

## Chapter 10

# Multivariate Linear Regression

This chapter will show that multivariate linear regression with  $m \geq 2$  response variables is nearly as easy to use, at least if  $m$  is small, as multiple linear regression which has 1 response variable. *For multivariate linear regression, at least one predictor variable is quantitative.* Plots for checking the model, including outlier detection, are given. Prediction regions that are robust to nonnormality are developed. For hypothesis testing, it is shown that the Wilks' lambda statistic, Hotelling Lawley trace statistic, and Pillai's trace statistic are robust to nonnormality.

### 10.1 Introduction

**Definition 10.1.** The **response variables** are the variables that you want to predict. The **predictor variables** are the variables used to predict the response variables.

**Definition 10.2.** The **multivariate linear regression model**

$$\mathbf{y}_i = \mathbf{B}^T \mathbf{x}_i + \epsilon_i$$

for  $i = 1, \dots, n$  has  $m \geq 2$  response variables  $Y_1, \dots, Y_m$  and  $p$  predictor variables  $x_1, x_2, \dots, x_p$  where  $x_1 \equiv 1$  is the trivial predictor. The  $i$ th case is  $(\mathbf{x}_i^T, \mathbf{y}_i^T) = (1, x_{i2}, \dots, x_{ip}, Y_{i1}, \dots, Y_{im})$  where the 1 could be omitted. The model is written in matrix form as  $\mathbf{Z} = \mathbf{X}\mathbf{B} + \mathbf{E}$  where the matrices are defined below. The model has  $E(\epsilon_k) = \mathbf{0}$  and  $\text{Cov}(\epsilon_k) = \mathbf{\Sigma}\epsilon = (\sigma_{ij})$  for  $k = 1, \dots, n$ . Then the  $p \times m$  coefficient matrix  $\mathbf{B} = [\beta_1 \beta_2 \dots \beta_m]$  and the  $m \times m$  covariance matrix  $\mathbf{\Sigma}\epsilon$  are to be estimated, and  $E(\mathbf{Z}) = \mathbf{X}\mathbf{B}$  while  $E(Y_{ij}) = \mathbf{x}_i^T \beta_j$ . The  $\epsilon_i$  are assumed to be iid. Multiple linear regression corresponds to  $m = 1$  response variable, and is written in matrix form as  $\mathbf{Y} = \mathbf{X}\beta + \mathbf{e}$ . Subscripts are needed for the  $m$  multiple linear regression

models  $\mathbf{Y}_j = \mathbf{X}\boldsymbol{\beta}_j + \mathbf{e}_j$  for  $j = 1, \dots, m$  where  $E(\mathbf{e}_j) = \mathbf{0}$ . For the multivariate linear regression model,  $\text{Cov}(\mathbf{e}_i, \mathbf{e}_j) = \sigma_{ij} \mathbf{I}_n$  for  $i, j = 1, \dots, m$  where  $\mathbf{I}_n$  is the  $n \times n$  identity matrix.

**Notation.** The **multiple linear regression model** uses  $m = 1$ . See Definition 1.9. The **multivariate linear model**  $\mathbf{y}_i = \mathbf{B}^T \mathbf{x}_i + \boldsymbol{\epsilon}_i$  for  $i = 1, \dots, n$  has  $m \geq 2$ , and multivariate linear regression and MANOVA models are special cases. See Definition 9.2. This chapter will use  $x_1 \equiv 1$  for the multivariate linear regression model. The **multivariate location and dispersion model** is the special case where  $\mathbf{X} = \mathbf{1}$  and  $p = 1$ .

The data matrix  $\mathbf{W} = [\mathbf{X} \ \mathbf{Z}]$  except usually the first column  $\mathbf{1}$  of  $\mathbf{X}$  is omitted for software. The  $n \times m$  matrix

$$\mathbf{Z} = \begin{bmatrix} Y_{1,1} & Y_{1,2} & \dots & Y_{1,m} \\ Y_{2,1} & Y_{2,2} & \dots & Y_{2,m} \\ \vdots & \vdots & \ddots & \vdots \\ Y_{n,1} & Y_{n,2} & \dots & Y_{n,m} \end{bmatrix} = [\mathbf{Y}_1 \ \mathbf{Y}_2 \ \dots \ \mathbf{Y}_m] = \begin{bmatrix} \mathbf{y}_1^T \\ \vdots \\ \mathbf{y}_n^T \end{bmatrix}.$$

The  $n \times p$  design matrix of predictor variables is

$$\mathbf{X} = \begin{bmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,p} \\ x_{2,1} & x_{2,2} & \dots & x_{2,p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n,1} & x_{n,2} & \dots & x_{n,p} \end{bmatrix} = [\mathbf{v}_1 \ \mathbf{v}_2 \ \dots \ \mathbf{v}_p] = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix}$$

where  $\mathbf{v}_1 = \mathbf{1}$ .

The  $p \times m$  matrix

$$\mathbf{B} = \begin{bmatrix} \beta_{1,1} & \beta_{1,2} & \dots & \beta_{1,m} \\ \beta_{2,1} & \beta_{2,2} & \dots & \beta_{2,m} \\ \vdots & \vdots & \ddots & \vdots \\ \beta_{p,1} & \beta_{p,2} & \dots & \beta_{p,m} \end{bmatrix} = [\boldsymbol{\beta}_1 \ \boldsymbol{\beta}_2 \ \dots \ \boldsymbol{\beta}_m].$$

The  $n \times m$  matrix

$$\mathbf{E} = \begin{bmatrix} \epsilon_{1,1} & \epsilon_{1,2} & \dots & \epsilon_{1,m} \\ \epsilon_{2,1} & \epsilon_{2,2} & \dots & \epsilon_{2,m} \\ \vdots & \vdots & \ddots & \vdots \\ \epsilon_{n,1} & \epsilon_{n,2} & \dots & \epsilon_{n,m} \end{bmatrix} = [\mathbf{e}_1 \ \mathbf{e}_2 \ \dots \ \mathbf{e}_m] = \begin{bmatrix} \boldsymbol{\epsilon}_1^T \\ \vdots \\ \boldsymbol{\epsilon}_n^T \end{bmatrix}.$$

Considering the  $i$ th row of  $\mathbf{Z}$ ,  $\mathbf{X}$ , and  $\mathbf{E}$  shows that  $\mathbf{y}_i^T = \mathbf{x}_i^T \mathbf{B} + \boldsymbol{\epsilon}_i^T$ .

Each response variable in a multivariate linear regression model follows a multiple linear regression model  $\mathbf{Y}_j = \mathbf{X}\boldsymbol{\beta}_j + \mathbf{e}_j$  for  $j = 1, \dots, m$  where it

is assumed that  $E(\mathbf{e}_j) = \mathbf{0}$  and  $\text{Cov}(\mathbf{e}_j) = \sigma_{jj}\mathbf{I}_n$ . Hence the errors corresponding to the  $j$ th response are uncorrelated with variance  $\sigma_j^2 = \sigma_{jj}$ . Notice that the **same design matrix**  $\mathbf{X}$  of predictors is used for each of the  $m$  models, but the  $j$ th response variable vector  $\mathbf{Y}_j$ , coefficient vector  $\boldsymbol{\beta}_j$ , and error vector  $\mathbf{e}_j$  change and thus depend on  $j$ .

Now consider the  $i$ th case  $(\mathbf{x}_i^T, \mathbf{y}_i^T)$  which corresponds to the  $i$ th row of  $\mathbf{Z}$  and the  $i$ th row of  $\mathbf{X}$ . Then

$$\begin{bmatrix} Y_{i1} = \beta_{11}x_{i1} + \cdots + \beta_{p1}x_{ip} + \epsilon_{i1} = \mathbf{x}_i^T \boldsymbol{\beta}_1 + \epsilon_{i1} \\ Y_{i2} = \beta_{12}x_{i1} + \cdots + \beta_{p2}x_{ip} + \epsilon_{i2} = \mathbf{x}_i^T \boldsymbol{\beta}_2 + \epsilon_{i2} \\ \vdots \\ Y_{im} = \beta_{1m}x_{i1} + \cdots + \beta_{pm}x_{ip} + \epsilon_{im} = \mathbf{x}_i^T \boldsymbol{\beta}_m + \epsilon_{im} \end{bmatrix}$$

or  $\mathbf{y}_i = \boldsymbol{\mu}_{\mathbf{x}_i} + \boldsymbol{\epsilon}_i = E(\mathbf{y}_i) + \boldsymbol{\epsilon}_i$  where

$$E(\mathbf{y}_i) = \boldsymbol{\mu}_{\mathbf{x}_i} = \mathbf{B}^T \mathbf{x}_i = \begin{bmatrix} \mathbf{x}_i^T \boldsymbol{\beta}_1 \\ \mathbf{x}_i^T \boldsymbol{\beta}_2 \\ \vdots \\ \mathbf{x}_i^T \boldsymbol{\beta}_m \end{bmatrix}.$$

The notation  $\mathbf{y}_i|\mathbf{x}_i$  and  $E(\mathbf{y}_i|\mathbf{x}_i)$  is more accurate, but usually the conditioning is suppressed. Taking  $\boldsymbol{\mu}_{\mathbf{x}_i}$  to be a constant (or condition on  $\mathbf{x}_i$  if the predictor variables are random variables),  $\mathbf{y}_i$  and  $\boldsymbol{\epsilon}_i$  have the same covariance matrix. In the multivariate regression model, this covariance matrix  $\boldsymbol{\Sigma}_{\boldsymbol{\epsilon}}$  does not depend on  $i$ . Observations from different cases are uncorrelated (often independent), but the  $m$  errors for the  $m$  different response variables for the *same case* are correlated. If  $\mathbf{X}$  is a random matrix, then assume  $\mathbf{X}$  and  $\mathbf{E}$  are independent and that expectations are conditional on  $\mathbf{X}$ .

**Example 10.1.** Suppose it is desired to predict the response variables  $Y_1 = \text{height}$  and  $Y_2 = \text{height at shoulder}$  of a person from partial skeletal remains. A model for prediction can be built from nearly complete skeletons or from living humans, depending on the population of interest (e.g. ancient Egyptians or modern US citizens). The predictor variables might be  $x_1 \equiv 1$ ,  $x_2 = \text{femur length}$ , and  $x_3 = \text{ulna length}$ . The two heights of individuals with  $x_2 = 200\text{mm}$  and  $x_3 = 140\text{mm}$  should be shorter on average than the two heights of individuals with  $x_2 = 500\text{mm}$  and  $x_3 = 350\text{mm}$ . In this example  $Y_1$ ,  $Y_2$ ,  $x_2$ , and  $x_3$  are quantitative variables. If  $x_4 = \text{gender}$  is a predictor variable, then gender (coded as male = 1 and female = 0) is qualitative.

**Definition 10.3.** Least squares is the classical method for fitting multivariate linear regression. The **least squares estimators** are

$$\hat{\mathbf{B}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Z} = [\hat{\boldsymbol{\beta}}_1 \hat{\boldsymbol{\beta}}_2 \cdots \hat{\boldsymbol{\beta}}_m].$$

The *predicted values* or *fitted values*

$$\hat{\mathbf{Z}} = \mathbf{X}\hat{\mathbf{B}} = [\hat{\mathbf{Y}}_1 \hat{\mathbf{Y}}_2 \dots \hat{\mathbf{Y}}_m] = \begin{bmatrix} \hat{Y}_{1,1} & \hat{Y}_{1,2} & \dots & \hat{Y}_{1,m} \\ \hat{Y}_{2,1} & \hat{Y}_{2,2} & \dots & \hat{Y}_{2,m} \\ \vdots & \vdots & \ddots & \vdots \\ \hat{Y}_{n,1} & \hat{Y}_{n,2} & \dots & \hat{Y}_{n,m} \end{bmatrix}.$$

The residuals  $\hat{\mathbf{E}} = \mathbf{Z} - \hat{\mathbf{Z}} = \mathbf{Z} - \mathbf{X}\hat{\mathbf{B}} =$

$$\begin{bmatrix} \hat{\epsilon}_1^T \\ \hat{\epsilon}_2^T \\ \vdots \\ \hat{\epsilon}_n^T \end{bmatrix} = [\mathbf{r}_1 \mathbf{r}_2 \dots \mathbf{r}_m] = \begin{bmatrix} \hat{\epsilon}_{1,1} & \hat{\epsilon}_{1,2} & \dots & \hat{\epsilon}_{1,m} \\ \hat{\epsilon}_{2,1} & \hat{\epsilon}_{2,2} & \dots & \hat{\epsilon}_{2,m} \\ \vdots & \vdots & \ddots & \vdots \\ \hat{\epsilon}_{n,1} & \hat{\epsilon}_{n,2} & \dots & \hat{\epsilon}_{n,m} \end{bmatrix}.$$

These quantities can be found from the  $m$  multiple linear regressions of  $\mathbf{Y}_j$  on the predictors:  $\hat{\boldsymbol{\beta}}_j = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}_j$ ,  $\hat{\mathbf{Y}}_j = \mathbf{X} \hat{\boldsymbol{\beta}}_j$ , and  $\mathbf{r}_j = \mathbf{Y}_j - \hat{\mathbf{Y}}_j$  for  $j = 1, \dots, m$ . Hence  $\hat{\epsilon}_{i,j} = Y_{i,j} - \hat{Y}_{i,j}$  where  $\hat{\mathbf{Y}}_j = (\hat{Y}_{1,j}, \dots, \hat{Y}_{n,j})^T$ . Finally,  $\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon},d} =$

$$\frac{(\mathbf{Z} - \hat{\mathbf{Z}})^T (\mathbf{Z} - \hat{\mathbf{Z}})}{n-d} = \frac{(\mathbf{Z} - \mathbf{X}\hat{\mathbf{B}})^T (\mathbf{Z} - \mathbf{X}\hat{\mathbf{B}})}{n-d} = \frac{\hat{\mathbf{E}}^T \hat{\mathbf{E}}}{n-d} = \frac{1}{n-d} \sum_{i=1}^n \hat{\epsilon}_i \hat{\epsilon}_i^T.$$

The choices  $d = 0$  and  $d = p$  are common. If  $d = 1$ , then  $\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon},d=1} = \mathbf{S}_r$ , the sample covariance matrix of the residual vectors  $\hat{\epsilon}_i$ , since the sample mean of the  $\hat{\epsilon}_i$  is  $\mathbf{0}$ . Let  $\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon}} = \hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon},p}$  be the unbiased estimator of  $\boldsymbol{\Sigma}_{\boldsymbol{\epsilon}}$ . Also,

$$\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon},d} = (n-d)^{-1} \mathbf{Z}^T [\mathbf{I} - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}] \mathbf{Z},$$

and

$$\hat{\mathbf{E}} = [\mathbf{I} - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}] \mathbf{Z}.$$

The following two theorems show that the least squares estimators are fairly good. Also see Theorem 10.7 in Section 10.4. Theorem 10.2 can also be used for  $\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon},d} = \frac{n-1}{n-d} \mathbf{S}_r$ .

**Theorem 10.1, Johnson and Wichern (1988, p. 304):** Suppose  $\mathbf{X}$  has full rank  $p < n$  and the covariance structure of Definition 10.2 holds. Then  $E(\hat{\mathbf{B}}) = \mathbf{B}$  so  $E(\hat{\boldsymbol{\beta}}_j) = \boldsymbol{\beta}_j$ ,  $\text{Cov}(\hat{\boldsymbol{\beta}}_j, \hat{\boldsymbol{\beta}}_k) = \sigma_{jk}(\mathbf{X}^T \mathbf{X})^{-1}$  for  $j, k = 1, \dots, p$ . Also  $\hat{\mathbf{E}}$  and  $\hat{\mathbf{B}}$  are uncorrelated,  $E(\hat{\mathbf{E}}) = \mathbf{0}$ , and

$$E(\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon}}) = E\left(\frac{\hat{\mathbf{E}}^T \hat{\mathbf{E}}}{n-p}\right) = \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}}.$$

**Theorem 10.2.**  $S_r = \Sigma_{\epsilon} + O_P(n^{-1/2})$  and  $\frac{1}{n} \sum_{i=1}^n \epsilon_i \epsilon_i^T = \Sigma_{\epsilon} + O_P(n^{-1/2})$  if the following three conditions hold:  $B - \hat{B} = O_P(n^{-1/2})$ ,  $\frac{1}{n} \sum_{i=1}^n \epsilon_i \mathbf{x}_i^T = O_P(1)$ , and  $\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T = O_P(n^{1/2})$ .

**Proof.** Note that  $\mathbf{y}_i = B^T \mathbf{x}_i + \epsilon_i = \hat{B}^T \mathbf{x}_i + \hat{\epsilon}_i$ . Hence  $\hat{\epsilon}_i = (B - \hat{B})^T \mathbf{x}_i + \epsilon_i$ . Thus

$$\begin{aligned} \sum_{i=1}^n \hat{\epsilon}_i \hat{\epsilon}_i^T &= \sum_{i=1}^n (\epsilon_i - \epsilon_i + \hat{\epsilon}_i)(\epsilon_i - \epsilon_i + \hat{\epsilon}_i)^T = \sum_{i=1}^n [\epsilon_i \epsilon_i^T + \epsilon_i (\hat{\epsilon}_i - \epsilon_i)^T + (\hat{\epsilon}_i - \epsilon_i) \hat{\epsilon}_i^T] \\ &= \sum_{i=1}^n \epsilon_i \epsilon_i^T + \left( \sum_{i=1}^n \epsilon_i \mathbf{x}_i^T \right) (B - \hat{B}) + (B - \hat{B})^T \left( \sum_{i=1}^n \mathbf{x}_i \epsilon_i^T \right) + \\ &\quad (B - \hat{B})^T \left( \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \right) (B - \hat{B}). \end{aligned}$$

Thus  $\frac{1}{n} \sum_{i=1}^n \hat{\epsilon}_i \hat{\epsilon}_i^T = \frac{1}{n} \sum_{i=1}^n \epsilon_i \epsilon_i^T +$

$$O_P(1)O_P(n^{-1/2}) + O_P(n^{-1/2})O_P(1) + O_P(n^{-1/2})O_P(n^{1/2})O_P(n^{-1/2}),$$

and the result follows since  $\frac{1}{n} \sum_{i=1}^n \epsilon_i \epsilon_i^T = \Sigma_{\epsilon} + O_P(n^{-1/2})$  and

$$S_r = \frac{n}{n-1} \frac{1}{n} \sum_{i=1}^n \hat{\epsilon}_i \hat{\epsilon}_i^T. \quad \square$$

$S_r$  and  $\hat{\Sigma}_{\epsilon}$  are also  $\sqrt{n}$  consistent estimators of  $\Sigma_{\epsilon}$  by Su and Cook (2012, p. 692). See Theorem 10.7.

## 10.2 Plots for the Multivariate Linear Regression Model

This section suggests using residual plots, response plots, and the DD plot to examine the multivariate linear model. The DD plot is used to examine the distribution of the iid error vectors. The residual plots are often used to check for lack of fit of the multivariate linear model. The response plots are used to check linearity and to detect influential cases for the linearity assumption. The response and residual plots are used exactly as in the  $m = 1$  case corresponding to multiple linear regression and experimental design models. See Olive (2010, 2017a), Olive et al. (2015), Olive and Hawkins (2005), and Cook and Weisberg (1999, p. 432).

**Notation.** Plots will be used to simplify the regression analysis, and in this text a plot of  $W$  versus  $Z$  uses  $W$  on the horizontal axis and  $Z$  on the vertical axis.

**Definition 10.4.** A **response plot** for the  $j$ th response variable is a plot of the fitted values  $\hat{Y}_{ij}$  versus the response  $Y_{ij}$ . The identity line with slope one and zero intercept is added to the plot as a visual aid. A **residual plot** corresponding to the  $j$ th response variable is a plot of  $\hat{Y}_{ij}$  versus  $r_{ij}$ .

**Remark 10.1.** Make the  $m$  response and residual plots for any multivariate linear regression. In a response plot, the vertical deviations from the identity line are the residuals  $r_{ij} = Y_{ij} - \hat{Y}_{ij}$ . Suppose the model is good, the  $j$ th error distribution is unimodal and not highly skewed for  $j = 1, \dots, m$ , and  $n \geq 10p$ . Then the plotted points should cluster about the identity line in each of the  $m$  response plots. If outliers are present or if the plot is not linear, then the current model or data need to be transformed or corrected. If the model is good, then each of the  $m$  residual plots should be ellipsoidal with no trend and should be centered about the  $r = 0$  line. There should not be any pattern in the residual plot: as a narrow vertical strip is moved from left to right, the behavior of the residuals within the strip should show little change. Outliers and patterns such as curvature or a fan shaped plot are bad.

**Rule of thumb 10.1.** Use multivariate linear regression if

$$n \geq \max((m + p)^2, mp + 30, 10p)$$

provided that the  $m$  response and residual plots all look good. Make the DD plot of the  $\hat{\epsilon}_i$ . If a residual plot would look good after several points have been deleted, and if these deleted points were not gross outliers (points far from the point cloud formed by the bulk of the data), then the residual plot is probably good. Beginners often find too many things wrong with a good model. For practice, use the computer to generate several multivariate linear regression data sets, and make the  $m$  response and residual plots for these data sets. This exercise will help show that the plots can have considerable variability even when the multivariate linear regression model is good. The *linmodpack* function `MLRsim` simulates response and residual plots for various distributions when  $m = 1$ .

**Rule of thumb 10.2.** If the plotted points in the residual plot look like a left or right opening megaphone, the first model violation to check is the assumption of nonconstant variance. (This is a rule of thumb because it is possible that such a residual plot results from another model violation such as nonlinearity, but nonconstant variance is much more common.)

**Remark 10.2.** Residual plots *magnify departures* from the model while the response plots emphasize *how well the multivariate linear regression model fits the data*.

**Definition 10.5.** An **RR plot** is a scatterplot matrix of the  $m$  sets of residuals  $\mathbf{r}_1, \dots, \mathbf{r}_m$ .

**Definition 10.6.** An **FF plot** is a scatterplot matrix of the  $m$  sets of fitted values of response variables  $\hat{Y}_1, \dots, \hat{Y}_m$ . The  $m$  response variables  $Y_1, \dots, Y_m$  can be added to the plot.

**Remark 10.3.** Some applications for multivariate linear regression need the  $m$  error vectors to be linearly related, and larger sample sizes may be needed if the error vectors are not linearly related. For example, the asymptotic optimality of the prediction regions of Section 10.3 needs the error vectors to be iid from an elliptically contoured distribution. Make the RR plot and a DD plot of the residual vectors  $\hat{\epsilon}_i$  to check that the error vectors are linearly related. Make a DD plot of the continuous predictor variables to check for  $\mathbf{x}$ -outliers. Make a DD plot of  $Y_1, \dots, Y_m$  to check for outliers, especially if it is assumed that the response variables come from an elliptically contoured distribution.

The RMVN DD plot of the residual vectors  $\hat{\epsilon}_i$  is used to check the error vector distribution, to detect outliers, and to display the nonparametric prediction region developed in Section 10.3. The DD plot suggests that the error vector distribution is elliptically contoured if the plotted points cluster tightly about a line through the origin as  $n \rightarrow \infty$ . The plot suggests that the error vector distribution is multivariate normal if the line is the identity line. If  $n$  is large and the plotted points do not cluster tightly about a line through the origin, then the error vector distribution may not be elliptically contoured. These applications of the DD plot for iid multivariate data are discussed in Olive (2002, 2008, 2013a, 2017b) and Chapter 7. The RMVN estimator has not yet been proven to be a consistent estimator when computed from residual vectors, but simulations suggest that the RMVN DD plot of the residual vectors is a useful diagnostic plot. The *linmodpack* function `mregdds` can be used to simulate the DD plots for various distributions.

Predictor transformations for the continuous predictors can be made exactly as in Section 1.2.

**Warning:** The log rule and other transformations do not always work. For example, the log rule may fail. If the relationships in the scatterplot matrix are already linear or if taking the transformation does not increase the linearity, then no transformation may be better than taking a transformation. For the Cook and Weisberg (1999) data set `evaporat.lsp` with  $m = 1$ , the log rule suggests transforming the response variable *Evap*, but no transformation works better.

Response transformations can also be made as in Section 1.2, but also make the response plot of  $\hat{Y}_j$  versus  $Y_j$ , and use the rules of Section 1.2 on  $Y_j$  to linearize the response plot for each of the  $m$  response variables  $Y_1, \dots, Y_m$ .

### 10.3 Asymptotically Optimal Prediction Regions

In this section, we will consider a more general multivariate regression model, and then consider the multivariate linear model as a special case. Given  $n$  cases of training or past data  $(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)$  and a vector of predictors  $\mathbf{x}_f$ , suppose it is desired to predict a future test vector  $\mathbf{y}_f$ .

**Definition 10.7.** A large sample  $100(1 - \delta)\%$  prediction region is a set  $\mathcal{A}_n$  such that  $P(\mathbf{y}_f \in \mathcal{A}_n) \rightarrow 1 - \delta$  as  $n \rightarrow \infty$ , and is *asymptotically optimal* if the volume of the region converges in probability to the volume of the population minimum volume covering region.

The classical large sample  $100(1 - \delta)\%$  prediction region for a future value  $\mathbf{x}_f$  given iid data  $\mathbf{x}_1, \dots, \mathbf{x}_n$  is  $\{\mathbf{x} : D_{\mathbf{x}}^2(\bar{\mathbf{x}}, \mathbf{S}) \leq \chi_{p,1-\delta}^2\}$ , while for multivariate linear regression, the classical large sample  $100(1 - \delta)\%$  prediction region for a future value  $\mathbf{y}_f$  given  $\mathbf{x}_f$  and past data  $(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)$  is  $\{\mathbf{y} : D_{\mathbf{y}}^2(\hat{\mathbf{y}}_f, \hat{\Sigma}\boldsymbol{\epsilon}) \leq \chi_{m,1-\delta}^2\}$ . See Johnson and Wichern (1988, pp. 134, 151, 312). By Equation (1.36), these regions may work for multivariate normal  $\mathbf{x}_i$  or  $\boldsymbol{\epsilon}_i$ , but otherwise tend to have undercoverage. Section 4.4 and Olive (2013a) replaced  $\chi_{p,1-\delta}^2$  by the order statistic  $D_{(U_n)}^2$  where  $U_n$  decreases to  $\lceil n(1 - \delta) \rceil$ . This section will use a similar technique from Olive (2018) to develop possibly the first practical large sample prediction region for the multivariate linear model with unknown error distribution. The following technical theorem will be needed to prove Theorem 10.4.

**Theorem 10.3.** Let  $a > 0$  and assume that  $(\hat{\boldsymbol{\mu}}_n, \hat{\Sigma}_n)$  is a consistent estimator of  $(\boldsymbol{\mu}, a\boldsymbol{\Sigma})$ .

a)  $D_{\mathbf{x}}^2(\hat{\boldsymbol{\mu}}_n, \hat{\Sigma}_n) - \frac{1}{a}D_{\mathbf{x}}^2(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = o_P(1)$ .

b) Let  $0 < \delta \leq 0.5$ . If  $(\hat{\boldsymbol{\mu}}_n, \hat{\Sigma}_n) - (\boldsymbol{\mu}, a\boldsymbol{\Sigma}) = O_P(n^{-\delta})$  and  $a\hat{\Sigma}_n^{-1} - \boldsymbol{\Sigma}^{-1} = O_P(n^{-\delta})$ , then

$$D_{\mathbf{x}}^2(\hat{\boldsymbol{\mu}}_n, \hat{\Sigma}_n) - \frac{1}{a}D_{\mathbf{x}}^2(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = O_P(n^{-\delta}).$$

**Proof.** Let  $B_n$  denote the subset of the sample space on which  $\hat{\Sigma}_n$  has an inverse. Then  $P(B_n) \rightarrow 1$  as  $n \rightarrow \infty$ . Now

$$\begin{aligned} D_{\mathbf{x}}^2(\hat{\boldsymbol{\mu}}_n, \hat{\Sigma}_n) &= (\mathbf{x} - \hat{\boldsymbol{\mu}}_n)^T \hat{\Sigma}_n^{-1} (\mathbf{x} - \hat{\boldsymbol{\mu}}_n) = \\ &= (\mathbf{x} - \hat{\boldsymbol{\mu}}_n)^T \left( \frac{\boldsymbol{\Sigma}^{-1}}{a} - \frac{\boldsymbol{\Sigma}^{-1}}{a} + \hat{\Sigma}_n^{-1} \right) (\mathbf{x} - \hat{\boldsymbol{\mu}}_n) = \\ &= (\mathbf{x} - \hat{\boldsymbol{\mu}}_n)^T \left( \frac{-\boldsymbol{\Sigma}^{-1}}{a} + \hat{\Sigma}_n^{-1} \right) (\mathbf{x} - \hat{\boldsymbol{\mu}}_n) + (\mathbf{x} - \hat{\boldsymbol{\mu}}_n)^T \left( \frac{\boldsymbol{\Sigma}^{-1}}{a} \right) (\mathbf{x} - \hat{\boldsymbol{\mu}}_n) = \end{aligned}$$

$$\begin{aligned}
& \frac{1}{a}(\mathbf{x} - \hat{\boldsymbol{\mu}}_n)^T(-\boldsymbol{\Sigma}^{-1} + a \hat{\boldsymbol{\Sigma}}_n^{-1})(\mathbf{x} - \hat{\boldsymbol{\mu}}_n) + \\
& (\mathbf{x} - \boldsymbol{\mu} + \boldsymbol{\mu} - \hat{\boldsymbol{\mu}}_n)^T \left( \frac{\boldsymbol{\Sigma}^{-1}}{a} \right) (\mathbf{x} - \boldsymbol{\mu} + \boldsymbol{\mu} - \hat{\boldsymbol{\mu}}_n) \\
& = \frac{1}{a}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) + \frac{2}{a}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}_n) + \\
& \frac{1}{a}(\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}_n)^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}_n) + \frac{1}{a}(\mathbf{x} - \hat{\boldsymbol{\mu}}_n)^T [a \hat{\boldsymbol{\Sigma}}_n^{-1} - \boldsymbol{\Sigma}^{-1}](\mathbf{x} - \hat{\boldsymbol{\mu}}_n)
\end{aligned}$$

on  $B_n$ , and the last three terms are  $o_P(1)$  under a) and  $O_P(n^{-\delta})$  under b).  
□

Now suppose a prediction region for an  $m \times 1$  random vector  $\mathbf{y}_f$  given a vector of predictors  $\mathbf{x}_f$  is desired for the multivariate linear model. If we had many cases  $\mathbf{z}_i = \mathbf{B}^T \mathbf{x}_f + \boldsymbol{\epsilon}_i$ , then we could use the multivariate prediction region for  $m$  variables from Section 2.2. Instead, Theorem 10.4 will use the prediction region from Section 4.4 on the pseudodata  $\hat{\mathbf{z}}_i = \hat{\mathbf{B}}^T \mathbf{x}_f + \hat{\boldsymbol{\epsilon}}_i = \hat{\mathbf{y}}_f + \hat{\boldsymbol{\epsilon}}_i$  for  $i = 1, \dots, n$ . This takes the data cloud of the  $n$  residual vectors  $\hat{\boldsymbol{\epsilon}}_i$  and centers the cloud at  $\hat{\mathbf{y}}_f$ . Note that  $\hat{\mathbf{z}}_i = (\mathbf{B} - \mathbf{B} + \hat{\mathbf{B}})^T \mathbf{x}_f + (\boldsymbol{\epsilon}_i - \boldsymbol{\epsilon}_i + \hat{\boldsymbol{\epsilon}}_i) = \mathbf{z}_i + (\hat{\mathbf{B}} - \mathbf{B})^T \mathbf{x}_f + \hat{\boldsymbol{\epsilon}}_i - \boldsymbol{\epsilon}_i = \mathbf{z}_i + (\hat{\mathbf{B}} - \mathbf{B})^T \mathbf{x}_f - (\hat{\mathbf{B}} - \mathbf{B})^T \mathbf{x}_i = \mathbf{z}_i + O_P(n^{-1/2})$ . Hence the distances based on the  $\mathbf{z}_i$  and the distances based on the  $\hat{\mathbf{z}}_i$  have the same quantiles, asymptotically (for quantiles that are continuity points of the distribution of  $\mathbf{z}_i$ ).

If the  $\boldsymbol{\epsilon}_i$  are iid from an  $EC_m(\mathbf{0}, \boldsymbol{\Sigma}, g)$  distribution with continuous decreasing  $g$  and nonsingular covariance matrix  $\boldsymbol{\Sigma}_{\boldsymbol{\epsilon}} = c\boldsymbol{\Sigma}$  for some constant  $c > 0$ , then the population asymptotically optimal prediction region is  $\{\mathbf{y} : D_{\mathbf{y}}(\mathbf{B}^T \mathbf{x}_f, \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}}) \leq D_{1-\delta}\}$  where  $P(D_{\mathbf{y}}(\mathbf{B}^T \mathbf{x}_f, \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}}) \leq D_{1-\delta}) = 1 - \delta$ . For example, if the iid  $\boldsymbol{\epsilon}_i \sim N_m(\mathbf{0}, \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}})$ , then  $D_{1-\delta} = \sqrt{\chi_{m,1-\delta}^2}$ . If the error distribution is not elliptically contoured, then the above region still has  $100(1 - \delta)\%$  coverage, but prediction regions with smaller volume may exist.

A natural way to make a large sample prediction region is to estimate the target population minimum volume covering region, but for moderate samples and many error distributions, the natural estimator that covers  $\lceil n(1 - \delta) \rceil$  of the cases tends to have undercoverage as high as  $\min(0.05, \delta/2)$ . This empirical result is not too surprising since it is well known that the performance of a prediction region on the training data is superior to the performance on future test data, due in part to the unknown variability of the estimator. To compensate for the undercoverage, let  $q_n$  be as in Theorem 10.4.

**Theorem 10.4.** Suppose  $\mathbf{y}_i = E(\mathbf{y}_i | \mathbf{x}_i) + \boldsymbol{\epsilon}_i = \hat{\mathbf{y}}_i + \hat{\boldsymbol{\epsilon}}_i$  where  $\text{Cov}(\boldsymbol{\epsilon}_i) = \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}} > 0$ , and where the zero mean  $\boldsymbol{\epsilon}_f$  and the  $\boldsymbol{\epsilon}_i$  are iid for  $i = 1, \dots, n$ . Given  $\mathbf{x}_f$ , suppose the fitted model produces  $\hat{\mathbf{y}}_f$  and nonsingular  $\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon}}$ . Let  $\hat{\mathbf{z}}_i = \hat{\mathbf{y}}_f + \hat{\boldsymbol{\epsilon}}_i$  and

$$D_i^2 \equiv D_i^2(\hat{\mathbf{y}}_f, \hat{\boldsymbol{\Sigma}}_\epsilon) = (\hat{\mathbf{z}}_i - \hat{\mathbf{y}}_f)^T \hat{\boldsymbol{\Sigma}}_\epsilon^{-1} (\hat{\mathbf{z}}_i - \hat{\mathbf{y}}_f)$$

for  $i = 1, \dots, n$ . Let  $q_n = \min(1 - \delta + 0.05, 1 - \delta + m/n)$  for  $\delta > 0.1$  and

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

If  $q_n < 1 - \delta + 0.001$ , set  $q_n = 1 - \delta$ . Let  $0 < \delta < 1$  and  $h = D_{(U_n)}$  where  $D_{(U_n)}$  is the 100  $q_n$ th sample quantile of the Mahalanobis distances  $D_i$ . Let the nominal  $100(1 - \delta)\%$  prediction region for  $\mathbf{y}_f$  be given by

$$\begin{aligned} \{\mathbf{z} : (\mathbf{z} - \hat{\mathbf{y}}_f)^T \hat{\boldsymbol{\Sigma}}_\epsilon^{-1} (\mathbf{z} - \hat{\mathbf{y}}_f) \leq D_{(U_n)}^2\} = \\ \{\mathbf{z} : D_{\mathbf{z}}^2(\hat{\mathbf{y}}_f, \hat{\boldsymbol{\Sigma}}_\epsilon) \leq D_{(U_n)}^2\} = \{\mathbf{z} : D_{\mathbf{z}}(\hat{\mathbf{y}}_f, \hat{\boldsymbol{\Sigma}}_\epsilon) \leq D_{(U_n)}\}. \end{aligned} \quad (10.1)$$

a) Consider the  $n$  prediction regions for the data where  $(\mathbf{y}_{f,i}, \mathbf{x}_{f,i}) = (\mathbf{y}_i, \mathbf{x}_i)$  for  $i = 1, \dots, n$ . If the order statistic  $D_{(U_n)}$  is unique, then  $U_n$  of the  $n$  prediction regions contain  $\mathbf{y}_i$  where  $U_n/n \rightarrow 1 - \delta$  as  $n \rightarrow \infty$ .

b) If  $(\hat{\mathbf{y}}_f, \hat{\boldsymbol{\Sigma}}_\epsilon)$  is a consistent estimator of  $(E(\mathbf{y}_f), \boldsymbol{\Sigma}_\epsilon)$ , then (10.1) is a large sample  $100(1 - \delta)\%$  prediction region for  $\mathbf{y}_f$ .

c) If  $(\hat{\mathbf{y}}_f, \hat{\boldsymbol{\Sigma}}_\epsilon)$  is a consistent estimator of  $(E(\mathbf{y}_f), \boldsymbol{\Sigma}_\epsilon)$ , and the  $\epsilon_i$  come from an elliptically contoured distribution such that the unique highest density region is  $\{\mathbf{z} : D_{\mathbf{z}}(\mathbf{0}, \boldsymbol{\Sigma}_\epsilon) \leq D_{1-\delta}\}$ , then the prediction region (10.1) is asymptotically optimal.

**Proof.** a) Suppose  $(\mathbf{x}_f, \mathbf{y}_f) = (\mathbf{x}_i, \mathbf{y}_i)$ . Then

$$D_{\mathbf{y}_i}^2(\hat{\mathbf{y}}_f, \hat{\boldsymbol{\Sigma}}_\epsilon) = (\mathbf{y}_i - \hat{\mathbf{y}}_f)^T \hat{\boldsymbol{\Sigma}}_\epsilon^{-1} (\mathbf{y}_i - \hat{\mathbf{y}}_f) = \hat{\epsilon}_i^T \hat{\boldsymbol{\Sigma}}_\epsilon^{-1} \hat{\epsilon}_i = D_{\hat{\epsilon}_i}^2(\mathbf{0}, \hat{\boldsymbol{\Sigma}}_\epsilon).$$

Hence  $\mathbf{y}_i$  is in the  $i$ th prediction region  $\{\mathbf{z} : D_{\mathbf{z}}(\hat{\mathbf{y}}_f, \hat{\boldsymbol{\Sigma}}_\epsilon) \leq D_{(U_n)}\}$  iff  $\hat{\epsilon}_i$  is in prediction region  $\{\mathbf{z} : D_{\mathbf{z}}(\mathbf{0}, \hat{\boldsymbol{\Sigma}}_\epsilon) \leq D_{(U_n)}\}$ , but exactly  $U_n$  of the  $\hat{\epsilon}_i$  are in the latter region by construction, if  $D_{(U_n)}$  is unique. Since  $D_{(U_n)}$  is the  $100(1 - \delta)$ th percentile of the  $D_i$  asymptotically,  $U_n/n \rightarrow 1 - \delta$ .

b) Let  $P[D_{\mathbf{z}}(E(\mathbf{y}_f), \boldsymbol{\Sigma}_\epsilon) \leq D_{1-\delta}(E(\mathbf{y}_f), \boldsymbol{\Sigma}_\epsilon)] = 1 - \delta$ . Since  $\boldsymbol{\Sigma}_\epsilon > 0$ , Theorem 10.3 shows that if  $(\hat{\mathbf{y}}_f, \hat{\boldsymbol{\Sigma}}_\epsilon) \xrightarrow{P} (E(\mathbf{y}_f), \boldsymbol{\Sigma}_\epsilon)$  then  $D(\hat{\mathbf{y}}_f, \hat{\boldsymbol{\Sigma}}_\epsilon) \xrightarrow{D} D_{\mathbf{z}}(E(\mathbf{y}_f), \boldsymbol{\Sigma}_\epsilon)$ . Hence the percentiles of the distances converge in distribution, and the probability that  $\mathbf{y}_f$  is in  $\{\mathbf{z} : D_{\mathbf{z}}(\hat{\mathbf{y}}_f, \hat{\boldsymbol{\Sigma}}_\epsilon) \leq D_{1-\delta}(\hat{\mathbf{y}}_f, \hat{\boldsymbol{\Sigma}}_\epsilon)\}$  converges to  $1 - \delta =$  the probability that  $\mathbf{y}_f$  is in  $\{\mathbf{z} : D_{\mathbf{z}}(E(\mathbf{y}_f), \boldsymbol{\Sigma}_\epsilon) \leq D_{1-\delta}(E(\mathbf{y}_f), \boldsymbol{\Sigma}_\epsilon)\}$  at continuity points  $D_{1-\delta}$  of the distribution of  $D(E(\mathbf{y}_f), \boldsymbol{\Sigma}_\epsilon)$ .

c) The asymptotically optimal prediction region is the region with the smallest volume (hence highest density) such that the coverage is  $1 - \delta$ , as  $n \rightarrow \infty$ . This region is  $\{\mathbf{z} : D_{\mathbf{z}}(E(\mathbf{y}_f), \boldsymbol{\Sigma}_\epsilon) \leq D_{1-\delta}(E(\mathbf{y}_f), \boldsymbol{\Sigma}_\epsilon)\}$  if the asymptotically optimal region for the  $\epsilon_i$  is  $\{\mathbf{z} : D_{\mathbf{z}}(\mathbf{0}, \boldsymbol{\Sigma}_\epsilon) \leq D_{1-\delta}(\mathbf{0}, \boldsymbol{\Sigma}_\epsilon)\}$ . Hence the result follows by b).  $\square$

Notice that if  $\hat{\Sigma}_{\epsilon}^{-1}$  exists, then  $100q_n\%$  of the  $n$  training data  $\mathbf{y}_i$  are in their corresponding prediction region with  $\mathbf{x}_f = \mathbf{x}_i$ , and  $q_n \rightarrow 1 - \delta$  even if  $(\hat{\mathbf{y}}_i, \hat{\Sigma}_{\epsilon})$  is not a good estimator or if the regression model is misspecified. 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  $(\hat{\mathbf{y}}_i, \hat{\Sigma}_{\epsilon})$  is used or if the  $\epsilon_i$  do not come from an elliptically contoured distribution. The response, residual, and DD plots can be used to check model assumptions. If the plotted points in the RMVN DD plot cluster tightly about some line through the origin and if  $n \geq \max[3(m+p)^2, mp+30]$ , we expect the volume of the prediction region may be fairly low for the least squares estimators.

If  $n$  is too small, then multivariate data is sparse and the covering ellipsoid for the training data may be far too small for future data, resulting in severe undercoverage. Also notice that  $q_n = 1 - \delta/2$  or  $q_n = 1 - \delta + 0.05$  for  $n \leq 20p$ . At the training data, the coverage  $q_n \geq 1 - \delta$ , and  $q_n$  converges to the nominal coverage  $1 - \delta$  as  $n \rightarrow \infty$ . Suppose  $n \leq 20p$ . Then the nominal 95% prediction region uses  $q_n = 0.975$  while the nominal 50% prediction region uses  $q_n = 0.55$ . Prediction distributions depend both on the error distribution and on the variability of the estimator  $(\hat{\mathbf{y}}_f, \hat{\Sigma}_{\epsilon})$ . This variability is typically unknown but converges to 0 as  $n \rightarrow \infty$ . Also, residuals tend to underestimate errors for small  $n$ . For moderate  $n$ , ignoring estimator variability and using  $q_n = 1 - \delta$  resulted in undercoverage as high as  $\min(0.05, \delta/2)$ . Letting the “coverage”  $q_n$  decrease to the nominal coverage  $1 - \delta$  inflates the volume of the prediction region for small  $n$ , compensating for the unknown variability of  $(\hat{\mathbf{y}}_f, \hat{\Sigma}_{\epsilon})$ .

Consider the multivariate linear regression model. Let  $\hat{\Sigma}_{\epsilon} = \hat{\Sigma}_{\epsilon, d=p}$ ,  $\hat{\mathbf{z}}_i = \hat{\mathbf{y}}_f + \hat{\epsilon}_i$ , and  $D_i^2(\hat{\mathbf{y}}_f, \mathbf{S}_r) = (\hat{\mathbf{z}}_i - \hat{\mathbf{y}}_f)^T \mathbf{S}_r^{-1} (\hat{\mathbf{z}}_i - \hat{\mathbf{y}}_f)$  for  $i = 1, \dots, n$ . Then the large sample nonparametric  $100(1 - \delta)\%$  prediction region is

$$\{\mathbf{z} : D_{\hat{\mathbf{z}}}^2(\hat{\mathbf{y}}_f, \mathbf{S}_r) \leq D_{(U_n)}^2\} = \{\mathbf{z} : D_{\mathbf{z}}(\hat{\mathbf{y}}_f, \mathbf{S}_r) \leq D_{(U_n)}\}. \quad (10.2)$$

Theorem 10.5 will show that this prediction region (10.2) can also be found by applying the nonparametric prediction region (2.24) on the  $\hat{\mathbf{z}}_i$ . Recall that  $\mathbf{S}_r$  defined in Definition 10.3 is the sample covariance matrix of the residual vectors  $\hat{\epsilon}_i$ . For the multivariate linear regression model, if  $D_{1-\delta}$  is a continuity point of the distribution of  $D$ , Assumption D1 above Theorem 10.7 holds, and the  $\epsilon_i$  have a nonsingular covariance matrix, then (10.2) is a large sample  $100(1 - \delta)\%$  prediction region for  $\mathbf{y}_f$ .

**Theorem 10.5.** For multivariate linear regression, when least squares is used to compute  $\hat{\mathbf{y}}_f$ ,  $\mathbf{S}_r$ , and the pseudodata  $\hat{\mathbf{z}}_i$ , prediction region (10.2) is the nonparametric prediction region (4.24) applied to the  $\hat{\mathbf{z}}_i$ .

**Proof.** Multivariate linear regression with least squares satisfies Theorem 10.4 by Su and Cook (2012). (See Theorem 10.7.) Let  $(T, \mathbf{C})$  be the sample mean and sample covariance matrix (see Definition 2.7) applied to the  $\hat{\mathbf{z}}_i$ . The sample mean and sample covariance matrix of the residual vectors is

$(\mathbf{0}, \mathbf{S}_r)$  since least squares was used. Hence the  $\hat{\mathbf{z}}_i = \hat{\mathbf{y}}_f + \hat{\boldsymbol{\epsilon}}_i$  have sample covariance matrix  $\mathbf{S}_r$ , and sample mean  $\hat{\mathbf{y}}_f$ . Hence  $(T, \mathbf{C}) = (\hat{\mathbf{y}}_f, \mathbf{S}_r)$ , and the  $D_i(\hat{\mathbf{y}}_f, \mathbf{S}_r)$  are used to compute  $D_{(U_n)}$ .  $\square$

The RMVN DD plot of the residual vectors will be used to display the prediction regions for multivariate linear regression. See Example 10.3. The nonparametric prediction region for multivariate linear regression of Theorem 10.5 uses  $(T, \mathbf{C}) = (\hat{\mathbf{y}}_f, \mathbf{S}_r)$  in (10.1), and has simple geometry. Let  $R_r$  be the nonparametric prediction region (10.2) applied to the residuals  $\hat{\boldsymbol{\epsilon}}_i$  with  $\hat{\mathbf{y}}_f = \mathbf{0}$ . Then  $R_r$  is a hyperellipsoid with center  $\mathbf{0}$ , and the nonparametric prediction region is the hyperellipsoid  $R_r$  translated to have center  $\hat{\mathbf{y}}_f$ . Hence in a DD plot, all points to the left of the line  $MD = D_{(U_n)}$  correspond to  $\mathbf{y}_i$  that are in their prediction region, while points to the right of the line are not in their prediction region.

The nonparametric prediction region has some interesting properties. This prediction region is asymptotically optimal if the  $\boldsymbol{\epsilon}_i$  are iid for a large class of elliptically contoured  $EC_m(\mathbf{0}, \boldsymbol{\Sigma}, g)$  distributions. Also, if there are 100 different values  $(\mathbf{x}_{jf}, \mathbf{y}_{jf})$  to be predicted, we only need to update  $\hat{\mathbf{y}}_{jf}$  for  $j = 1, \dots, 100$ , we do not need to update the covariance matrix  $\mathbf{S}_r$ .

It is common practice to examine how well the prediction regions work on the training data. That is, for  $i = 1, \dots, n$ , set  $\mathbf{x}_f = \mathbf{x}_i$  and see if  $\mathbf{y}_i$  is in the region with probability near to  $1 - \delta$  with a simulation study. Note that  $\hat{\mathbf{y}}_f = \hat{\mathbf{y}}_i$  if  $\mathbf{x}_f = \mathbf{x}_i$ . Simulation is not needed for the nonparametric prediction region (10.2) for the data since the prediction region (10.2) centered at  $\hat{\mathbf{y}}_i$  contains  $\mathbf{y}_i$  iff  $R_r$ , the prediction region centered at  $\mathbf{0}$ , contains  $\hat{\boldsymbol{\epsilon}}_i$  since  $\hat{\boldsymbol{\epsilon}}_i = \mathbf{y}_i - \hat{\mathbf{y}}_i$ . Thus  $100q_n\%$  of prediction regions corresponding to the data  $(\mathbf{y}_i, \mathbf{x}_i)$  contain  $\mathbf{y}_i$ , and  $100q_n\% \rightarrow 100(1 - \delta)\%$ . Hence the prediction regions work well on the training data and should work well on  $(\mathbf{x}_f, \mathbf{y}_f)$  similar to the training data. Of course simulation should be done for test data  $(\mathbf{x}_f, \mathbf{y}_f)$  that are not equal to training data cases. See Problem 10.11.

This training data result holds provided that the multivariate linear regression using least squares is such that the sample covariance matrix  $\mathbf{S}_r$  of the residual vectors is nonsingular, **the multivariate regression model need not be correct**. Hence the coverage at the  $n$  training data cases  $(\mathbf{x}_i, \mathbf{y}_i)$  is robust to model misspecification. Of course, the prediction regions may be very large if the model is severely misspecified, but severity of misspecification can be checked with the response and residual plots. Coverage for a future value  $\mathbf{y}_f$  can also be arbitrarily bad if there is extrapolation or if  $(\mathbf{x}_f, \mathbf{y}_f)$  comes from a different population than that of the data.

## 10.4 Testing Hypotheses

This section considers testing a linear hypothesis  $H_0 : \mathbf{LB} = \mathbf{0}$  versus  $H_1 : \mathbf{LB} \neq \mathbf{0}$  where  $\mathbf{L}$  is a full rank  $r \times p$  matrix.

**Definition 10.8.** Assume  $\text{rank}(\mathbf{X}) = p$ . The *total corrected (for the mean) sum of squares and cross products matrix* is

$$\mathbf{T} = \mathbf{R} + \mathbf{W}_e = \mathbf{Z}^T \left( \mathbf{I}_n - \frac{1}{n} \mathbf{1}\mathbf{1}^T \right) \mathbf{Z}.$$

Note that  $\mathbf{T}/(n-1)$  is the usual sample covariance matrix  $\hat{\Sigma}_{\mathbf{y}}$  if all  $n$  of the  $\mathbf{y}_i$  are iid, e.g. if  $\mathbf{B} = \mathbf{0}$ . The *regression sum of squares and cross products matrix* is

$$\mathbf{R} = \mathbf{Z}^T \left[ \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T - \frac{1}{n} \mathbf{1}\mathbf{1}^T \right] \mathbf{Z} = \mathbf{Z}^T \mathbf{X} \hat{\mathbf{B}} - \frac{1}{n} \mathbf{Z}^T \mathbf{1}\mathbf{1}^T \mathbf{Z}.$$

Let  $\mathbf{H} = \hat{\mathbf{B}}^T \mathbf{L}^T [\mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T]^{-1} \mathbf{L} \hat{\mathbf{B}}$ . The *error or residual sum of squares and cross products matrix* is

$$\mathbf{W}_e = (\mathbf{Z} - \hat{\mathbf{Z}})^T (\mathbf{Z} - \hat{\mathbf{Z}}) = \mathbf{Z}^T \mathbf{Z} - \mathbf{Z}^T \mathbf{X} \hat{\mathbf{B}} = \mathbf{Z}^T [\mathbf{I}_n - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T] \mathbf{Z}.$$

Note that  $\mathbf{W}_e = \hat{\mathbf{E}}^T \hat{\mathbf{E}}$  and  $\mathbf{W}_e/(n-p) = \hat{\Sigma}_{\epsilon}$ .

**Warning:** SAS output uses  $\mathbf{E}$  instead of  $\mathbf{W}_e$ .

The MANOVA table is shown below.

Summary MANOVA Table

Source	matrix	df
Regression or Treatment	$\mathbf{R}$	$p-1$
Error or Residual	$\mathbf{W}_e$	$n-p$
Total (corrected)	$\mathbf{T}$	$n-1$

**Definition 10.9.** Let  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m$  be the ordered eigenvalues of  $\mathbf{W}_e^{-1} \mathbf{H}$ . Then there are four commonly used test statistics.

The *Roy's maximum root statistic* is  $\lambda_{\max}(\mathbf{L}) = \lambda_1$ .

The *Wilks'  $\Lambda$  statistic* is  $\Lambda(\mathbf{L}) = |(\mathbf{H} + \mathbf{W}_e)^{-1} \mathbf{W}_e| = |\mathbf{W}_e^{-1} \mathbf{H} + \mathbf{I}|^{-1} = \prod_{i=1}^m (1 + \lambda_i)^{-1}$ .

The *Pillai's trace statistic* is  $V(\mathbf{L}) = \text{tr}[(\mathbf{H} + \mathbf{W}_e)^{-1} \mathbf{H}] = \sum_{i=1}^m \frac{\lambda_i}{1 + \lambda_i}$ .

The *Hotelling-Lawley trace statistic* is  $U(\mathbf{L}) = \text{tr}[\mathbf{W}_e^{-1}\mathbf{H}] = \sum_{i=1}^m \lambda_i$ .

Typically some function of one of the four above statistics is used to get pval, the estimated pvalue. Output often gives the pvals for all four test statistics. Be cautious about inference if the last three test statistics do not lead to the same conclusions (Roy's test may not be trustworthy for  $r > 1$ ). Theory and simulations developed below for the four statistics will provide more information about the sample sizes needed to use the four test statistics. See the paragraphs after the following theorem for the notation used in that theorem.

**Theorem 10.6.** *The Hotelling-Lawley trace statistic*

$$U(\mathbf{L}) = \frac{1}{n-p} [\text{vec}(\mathbf{L}\hat{\mathbf{B}})]^T [\hat{\Sigma}_\epsilon^{-1} \otimes (\mathbf{L}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{L}^T)^{-1}] [\text{vec}(\mathbf{L}\hat{\mathbf{B}})]. \quad (10.3)$$

**Proof.** Using the Searle (1982, p. 333) identity  $\text{tr}(\mathbf{A}\mathbf{G}^T\mathbf{D}\mathbf{G}\mathbf{C}) = [\text{vec}(\mathbf{G})]^T [\mathbf{C}\mathbf{A} \otimes \mathbf{D}^T] [\text{vec}(\mathbf{G})]$ , it follows that  $(n-p)U(\mathbf{L}) = \text{tr}[\hat{\Sigma}_\epsilon^{-1}\hat{\mathbf{B}}^T\mathbf{L}^T[\mathbf{L}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{L}^T]^{-1}\mathbf{L}\hat{\mathbf{B}}] = [\text{vec}(\mathbf{L}\hat{\mathbf{B}})]^T [\hat{\Sigma}_\epsilon^{-1} \otimes (\mathbf{L}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{L}^T)^{-1}] [\text{vec}(\mathbf{L}\hat{\mathbf{B}})] = T$  where  $\mathbf{A} = \hat{\Sigma}_\epsilon^{-1}$ ,  $\mathbf{G} = \mathbf{L}\hat{\mathbf{B}}$ ,  $\mathbf{D} = [\mathbf{L}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{L}^T]^{-1}$ , and  $\mathbf{C} = \mathbf{I}$ . Hence (10.3) holds.  $\square$

Some notation is useful to show (10.3) and to show that  $(n-p)U(\mathbf{L}) \xrightarrow{D} \chi_{rm}^2$  under mild conditions if  $H_0$  is true. Following Henderson and Searle (1979), let matrix  $\mathbf{A} = [\mathbf{a}_1 \ \mathbf{a}_2 \ \dots \ \mathbf{a}_p]$ . Then the vec operator stacks the columns of  $\mathbf{A}$  on top of one another so

$$\text{vec}(\mathbf{A}) = \begin{pmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \\ \vdots \\ \mathbf{a}_p \end{pmatrix}.$$

Let  $\mathbf{A} = (a_{ij})$  be an  $m \times n$  matrix and  $\mathbf{B}$  a  $p \times q$  matrix. Then the Kronecker product of  $\mathbf{A}$  and  $\mathbf{B}$  is the  $mp \times nq$  matrix

$$\mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} a_{11}\mathbf{B} & a_{12}\mathbf{B} & \cdots & a_{1n}\mathbf{B} \\ a_{21}\mathbf{B} & a_{22}\mathbf{B} & \cdots & a_{2n}\mathbf{B} \\ \vdots & \vdots & \cdots & \vdots \\ a_{m1}\mathbf{B} & a_{m2}\mathbf{B} & \cdots & a_{mn}\mathbf{B} \end{bmatrix}.$$

An important fact is that if  $\mathbf{A}$  and  $\mathbf{B}$  are nonsingular square matrices, then  $[\mathbf{A} \otimes \mathbf{B}]^{-1} = \mathbf{A}^{-1} \otimes \mathbf{B}^{-1}$ . The following assumption is important.

**Assumption D1:** Let  $h_i$  be the  $i$ th diagonal element of  $\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$ . Assume  $\max_{1 \leq i \leq n} h_i \xrightarrow{P} 0$  as  $n \rightarrow \infty$ , assume that the zero mean iid error vectors have finite fourth moments, and assume that  $\frac{1}{n}\mathbf{X}^T\mathbf{X} \xrightarrow{P} \mathbf{W}^{-1}$ .

Su and Cook (2012) proved a central limit type theorem for  $\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon}}$  and  $\hat{\mathbf{B}}$  for the partial envelopes estimator, and the least squares estimator is a special case. These results prove the following theorem. Their theorem also shows that for multiple linear regression ( $m = 1$ ),  $\hat{\sigma}^2 = MSE$  is a  $\sqrt{n}$  consistent estimator of  $\sigma^2$ .

**Theorem 10.7: Multivariate Least Squares Central Limit Theorem (MLS CLT).** For the least squares estimator, if assumption D1 holds, then  $\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon}}$  is a  $\sqrt{n}$  consistent estimator of  $\boldsymbol{\Sigma}_{\boldsymbol{\epsilon}}$  and

$$\sqrt{n} \operatorname{vec}(\hat{\mathbf{B}} - \mathbf{B}) \xrightarrow{D} N_{pm}(\mathbf{0}, \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}} \otimes \mathbf{W}).$$

**Theorem 10.8.** If assumption D1 holds and if  $H_0$  is true, then  $(n-p)U(\mathbf{L}) \xrightarrow{D} \chi_{rm}^2$ .

**Proof.** By Theorem 10.7,  $\sqrt{n} \operatorname{vec}(\hat{\mathbf{B}} - \mathbf{B}) \xrightarrow{D} N_{pm}(\mathbf{0}, \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}} \otimes \mathbf{W})$ . Then under  $H_0$ ,  $\sqrt{n} \operatorname{vec}(\mathbf{L}\hat{\mathbf{B}}) \xrightarrow{D} N_{rm}(\mathbf{0}, \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}} \otimes \mathbf{L}\mathbf{W}\mathbf{L}^T)$ , and  $n [\operatorname{vec}(\mathbf{L}\hat{\mathbf{B}})]^T [\boldsymbol{\Sigma}_{\boldsymbol{\epsilon}}^{-1} \otimes (\mathbf{L}\mathbf{W}\mathbf{L}^T)^{-1}] [\operatorname{vec}(\mathbf{L}\hat{\mathbf{B}})] \xrightarrow{D} \chi_{rm}^2$ . This result also holds if  $\mathbf{W}$  and  $\boldsymbol{\Sigma}_{\boldsymbol{\epsilon}}$  are replaced by  $\hat{\mathbf{W}} = n(\mathbf{X}^T\mathbf{X})^{-1}$  and  $\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon}}$ . Hence under  $H_0$  and using the proof of Theorem 10.6,

$$T = (n-p)U(\mathbf{L}) = [\operatorname{vec}(\mathbf{L}\hat{\mathbf{B}})]^T [\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon}}^{-1} \otimes (\mathbf{L}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{L}^T)^{-1}] [\operatorname{vec}(\mathbf{L}\hat{\mathbf{B}})] \xrightarrow{D} \chi_{rm}^2.$$

□

Some more details on the above results may be useful. Consider testing a linear hypothesis  $H_0 : \mathbf{L}\mathbf{B} = \mathbf{0}$  versus  $H_1 : \mathbf{L}\mathbf{B} \neq \mathbf{0}$  where  $\mathbf{L}$  is a full rank  $r \times p$  matrix. For now assume the error distribution is multivariate normal  $N_m(\mathbf{0}, \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}})$ . Then

$$\operatorname{vec}(\hat{\mathbf{B}} - \mathbf{B}) = \begin{pmatrix} \hat{\boldsymbol{\beta}}_1 - \boldsymbol{\beta}_1 \\ \hat{\boldsymbol{\beta}}_2 - \boldsymbol{\beta}_2 \\ \vdots \\ \hat{\boldsymbol{\beta}}_m - \boldsymbol{\beta}_m \end{pmatrix} \sim N_{pm}(\mathbf{0}, \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}} \otimes (\mathbf{X}^T\mathbf{X})^{-1})$$

where

$$\mathbf{C} = \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}} \otimes (\mathbf{X}^T \mathbf{X})^{-1} = \begin{bmatrix} \sigma_{11}(\mathbf{X}^T \mathbf{X})^{-1} & \sigma_{12}(\mathbf{X}^T \mathbf{X})^{-1} & \cdots & \sigma_{1m}(\mathbf{X}^T \mathbf{X})^{-1} \\ \sigma_{21}(\mathbf{X}^T \mathbf{X})^{-1} & \sigma_{22}(\mathbf{X}^T \mathbf{X})^{-1} & \cdots & \sigma_{2m}(\mathbf{X}^T \mathbf{X})^{-1} \\ \vdots & \vdots & \cdots & \vdots \\ \sigma_{m1}(\mathbf{X}^T \mathbf{X})^{-1} & \sigma_{m2}(\mathbf{X}^T \mathbf{X})^{-1} & \cdots & \sigma_{mm}(\mathbf{X}^T \mathbf{X})^{-1} \end{bmatrix}.$$

Now let  $\mathbf{A}$  be an  $rm \times pm$  block diagonal matrix:  $\mathbf{A} = \text{diag}(\mathbf{L}, \dots, \mathbf{L})$ . Then  $\mathbf{A} \text{vec}(\hat{\mathbf{B}} - \mathbf{B}) = \text{vec}(\mathbf{L}(\hat{\mathbf{B}} - \mathbf{B})) =$

$$\begin{pmatrix} \mathbf{L}(\hat{\boldsymbol{\beta}}_1 - \boldsymbol{\beta}_1) \\ \mathbf{L}(\hat{\boldsymbol{\beta}}_2 - \boldsymbol{\beta}_2) \\ \vdots \\ \mathbf{L}(\hat{\boldsymbol{\beta}}_m - \boldsymbol{\beta}_m) \end{pmatrix} \sim N_{rm}(\mathbf{0}, \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}} \otimes \mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T)$$

where  $\mathbf{D} = \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}} \otimes \mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T = \mathbf{A} \mathbf{C} \mathbf{A}^T =$

$$\begin{bmatrix} \sigma_{11} \mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T & \sigma_{12} \mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T & \cdots & \sigma_{1m} \mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T \\ \sigma_{21} \mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T & \sigma_{22} \mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T & \cdots & \sigma_{2m} \mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T \\ \vdots & \vdots & \cdots & \vdots \\ \sigma_{m1} \mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T & \sigma_{m2} \mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T & \cdots & \sigma_{mm} \mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T \end{bmatrix}.$$

Under  $H_0$ ,  $\text{vec}(\mathbf{L}\mathbf{B}) = \mathbf{A} \text{vec}(\mathbf{B}) = \mathbf{0}$ , and

$$\text{vec}(\mathbf{L}\hat{\mathbf{B}}) = \begin{pmatrix} \mathbf{L}\hat{\boldsymbol{\beta}}_1 \\ \mathbf{L}\hat{\boldsymbol{\beta}}_2 \\ \vdots \\ \mathbf{L}\hat{\boldsymbol{\beta}}_m \end{pmatrix} \sim N_{rm}(\mathbf{0}, \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}} \otimes \mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T).$$

Hence under  $H_0$ ,

$$[\text{vec}(\mathbf{L}\hat{\mathbf{B}})]^T [\boldsymbol{\Sigma}_{\boldsymbol{\epsilon}}^{-1} \otimes (\mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T)^{-1}] [\text{vec}(\mathbf{L}\hat{\mathbf{B}})] \sim \chi_{rm}^2,$$

and

$$T = [\text{vec}(\mathbf{L}\hat{\mathbf{B}})]^T [\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon}}^{-1} \otimes (\mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T)^{-1}] [\text{vec}(\mathbf{L}\hat{\mathbf{B}})] \xrightarrow{D} \chi_{rm}^2. \quad (10.4)$$

A large sample level  $\delta$  test will reject  $H_0$  if  $pval \leq \delta$  where

$$pval = P\left(\frac{T}{rm} < F_{rm, n-mp}\right). \quad (10.5)$$

Since least squares estimators are asymptotically normal, if the  $\boldsymbol{\epsilon}_i$  are iid for a large class of distributions,

$$\sqrt{n} \operatorname{vec}(\hat{\mathbf{B}} - \mathbf{B}) = \sqrt{n} \begin{pmatrix} \hat{\boldsymbol{\beta}}_1 - \boldsymbol{\beta}_1 \\ \hat{\boldsymbol{\beta}}_2 - \boldsymbol{\beta}_2 \\ \vdots \\ \hat{\boldsymbol{\beta}}_m - \boldsymbol{\beta}_m \end{pmatrix} \xrightarrow{D} N_{pm}(\mathbf{0}, \boldsymbol{\Sigma}_\epsilon \otimes \mathbf{W})$$

where

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

Then under  $H_0$ ,

$$\sqrt{n} \operatorname{vec}(\mathbf{L}\hat{\mathbf{B}}) = \sqrt{n} \begin{pmatrix} \mathbf{L}\hat{\boldsymbol{\beta}}_1 \\ \mathbf{L}\hat{\boldsymbol{\beta}}_2 \\ \vdots \\ \mathbf{L}\hat{\boldsymbol{\beta}}_m \end{pmatrix} \xrightarrow{D} N_{rm}(\mathbf{0}, \boldsymbol{\Sigma}_\epsilon \otimes \mathbf{LW}\mathbf{L}^T),$$

and

$$n [\operatorname{vec}(\mathbf{L}\hat{\mathbf{B}})]^T [\boldsymbol{\Sigma}_\epsilon^{-1} \otimes (\mathbf{LW}\mathbf{L}^T)^{-1}] [\operatorname{vec}(\mathbf{L}\hat{\mathbf{B}})] \xrightarrow{D} \chi_{rm}^2.$$

Hence (10.4) holds, and (10.5) gives a large sample level  $\delta$  test if the least squares estimators are asymptotically normal.

Kakizawa (2009) showed, under stronger assumptions than Theorem 10.8, that for a large class of iid error distributions, the following test statistics have the same  $\chi_{rm}^2$  limiting distribution when  $H_0$  is true, and the same non-central  $\chi_{rm}^2(\omega^2)$  limiting distribution with noncentrality parameter  $\omega^2$  when  $H_0$  is false under a local alternative. Hence the three tests are robust to the assumption of normality. The limiting null distribution is well known when the zero mean errors are iid from a multivariate normal distribution. See Khattree and Naik (1999, p. 68):  $(n-p)U(\mathbf{L}) \xrightarrow{D} \chi_{rm}^2$ ,  $(n-p)V(\mathbf{L}) \xrightarrow{D} \chi_{rm}^2$ , and  $-[n-p-0.5(m-r+3)] \log(\Lambda(\mathbf{L})) \xrightarrow{D} \chi_{rm}^2$ . Results from Kshirsagar (1972, p. 301) suggest that the third chi-square approximation is very good if  $n \geq 3(m+p)^2$  for multivariate normal error vectors.

Theorems 10.6 and 10.8 are useful for relating multivariate tests with the partial  $F$  test for multiple linear regression that tests whether a reduced model that omits some of the predictors can be used instead of the full model that uses all  $p$  predictors. The partial  $F$  test statistic is

$$F_R = \left[ \frac{SSE(R) - SSE(F)}{df_R - df_F} \right] / MSE(F)$$

where the residual sums of squares  $SSE(F)$  and  $SSE(R)$  and degrees of freedom  $df_F$  and  $df_r$  are for the full and reduced model while the mean square error  $MSE(F)$  is for the full model. Let the null hypothesis for the partial  $F$  test be  $H_0 : \mathbf{L}\boldsymbol{\beta} = \mathbf{0}$  where  $\mathbf{L}$  sets the coefficients of the predictors in the full model but not in the reduced model to 0. Seber and Lee (2003, p. 100) shows that

$$F_R = \frac{[\mathbf{L}\hat{\boldsymbol{\beta}}]^T (\mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T)^{-1} [\mathbf{L}\hat{\boldsymbol{\beta}}]}{r\hat{\sigma}^2}$$

is distributed as  $F_{r,n-p}$  if  $H_0$  is true and the errors are iid  $N(0, \sigma^2)$ . Note that for multiple linear regression with  $m = 1$ ,  $F_R = (n-p)U(\mathbf{L})/r$  since  $\hat{\boldsymbol{\Sigma}}_\epsilon^{-1} = 1/\hat{\sigma}^2$ . Hence the scaled Hotelling Lawley test statistic is the partial  $F$  test statistic extended to  $m > 1$  predictor variables by Theorem 10.6.

By Theorem 10.8, for example,  $rF_R \xrightarrow{D} \chi_r^2$  for a large class of nonnormal error distributions. If  $Z_n \sim F_{k,d_n}$ , then  $Z_n \xrightarrow{D} \chi_k^2/k$  as  $d_n \rightarrow \infty$ . Hence using the  $F_{r,n-p}$  approximation gives a large sample test with correct asymptotic level, and the partial  $F$  test is robust to nonnormality.

Similarly, using an  $F_{rm,n-pm}$  approximation for the following test statistics gives large sample tests with correct asymptotic level by Kakizawa (2009) and similar power for large  $n$ . The large sample test will have correct asymptotic level as long as the denominator degrees of freedom  $d_n \rightarrow \infty$  as  $n \rightarrow \infty$ , and  $d_n = n - pm$  reduces to the partial  $F$  test if  $m = 1$  and  $U(\mathbf{L})$  is used. Then the three test statistics are

$$\frac{-[n-p-0.5(m-r+3)]}{rm} \log(\Lambda(\mathbf{L})), \quad \frac{n-p}{rm} V(\mathbf{L}), \quad \text{and} \quad \frac{n-p}{rm} U(\mathbf{L}).$$

By Berndt and Savin (1977) and Anderson (1984, pp. 333, 371),

$$V(\mathbf{L}) \leq -\log(\Lambda(\mathbf{L})) \leq U(\mathbf{L}).$$

Hence the Hotelling Lawley test will have the most power and Pillai's test will have the least power.

Following Khattree and Naik (1999, pp. 67-68), there are several approximations used by the SAS software. For the Roy's largest root test, if  $h = \max(r, m)$ , use

$$\frac{n-p-h+r}{h} \lambda_{\max}(\mathbf{L}) \approx F(h, n-p-h+r).$$

The simulations in Section 10.5 suggest that this approximation is good for  $r = 1$  but poor for  $r > 1$ . Anderson (1984, p. 333) stated that Roy's largest root test has the greatest power if  $r = 1$  but is an inferior test for  $r > 1$ . Let  $g = n-p-(m-r+1)/2$ ,  $u = (rm-2)/4$  and  $t = \sqrt{r^2m^2-4}/\sqrt{m^2+r^2-5}$  for  $m^2+r^2-5 > 0$  and  $t = 1$ , otherwise. Assume  $H_0$  is true. Thus  $U \xrightarrow{P} 0$ ,  $V \xrightarrow{P} 0$ , and  $\Lambda \xrightarrow{P} 1$  as  $n \rightarrow \infty$ . Then

$$\frac{gt-2u}{rm} \frac{1-\Lambda^{1/t}}{\Lambda^{1/t}} \approx F(rm, gt-2u) \quad \text{or} \quad (n-p)t \frac{1-\Lambda^{1/t}}{\Lambda^{1/t}} \approx \chi_{rm}^2.$$

For large  $n$  and  $t > 0$ ,  $-\log(\Lambda) = -t \log(\Lambda^{1/t}) = -t \log(1 + \Lambda^{1/t} - 1) \approx t(1 - \Lambda^{1/t}) \approx t(1 - \Lambda^{1/t})/\Lambda^{1/t}$ . If it can not be shown that

$$(n-p)[- \log(\Lambda) - t(1 - \Lambda^{1/t})/\Lambda^{1/t}] \xrightarrow{P} 0 \text{ as } n \rightarrow \infty,$$

then it is possible that the approximate  $\chi_{rm}^2$  distribution may be the limiting distribution for only a small class of iid error distributions. When the  $\epsilon_i$  are iid  $N_m(\mathbf{0}, \boldsymbol{\Sigma}_\epsilon)$ , there are some exact results. For  $r = 1$ ,

$$\frac{n-p-m+1}{m} \frac{1-\Lambda}{\Lambda} \sim F(m, n-p-m+1).$$

For  $r = 2$ ,

$$\frac{2(n-p-m+1)}{2m} \frac{1-\Lambda^{1/2}}{\Lambda^{1/2}} \sim F(2m, 2(n-p-m+1)).$$

For  $m = 2$ ,

$$\frac{2(n-p)}{2r} \frac{1-\Lambda^{1/2}}{\Lambda^{1/2}} \sim F(2r, 2(n-p)).$$

Let  $s = \min(r, m)$ ,  $m_1 = (|r-m| - 1)/2$  and  $m_2 = (n-p-m-1)/2$ . Note that  $s(|r-m|+s) = \min(r, m) \max(r, m) = rm$ . Then

$$\frac{n-p}{rm} \frac{V}{1-V/s} = \frac{n-p}{s(|r-m|+s)} \frac{V}{1-V/s} \approx \frac{2m_2+s+1}{2m_1+s+1} \frac{V}{s-V} \approx$$

$$F(s(2m_1+s+1), s(2m_2+s+1)) \approx F(s(|r-m|+s), s(n-p)) = F(rm, s(n-p)).$$

This approximation is asymptotically correct by Slutsky's theorem since  $1 - V/s \xrightarrow{P} 1$ . Finally,  $\frac{n-p}{rm} U =$

$$\begin{aligned} \frac{n-p}{s(|r-m|+s)} U &\approx \frac{2(sm_2+1)}{s^2(2m_1+s+1)} U \approx F(s(2m_1+s+1), 2(sm_2+1)) \\ &\approx F(s(|r-m|+s), s(n-p)) = F(rm, s(n-p)). \end{aligned}$$

This approximation is asymptotically correct for a wide range of iid error distributions.

Multivariate analogs of tests for multiple linear regression can be derived with appropriate choice of  $\mathbf{L}$ . Assume a constant  $x_1 = 1$  is in the model. As a textbook convention, use  $\delta = 0.05$  if  $\delta$  is not given.

The four step MANOVA test of linear hypotheses is useful.

- i) State the hypotheses  $H_0 : \mathbf{LB} = \mathbf{0}$  and  $H_1 : \mathbf{LB} \neq \mathbf{0}$ .
- ii) Get test statistic from output.
- iii) Get pval from output.
- iv) State whether you reject  $H_0$  or fail to reject  $H_0$ . If  $\text{pval} \leq \delta$ , reject  $H_0$  and conclude that  $\mathbf{LB} \neq \mathbf{0}$ . If  $\text{pval} > \delta$ , fail to reject  $H_0$  and conclude that  $\mathbf{LB} = \mathbf{0}$  or that there is not enough evidence to conclude that  $\mathbf{LB} \neq \mathbf{0}$ .

The MANOVA test of  $H_0 : \mathbf{B} = \mathbf{0}$  versus  $H_1 : \mathbf{B} \neq \mathbf{0}$  is the special case corresponding to  $\mathbf{L} = \mathbf{I}$  and  $\mathbf{H} = \hat{\mathbf{B}}^T \mathbf{X}^T \mathbf{X} \hat{\mathbf{B}} = \hat{\mathbf{Z}}^T \hat{\mathbf{Z}}$ , but is usually not a test of interest.

The analog of the ANOVA  $F$  test for multiple linear regression is the MANOVA  $F$  test that uses  $\mathbf{L} = [\mathbf{0} \ \mathbf{I}_{p-1}]$  to test whether the nontrivial predictors are needed in the model. This test should reject  $H_0$  if the response and residual plots look good,  $n$  is large enough, and at least one response plot does not look like the corresponding residual plot. A response plot for  $Y_j$  will look like a residual plot if the identity line appears almost horizontal, hence the range of  $\hat{Y}_j$  is small. Response and residual plots are often useful for  $n \geq 10p$ .

The 4 step **MANOVA  $F$  test** of hypotheses uses  $\mathbf{L} = [\mathbf{0} \ \mathbf{I}_{p-1}]$ .

- i) State the hypotheses  $H_0$ : the nontrivial predictors are not needed in the mreg model  $H_1$ : at least one of the nontrivial predictors is needed.
- ii) Find the test statistic  $F_0$  from output.
- iii) Find the pval from output.
- iv) If  $\text{pval} \leq \delta$ , reject  $H_0$ . If  $\text{pval} > \delta$ , fail to reject  $H_0$ . If  $H_0$  is rejected, conclude that there is a mreg relationship between the response variables  $Y_1, \dots, Y_m$  and the predictors  $x_2, \dots, x_p$ . If you fail to reject  $H_0$ , conclude that there is a not a mreg relationship between  $Y_1, \dots, Y_m$  and the predictors  $x_2, \dots, x_p$ . (Or there is not enough evidence to conclude that there is a mreg relationship between the response variables and the predictors. Get the variable names from the story problem.)

The  $F_j$  test of hypotheses uses  $\mathbf{L}_j = [0, \dots, 0, 1, 0, \dots, 0]$ , where the 1 is in the  $j$ th position, to test whether the  $j$ th predictor  $x_j$  is needed in the model given that the other  $p - 1$  predictors are in the model. This test is an analog of the  $t$  tests for multiple linear regression. Note that  $x_j$  is not needed in the model corresponds to  $H_0 : \mathbf{B}_j = \mathbf{0}$  while  $x_j$  needed in the model corresponds to  $H_1 : \mathbf{B}_j \neq \mathbf{0}$  where  $\mathbf{B}_j^T$  is the  $j$ th row of  $\mathbf{B}$ .

The 4 step  $F_j$  test of hypotheses uses  $\mathbf{L}_j = [0, \dots, 0, 1, 0, \dots, 0]$  where the 1 is in the  $j$ th position.

- i) State the hypotheses  $H_0 : x_j$  is not needed in the model  $H_1 : x_j$  is needed.
- ii) Find the test statistic  $F_j$  from output.
- iii) Find pval from output.
- iv) If  $\text{pval} \leq \delta$ , reject  $H_0$ . If  $\text{pval} > \delta$ , fail to reject  $H_0$ . Give a nontechnical sentence restating your conclusion in terms of the story problem. If  $H_0$  is rejected, then conclude that  $x_j$  is needed in the mreg model for  $Y_1, \dots, Y_m$  given that the other predictors are in the model. If you fail to reject  $H_0$ , then conclude that  $x_j$  is not needed in the mreg model for  $Y_1, \dots, Y_m$  given that the other predictors are in the model. (Or there is not enough evidence to conclude that  $x_j$  is needed in the model. Get the variable names from the story problem.)

The Hotelling Lawley statistic

$$F_j = \frac{1}{d_j} \hat{\mathbf{B}}_j^T \hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon}}^{-1} \hat{\mathbf{B}}_j = \frac{1}{d_j} (\hat{\beta}_{j1}, \hat{\beta}_{j2}, \dots, \hat{\beta}_{jm}) \hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon}}^{-1} \begin{pmatrix} \hat{\beta}_{j1} \\ \hat{\beta}_{j2} \\ \vdots \\ \hat{\beta}_{jm} \end{pmatrix}$$

where  $\hat{\mathbf{B}}_j^T$  is the  $j$ th row of  $\hat{\mathbf{B}}$  and  $d_j = (\mathbf{X}^T \mathbf{X})_{jj}^{-1}$ , the  $j$ th diagonal entry of  $(\mathbf{X}^T \mathbf{X})^{-1}$ . The statistic  $F_j$  could be used for forward selection and backward elimination in variable selection.

The 4 step **MANOVA partial F test** of hypotheses has a full model using all of the variables and a reduced model where  $r$  of the variables are deleted. The  $i$ th row of  $\mathbf{L}$  has a 1 in the position corresponding to the  $i$ th variable to be deleted. Omitting the  $j$ th variable corresponds to the  $F_j$  test while omitting variables  $x_2, \dots, x_p$  corresponds to the MANOVA  $F$  test. Using  $\mathbf{L} = [\mathbf{0} \ \mathbf{I}_k]$  tests whether the last  $k$  predictors are needed in the multivariate linear regression model given that the remaining predictors are in the model.

- i) State the hypotheses  $H_0$ : the reduced model is good  $H_1$ : use the full model.
- ii) Find the test statistic  $F_R$  from output.
- iii) Find the pval from output.
- iv) If  $\text{pval} \leq \delta$ , reject  $H_0$  and conclude that the full model should be used. If  $\text{pval} > \delta$ , fail to reject  $H_0$  and conclude that the reduced model is good.

The *linmodpack* function `mltreg` produces the  $m$  response and residual plots, gives  $\hat{\mathbf{B}}$ ,  $\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon}}$ , the MANOVA partial  $F$  test statistic and pval corresponding to the reduced model that leaves out the variables given by indices (so  $x_2$  and  $x_4$  in the output below with  $F = 0.77$  and  $\text{pval} = 0.614$ ),  $F_j$  and the pval for the  $F_j$  test for variables 1, 2, ...,  $p$  (where  $p = 4$  in the output below so  $F_2 = 1.51$  with  $\text{pval} = 0.284$ ), and  $F_0$  and pval for the MANOVA  $F$  test (in the output below  $F_0 = 3.15$  and  $\text{pval} = 0.06$ ). Right click `stop` on the plots  $m$  times to advance the plots and to get the cursor back on the command line in *R*.

The command `out <- mltreg(x, y, indices=c(2))` would produce a MANOVA partial  $F$  test corresponding to the  $F_2$  test while the command `out <- mltreg(x, y, indices=c(2, 3, 4))` would produce a MANOVA partial  $F$  test corresponding to the MANOVA  $F$  test for a data set with  $p = 4$  predictor variables. The Hotelling Lawley trace statistic is used in the tests.

```
out <- mltreg(x, y, indices=c(2, 4))
$Bhat
      [,1]      [,2]      [,3]
[1,] 47.96841291 623.2817463 179.8867890
```

```

[2,] 0.07884384 0.7276600 -0.5378649
[3,] -1.45584256 -17.3872206 0.2337900
[4,] -0.01895002 0.1393189 -0.3885967
$Covhat
      [,1]      [,2]      [,3]
[1,] 21.91591 123.2557 132.339
[2,] 123.25566 2619.4996 2145.780
[3,] 132.33902 2145.7797 2954.082
$partial
      partialF      Pval
[1,] 0.7703294 0.6141573

$Ftable
      Fj      pvals
[1,] 6.30355375 0.01677169
[2,] 1.51013090 0.28449166
[3,] 5.61329324 0.02279833
[4,] 0.06482555 0.97701447

$MANOVA
      MANOVAF      pval
[1,] 3.150118 0.06038742

#Output for Example 10.2
y<-marry[,c(2,3)]; x<-marry[,-c(2,3)];
mltreg(x,y,indices=c(3,4))
$partial
      partialF      Pval
[1,] 0.2001622 0.9349877

$Ftable
      Fj      pvals
[1,] 4.35326807 0.02870083
[2,] 600.57002201 0.00000000
[3,] 0.08819810 0.91597268
[4,] 0.06531531 0.93699302

$MANOVA
      MANOVAF      pval
[1,] 295.071 1.110223e-16

```

**Example 10.2.** The above output is for the Hebbler (1847) data from the 1843 Prussia census. Sometimes if the wife or husband was not at the household, then s/he would not be counted.  $Y_1$  = number of married civilian men in the district,  $Y_2$  = number of women married to civilians in the district,  $x_2$  = population of the district in 1843,  $x_3$  = number of married military men

in the district, and  $x_4$  = number of women married to military men in the district. The reduced model deletes  $x_3$  and  $x_4$ . The constant uses  $x_1 = 1$ .

- a) Do the MANOVA  $F$  test.
- b) Do the  $F_2$  test.
- c) Do the  $F_4$  test.
- d) Do an appropriate 4 step test for the reduced model that deletes  $x_3$  and  $x_4$ .
- e) The output for the reduced model that deletes  $x_1$  and  $x_2$  is shown below. Do an appropriate 4 step test.

```
$partial
      partialF Pval
[1,] 569.6429    0
```

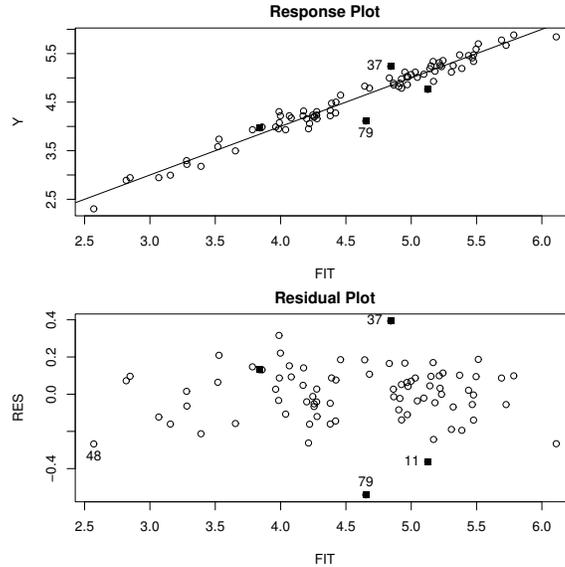
**Solution:**

- a) i)  $H_0$ : the nontrivial predictors are not needed in the mreg model  
 $H_1$ : at least one of the nontrivial predictors is needed
  - ii)  $F_0 = 295.071$
  - iii)  $pval = 0$
  - iv) Reject  $H_0$ , the nontrivial predictors are needed in the mreg model.
- b) i)  $H_0$ :  $x_2$  is not needed in the model  $H_1$ :  $x_2$  is needed
  - ii)  $F_2 = 600.57$
  - iii)  $pval = 0$
  - iv) Reject  $H_0$ , *population of the district* is needed in the model.
- c) i)  $H_0$ :  $x_4$  is not needed in the model  $H_1$ :  $x_4$  is needed
  - ii)  $F_4 = 0.065$
  - iii)  $pval = 0.937$
  - iv) Fail to reject  $H_0$ , *number of women married to military men* is not needed in the model given that the other predictors are in the model.
- d) i)  $H_0$ : the reduced model is good  $H_1$ : use the full model.
  - ii)  $F_R = 0.200$
  - iii)  $pval = 0.935$
  - iv) Fail to reject  $H_0$ , so the reduced model is good.
- e) i)  $H_0$ : the reduced model is good  $H_1$ : use the full model.
  - ii)  $F_R = 569.6$
  - iii)  $pval = 0.00$
  - iv) Reject  $H_0$ , so use the full model.

## 10.5 An Example and Simulations

In the DD plot, cases to the left of the vertical line are in their nonparametric prediction region. The long horizontal line corresponds to a similar cutoff based on the RD. The shorter horizontal line that ends at the identity line

is the parametric MVN prediction region from Section 4.4 applied to the  $\hat{z}_i$ . Points below these two lines are only conjectured to be large sample prediction regions, but are added to the DD plot as visual aids. Note that  $\hat{z}_i = \hat{\mathbf{y}}_f + \hat{\boldsymbol{\epsilon}}_i$ , and adding a constant  $\hat{\mathbf{y}}_f$  to all of the residual vectors does not change the Mahalanobis distances, so the DD plot of the residual vectors can be used to display the prediction regions.



**Fig. 10.1** Plots for  $Y_1 = \log(S)$ .

**Example 10.3.** Cook and Weisberg (1999, pp. 351, 433, 447) gave a data set on 82 mussels sampled off the coast of New Zealand. Let  $Y_1 = \log(S)$  and  $Y_2 = \log(M)$  where  $S$  is the shell mass and  $M$  is the muscle mass. The predictors are  $X_2 = L$ ,  $X_3 = \log(W)$ , and  $X_4 = H$ : the shell length,  $\log(\text{width})$ , and height. To check linearity of the multivariate linear regression model, Figures 10.1 and 10.2 give the response and residual plots for  $Y_1$  and  $Y_2$ . The response plots show strong linear relationships. For  $Y_1$ , case 79 sticks out while for  $Y_2$ , cases 8, 25, and 48 are not fit well. Highlighted cases had Cook's distance  $> \min(0.5, 2p/n)$ . See Cook (1977).

To check the error vector distribution, the DD plot should be used instead of univariate residual plots, which do not take into account the correlations of the random variables  $\epsilon_1, \dots, \epsilon_m$  in the error vector  $\boldsymbol{\epsilon}$ . A residual vector  $\hat{\boldsymbol{\epsilon}} = (\hat{\boldsymbol{\epsilon}} - \boldsymbol{\epsilon}) + \boldsymbol{\epsilon}$  is a combination of  $\boldsymbol{\epsilon}$  and a discrepancy  $\hat{\boldsymbol{\epsilon}} - \boldsymbol{\epsilon}$  that tends to have an approximate multivariate normal distribution. The  $\hat{\boldsymbol{\epsilon}} - \boldsymbol{\epsilon}$  term can dominate for small to moderate  $n$  when  $\boldsymbol{\epsilon}$  is not multivariate normal,

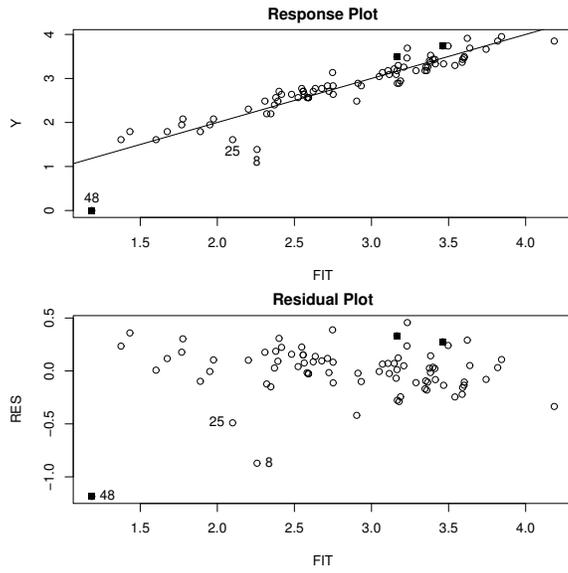


Fig. 10.2 Plots for  $Y_2 = \log(M)$ .

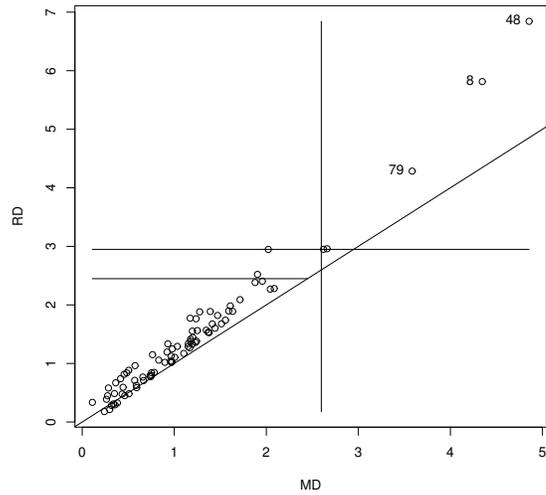
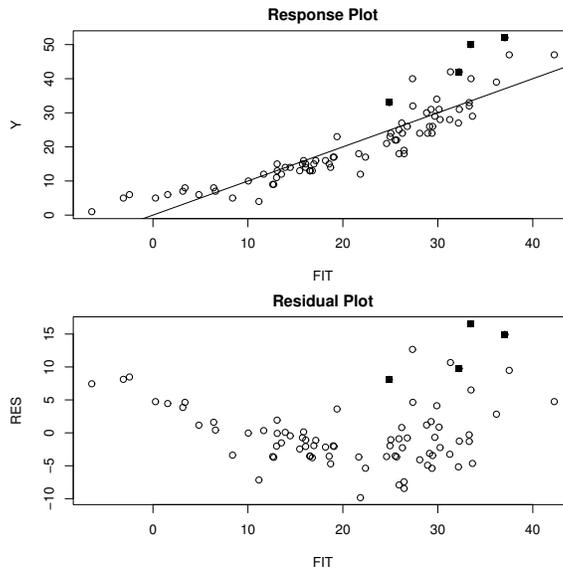


Fig. 10.3 DD Plot of the Residual Vectors for the Mussels Data.

incorrectly suggesting that the distribution of the error vector  $\epsilon$  is closer to a multivariate normal distribution than is actually the case. Figure 10.3 shows the DD plot of the residual vectors. The plotted points are highly correlated but do not cover the identity line, suggesting an elliptically contoured error distribution that is not multivariate normal. The nonparametric 90% prediction region for the residuals consists of the points to the left of the vertical line  $MD = 2.60$ . Cases 8, 48, and 79 have especially large distances.

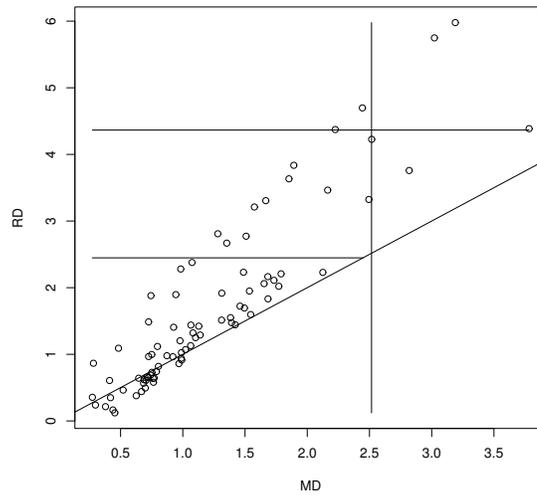
The four Hotelling Lawley  $F_j$  statistics were greater than 5.77 with pvalues less than 0.005, and the MANOVA  $F$  statistic was 337.8 with pvalue  $\approx 0$ .

The response, residual, and DD plots are effective for finding influential cases, for checking linearity, for checking whether the error distribution is multivariate normal or some other elliptically contoured distribution, and for displaying the nonparametric prediction region. Note that cases to the right of the vertical line correspond to cases with  $\mathbf{y}_i$  that are not in their prediction region. These are the cases corresponding to residual vectors with large Mahalanobis distances. Adding a constant does not change the distance, so the DD plot for the residual vectors is the same as the DD plot for the  $\hat{\mathbf{z}}_i$ .



**Fig. 10.4** Plots for  $Y_2 = M$ .

c) Now suppose the same model is used except  $Y_2 = M$ . Then the response and residual plots for  $Y_1$  remain the same, but the plots shown in Figure 10.4 show curvature about the identity and  $r = 0$  lines. Hence the linearity condition is violated. Figure 10.5 shows that the plotted points in the DD plot have correlation well less than one, suggesting that the error vector distribution



**Fig. 10.5** DD Plot When  $Y_2 = M$ .

is no longer elliptically contoured. The nonparametric 90% prediction region for the residual vectors consists of the points to the left of the vertical line  $MD = 2.52$ , and contains 95% of the training data. Note that the plots can be used to quickly assess whether power transformations have resulted in a linear model, and whether influential cases are present. *R* code for producing the five figures is shown below.

```

y <- log(mussels)[,4:5]
x <- mussels[,1:3]
x[,2] <- log(x[,2])
z<-cbind(x,y) #scatterplot matrix
pairs(z, labels=c("L","log(W)","H","log(S)","log(M)"))
ddplot4(z) #right click Stop, DD plot of MLD model
out <- mltreg(x,y) #right click Stop 4 times, Fig. 10.1, 10.2
ddplot4(out$res) #right click Stop, Fig. 10.3
y[,2] <- mussels[,5]
tem <- mltreg(x,y) #right click Stop 4 times, Fig. 10.4
ddplot4(tem$res) #right click Stop, Fig. 10.5

```

### 10.5.1 Simulations for Testing

A small simulation was used to study the Wilks'  $\Lambda$  test, the Pillai's trace test, the Hotelling Lawley trace test, and the Roy's largest root test for the  $F_j$  tests and the MANOVA  $F$  test for multivariate linear regression. The first row of  $\mathbf{B}$  was always  $\mathbf{1}^T$  and the last row of  $\mathbf{B}$  was always  $\mathbf{0}^T$ . When the null hypothesis for the MANOVA  $F$  test is true, all but the first row corresponding to the constant are equal to  $\mathbf{0}^T$ . When  $p \geq 3$  and the null hypothesis for the MANOVA  $F$  test is false, then the second to last row of  $\mathbf{B}$  is  $(1, 0, \dots, 0)$ , the third to last row is  $(1, 1, 0, \dots, 0)$  et cetera as long as the first row is not changed from  $\mathbf{1}^T$ . First  $m \times 1$  error vectors  $\mathbf{w}_i$  were generated such that the  $m$  random variables in the vector  $\mathbf{w}_i$  are iid with variance  $\sigma^2$ . Let the  $m \times m$  matrix  $\mathbf{A} = (a_{ij})$  with  $a_{ii} = 1$  and  $a_{ij} = \psi$  where  $0 \leq \psi < 1$  for  $i \neq j$ . Then  $\boldsymbol{\epsilon}_i = \mathbf{A}\mathbf{w}_i$  so that  $\boldsymbol{\Sigma}_{\boldsymbol{\epsilon}} = \sigma^2 \mathbf{A}\mathbf{A}^T = (\sigma_{ij})$  where the diagonal entries  $\sigma_{ii} = \sigma^2[1 + (m-1)\psi^2]$  and the off diagonal entries  $\sigma_{ij} = \sigma^2[2\psi + (m-2)\psi^2]$  where  $\psi = 0.10$ . Hence the correlations are  $(2\psi + (m-2)\psi^2)/(1 + (m-1)\psi^2)$ . As  $\psi$  gets close to 1, the error vectors cluster about the line in the direction of  $(1, \dots, 1)^T$ . We used  $\mathbf{w}_i \sim N_m(\mathbf{0}, \mathbf{I})$ ,  $\mathbf{w}_i \sim (1 - \tau)N_m(\mathbf{0}, \mathbf{I}) + \tau N_m(\mathbf{0}, 25\mathbf{I})$  with  $0 < \tau < 1$  and  $\tau = 0.25$  in the simulation,  $\mathbf{w}_i \sim$  multivariate  $t_d$  with  $d = 7$  degrees of freedom, or  $\mathbf{w}_i \sim$  lognormal - E(lognormal): where the  $m$  components of  $\mathbf{w}_i$  were iid with distribution  $e^z - E(e^z)$  where  $z \sim N(0, 1)$ . Only the lognormal distribution is not elliptically contoured.

**Table 10.1** Test Coverages: MANOVA  $F$   $H_0$  is True.

$\mathbf{w}$ dist	$n$	test	$F_1$	$F_2$	$F_{p-1}$	$F_p$	$F_M$
MVN 300	W	1	0.043	0.042	0.041	0.018	
MVN 300	P	1	0.040	0.038	0.038	0.007	
MVN 300	HL	1	0.059	0.058	0.057	0.045	
MVN 300	R	1	0.051	0.049	0.048	0.993	
MVN 600	W	1	0.048	0.043	0.043	0.034	
MVN 600	P	1	0.046	0.042	0.041	0.026	
MVN 600	HL	1	0.055	0.052	0.050	0.052	
MVN 600	R	1	0.052	0.048	0.047	0.994	
MIX 300	W	1	0.042	0.043	0.044	0.017	
MIX 300	P	1	0.039	0.040	0.042	0.008	
MIX 300	HL	1	0.057	0.059	0.058	0.039	
MIX 300	R	1	0.050	0.050	0.051	0.993	
MVT(7) 300	W	1	0.048	0.036	0.045	0.020	
MVT(7) 300	P	1	0.046	0.032	0.042	0.011	
MVT(7) 300	HL	1	0.064	0.049	0.058	0.045	
MVT(7) 300	R	1	0.055	0.043	0.051	0.993	
LN 300	W	1	0.043	0.047	0.040	0.020	
LN 300	P	1	0.039	0.045	0.037	0.009	
LN 300	HL	1	0.057	0.061	0.058	0.041	
LN 300	R	1	0.049	0.055	0.050	0.994	

**Table 10.2** Test Coverages: MANOVA  $F$   $H_0$  is False.

$n$	$m = p$	test	$F_1$	$F_2$	$F_{p-1}$	$F_p$	$F_M$
30	5	W	0.012	0.222	0.058	0.000	0.006
30	5	P	0.000	0.000	0.000	0.000	0.000
30	5	HL	0.382	0.694	0.322	0.007	0.579
30	5	R	0.799	0.871	0.549	0.047	0.997
50	5	W	0.984	0.955	0.644	0.017	0.963
50	5	P	0.971	0.940	0.598	0.012	0.871
50	5	HL	0.997	0.979	0.756	0.053	0.991
50	5	R	0.996	0.978	0.744	0.049	1
105	10	W	0.650	0.970	0.191	0.000	0.633
105	10	P	0.109	0.812	0.050	0.000	0.000
105	10	HL	0.964	0.997	0.428	0.000	1
105	10	R	1	1	0.892	0.052	1
150	10	W	1	1	0.948	0.032	1
150	10	P	1	1	0.941	0.025	1
150	10	HL	1	1	0.966	0.060	1
150	10	R	1	1	0.965	0.057	1
450	20	W	1	1	0.999	0.020	1
450	20	P	1	1	0.999	0.016	1
450	20	HL	1	1	0.999	0.035	1
450	20	R	1	1	0.999	0.056	1

The simulation used 5000 runs, and  $H_0$  was rejected if the  $F$  statistic was greater than  $F_{d_1, d_2}(0.95)$  where  $P(F_{d_1, d_2} < F_{d_1, d_2}(0.95)) = 0.95$  with  $d_1 = rm$  and  $d_2 = n - mp$  for the test statistics

$$\frac{-(n - p - 0.5(m - r + 3))}{rm} \log(\Lambda(\mathbf{L})), \quad \frac{n - p}{rm} V(\mathbf{L}), \quad \text{and} \quad \frac{n - p}{rm} U(\mathbf{L}),$$

while  $d_1 = h = \max(r, m)$  and  $d_2 = n - p - h + r$  for the test statistic

$$\frac{n - p - h + r}{h} \lambda_{max}(\mathbf{L}).$$

Denote these statistics by  $W$ ,  $P$ ,  $HL$ , and  $R$ . Let the coverage be the proportion of times that  $H_0$  is rejected. We want coverage near 0.05 when  $H_0$  is true and coverage close to 1 for good power when  $H_0$  is false. With 5000 runs, coverage outside of (0.04, 0.06) suggests that the true coverage is not 0.05. Coverages are tabled for the  $F_1, F_2, F_{p-1}$ , and  $F_p$  test and for the MANOVA  $F$  test denoted by  $F_M$ . The null hypothesis  $H_0$  was always true for the  $F_p$  test and always false for the  $F_1$  test. When the MANOVA  $F$  test was true,  $H_0$  was true for the  $F_j$  tests with  $j \neq 1$ . When the MANOVA  $F$  test was false,  $H_0$  was false for the  $F_j$  tests with  $j \neq p$ , but the  $F_{p-1}$  test should be hardest to reject for  $j \neq p$  by construction of  $\mathbf{B}$  and the error vectors.

When the null hypothesis  $H_0$  was true, simulated values started to get close to nominal levels for  $n \geq 0.8(m+p)^2$ , and were fairly good for  $n \geq 1.5(m+p)^2$ . The exception was Roy's test which rejects  $H_0$  far too often if  $r > 1$ . See

Table 10.1 where we want values for the  $F_1$  test to be close to 1 since  $H_0$  is false for the  $F_1$  test, and we want values close to 0.05, otherwise. Roy's test was very good for the  $F_j$  tests but very poor for the MANOVA  $F$  test. Results are shown for  $m = p = 10$ . As expected from Berndt and Savin (1977), Pillai's test rejected  $H_0$  less often than Wilks' test which rejected  $H_0$  less often than the Hotelling Lawley test. Based on a much larger simulation study, using the four types of error vector distributions and  $m = p$ , the tests had approximately correct level if  $n \geq 0.83(m+p)^2$  for the Hotelling Lawley test, if  $n \geq 2.80(m+p)^2$  for the Wilks' test (agreeing with Kshirsagar (1972)  $n \geq 3(m+p)^2$  for multivariate normal data), and if  $n \geq 4.2(m+p)^2$  for Pillai's test.

In Table 10.2,  $H_0$  is only true for the  $F_p$  test where  $p = m$ , and we want values in the  $F_p$  column near 0.05. We want values near 1 for high power otherwise. If  $H_0$  is false, often  $H_0$  will be rejected for small  $n$ . For example, if  $n \geq 10p$ , then the  $m$  residual plots should start to look good, and the MANOVA  $F$  test should be rejected. For the simulated data, the test had fair power for  $n$  not much larger than  $mp$ . Results are shown for the lognormal distribution.

Some  $R$  output for reproducing the simulation is shown below. The *linmod-pack* function is `mregsim` and `etype = 1` uses data from a MVN distribution. The `fcov` line computed the Hotelling Lawley statistic using Equation (10.3) while the `hotlawcov` line used Definition 10.9. The `mnull=T` part of the command means we want the first value near 1 for high power and the next three numbers near the nominal level 0.05 except for `mancv` where we want all of the MANOVA  $F$  test statistics to be near the nominal level of 0.05. The `mnull=F` part of the command means want all values near 1 for high power except for the last column (for the terms other than `mancv`) corresponding to the  $F_p$  test where  $H_0$  is true so we want values near the nominal level of 0.05. The "coverage" is the proportion of times that  $H_0$  is rejected, so "coverage" is short for "power" and "level": we want the coverage near 1 for high power when  $H_0$  is false and we want the coverage near the nominal level 0.05 when  $H_0$  is true. Also see Problem 10.10.

```
mregsim(nruns=5000,etype=1,mnull=T)
$wilkcov
[1] 1.0000 0.0450 0.0462 0.0430
$pilcov
[1] 1.0000 0.0414 0.0432 0.0400
$hotlawcov
[1] 1.0000 0.0522 0.0516 0.0490
$roycov
[1] 1.0000 0.0512 0.0500 0.0480
$fcov
[1] 1.0000 0.0522 0.0516 0.0490
$mancv
      wcv   pcv  hlcw   rcv   fcw
```

```
[1,] 0.0406 0.0332 0.049 0.1526 0.049

mregsim(nruns=5000, etype=2, mnull=F)

$wilkcov
[1] 0.9834 0.9814 0.9104 0.0408
$pilecov
[1] 0.9824 0.9804 0.9064 0.0372
$shotlawcov
[1] 0.9856 0.9838 0.9162 0.0480
$roycov
[1] 0.9848 0.9834 0.9156 0.0462
$fcov
[1] 0.9856 0.9838 0.9162 0.0480
$mancv
      wcv      pcv      hlcv      rcv      fcv
[1,] 0.993 0.9918 0.9942 0.9978 0.9942
```

See Olive (2017b, § 12.5.2) for simulations for the prediction region. Also see Problem 10.11.

## 10.6 The Robust `rmreg2` Estimator

The robust multivariate linear regression estimator `rmreg2` is the classical multivariate linear regression estimator applied to the RMVN set when RMVN is computed from the vectors  $\mathbf{u}_i = (x_{i2}, \dots, x_{ip}, Y_{i1}, \dots, Y_{im})^T$  for  $i = 1, \dots, n$ . Hence  $\mathbf{u}_i$  is the  $i$ th case with  $x_{i1} = 1$  deleted. This regression estimator has considerable outlier resistance, and is one of the most outlier resistant practical robust regression estimator for the  $m = 1$  multiple linear regression case. See Chapter 7. The `rmreg2` estimator has been shown to be consistent if the  $\mathbf{u}_i$  are iid from a large class of elliptically contoured distributions, which is a much stronger assumption than having iid error vectors  $\epsilon_i$ .

Theorem 2.20 gave a second way to compute  $\hat{\boldsymbol{\beta}}$ , and there is a similar result for multivariate linear regression. Let  $\mathbf{x} = (1, \mathbf{u}^T)^T$  and let  $\boldsymbol{\beta} = (\beta_1, \boldsymbol{\beta}_2^T)^T = (\alpha, \boldsymbol{\eta}^T)^T$ . Now for multivariate linear regression,  $\hat{\boldsymbol{\beta}}_j = (\hat{\alpha}_j, \hat{\boldsymbol{\eta}}_j^T)^T$  where  $\hat{\alpha}_j = \bar{Y}_j - \hat{\boldsymbol{\eta}}_j^T \bar{\mathbf{u}}$  and  $\hat{\boldsymbol{\eta}}_j = \hat{\boldsymbol{\Sigma}}_{\mathbf{u}}^{-1} \hat{\boldsymbol{\Sigma}}_{\mathbf{u}Y_j}$  by Theorem 2.20. Let  $\hat{\boldsymbol{\Sigma}}_{\mathbf{u}Y_j} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{w}_i - \bar{\mathbf{w}})(\mathbf{y}_i - \bar{\mathbf{y}})^T$  which has  $j$ th column  $\hat{\boldsymbol{\Sigma}}_{\mathbf{w}Y_j}$  for  $j = 1, \dots, m$ . Let

$$\mathbf{v} = \begin{pmatrix} \mathbf{u} \\ \mathbf{y} \end{pmatrix}, \quad E(\mathbf{v}) = \boldsymbol{\mu}_v = \begin{pmatrix} E(\mathbf{u}) \\ E(\mathbf{y}) \end{pmatrix} = \begin{pmatrix} \boldsymbol{\mu}_u \\ \boldsymbol{\mu}_y \end{pmatrix}, \quad \text{and} \quad \text{Cov}(\mathbf{v}) = \boldsymbol{\Sigma}_v =$$

$$\begin{pmatrix} \Sigma_{uu} & \Sigma_{uy} \\ \Sigma_{yu} & \Sigma_{yy} \end{pmatrix}.$$

Let the vector of constants be  $\alpha^T = (\alpha_1, \dots, \alpha_m)$  and the matrix of slope vectors  $B_S = [\eta_1 \ \eta_2 \ \dots \ \eta_m]$ . Then the population least squares coefficient matrix is

$$B = \begin{pmatrix} \alpha^T \\ B_S \end{pmatrix}$$

where  $\alpha = \mu_y - B_S^T \mu_u$  and  $B_S = \Sigma_u^{-1} \Sigma_{uy}$  where  $\Sigma_u = \Sigma_{uu}$ .

If the  $u_i$  are iid with nonsingular covariance matrix  $\text{Cov}(u)$ , the least squares estimator

$$\hat{B} = \begin{pmatrix} \hat{\alpha}^T \\ \hat{B}_S \end{pmatrix}$$

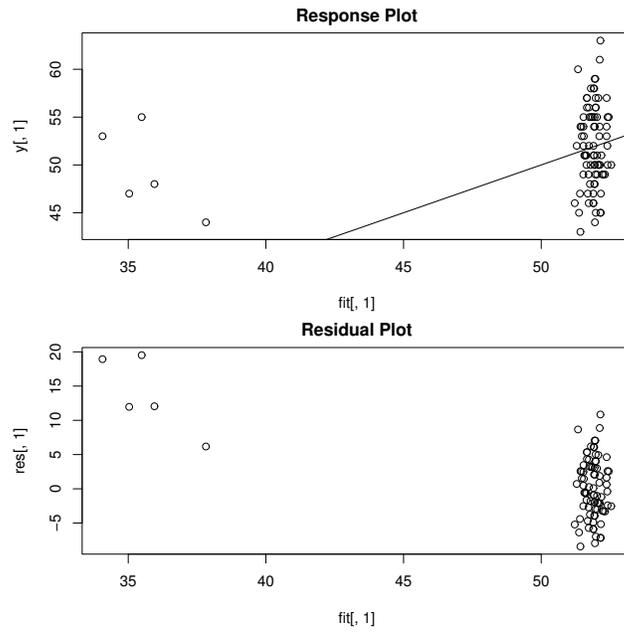
where  $\hat{\alpha} = \bar{y} - \hat{B}_S^T \bar{u}$  and  $\hat{B}_S = \hat{\Sigma}_u^{-1} \hat{\Sigma}_{uy}$ . The least squares multivariate linear regression estimator can be calculated by computing the classical estimator  $(\bar{v}, S_v) = (\bar{v}, \hat{\Sigma}_v)$  of multivariate location and dispersion on the  $v_i$ , and then plug in the results into the formulas for  $\hat{\alpha}$  and  $\hat{B}_S$ .

Let  $(T, C) = (\tilde{\mu}_v, \tilde{\Sigma}_v)$  be a robust estimator of multivariate location and dispersion. If  $\tilde{\mu}_v$  is a consistent estimator of  $\mu_v$  and  $\tilde{\Sigma}_v$  is a consistent estimator of  $c \Sigma_v$  for some constant  $c > 0$ , then a robust estimator of multivariate linear regression is the plug in estimator  $\tilde{\alpha} = \tilde{\mu}_y - \tilde{B}_S^T \tilde{\mu}_u$  and  $\tilde{B}_S = \tilde{\Sigma}_u^{-1} \tilde{\Sigma}_{uy}$ .

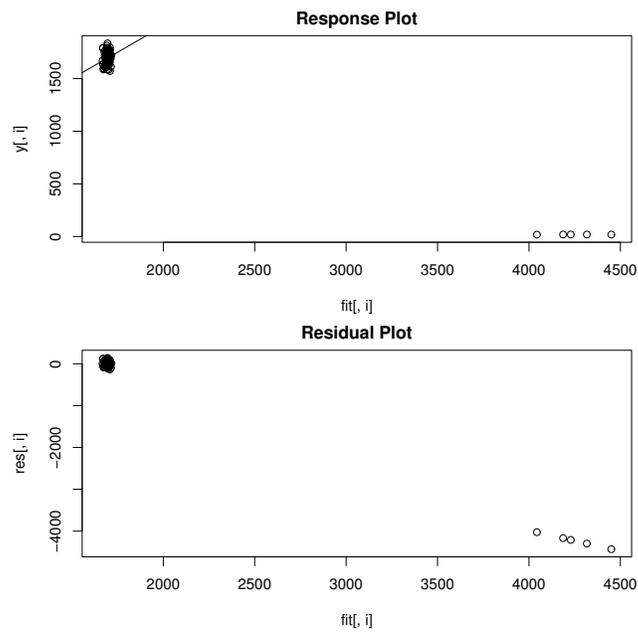
For the `rmreg2` estimator,  $(T, C)$  is the classical estimator applied to the RMVN set when RMVN is applied to vectors  $v_i$  for  $i = 1, \dots, n$  (could use  $(T, C) = \text{RMVN}$  estimator since the scaling does not matter for this application). Then  $(T, C)$  is a  $\sqrt{n}$  consistent estimator of  $(\mu_v, c \Sigma_v)$  if the  $v_i$  are iid from a large class of  $EC_d(\mu_v, \Sigma_v, g)$  distributions where  $d = m + p - 1$ . Thus the classical and robust estimators of multivariate linear regression are both  $\sqrt{n}$  consistent estimators of  $B$  if the  $v_i$  are iid from a large class of elliptically contoured distributions. This assumption is quite strong, but the robust estimator is useful for detecting outliers. When there are categorical predictors or the joint distribution of  $v$  is not elliptically contoured, it is possible that the robust estimator is bad and very different from the good classical least squares estimator. The `linmodpack` function `rmreg2` computes the `rmreg2` estimator and produces the response and residual plots.

**Example 10.4.** Buxton (1920) gave various measurements of 88 men. Let  $Y_1 = \text{nasal height}$  and  $Y_2 = \text{height}$  with  $x_2 = \text{head length}$ ,  $x_3 = \text{bigonal breadth}$ , and  $x_4 = \text{cephalic index}$ . Five individuals, numbers 62–66, were reported to be about 0.75 inches tall with head lengths well over five feet! Thus  $Y_2$  and  $x_2$  have massive outliers. Figures 10.6 and 10.7 show that the response and residual plots corresponding to `rmreg2` do not have fits that pass through the outliers.

These figures can be made with the following *R* commands.



**Fig. 10.6** Plots for  $Y_1 = \text{nasal height}$  using `rmreg2`.



**Fig. 10.7** Plots for  $Y_2 = \text{height}$  using `rmreg2`.

```

ht <- buxy; z <- cbind(buxx,ht);
y <- z[,c(2,5)]; x <- z[,c(1,3,4)]
# compare mltrreg(x,y) #right click Stop 4 times
out <- rmreg2(x,y) #right click Stop 4 times
# try ddplot4(out$res) #right click Stop

```

The residual bootstrap for the test  $H_0 : \mathbf{LB} = \mathbf{0}$  may be useful. Take a sample of size  $n$  with replacement from the residual vectors to form  $\mathbf{Z}_1^*$  with  $i$ th row  $\mathbf{y}_i^{*T}$  where  $\mathbf{y}_i^* = \hat{\mathbf{y}}_i + \boldsymbol{\epsilon}_i^*$ . The function `rmreg3` gets the `rmreg2` estimator without the plots. Using `rmreg3`, regress  $\mathbf{Z}$  on  $\mathbf{X}$  to get  $\text{vec}(\mathbf{L}\hat{\mathbf{B}}_1^*)$ . Repeat  $B$  times to get a bootstrap sample  $\mathbf{w}_1, \dots, \mathbf{w}_B$  where  $\mathbf{w}_i = \text{vec}(\mathbf{L}\hat{\mathbf{B}}_i^*)$ . The nonparametric bootstrap uses  $n$  cases drawn with replacement, and may also be useful. Apply the nonparametric prediction region to the  $\mathbf{w}_i$  and see if  $\mathbf{0}$  is in the region. If  $\mathbf{L}$  is  $r \times p$ , then  $\mathbf{w}$  is  $rp \times 1$ , and we likely need  $n \geq \max[50rp, 3(m+p)^2]$ .

## 10.7 Bootstrap

### 10.7.1 Parametric Bootstrap

The parametric bootstrap for the multivariate linear regression model uses  $\mathbf{y}_i^* \sim N_m(\hat{\mathbf{B}}^T \mathbf{x}_i, \hat{\boldsymbol{\Sigma}}\boldsymbol{\epsilon})$  for  $i = 1, \dots, n$  where **we are not assuming** that the  $\boldsymbol{\epsilon}_i \sim N_m(\mathbf{0}, \boldsymbol{\Sigma}\boldsymbol{\epsilon})$ . Let  $\mathbf{Z}_j^*$  have  $i$ th row  $\mathbf{y}_i^{*T}$  and regress  $\mathbf{Z}_j^*$  on  $\mathbf{X}$  to obtain  $\hat{\mathbf{B}}_j^*$  for  $j = 1, \dots, B$ . Let  $S \subseteq I$ , let  $\hat{\mathbf{B}}_I = (\mathbf{X}_I^T \mathbf{X}_I)^{-1} \mathbf{X}_I^T \mathbf{Z}^*$ , and assume  $n(\mathbf{X}_I^T \mathbf{X}_I)^{-1} \xrightarrow{P} \mathbf{W}_I$  for any  $I$  such that  $S \subseteq I$ . Then with calculations similar to those for the multiple linear regression model parametric bootstrap of Section 4.6.1,  $E(\hat{\mathbf{B}}_I^*) = \hat{\mathbf{B}}_I$ ,

$$\sqrt{n} \text{vec}(\hat{\mathbf{B}}_I - \mathbf{B}_I) \xrightarrow{D} N_{arm}(\mathbf{0}, \boldsymbol{\Sigma}\boldsymbol{\epsilon} \otimes \mathbf{W}_I),$$

and  $\sqrt{n} \text{vec}(\hat{\mathbf{B}}_I^* - \hat{\mathbf{B}}_I) \sim N_{arm}(\mathbf{0}, \hat{\boldsymbol{\Sigma}}\boldsymbol{\epsilon} \otimes n(\mathbf{X}_I^T \mathbf{X}_I)^{-1}) \xrightarrow{D} N_{arm}(\mathbf{0}, \boldsymbol{\Sigma}\boldsymbol{\epsilon} \otimes \mathbf{W}_I)$

as  $n, B \rightarrow \infty$  if  $S \subseteq I$ . Let  $\hat{\mathbf{B}}_{I,0}^*$  be formed from  $\hat{\mathbf{B}}_I^*$  by adding rows of zeros corresponding to omitted variables.

### 10.7.2 Residual Bootstrap

The residual bootstrap uses the multivariate linear regression model

$$\mathbf{Z}^* = \mathbf{X}\hat{\mathbf{B}} + \hat{\mathbf{E}}^W$$

where the rows of  $\hat{\mathbf{E}}^W$  are sampled with replacement from the rows of  $\hat{\mathbf{E}}^W$ . Regress  $\mathbf{Z}^*$  of  $\mathbf{X}$  and repeat to get the bootstrap sample  $\hat{\mathbf{B}}_1^*, \dots, \hat{\mathbf{B}}_B^*$ .

### 10.7.3 Nonparametric Bootstrap

The nonparametric bootstrap samples cases  $(\mathbf{y}_i^T, \mathbf{x}_i^T)^T$  with replacement to form  $(\mathbf{Z}_j^*, \mathbf{X}_j^*)$ , and regresses  $\mathbf{Z}_j^*$  on  $\mathbf{X}_j^*$  to get  $\hat{\mathbf{B}}_j^*$  for  $j = 1, \dots, B$ . The nonparametric bootstrap can be useful even if heteroscedasticity or overdispersion is present, if the cases are an iid sample from some population, a very strong assumption. See Eck (2018) for using the residual bootstrap and nonparametric bootstrap to bootstrap multivariate linear regression.

## 10.8 Data Splitting

The theory for multivariate linear regression assumes that the model is known before gathering data. If variable selection and response transformations are performed to build a model, then the estimators are biased and results for inference fail to hold in that pvalues and coverage of confidence and prediction regions will be wrong.

Data splitting can be used in a manner similar to how data splitting is used for MLR and other regression models. A pilot study is an alternative to data splitting.

## 10.9 Ridge Regression, PCR, and Other High Dimensional Methods

### 10.10 Summary

1) The multivariate linear regression model is a special case of the multivariate linear model where at least one predictor variable  $x_j$  is continuous. The MANOVA model in Chapter 9 is a multivariate linear model where all of the predictors are categorical variables so the  $x_j$  are coded and are often indicator variables.

2) The **multivariate linear regression model**  $\mathbf{y}_i = \mathbf{B}^T \mathbf{x}_i + \epsilon_i$  for  $i = 1, \dots, n$  has  $m \geq 2$  response variables  $Y_1, \dots, Y_m$  and  $p$  predictor variables  $x_1, x_2, \dots, x_p$ . The  $i$ th case is  $(\mathbf{x}_i^T, \mathbf{y}_i^T) = (x_{i1}, x_{i2}, \dots, x_{ip}, Y_{i1}, \dots, Y_{im})$ . The constant  $x_{i1} = 1$  is in the model, and is often omitted from the case and the data matrix. The model is written in matrix form as  $\mathbf{Z} = \mathbf{X}\mathbf{B} + \mathbf{E}$ .

The model has  $E(\boldsymbol{\epsilon}_k) = \mathbf{0}$  and  $\text{Cov}(\boldsymbol{\epsilon}_k) = \boldsymbol{\Sigma}_\epsilon = (\sigma_{ij})$  for  $k = 1, \dots, n$ . Also  $E(\mathbf{e}_i) = \mathbf{0}$  while  $\text{Cov}(\mathbf{e}_i, \mathbf{e}_j) = \sigma_{ij}\mathbf{I}_n$  for  $i, j = 1, \dots, m$ . Then  $\mathbf{B}$  and  $\boldsymbol{\Sigma}_\epsilon$  are unknown matrices of parameters to be estimated, and  $E(\mathbf{Z}) = \mathbf{X}\mathbf{B}$  while  $E(Y_{ij}) = \mathbf{x}_i^T \boldsymbol{\beta}_j$ .

3) Each response variable in a multivariate linear regression model follows a multiple linear regression model  $\mathbf{Y}_j = \mathbf{X}\boldsymbol{\beta}_j + \mathbf{e}_j$  for  $j = 1, \dots, m$  where it is assumed that  $E(\mathbf{e}_j) = \mathbf{0}$  and  $\text{Cov}(\mathbf{e}_j) = \sigma_{jj}\mathbf{I}_n$ .

4) For each variable  $Y_k$  make a response plot of  $\hat{Y}_{ik}$  versus  $Y_{ik}$  and a residual plot of  $\hat{Y}_{ik}$  versus  $r_{ik} = Y_{ik} - \hat{Y}_{ik}$ . If the multivariate linear regression model is appropriate, then the plotted points should cluster about the identity line in each of the  $m$  response plots. If outliers are present or if the plot is not linear, then the current model or data need to be transformed or corrected. If the model is good, then each of the  $m$  residual plots should be ellipsoidal with no trend and should be centered about the  $r = 0$  line. There should not be any pattern in the residual plot: as a narrow vertical strip is moved from left to right, the behavior of the residuals within the strip should show little change. Outliers and patterns such as curvature or a fan shaped plot are bad.

5) Make a scatterplot matrix of  $Y_1, \dots, Y_m$  and of the continuous predictors. Use power transformations to remove strong nonlinearities.

6) Consider testing  $\mathbf{L}\mathbf{B} = \mathbf{0}$  where  $\mathbf{L}$  is an  $r \times p$  full rank matrix. Let  $\mathbf{W}_e = \hat{\mathbf{E}}^T \hat{\mathbf{E}}$  and  $\mathbf{W}_e/(n-p) = \hat{\boldsymbol{\Sigma}}_\epsilon$ . Let  $\mathbf{H} = \hat{\mathbf{B}}^T \mathbf{L}^T [\mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T]^{-1} \mathbf{L} \hat{\mathbf{B}}$ . Let  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m$  be the ordered eigenvalues of  $\mathbf{W}_e^{-1} \mathbf{H}$ . Then there are four commonly used test statistics.

The Wilks'  $\Lambda$  statistic is  $\Lambda(\mathbf{L}) = |(\mathbf{H} + \mathbf{W}_e)^{-1} \mathbf{W}_e| = |\mathbf{W}_e^{-1} \mathbf{H} + \mathbf{I}|^{-1} = \prod_{i=1}^m (1 + \lambda_i)^{-1}$ .

The Pillai's trace statistic is  $V(\mathbf{L}) = \text{tr}[(\mathbf{H} + \mathbf{W}_e)^{-1} \mathbf{H}] = \sum_{i=1}^m \frac{\lambda_i}{1 + \lambda_i}$ .

The Hotelling-Lawley trace statistic is  $U(\mathbf{L}) = \text{tr}[\mathbf{W}_e^{-1} \mathbf{H}] = \sum_{i=1}^m \lambda_i$ .

The Roy's maximum root statistic is  $\lambda_{\max}(\mathbf{L}) = \lambda_1$ .

7) **Theorem:** The Hotelling-Lawley trace statistic

$$U(\mathbf{L}) = \frac{1}{n-p} [\text{vec}(\mathbf{L}\hat{\mathbf{B}})]^T [\hat{\boldsymbol{\Sigma}}_\epsilon^{-1} \otimes (\mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T)^{-1}] [\text{vec}(\mathbf{L}\hat{\mathbf{B}})].$$

8) **Assumption D1:** Let  $h_i$  be the  $i$ th diagonal element of  $\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ . Assume  $\max(h_1, \dots, h_n) \xrightarrow{P} 0$  as  $n \rightarrow \infty$ , assume that the zero mean iid error vectors have finite fourth moments, and assume that  $\frac{1}{n} \mathbf{X}^T \mathbf{X} \xrightarrow{P} \mathbf{W}^{-1}$ .

9) **Multivariate Least Squares Central Limit Theorem (MLS CLT):** For the least squares estimator, if assumption D1 holds, then  $\hat{\boldsymbol{\Sigma}}_\epsilon$  is a  $\sqrt{n}$  consistent estimator of  $\boldsymbol{\Sigma}_\epsilon$ , and  $\sqrt{n} \text{vec}(\hat{\mathbf{B}} - \mathbf{B}) \xrightarrow{D} N_{pm}(\mathbf{0}, \boldsymbol{\Sigma}_\epsilon \otimes \mathbf{W})$ .

10) **Theorem:** If assumption D1 holds and if  $H_0$  is true, then

$$(n-p)U(\mathbf{L}) \xrightarrow{D} \chi_{rm}^2.$$

11) Under regularity conditions,  $-[n-p+1-0.5(m-r+3)] \log(\Lambda(\mathbf{L})) \xrightarrow{D} \chi_{rm}^2$ ,  $(n-p)V(\mathbf{L}) \xrightarrow{D} \chi_{rm}^2$ , and  $(n-p)U(\mathbf{L}) \xrightarrow{D} \chi_{rm}^2$ .

These statistics are robust against nonnormality.

12) For the Wilks' Lambda test,

$$pval = P\left(\frac{-[n-p+1-0.5(m-r+3)]}{rm} \log(\Lambda(\mathbf{L})) < F_{rm, n-rm}\right).$$

$$\text{For the Pillai's trace test, } pval = P\left(\frac{n-p}{rm} V(\mathbf{L}) < F_{rm, n-rm}\right).$$

$$\text{For the Hotelling Lawley trace test, } pval = P\left(\frac{n-p}{rm} U(\mathbf{L}) < F_{rm, n-rm}\right).$$

The above three tests are large sample tests,  $P(\text{reject } H_0 | H_0 \text{ is true}) \rightarrow \delta$  as  $n \rightarrow \infty$ , under regularity conditions.

13) The 4 step MANOVA  $F$  test of hypotheses uses  $\mathbf{L} = [\mathbf{0} \ \mathbf{I}_{p-1}]$ .

i) State the hypotheses  $H_0$ : the nontrivial predictors are not needed in the mreg model  $H_1$ : at least one of the nontrivial predictors is needed.

ii) Find the test statistic  $F_o$  from output.

iii) Find the pval from output.

iv) If  $pval \leq \delta$ , reject  $H_0$ . If  $pval > \delta$ , fail to reject  $H_0$ . If  $H_0$  is rejected, conclude that there is a mreg relationship between the response variables  $Y_1, \dots, Y_m$  and the predictors  $x_2, \dots, x_p$ . If you fail to reject  $H_0$ , conclude that there is a not a mreg relationship between  $Y_1, \dots, Y_m$  and the predictors  $x_2, \dots, x_p$ . (Get the variable names from the story problem.)

14) The 4 step  $F_j$  test of hypotheses uses  $\mathbf{L}_j = [0, \dots, 0, 1, 0, \dots, 0]$  where the 1 is in the  $j$ th position. Let  $\mathbf{B}_j^T$  be the  $j$ th row of  $\mathbf{B}$ . The hypotheses are equivalent to  $H_0: \mathbf{B}_j^T = \mathbf{0}$   $H_1: \mathbf{B}_j^T \neq \mathbf{0}$ . i) State the hypotheses  $H_0$ :  $x_j$  is not needed in the model  $H_1$ :  $x_j$  is needed in the model.

ii) Find the test statistic  $F_j$  from output.

iii) Find pval from output.

iv) If  $pval \leq \delta$ , reject  $H_0$ . If  $pval > \delta$ , fail to reject  $H_0$ . Give a nontechnical sentence restating your conclusion in terms of the story problem. If  $H_0$  is rejected, then conclude that  $x_j$  is needed in the mreg model for  $Y_1, \dots, Y_m$ . If you fail to reject  $H_0$ , then conclude that  $x_j$  is not needed in the mreg model for  $Y_1, \dots, Y_m$  given that the other predictors are in the model.

15) The 4 step **MANOVA partial F test** of hypotheses has a full model using all of the variables and a reduced model where  $r$  of the variables are deleted. The  $i$ th row of  $\mathbf{L}$  has a 1 in the position corresponding to the  $i$ th variable to be deleted. Omitting the  $j$ th variable corresponds to the  $F_j$  test while omitting variables  $x_2, \dots, x_p$  corresponds to the MANOVA  $F$  test.

i) State the hypotheses  $H_0$ : the reduced model is good

$H_1$ : use the full model.

ii) Find the test statistic  $F_R$  from output.

iii) Find the pval from output.

iv) If  $pval \leq \delta$ , reject  $H_0$  and conclude that the full model should be used. If  $pval > \delta$ , fail to reject  $H_0$  and conclude that the reduced model is good.

16) The 4 step MANOVA  $F$  test should reject  $H_0$  if the response and residual plots look good,  $n$  is large enough, and at least one response plot does not look like the corresponding residual plot. A response plot for  $Y_j$  will look like a residual plot if the identity line appears almost horizontal, hence the range of  $\hat{Y}_j$  is small.

17) The *linmodpack* function `mltreg` produces the  $m$  response and residual plots, gives  $\hat{\mathbf{B}}$ ,  $\hat{\Sigma}_\epsilon$ , the MANOVA partial  $F$  test statistic and pval corresponding to the reduced model that leaves out the variables given by indices (so  $x_2$  and  $x_4$  in the output below with  $F = 0.77$  and  $pval = 0.614$ ),  $F_j$  and the pval for the  $F_j$  test for variables 1, 2, ...,  $p$  (where  $p = 4$  in the output below so  $F_2 = 1.51$  with  $pval = 0.284$ ), and  $F_0$  and pval for the MANOVA  $F$  test (in the output below  $F_0 = 3.15$  and  $pval = 0.06$ ). The command `out <- mltreg(x,y,indices=c(2))` would produce a MANOVA partial  $F$  test corresponding to the  $F_2$  test while the command `out <- mltreg(x,y,indices=c(2,3,4))` would produce a MANOVA partial  $F$  test corresponding to the MANOVA  $F$  test for a data set with  $p = 4$  predictor variables. The Hotelling Lawley trace statistic is used in the tests.

```

out <- mltreg(x,y,indices=c(2,4))
$Bhat      [,1]      [,2]      [,3]
[1,] 47.96841291 623.2817463 179.8867890
[2,]  0.07884384  0.7276600 -0.5378649
[3,] -1.45584256 -17.3872206  0.2337900
[4,] -0.01895002  0.1393189 -0.3885967
$Covhat
      [,1]      [,2]      [,3]
[1,] 21.91591 123.2557 132.339
[2,] 123.25566 2619.4996 2145.780
[3,] 132.33902 2145.7797 2954.082
$partial
      partialF      Pval
[1,] 0.7703294 0.6141573
$Ftable
      Fj      pvals
[1,] 6.30355375 0.01677169
[2,] 1.51013090 0.28449166
[3,] 5.61329324 0.02279833
[4,] 0.06482555 0.97701447
$MANOVA
      MANOVAF      pval
[1,] 3.150118 0.06038742

```

18) Given  $\hat{\mathbf{B}} = [\hat{\beta}_1 \ \hat{\beta}_2 \ \cdots \ \hat{\beta}_m]$  and  $\mathbf{x}_f$ , find  $\hat{\mathbf{y}}_f = (\hat{y}_1, \dots, \hat{y}_m)^T$  where  $\hat{y}_i = \hat{\beta}_i^T \mathbf{x}_f$ .

19)  $\hat{\Sigma}_\epsilon = \frac{\hat{\mathbf{E}}^T \hat{\mathbf{E}}}{n-p} = \frac{1}{n-p} \sum_{i=1}^n \hat{\epsilon}_i \hat{\epsilon}_i^T$  while the sample covariance matrix of

the residuals is  $\mathbf{S}_r = \frac{n-p}{n-1} \hat{\Sigma}_\epsilon = \frac{\hat{\mathbf{E}}^T \hat{\mathbf{E}}}{n-1}$ . Both  $\hat{\Sigma}_\epsilon$  and  $\mathbf{S}_r$  are  $\sqrt{n}$  consistent estimators of  $\Sigma_\epsilon$  for a large class of distributions for the error vectors  $\epsilon_i$ .

20) The  $100(1-\delta)\%$  nonparametric prediction region for  $\mathbf{y}_f$  given  $\mathbf{x}_f$  is the nonparametric prediction region from § 2.2 applied to  $\hat{\mathbf{z}}_i = \hat{\mathbf{y}}_f + \hat{\epsilon}_i = \hat{\mathbf{B}}^T \mathbf{x}_f + \hat{\epsilon}_i$  for  $i = 1, \dots, n$ . This takes the data cloud of the  $n$  residual vectors  $\hat{\epsilon}_i$  and centers the cloud at  $\hat{\mathbf{y}}_f$ . Let

$$D_i^2(\hat{\mathbf{y}}_f, \mathbf{S}_r) = (\hat{\mathbf{z}}_i - \hat{\mathbf{y}}_f)^T \mathbf{S}_r^{-1} (\hat{\mathbf{z}}_i - \hat{\mathbf{y}}_f)$$

for  $i = 1, \dots, n$ . Let  $q_n = \min(1 - \delta + 0.05, 1 - \delta + m/n)$  for  $\delta > 0.1$  and

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

If  $q_n < 1 - \delta + 0.001$ , set  $q_n = 1 - \delta$ . Let  $0 < \delta < 1$  and  $h = D_{(U_n)}$  where  $D_{(U_n)}$  is the  $q_n$ th sample quantile of the  $D_i$ . The  $100(1-\delta)\%$  nonparametric prediction region for  $\mathbf{y}_f$  is

$$\{\mathbf{y} : (\mathbf{y} - \hat{\mathbf{y}}_f)^T \mathbf{S}_r^{-1} (\mathbf{y} - \hat{\mathbf{y}}_f) \leq D_{(U_n)}^2\} = \{\mathbf{y} : D_{\mathbf{y}}(\hat{\mathbf{y}}_f, \mathbf{S}_r) \leq D_{(U_n)}\}.$$

a) Consider the  $n$  prediction regions for the data where  $(\mathbf{y}_{f,i}, \mathbf{x}_{f,i}) = (\mathbf{y}_i, \mathbf{x}_i)$  for  $i = 1, \dots, n$ . If the order statistic  $D_{(U_n)}$  is unique, then  $U_n$  of the  $n$  prediction regions contain  $\mathbf{y}_i$  where  $U_n/n \rightarrow 1 - \delta$  as  $n \rightarrow \infty$ .

b) If  $(\hat{\mathbf{y}}_f, \mathbf{S}_r)$  is a consistent estimator of  $(E(\mathbf{y}_f), \Sigma_\epsilon)$  then the nonparametric prediction region is a large sample  $100(1-\delta)\%$  prediction region for  $\mathbf{y}_f$ .

c) If  $(\hat{\mathbf{y}}_f, \mathbf{S}_r)$  is a consistent estimator of  $(E(\mathbf{y}_f), \Sigma_\epsilon)$ , and the  $\epsilon_i$  come from an elliptically contoured distribution such that the unique highest density region is  $\{\mathbf{y} : D_{\mathbf{y}}(\mathbf{0}, \Sigma_\epsilon) \leq D_{1-\delta}\}$ , then the nonparametric prediction region is asymptotically optimal.

21) On the DD plot for the residual vectors, the cases to the left of the vertical line correspond to cases that would have  $\mathbf{y}_f = \mathbf{y}_i$  in the nonparametric prediction region if  $\mathbf{x}_f = \mathbf{x}_i$ , while the cases to the right of the line would not have  $\mathbf{y}_f = \mathbf{y}_i$  in the nonparametric prediction region.

22) The DD plot for the residual vectors is interpreted almost exactly as a DD plot for iid multivariate data is interpreted. Plotted points clustering about the identity line suggests that the  $\epsilon_i$  may be iid from a multivariate normal distribution, while plotted points that cluster about a line through the origin with slope greater than 1 suggests that the  $\epsilon_i$  may be iid from an

elliptically contoured distribution that is not MVN. Points to the left of the vertical line corresponds to the cases that are in their nonparametric prediction region. Robust distances have not been shown to be consistent estimators of the population distances, but are useful for a graphical diagnostic.

23)	Multiple Linear Regression	Multivariate Linear Regression
	$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$	$\mathbf{Z} = \mathbf{X}\mathbf{B} + \mathbf{E}$
1)	$E(\mathbf{Y}) = \mathbf{X}\boldsymbol{\beta}$	$E[\mathbf{Z}] = \mathbf{X}\mathbf{B}$
2)	$Y_i = \mathbf{x}_i^T \boldsymbol{\beta} + e_i$	$\mathbf{y}_i = \mathbf{B}^T \mathbf{x}_i + \boldsymbol{\epsilon}_i$
3)	$E(\mathbf{e}) = \mathbf{0}$	$E[\mathbf{E}] = \mathbf{0}$
4)	$\mathbf{H} = \mathbf{P} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$	$\mathbf{H} = \mathbf{P} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$
5)	$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$	$\hat{\mathbf{B}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Z}$
6)	$\hat{\mathbf{Y}} = \mathbf{P}\mathbf{Y}$	$\hat{\mathbf{Z}} = \mathbf{P}\mathbf{Z}$
7)	$\mathbf{r} = \hat{\mathbf{e}} = (\mathbf{I} - \mathbf{P})\mathbf{Y}$	$\hat{\mathbf{E}} = (\mathbf{I} - \mathbf{P})\mathbf{Z}$
8)	$E[\hat{\boldsymbol{\beta}}] = \boldsymbol{\beta}$	$E[\hat{\mathbf{B}}] = \mathbf{B}$
9)	$E(\hat{\mathbf{Y}}) = E(\mathbf{Y}) = \mathbf{X}\boldsymbol{\beta}$	$E[\hat{\mathbf{Z}}] = \mathbf{X}\mathbf{B}$
10)	$\hat{\sigma}^2 = \frac{\mathbf{r}^T \mathbf{r}}{n-p}$	$\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon}} = \frac{\hat{\mathbf{E}}^T \hat{\mathbf{E}}}{n-p}$
11)	$V(e_i) = \sigma^2$	$\text{Cov}(\boldsymbol{\epsilon}_i) = \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}}$
12)	$E(Y_i) = \boldsymbol{\beta}^T \mathbf{x}_i$	$E[\mathbf{y}_i] = \mathbf{B}^T \mathbf{x}_i$
13)	$H_0 : \mathbf{L}\boldsymbol{\beta} = \mathbf{0}$ $rF_R \xrightarrow{D} \chi_r^2$	$H_0 : \mathbf{L}\mathbf{B} = \mathbf{0}$ $(n-p)U(\mathbf{L}) \xrightarrow{D} \chi_{rm}^2$
14)	LS CLT $\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \xrightarrow{D} N_p(\mathbf{0}, \sigma^2 \mathbf{W})$	MLS CLT $\sqrt{n} \text{vec}(\hat{\mathbf{B}} - \mathbf{B}) \xrightarrow{D} N_{pm}(\mathbf{0}, \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}} \otimes \mathbf{W})$

23) The table on the previous page compares MLR and MREG.

24) The robust multivariate linear regression method `rmreg2` computes the classical estimator on the RMVN set where RMVN is computed from the  $n$  cases  $\mathbf{v}_i = (x_{i2}, \dots, x_{pi}, Y_{i1}, \dots, Y_{im})^T$ . This estimator has considerable outlier resistance but theory currently needs very strong assumptions. The response and residual plots and DD plot of the residuals from this estimator are useful for outlier detection. The `rmreg2` estimator is superior to the `rmreg` estimator for outlier detection.

## 10.11 Complements

This chapter followed Olive (2017b, ch. 12) closely. Multivariate linear regression is a semiparametric method that is nearly as easy to use as multiple linear regression if  $m$  is small. Section 10.3 followed Olive (2018) closely. The material on plots and testing followed Olive et al. (2015) closely. The  $m$  response and residual plots should be made as well as the DD plot, and the response and residual plots are very useful for the  $m = 1$  case of multiple linear regression and experimental design. These plots speed up the model building process for multivariate linear models since the success of power transformations achieving linearity can be quickly assessed, and influential cases can be quickly detected. See Cook and Olive (2001).

Work is needed on variable selection and on determining the sample sizes for when the tests and prediction regions start to work well. Response and residual plots can look good for  $n \geq 10p$ , but for testing and prediction regions, we may need  $n \geq a(m+p)^2$  where  $0.8 \leq a \leq 5$  even for well behaved elliptically contoured error distributions. Variable selection for multivariate linear regression is discussed in Fujikoshi et al. (2014).  $R$  programs are needed to make variable selection easy. Forward selection would be especially useful.

Often observations  $(Y_1, \dots, Y_m, x_2, \dots, x_p)$  are collected on the same person or thing and hence are correlated. If transformations can be found such that the DD plot and the  $m$  response plots and residual plots look good, and  $n$  is large ( $n \geq \max[(m+p)^2, mp+30]$  starts to give good results), then multivariate linear regression can be used to efficiently analyze the data. Examining  $m$  multiple linear regressions is an incorrect method for analyzing the data.

In addition to robust estimators and seemingly unrelated regressions, envelope estimators and partial least squares (PLS) are competing methods for multivariate linear regression. See recent work by Cook such as Cook (2018), Cook and Su (2013), Cook et al. (2013), and Su and Cook (2012). Methods like ridge regression and lasso can also be extended to multivariate linear regression. See, for example, Obozinski et al. (2011). Relaxed lasso extensions are likely useful. Prediction regions for alternative methods with  $n \gg p$  could be made following Section 10.3.

Plugging in robust dispersion estimators in place of the covariance matrices, as done in Section 10.6, is not a new idea. Maronna and Morgenthaler (1986) used  $M$ -estimators when  $m = 1$ . Problems can occur if the error distribution is not elliptically contoured. See Nordhausen and Tyler (2015).

Khattree and Naik (1999, pp. 91-98) discussed testing  $H_0 : \mathbf{LBM} = \mathbf{0}$  versus  $H_1 : \mathbf{LBM} \neq \mathbf{0}$  where  $\mathbf{M} = \mathbf{I}$  gives a linear test of hypotheses. Johnstone and Nadler (2017) gave useful approximations for Roy's largest root test when the error vector distribution is multivariate normal.

## 10.12 Problems

**PROBLEMS WITH AN ASTERISK \* ARE ESPECIALLY USEFUL.**

**10.1\***. Consider the Hotelling Lawley test statistic. Let

$$T(\mathbf{W}) = n [\text{vec}(\mathbf{L}\hat{\mathbf{B}})]^T [\hat{\Sigma}_{\epsilon}^{-1} \otimes (\mathbf{L}\mathbf{W}\mathbf{L}^T)^{-1}] [\text{vec}(\mathbf{L}\hat{\mathbf{B}})].$$

Let

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

Show  $T(\hat{\mathbf{W}}) = [\text{vec}(\mathbf{L}\hat{\mathbf{B}})]^T [\hat{\Sigma}_{\epsilon}^{-1} \otimes (\mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T)^{-1}] [\text{vec}(\mathbf{L}\hat{\mathbf{B}})]$ .

**10.2.** Consider the Hotelling Lawley test statistic. Let  $T =$

$$[\text{vec}(\mathbf{L}\hat{\mathbf{B}})]^T [\hat{\Sigma}_{\epsilon}^{-1} \otimes (\mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T)^{-1}] [\text{vec}(\mathbf{L}\hat{\mathbf{B}})].$$

Let  $\mathbf{L} = \mathbf{L}_j = [0, \dots, 0, 1, 0, \dots, 0]$  have a 1 in the  $j$ th position. Let  $\hat{\mathbf{b}}_j^T = \mathbf{L}\hat{\mathbf{B}}$  be the  $j$ th row of  $\hat{\mathbf{B}}$ . Let  $d_j = \mathbf{L}_j(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}_j^T = (\mathbf{X}^T \mathbf{X})_{jj}^{-1}$ , the  $j$ th diagonal entry of  $(\mathbf{X}^T \mathbf{X})^{-1}$ . Then  $T_j = \frac{1}{d_j} \hat{\mathbf{b}}_j^T \hat{\Sigma}_{\epsilon}^{-1} \hat{\mathbf{b}}_j$ . The Hotelling Lawley statistic

$$U = \text{tr}([(n-p)\hat{\Sigma}_{\epsilon}]^{-1} \hat{\mathbf{B}}^T \mathbf{L}^T [\mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T]^{-1} \mathbf{L}\hat{\mathbf{B}}).$$

Hence if  $\mathbf{L} = \mathbf{L}_j$ , then  $U_j = \frac{1}{d_j(n-p)} \text{tr}(\hat{\Sigma}_{\epsilon}^{-1} \hat{\mathbf{b}}_j \hat{\mathbf{b}}_j^T)$ .

Using  $\text{tr}(\mathbf{ABC}) = \text{tr}(\mathbf{CAB})$  and  $\text{tr}(a) = a$  for scalar  $a$ , show that  $(n-p)U_j = T_j$ .

**10.3.** Consider the Hotelling Lawley test statistic. Using the Searle (1982, p. 333) identity

$$\text{tr}(\mathbf{AG}^T \mathbf{DGC}) = [\text{vec}(\mathbf{G})]^T [\mathbf{CA} \otimes \mathbf{D}^T] [\text{vec}(\mathbf{G})],$$

show  $(n-p)U(\mathbf{L}) = \text{tr}[\hat{\Sigma}_{\epsilon}^{-1} \hat{\mathbf{B}}^T \mathbf{L}^T [\mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T]^{-1} \mathbf{L} \hat{\mathbf{B}}]$   
 $= [\text{vec}(\mathbf{L} \hat{\mathbf{B}})]^T [\hat{\Sigma}_{\epsilon}^{-1} \otimes (\mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T)^{-1}] [\text{vec}(\mathbf{L} \hat{\mathbf{B}})]$  by identifying  $\mathbf{A}$ ,  $\mathbf{G}$ ,  $\mathbf{D}$ ,  
and  $\mathbf{C}$ .

```
$Ftable      Fj          pvals  #Output for problem 10.4.
[1, ] 82.147221 0.000000e+00
[2, ] 58.448961 0.000000e+00
[3, ] 15.700326 4.258563e-09
[4, ]  9.072358 1.281220e-05
[5, ] 45.364862 0.000000e+00
```

```
$MANOVA
      MANOVAF pval
[1, ] 67.80145    0
```

**10.4.** The output above is for the *R* Seatbelts data set where  $Y_1 = \text{drivers}$  = number of drivers killed or seriously injured,  $Y_2 = \text{front}$  = number of front seat passengers killed or seriously injured, and  $Y_3 = \text{back}$  = number of back seat passengers killed or seriously injured. The predictors were  $x_2 = \text{kms}$  = distance driven,  $x_3 = \text{price}$  = petrol price,  $x_4 = \text{van}$  = number of van drivers killed, and  $x_5 = \text{law}$  = 0 if the law was in effect that month and 1 otherwise. The data consists of 192 monthly totals in Great Britain from January 1969 to December 1984, and the compulsory wearing of seat belts law was introduced in February 1983.

- Do the MANOVA  $F$  test.
- Do the  $F_4$  test.

**10.5.** a) Sketch a DD plot of the residual vectors  $\hat{\epsilon}_i$  for the multivariate linear regression model if the error vectors  $\epsilon_i$  are iid from a multivariate normal distribution. b) Does the DD plot change if the one way MANOVA model is used instead of the multivariate linear regression model?

**10.6.** The output below is for the *R* judge ratings data set consisting of lawyer ratings for  $n = 43$  judges.  $Y_1 = \text{oral}$  = sound oral rulings,  $Y_2 = \text{writ}$  = sound written rulings, and  $Y_3 = \text{rten}$  = worthy of retention. The predictors were  $x_2 = \text{cont}$  = number of contacts of lawyer with judge,  $x_3 = \text{intg}$  = judicial integrity,  $x_4 = \text{dmnr}$  = demeanor,  $x_5 = \text{dilig}$  = diligence,  $x_6 = \text{cfmg}$  = case flow managing,  $x_7 = \text{deci}$  = prompt decisions,  $x_8 = \text{prep}$  = preparation for trial,  $x_9 = \text{fami}$  = familiarity with law, and  $x_{10} = \text{phys}$  = physical ability.

- Do the MANOVA  $F$  test.
- Do the MANOVA partial  $F$  test for the reduced model that deletes  $x_2, x_5, x_6, x_7$ , and  $x_8$ .

```
y<-USJudgeRatings[,c(9,10,12)] #See problem 8.6.
```

```

x<-USJudgeRatings[, -c(9, 10, 12)]
mltreg(x, y, indices=c(2, 5, 6, 7, 8))
$partial
      partialF      Pval
[1,] 1.649415 0.1855314

$MANOVA
      MANOVAF      pval
[1,] 340.1018 1.121325e-14

```

**10.7.** Let  $\beta_i$  be  $p \times 1$  and suppose

$$\begin{pmatrix} \hat{\beta}_1 - \beta_1 \\ \hat{\beta}_2 - \beta_2 \end{pmatrix} \sim N_{2p} \left( \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \end{pmatrix}, \begin{bmatrix} \sigma_{11}(\mathbf{X}^T \mathbf{X})^{-1} & \sigma_{12}(\mathbf{X}^T \mathbf{X})^{-1} \\ \sigma_{21}(\mathbf{X}^T \mathbf{X})^{-1} & \sigma_{22}(\mathbf{X}^T \mathbf{X})^{-1} \end{bmatrix} \right).$$

Find the distribution of

$$[\mathbf{L} \ \mathbf{0}] \begin{pmatrix} \hat{\beta}_1 - \beta_1 \\ \hat{\beta}_2 - \beta_2 \end{pmatrix} = \mathbf{L} \hat{\beta}_1$$

where  $\mathbf{L} \beta_1 = \mathbf{0}$  and  $\mathbf{L}$  is  $r \times p$  with  $r \leq p$ . Simplify.

**10.8.** Let  $\mathbf{y} = \mathbf{B}^T \mathbf{x} + \epsilon$ . Suppose  $\mathbf{x} = (1, x_2, \dots, x_p)^T = (1 \ \mathbf{w}^T)^T$  where  $\mathbf{w} = (x_2, \dots, x_p)^T$ . Let

$$\mathbf{B} = \begin{pmatrix} \boldsymbol{\alpha}^T \\ \mathbf{B}_S \end{pmatrix}.$$

Suppose

$$\begin{pmatrix} \mathbf{y} \\ \mathbf{w} \end{pmatrix} \sim N_{m+p-1} \left[ \begin{pmatrix} \boldsymbol{\mu}_y \\ \boldsymbol{\mu}_w \end{pmatrix}, \begin{pmatrix} \boldsymbol{\Sigma}_{yy} & \boldsymbol{\Sigma}_{yw} \\ \boldsymbol{\Sigma}_{wy} & \boldsymbol{\Sigma}_{ww} \end{pmatrix} \right].$$

Then  $\mathbf{y}|\mathbf{w} \sim N_m(\boldsymbol{\mu}_y + \boldsymbol{\Sigma}_{yw} \boldsymbol{\Sigma}_{ww}^{-1}(\mathbf{w} - \boldsymbol{\mu}_w), \boldsymbol{\Sigma}_{yy} - \boldsymbol{\Sigma}_{yw} \boldsymbol{\Sigma}_{ww}^{-1} \boldsymbol{\Sigma}_{wy})$ , and  $\epsilon \sim N_m(\mathbf{0}, \boldsymbol{\Sigma}_{yy} - \boldsymbol{\Sigma}_{yw} \boldsymbol{\Sigma}_{ww}^{-1} \boldsymbol{\Sigma}_{wy}) = N_m(\mathbf{0}, \boldsymbol{\Sigma}_\epsilon)$ .

Now

$$\mathbf{y}|\mathbf{x} = \mathbf{y} \begin{pmatrix} 1 \\ \mathbf{w} \end{pmatrix} = \mathbf{B}^T \mathbf{x} + \epsilon,$$

and

$$\mathbf{y}|\mathbf{w} = \mathbf{B}^T \mathbf{x} + \epsilon = \begin{pmatrix} \boldsymbol{\alpha}^T \\ \mathbf{B}_S \end{pmatrix}^T \begin{pmatrix} 1 \\ \mathbf{w} \end{pmatrix} + \epsilon = (\boldsymbol{\alpha} \ \mathbf{B}_S^T) \begin{pmatrix} 1 \\ \mathbf{w} \end{pmatrix} + \epsilon = \boldsymbol{\alpha} + \mathbf{B}_S^T \mathbf{w} + \epsilon.$$

Hence  $E(\mathbf{y}|\mathbf{w}) = \boldsymbol{\mu}_y + \boldsymbol{\Sigma}_{yw} \boldsymbol{\Sigma}_{ww}^{-1}(\mathbf{w} - \boldsymbol{\mu}_w) = \boldsymbol{\alpha} + \mathbf{B}_S^T \mathbf{w}$ .

a) Show  $\boldsymbol{\alpha} = \boldsymbol{\mu}_y - \mathbf{B}_S^T \boldsymbol{\mu}_w$ .

b) Show  $\mathbf{B}_S = \boldsymbol{\Sigma}_w^{-1} \boldsymbol{\Sigma}_{wy}$  where  $\boldsymbol{\Sigma}_w = \boldsymbol{\Sigma}_{ww}$ .

(Hence  $\mathbf{B}_S^T = \boldsymbol{\Sigma}_{yw} \boldsymbol{\Sigma}_w^{-1}$ .)

**R Problems**

**Warning: Use the command `source("G:/linmodpack.txt")` to download the programs. See Preface or Section 11.1.** Typing the name of the `mpack` function, e.g. `ddplot`, will display the code for the function. Use the `args` command, e.g. `args(ddplot)`, 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/linmodrhw.txt>) into *R*.

**10.9.** This problem examines multivariate linear regression on the Cook and Weisberg (1999) mussels data with  $Y_1 = \log(S)$  and  $Y_2 = \log(M)$  where  $S$  is the shell mass and  $M$  is the muscle mass. The predictors are  $X_2 = L$ ,  $X_3 = \log(W)$ , and  $X_4 = H$ : the shell length,  $\log(\text{width})$ , and height.

a) The *R* command for this part makes the response and residual plots for each of the two response variables. Click the rightmost mouse button and highlight *Stop* to advance the plot. When you have the response and residual plots for one variable on the screen, copy and paste the two plots into *Word*. Do this two times, once for each response variable. The plotted points fall in roughly evenly populated bands about the identity or  $r = 0$  line.

b) Copy and paste the output produced from the *R* command for this part from \$partial on. This gives the output needed to do the MANOVA  $F$  test, MANOVA partial  $F$  test, and the  $F_j$  tests.

c) The *R* command for this part makes a DD plot of the residual vectors and adds the lines corresponding to those in Figure 10.3. Place the plot in *Word*. Do the residual vectors appear to follow a multivariate normal distribution? (Right click *Stop* once.)

d) Do the MANOVA partial  $F$  test where the reduced model deletes  $X_3$  and  $X_4$ .

e) Do the  $F_2$  test.

f) Do the MANOVA  $F$  test.

**10.10.** This problem examines multivariate linear regression on the SAS Institute (1985, p. 146) Fitness Club Data with  $Y_1 = \text{chinups}$ ,  $Y_2 = \text{situps}$ , and  $Y_3 = \text{jumps}$ . The predictors are  $X_2 = \text{weight}$ ,  $X_3 = \text{waist}$ , and  $X_4 = \text{pulse}$ .

a) The *R* command for this part makes the response and residual plots for each of the three variables. Click the rightmost mouse button and highlight *Stop* to advance the plot. When you have the response and residual plots for one variable on the screen, copy and paste the three plots into *Word*. Do this three times, once for each response variable. Are there any outliers?

b) The *R* command for this part makes a DD plot of the residual vectors and adds the lines corresponding to those in Figure 10.3. Place the plot in *Word*. Are there any outliers? (Right click *Stop* once.)

**10.11.** This problem uses the *linmodpack* function `mregsim` to simulate the Wilks'  $A$  test, Pillai's trace test, Hotelling Lawley trace test, and Roy's largest root test for the  $F_j$  tests and the MANOVA  $F$  test for multivariate linear regression. When `mnull = T` the first row of  $\mathbf{B}$  is  $\mathbf{1}^T$  while the re-

maining rows are equal to  $\mathbf{0}^T$ . Hence the null hypothesis for the MANOVA  $F$  test is true. When `mnull = F` the null hypothesis is true for  $p = 2$ , but false for  $p > 2$ . Now the first row of  $\mathbf{B}$  is  $\mathbf{1}^T$  and the last row of  $\mathbf{B}$  is  $\mathbf{0}^T$ . If  $p > 2$ , then the second to last row of  $\mathbf{B}$  is  $(1, 0, \dots, 0)$ , the third to last row is  $(1, 1, 0, \dots, 0)$  et cetera as long as the first row is not changed from  $\mathbf{1}^T$ . First  $m$  iid errors  $\mathbf{z}_i$  are generated such that the  $m$  errors are iid with variance  $\sigma^2$ . Then  $\boldsymbol{\epsilon}_i = \mathbf{A}\mathbf{z}_i$  so that  $\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon}} = \sigma^2 \mathbf{A}\mathbf{A}^T = (\sigma_{ij})$  where the diagonal entries  $\sigma_{ii} = \sigma^2[1 + (m-1)\psi^2]$  and the off diagonal entries  $\sigma_{ij} = \sigma^2[2\psi + (m-2)\psi^2]$  where  $\psi = 0.10$ . Terms like `Wilkcov` give the percentage of times the Wilks' test rejected the  $F_1, F_2, \dots, F_p$  tests. The `$mancv wcv pcv hlcov rcv fcov` output gives the percentage of times the 4 test statistics reject the MANOVA  $F$  test. Here `hlcov` and `fcov` both correspond to the Hotelling Lawley test using the formulas in Problem 10.3.

5000 runs will be used so the simulation may take several minutes. Sample sizes  $n = (m+p)^2$ ,  $n = 3(m+p)^2$ , and  $n = 4(m+p)^2$  were interesting. We want coverage near 0.05 when  $H_0$  is true and coverage close to 1 for good power when  $H_0$  is false. Multivariate normal errors were used in a) and b) below.

a) Copy the coverage parts of the output produced by the `R` commands for this part where  $n = 20, m = 2$ , and  $p = 4$ . Here  $H_0$  is true except for the  $F_1$  test. Wilks' and Pillai's tests had low coverage  $< 0.05$  when  $H_0$  was false. Roy's test was good for the  $F_j$  tests, but why was Roy's test bad for the MANOVA  $F$  test?

b) Copy the coverage parts of the output produced by the `R` commands for this part where  $n = 20, m = 2$ , and  $p = 4$ . Here  $H_0$  is false except for the  $F_4$  test. Which two tests seem to be the best for this part?

**10.12.** This problem uses the `linmodpack` function `mpredsim` to simulate the prediction regions for  $\mathbf{y}_f$  given  $\mathbf{x}_f$  for multivariate regression. With 5000 runs this simulation may take several minutes. The `R` command for this problem generates iid lognormal errors then subtracts the mean, producing  $\mathbf{z}_i$ . Then the  $\boldsymbol{\epsilon}_i = \mathbf{A}\mathbf{z}_i$  are generated as in Problem 10.11 with  $n=100, m=2$ , and  $p=4$ . The nominal coverage of the prediction region is 90%, and 92% of the training data is covered. The `ncvr` output gives the coverage of the nonparametric region. What was `ncvr`?

