$ prediction region, but regions with smaller volume may exist.

Notice that for the training data $\mathbf{x}_1, \dots, \mathbf{x}_n$, if \mathbf{C}^{-1} exists, then $c \approx 100q_n\%$ of the n cases are in the prediction regions for $\mathbf{x}_f = \mathbf{x}_i$, and $q_n \rightarrow 1 - \delta$ even if (T, \mathbf{C}) is not a good estimator. Hence the coverage q_n of the training data is robust to model assumptions. Of course the volume of the prediction region could be large if a poor estimator (T, \mathbf{C}) is used or if the \mathbf{x}_i do not come from an elliptically contoured distribution. Also notice that $q_n = 1 - \delta/2$ or $q_n = 1 - \delta + 0.05$ for $n \leq 20p$ and $q_n \rightarrow 1 - \delta$ as $n \rightarrow \infty$. If $q_n \equiv 1 - \delta$ and (T, \mathbf{C}) is a consistent estimator of $(\boldsymbol{\mu}, d\Sigma)$ where $d > 0$ and Σ is nonsingular, then (4.7) is a large sample prediction region, but taking q_n given by (4.4) 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 small n the resulting prediction region tended to have undercoverage as high as $\min(0.05, \delta/2)$. Using (4.4) helped reduce undercoverage for small n due to the unknown variability of (T, \mathbf{C}) .

Example 4.1. An artificial data set consisting of 100 iid cases from a

$$N_2 \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1.49 & 1.4 \\ 1.4 & 1.49 \end{pmatrix} \right)$$

distribution and 40 iid cases from a bivariate normal distribution with mean $(0, -3)^T$ and covariance \mathbf{I}_2 . Figure 4.1 shows the classical ellipsoid (with $MD \leq \sqrt{\chi_{2, 0.95}^2}$) that uses $(T, \mathbf{C}) = (\bar{\mathbf{x}}, \mathbf{S})$. The symbol “1” denotes the data while the symbol “2” is on the border of the covering ellipse. There is an *R* package that makes an ellipse. Notice that the classical parametric ellipsoid covers almost all of the data. Figure 4.2 displays the robust ellipsoid (using $RD \leq \sqrt{\chi_{2, 0.95}^2}$) which contains most of the 100 “clean” cases and excludes the 40 outliers. Problem 4.5 recreates similar figures with the classical and RMVN estimators using $q_n = 0.95$.

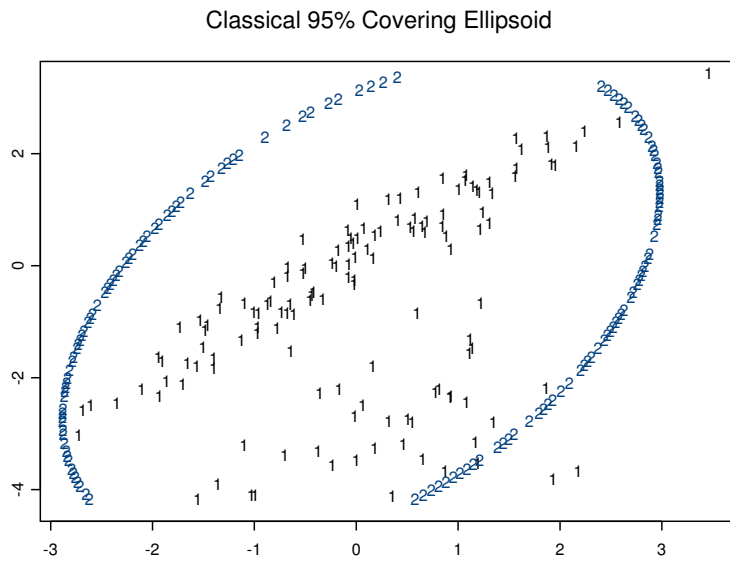


Fig. 4.1 Artificial Bivariate Data

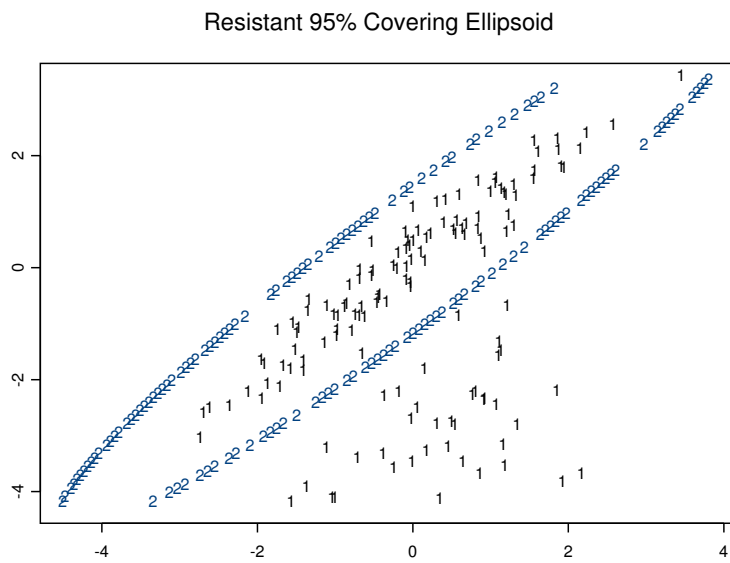


Fig. 4.2 Artificial Data

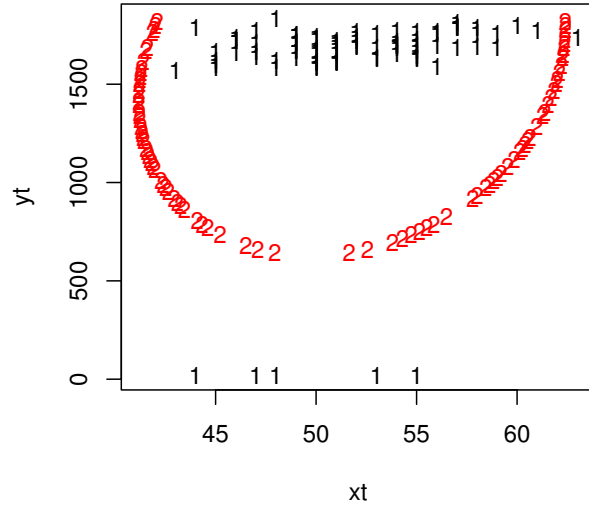


Fig. 4.3 Ellipsoid is Inflated by Outliers

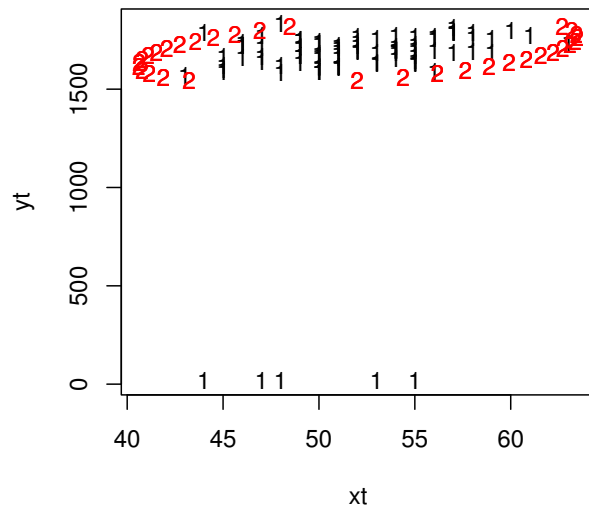


Fig. 4.4 Ellipsoid Ignores Outliers

Example 4.2. Buxton (1920) gave various measurements on 87 men including *height*, *head length*, *nasal height*, *bigonal breadth*, and *cephalic index*. Five *heights* were recorded to be about 19mm (and the actual heights for these cases were recorded as the head lengths) and are massive outliers. First *height* and *nasal height* were used with $q_n = 0.95$. Figure 4.3 shows that the classical parametric prediction region (using $MD \leq \sqrt{\chi_{2,.95}^2}$) is quite large but does not include any of the outliers. Figure 4.4 shows that the parametric MVN prediction region (using $RD \leq \sqrt{\chi_{2,.95}^2}$) is not inflated by the outliers.

Next all 87 cases and 5 predictors were used. Figure 4.5 shows the RMVN DD plot with the identity line added as a visual aid. Points to the left of the vertical line are in the nonparametric large sample 90% prediction region. Points below the horizontal line are in the semiparametric region. The horizontal line at $RD = 3.33$ corresponding to the parametric MVN 90% region is obscured by the identity line. This region contains 78 of the cases. Since $n = 87$, the nonparametric and semiparametric regions used the 95th quantile. Since there were 5 outliers, this quantile was a linear combination of the largest clean distance and the smallest outlier distance. The nonparametric and semiparametric 90% regions blow up unless the outlier proportion is small.

Figure 4.5 can be made with the following *R* commands, assuming source commands for *pack* and *robddata* have been performed. See the Preface or Section 11.2. Right click *Stop* to get the cursor.

```
x <- cbind(buxy, buxx)
ddplot4(x) #right click Stop
```

Figure 4.6 shows the DD plot and 3 prediction regions after the 5 outliers were removed. The classical and robust distances cluster about the identity line and the three regions are similar, with the parametric MVN region cutoff again at 3.33, slightly below the semiparametric region cutoff of 3.44. Cases to the left of the vertical line $MD = 3.33$ (not shown since you can mentally drop down a vertical line where the horizontal line ends at the identity line), correspond to a (modified) classical prediction region.

Figure 4.6 can be made with the following *R* commands. Right click *Stop* to get the cursor and the output following the two commands.

```
zx <- x[-c(61:65), ]
ddplot4(zx) #right click Stop
$cuplim
  95%
3.086005
$ruplim
  95%
3.438821
$mvnlim
[1] 3.327236
```

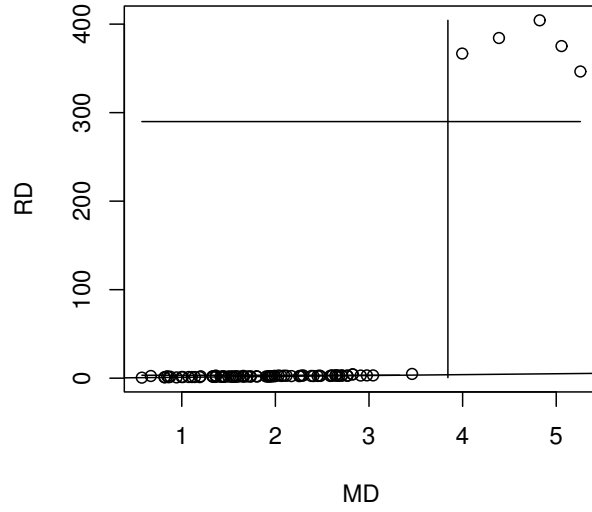


Fig. 4.5 Prediction Regions for Buxton Data

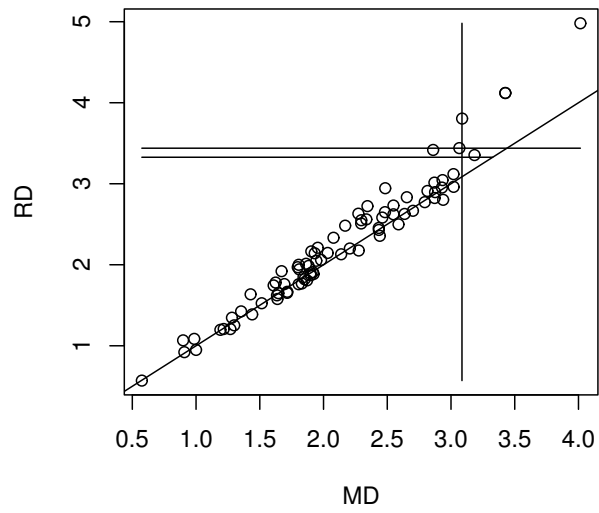


Fig. 4.6 Prediction Regions for Buxton Data without Outliers

Simulations for the prediction regions used $\mathbf{x} = \mathbf{A}\mathbf{w}$ where $\mathbf{A} = \text{diag}(\sqrt{1}, \dots, \sqrt{p})$, $\mathbf{w} \sim N_p(\mathbf{0}, \mathbf{I}_p)$ (MVN), $\mathbf{w} \sim LN(\mathbf{0}, \mathbf{I}_p)$ where the marginals are iid lognormal(0,1), or $\mathbf{w} \sim MVT_p(1)$, a multivariate t distribution with 1 degree of freedom so the marginals are iid Cauchy(0,1). All simulations used 5000 runs and $\delta = 0.1$.

Often the coverage for the semiparametric region was better than that of the nonparametric region for n near $10p$. The nonparametric covering region $\{\mathbf{z} : (\mathbf{z} - \bar{\mathbf{x}})^T \mathbf{S}^{-1}(\mathbf{z} - \bar{\mathbf{x}}) \leq D_{(n)}^2(\bar{\mathbf{x}}, \mathbf{S})\}$ uses all of the data, but for small n , data is sparse, and the covering region overfits and hence the volume is too small. The nonparametric prediction region is a hyperellipsoid that is concentric with the covering region (that replaces $D_{(U_n)}^2$ with $D_{(n)}^2$). The semiparametric region is based on the RMVN half set of data. This region is not a good estimator of the population 50% covering region for small n . Hence when it is blown up to cover 95% of the training data, the region is quite large, so it is likely that a future \mathbf{x}_f is in the region.

For large n , the semiparametric and nonparametric regions are likely to have coverage near 0.90 because the coverage on the training sample is slightly larger than 0.9 and \mathbf{x}_f comes from the same distribution as the \mathbf{x}_i . For $n = 10p$ and $2 \leq p \leq 40$, the semiparametric region had coverage near 0.9. The ratio of the volumes

$$\frac{h_i^p \sqrt{\det(\mathbf{C}_i)}}{h_2^p \sqrt{\det(\mathbf{C}_2)}}$$

was recorded where $i = 1$ was the nonparametric region, $i = 2$ was the semiparametric region, and $i = 3$ was the parametric MVN region. The volume ratio converges in probability to 1 for $N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ data, and the ratio converges to 1 for $i = 1$ if Assumption (E1) holds. The parametric MVN region often had coverage much lower than 0.9 with a volume ratio near 0, recorded as 0+. The volume ratio tends to be tiny when the coverage is much less than the nominal value 0.9. For $10p \leq n \leq 20p$, the nonparametric region often had good coverage (and volume ratio near 0.5 for MVN data).

Table 4.1 Coverages for 90% Prediction Regions

\mathbf{w} dist	n	p	ncov	scov	mcov	voln	volm
MVN	600	30	0.906	0.919	0.902	0.503	0.512
MVN	1500	30	0.899	0.899	0.900	1.014	1.027
LN	1000	10	0.903	0.906	0.567	0.659	0+
MVT(1)	1000	10	0.914	0.914	0.541	22634.3	0+

Simulations and Table 4.1 suggest that for MVN data, the coverages (ncov, scov, and mcov) for the 3 regions are near 90% for $n = 20p$ and that the volume ratios voln and volm are near 1 for $n = 50p$. With fewer than 5000 runs, this result held for $2 \leq p \leq 80$. For the non-elliptically contoured LN

data, the nonparametric region had volm well under 1, but the volume ratio blew up for $\mathbf{w} \sim MVT_p(1)$.

4.2 Bootstrap Confidence Regions

This section shows that, under regularity conditions, applying the nonparametric prediction region of Section 4.1 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. 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 $n \rightarrow \infty$.

Definition 4.5. 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 $n \rightarrow \infty$.

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 4.6. 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 4.3. Let \mathbf{w} be a random variable having the empirical distribution given by Definition 4.6. 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 2.8 was similar to Example 4.3, and see Example 2.9 for the empirical cdf of a random variable. 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 $j = 1, \dots, B$ where the n is often suppressed.

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

$$\bar{T}^* = \frac{1}{B} \sum_{i=1}^B T_i^* \quad \text{and} \quad \mathbf{S}_T^* = \frac{1}{B-1} \sum_{i=1}^B (T_i^* - \bar{T}^*)(T_i^* - \bar{T}^*)^T \quad (4.9)$$

be the sample mean and sample covariance matrix of the bootstrap sample T_1^*, \dots, T_B^* where $T_i^* = T_{i,n}^*$.

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 $k_B = \lceil B(1-\delta) \rceil$.

Remark 4.4. A useful fact for the F and chi-square distributions 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 $X \sim F_{g,d_n}$.

Definition 4.7. a) The standard bootstrap large sample $100(1-\delta)\%$ confidence region for $\boldsymbol{\theta}$ is $\{\mathbf{w} : (\mathbf{w} - T_n)^T [\mathbf{S}_T^*]^{-1} (\mathbf{w} - T_n) \leq D_{1-\delta}^2\} =$

$$\{\mathbf{w} : D_{\mathbf{w}}^2(T_n, \mathbf{S}_T^*) \leq D_{1-\delta}^2\} \quad (4.10)$$

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

$$\{\mathbf{w} : D_{\mathbf{w}}^2(T_n, \hat{\Sigma}_A/n) \leq D_{(k_B, T)}^2\} \quad (4.11)$$

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

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

Confidence region (4.10) needs $\sqrt{n}(T_n - \theta) \xrightarrow{D} N_g(\mathbf{0}, \Sigma_A)$ and $n\mathbf{S}_T^* \xrightarrow{P} \Sigma_A > 0$ as $n, B \rightarrow \infty$. See Machado and Parente (2005) for regularity conditions for this assumption.

The following three confidence regions will be used for inference after variable selection. The Olive (2017ab, 2018) prediction region method applies the nonparametric prediction region (4.5) to the bootstrap sample. Olive (2017ab, 2018) also gave the modified Bickel and Ren confidence region that uses $\hat{\Sigma}_A = n\mathbf{S}_T^*$. The hybrid confidence region is due to Pelawa Watagoda and Olive (2019). Let $q_B = \min(1 - \delta + 0.05, 1 - \delta + g/B)$ for $\delta > 0.1$ and

$$q_B = \min(1 - \delta/2, 1 - \delta + 10\delta g/B), \quad \text{otherwise.} \quad (4.12)$$

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

Definition 4.8. a) The prediction region method large sample $100(1 - \delta)\%$ confidence region for θ is $\{\mathbf{w} : (\mathbf{w} - \bar{T}^*)^T [\mathbf{S}_T^*]^{-1} (\mathbf{w} - \bar{T}^*) \leq D_{(U_B)}^2\} =$

$$\{\mathbf{w} : D_{\mathbf{w}}^2(\bar{T}^*, \mathbf{S}_T^*) \leq D_{(U_B)}^2\} \quad (4.13)$$

where $D_{(U_B)}^2$ is computed from $D_i^2 = (T_i^* - \bar{T}^*)^T [\mathbf{S}_T^*]^{-1} (T_i^* - \bar{T}^*)$ for $i = 1, \dots, B$. Note that the corresponding test for $H_0 : \theta = \theta_0$ rejects H_0 if $(\bar{T}^* - \theta_0)^T [\mathbf{S}_T^*]^{-1} (\bar{T}^* - \theta_0) > D_{(U_B)}^2$. (This procedure is basically the one sample Hotelling's T^2 test applied to the T_i^* using \mathbf{S}_T^* as the estimated covariance matrix and replacing the $\chi_{g, 1-\delta}^2$ cutoff by $D_{(U_B)}^2$.) b) The modified Bickel and Ren (2001) large sample $100(1 - \delta)\%$ confidence region is $\{\mathbf{w} : (\mathbf{w} - T_n)^T [\mathbf{S}_T^*]^{-1} (\mathbf{w} - T_n) \leq D_{(U_B, T)}^2\} =$

$$\{\mathbf{w} : D_{\mathbf{w}}^2(T_n, \mathbf{S}_T^*) \leq D_{(U_B, T)}^2\} \quad (4.14)$$

where the cutoff $D_{(U_B, T)}^2$ is the $100q_B$ th sample quantile of the $D_i^2 = (T_i^* - T_n)^T [\mathbf{S}_T^*]^{-1} (T_i^* - T_n)$. Note that the corresponding test for $H_0 : \theta = \theta_0$ rejects H_0 if $(T_n - \theta_0)^T [\mathbf{S}_T^*]^{-1} (T_n - \theta_0) > D_{(U_B, T)}^2$. c) Shift region (4.13) to have center T_n , or equivalently, change the cutoff of region (4.14) to $D_{(U_B)}^2$ to get the hybrid large sample $100(1 - \delta)\%$ confidence region: $\{\mathbf{w} : (\mathbf{w} - T_n)^T [\mathbf{S}_T^*]^{-1} (\mathbf{w} - T_n) \leq D_{(U_B)}^2\} =$

$$\{\mathbf{w} : D_{\mathbf{w}}^2(T_n, \mathbf{S}_T^*) \leq D_{(U_B)}^2\}. \quad (4.15)$$

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

Hyperellipsoids (4.13) and (4.15) have the same volume since they are the same region shifted to have a different center. The ratio of the volumes of regions (4.13) and (4.14) is

$$\frac{|\mathbf{S}_T^*|^{1/2}}{|\mathbf{S}_T^*|^{1/2}} \left(\frac{D_{(U_B)}}{D_{(U_B, T)}} \right)^g = \left(\frac{D_{(U_B)}}{D_{(U_B, T)}} \right)^g. \quad (4.16)$$

The volume of confidence region (4.14) tends to be greater than that of (4.13) since the T_i^* are closer to \bar{T}^* than T_n on average.

Next we review the Section 2.5 confidence intervals corresponding to the three confidence regions if $g = 1$. 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 $g = 1$. Efron (2014) used a similar large sample $100(1 - \delta)\%$ confidence interval assuming that \bar{T}^* is asymptotically normal. The CI corresponding to (4.14) is defined similarly, and $[T_n - a_{(U_B)}, T_n + a_{(U_B)}]$ is the CI for (4.15). Note that the three CIs corresponding to (4.13)–(4.15) can be computed without finding S_T^* or $D_{(U_B)}$ even if $S_T^* = 0$. The shorth(c) CI (2.13) computed from the T_i^* can be much shorter than the Efron (2014) or prediction region method confidence intervals. See Remark 2.5 for some theory for bootstrap CIs.

Remark 4.5. Suppose the $p \times 1$ vector $\hat{\boldsymbol{\beta}}$, and $\boldsymbol{\theta} = \mathbf{A}\hat{\boldsymbol{\beta}}$ is $g \times 1$. We will often want $n \geq 20p$ and $B \geq \max(100, n, 50p)$. If $T_n = \mathbf{A}\hat{\boldsymbol{\beta}}$ is $g \times 1$, we might replace p by g or replace p by d if d is the model degrees of freedom. Sometimes much larger n is needed to avoid undercoverage. We want $B \geq 50g$ so that \mathbf{S}_T^* is a good estimator of $Cov(T_n^*)$. Prediction region theory uses correction factors like (4.4) to compensate for finite n . The bootstrap confidence regions (4.13)–(4.15) and the shorth CI use the correction factors (4.12) and (2.14) 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.

4.3 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 (4.14) is a large sample confidence region by Bickel and Ren (2001), (4.13) and (4.15) 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 (2019) proof is simpler. Remark 2.5 gave theory for bootstrap confidence intervals.

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

$$\sqrt{n}(T_n - \boldsymbol{\theta}) \xrightarrow{D} \mathbf{u} \quad (4.17)$$

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

$$(T_n - \boldsymbol{\theta})^T [\mathbf{S}_T^*]^{-1} (T_n - \boldsymbol{\theta}) \xrightarrow{D} D^2 \quad (4.18)$$

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 (4.10) $\{\mathbf{w} : (\mathbf{w} - T_n)^T [\mathbf{S}_T^*]^{-1} (\mathbf{w} - T_n) \leq D_{1-\delta}^2\}$. 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 4.6. Under reasonable conditions, i) $\sqrt{n}(T_n - \boldsymbol{\theta}) \xrightarrow{D} \mathbf{u}$, ii) $\sqrt{n}(T_i^* - T_n) \xrightarrow{D} \mathbf{u}$, iii) $\sqrt{n}(\bar{T}^* - \boldsymbol{\theta}) \xrightarrow{D} \mathbf{u}$, and iv) $\sqrt{n}(T_i^* - \bar{T}^*) \xrightarrow{D} \mathbf{u}$. Then

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

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

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

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

are well behaved. If $(n\mathbf{S}_T^*)^{-1} \xrightarrow{P} \boldsymbol{\Sigma}_A^{-1}$, then $D_j^2 \xrightarrow{D} D^2 = \mathbf{u}^T \boldsymbol{\Sigma}_A^{-1} \mathbf{u}$. 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 (4.13), (4.14), and (4.15) will have coverage near $1 - \delta$. The regularity conditions for (4.13)–(4.15) are weaker when $g = 1$, since \mathbf{S}_T^* does not need to be computed.

The following Pelawa Watagoda and Olive (2019) theorem is very useful. Let $D_{(U_B)}^2$ be the cutoff for the nonparametric prediction region (4.5) computed from the $D_i^2(\bar{T}, \mathbf{S}_T)$ for $i = 1, \dots, B$. Hence n is replaced by B . Since T_n depends on the sample size n , we need $(n\mathbf{S}_T)^{-1}$ to be fairly well behaved (“not too ill conditioned”) for each $n \geq 20g$, say. This condition is weaker than $(n\mathbf{S}_T)^{-1} \xrightarrow{P} \boldsymbol{\Sigma}_A^{-1}$. Note that $T_i = T_{in}$. In the following theorem, note that we can replace \sqrt{n} by n^δ where $0 < \delta \leq 1$.

Theorem 4.1: 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 $\boldsymbol{\Sigma}_{T_n}$. Assume $(n\mathbf{S}_T)^{-1} \xrightarrow{P} \boldsymbol{\Sigma}_A^{-1}$. Then the large sample

100(1 - δ)% prediction region $R_p = \{\mathbf{w} : D_{\mathbf{w}}^2(\bar{\mathbf{T}}, \mathbf{S}_T) \leq D_{(U_B)}^2\}$ centered at $\bar{\mathbf{T}}$ 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 - δ)% 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 $\bar{\mathbf{T}}$ 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 11.26 and 11.27, and see Example 11.12.) For fixed B , the average of these random vectors is

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

by Theorem 11.29. Hence $(\bar{\mathbf{T}} - \boldsymbol{\theta}) = O_P((nB)^{-1/2})$, and $\bar{\mathbf{T}}$ gets arbitrarily close to $\boldsymbol{\theta}$ compared to T_n as $B \rightarrow \infty$. Thus R_c is a large sample 100(1 - δ)% confidence region for $\boldsymbol{\theta}$ as $n, B \rightarrow \infty$. \square

Remark 4.7. Theorem 4.1 is useful for explaining why a correction factor is needed. R_c contains $\bar{\mathbf{T}}$ with probability $1 - \delta_B$ and there is a hyperellipsoid $R_{\bar{\mathbf{T}}}$ about $\bar{\mathbf{T}}$ that contains $\boldsymbol{\theta}$ with high probability where the volume of $R_{\bar{\mathbf{T}}}$ goes to 0 as $B \rightarrow \infty$. For finite $B \approx 50g$, the volume of the hyperellipsoid $R_{\bar{\mathbf{T}}}$ is small compared to that of R_c but is not zero. As B increases, covering $\bar{\mathbf{T}}$ also covers most of $R_{\bar{\mathbf{T}}}$, and the confidence region coverage gets near the nominal level $1 - \delta$. When B is near $50g$, covering $\bar{\mathbf{T}}$ does not necessarily cover most of $R_{\bar{\mathbf{T}}}$, and hence there may be undercoverage for $\boldsymbol{\theta}$ if $U_B = \lceil B(1 - \delta) \rceil$. Using the correction factor (4.12) increases the coverage of $\bar{\mathbf{T}}$, $R_{\bar{\mathbf{T}}}$, and $\boldsymbol{\theta}$ when B is near $50g$.

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 (4.13) applies the prediction region to the bootstrap. Then the hybrid region (4.15) centers that region at T_n . Hence (4.15) is a confidence region by the geometric argument, and (4.13) is a confidence region if $\sqrt{n}(\bar{\mathbf{T}}^* - T_n) \xrightarrow{P} \mathbf{0}$. Since the T_i^* are closer to $\bar{\mathbf{T}}^*$ 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 (4.14) tend to be at least as large as the coverage and volume of (4.15).

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 3.1 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 3.1 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 4.1, 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 3.1. 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 4.8. 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 4.1 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 . 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 .

4.4 Data Splitting

Data splitting can be used to get prediction regions using estimators such as $(T, \mathbf{C}) = (T_{RMV N}, \mathbf{C}_{RMV N})$ or $(T, \mathbf{C}) = (\text{MED}(\mathbf{W}), \mathbf{I}_p)$ if $\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{x}_f$ are iid. Data splitting divides the training data $\mathbf{x}_1, \dots, \mathbf{x}_n$ into two sets H and V where H has $n_H \geq n/2$ of the cases and V has the remaining $n_V = n - n_H$ cases i_1, \dots, i_{n_V} . The estimator (T_H, \mathbf{C}_H) is computed using the data set H . Then the squared validation distances $D_j^2 = D_{\mathbf{x}_{i_j}}^2(T_H, \mathbf{C}_H) = (\mathbf{x}_{i_j} - T_H)^T \mathbf{C}_H^{-1} (\mathbf{x}_{i_j} - T_H)$ are computed for the $j = 1, \dots, n_V$ cases in the validation set V . Let $D_{(U_V)}^2$ be the U_V th order statistic of the D_j^2 where

$$U_V = \min(n_V, \lceil (n_V + 1)(1 - \delta) \rceil). \quad (4.19)$$

Definition 4.9. The large sample $100(1 - \delta)\%$ data splitting prediction region for \mathbf{x}_f is

$$\{\mathbf{z} : D_{\mathbf{z}}^2(T_H, \mathbf{C}_H) \leq D_{(U_V)}^2\}. \quad (4.20)$$

To show that (4.20) is a prediction region, suppose the \mathbf{x}_i are iid for $i = 1, \dots, n, n+1$ where $\mathbf{x}_f = \mathbf{x}_{n+1}$. Compute (T_H, \mathbf{C}_H) from the cases in H . Consider the squared validation distances D_k^2 for $k = 1, \dots, n_V$ and the squared validation distance $D_{n_V+1}^2$ for case \mathbf{x}_f . Since these $n_V + 1$ cases are iid, the probability that D_t^2 has rank j for $j = 1, \dots, n_V + 1$ is $1/(n_V + 1)$ for each t , i.e., the ranks follow the discrete uniform distribution. Let $t = n_V + 1$ and let the $D_{(j)}^2$ be the ordered squared validation distances using $j = 1, \dots, n_V$. That is, get the order statistics without using the unknown squared validation distance $D_{n_V+1}^2$. Then $D_{(i)}^2$ has rank i if $D_{(i)}^2 < D_{n_V+1}^2$ but rank $i + 1$ if $D_{(i)}^2 > D_{n_V+1}^2$. Thus $D_{(U_V)}^2$ has rank $U_V + 1$ if $D_{\mathbf{x}_f}^2 < D_{(U_V)}^2$ and

$$P(\mathbf{x}_f \in \{\mathbf{z} : D_{\mathbf{z}}^2(T_H, \mathbf{C}_H) \leq D_{(U_V)}^2\}) = P(D_{\mathbf{x}_f}^2 \leq D_{(U_V)}^2) \geq U_V / (1 + n_V) \rightarrow$$

$1 - \delta$ as $n_V \rightarrow \infty$. If there are no tied ranks, then

$$P(D_{\mathbf{x}_f}^2 \leq D_{(U_V)}^2) = P(D_{\mathbf{x}_f}^2 < D_{(U_V)}^2) = P(\text{rank of } D_{\mathbf{x}_f}^2 \leq U_V) = U_V / (1 + n_V).$$

Note that we can get coverage close to $1 - \delta$ for $n_V \geq 20$ for $\delta = 0.05$ even if (T_H, \mathbf{C}_H) is a bad estimator. The volume of the prediction region tends to be much larger than that of the highest density region, even if \mathbf{C}_H is well conditioned. We likely need $U_V \geq 50$ for $D_{(U_V)}^2$ to approximate the population percentile of $D_j^2 = (\mathbf{x}_{i_j} - T_H)^T \mathbf{C}_H^{-1} (\mathbf{x}_{i_j} - T_H)$.

As an example, consider using $(T, \mathbf{C}) = (\text{MED}(\mathbf{W}), \mathbf{I}_p)$. Then the prediction region is a hypersphere centered at the coordinatewise median. The prediction region is good if the iid $\mathbf{x}_i \sim N_p(\boldsymbol{\mu}, \sigma^2 \mathbf{I}_p)$, but if the $\mathbf{x}_i \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ such that highest density region is a hyperellipsoid tightly clustered about a vector in the direction of $\mathbf{1}$, then the prediction region (4.20) has huge volume compared to the highest density region.

If $p > n$, prediction region (4.20) can be used as long as \mathbf{C} is nonsingular. Then $\mathbf{C} = \mathbf{I}_p$, $\mathbf{C} = \text{diag}(S_1^2, \dots, S_p^2)$, or

$$\mathbf{C} = \text{diag}([\text{MAD}(x_{11}, \dots, x_{n1})]^2, \dots, [\text{MAD}(x_{1p}, \dots, x_{np})]^2)$$

could be used. Regularized covariance matrices or precision matrices could also be used.

If $n \geq 20p$, using $(T, \mathbf{C}) = (T_{RMV N}, \mathbf{C}_{RMV N})$ might result in a prediction region with smaller volume than using $(T, \mathbf{C}) = (\bar{\mathbf{x}}, \mathbf{S})$ since the robust estimator attempts to estimate a small volume hyperellipsoid. Also, if $D_{(U_V)}^2 \approx D_{(U_n)}^2$ in Definition 4.4, then the semiparametric region using all n cases should have good coverage.

4.5 Summary

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 4.1: 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 (4.24) to the iid data T_1, \dots, T_B results in the $100(1 - \delta)\%$ confidence region $\{\mathbf{w} : (\mathbf{w} - T_n)^T \mathbf{S}_T^{-1} (\mathbf{w} - T_n) \leq D_{(U_B)}^2(T_n, \mathbf{S}_T)\}$ where $D_{(U_B)}^2(T_n, \mathbf{S}_T)$ is computed from the $(T_i - T_n)^T \mathbf{S}_T^{-1} (T_i - T_n)$ provided the $\mathbf{S}_T = \mathbf{S}_{T_n}$ are “not too ill conditioned.” For OLS variable selection, assume there are two or more component clouds. The bootstrap component data clouds have the same asymptotic covariance matrix as the iid component data clouds, which are centered at $\boldsymbol{\theta}$. The j th bootstrap component data cloud is centered at $E(T_{ij}^*)$ and often $E(T_{jn}^*) = T_{jn}$. Confidence region (4.32) is the prediction region (4.24) applied to the bootstrap sample, and (4.32) is slightly larger in volume than (4.24) applied to the iid sample, asymptotically. The hybrid region (4.34) shifts (4.32) to be centered at T_n . Shifting the component clouds slightly and computing (4.24) does not change the axes of the prediction region (4.24) much compared to not shifting the component clouds. Hence by the geometric argument, we expect (4.34) to have coverage at least as high as the nominal, asymptotically, provided the \mathbf{S}_T^* are “not too ill conditioned.” The Bickel and Ren confidence region (4.33) tends to have higher coverage and volume than (4.34). Since \bar{T}^* tends to be closer to $\boldsymbol{\theta}$ than T_n , (4.32) tends to have good coverage.

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

4.6 Complements

There are few practical competitors for the prediction regions in Sections 4.1 and 4.3. Parametric regions such as the classical region for multivariate normal data tend to have severe undercoverage because the data rarely follows the parametric distribution. Procedures that use brand name high breakdown multivariate location and dispersion estimators take too long to compute for $p > 2$. An interesting idea is to estimate the pdf of the data, then use the pdf to find small prediction regions. The problem with these regions is that nonparametric pdf estimators do not work well for $p > 4$. See Lei et al. (2013). A useful application of prediction regions is Mykland (2003).

Bickel and Ren (2001) have interesting sufficient conditions for (4.11) to be a confidence region when $\hat{\Sigma}_A$ is a consistent estimator of positive definite Σ_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 $\sqrt{n}(T_i^* - T_n) \xrightarrow{D} \lim_{n \rightarrow \infty} \dot{T}(F_n)\mathbf{z}_{F_n} = \dot{T}(F)\mathbf{z}_F \sim \mathbf{u}$.

Good references for the bootstrap include Efron (1979, 1982), Efron and Hastie (2016, ch. 10–11), and Efron and Tibshirani (1993). Also see Chen (2016), Hesterberg (2014), and Rajapaksha and Olive (2021). One of the sufficient conditions for the bootstrap confidence region is that T has a well behaved Hadamard derivative. Fréchet differentiability implies Hadamard differentiability, and many statistics are shown to be Hadamard differentiable in Bickel and Ren (2001), Clarke (1986b, 2000), Fernholtz (1983), Gill (1989),

Ren (1991), and Ren and Sen (1995). Bickel and Ren (2001) showed that their method can work when Hadamard differentiability fails.

For bootstrapping robust estimators, see Olive (2017b), Rupasinghe Arachchige Don and Olive D.J. (2019) and Rupasinghe Arachchige Don and Pelawa Watagoda (2018).

4.7 Problems

R Problems Use the command `source("G:/rpack.txt")` to download the functions and the command `source("G:/robddata.txt")` to download the data. See Preface or Section 11.2. Typing the name of the `rpack` function, e.g. `covmba`, will display the code for the function. Use the `args` command, e.g. `args(covmba)`, to display the needed arguments for the function. For some of the following problems, the *R* commands can be copied and pasted from (<http://parker.ad.siu.edu/Olive/robRhw.txt>) into *R*.

4.1. Use the *R* source commands and then type `ddplot4(buwx, alpha=0.2)` and put the plot in *Word*. The Buxton data has 5 outliers, $p = 4$, and $n = 87$, so the 80% prediction region uses the $100(1 - \delta + p/n) = 84.6$ th percentile. The output shows that the cutoffs are 2.527, 2.734, and 2.583 for the nonparametric, semiparametric, and robust parametric prediction regions. The two horizontal lines that correspond to the robust distances are obscured by the identity line. (Right click *Stop* once on the plot.)

4.2. 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 $n \rightarrow \infty$. Were the three coverages near 90%?

4.3. The function `predsim2` computes the data splitting prediction region. The output gives `cvr` = observed coverage, `up` \approx actual coverage, and `mnhsq` = mean cutoff $D_{(U_V)}^2$. With 5000 runs, expect observed coverage $\in [0.94, 0.96]$ if the actual coverage is close to 0.95.

a) When `xtype=3` and `dtype=1`, $(T, C) = (\bar{\mathbf{x}}, \mathbf{I}_p)$ where $\mathbf{x}_i \sim N_p(\mathbf{0}, \mathbf{I}_p)$. If $n \geq \max(20p, 200)$ and $n_V = 100$, then $D_{(U_V)}^2$ should estimate the population percentile $\chi_{p,0.95}^2$. Copy and paste the commands for this problem into *R*. Include the output in *Word*.

- i) Was the observed coverage near the actual coverage?
- ii) Was the `mnhsq` near 18.3?

b) When $xtype = 1$, $\mathbf{x}_i \sim N_p(\mathbf{0}, \text{diag}(1, \dots, p))$ and the χ^2 approximation no longer holds. Copy and paste the commands for this problem into *R*. Include the output in *Word*.

i) Was the observed coverage near the actual coverage?

ii) Was the `mnhsq` a lot larger than 18.3? (If so, then the volume of the prediction region is much larger than that in a.)

c) Copy and paste the commands for this problem into *R*. Include the output in *Word*. Now $p > n$. Were the observed and actual coverages close?