

# Chapter 10

## Graphical Diagnostics

### 10.1 1D Regression

From Chapter 6, in a **1D regression model**,  $Y$  is conditionally independent of  $\mathbf{x}$  given the **sufficient predictor**  $SP = h(\mathbf{x})$ , written

$$Y \perp\!\!\!\perp \mathbf{x} | SP \text{ or } Y \perp\!\!\!\perp \mathbf{x} | h(\mathbf{x}), \quad (10.1)$$

where the real valued function  $h : \mathbb{R}^p \rightarrow \mathbb{R}$ . The **estimated sufficient predictor**  $ESP = \hat{h}(\mathbf{x})$ .

**Definition 10.1.** A **response plot** is a plot of the ESP versus  $Y$ . A *residual plot* is a plot of the ESP versus the residuals.

A response plot is also called an *estimated sufficient summary plot* (ESSP). A sufficient summary plot is a plot of SP versus  $Y$ . Hence if the ESP is a consistent estimator of the SP, then the response plot estimates the sufficient summary plot.

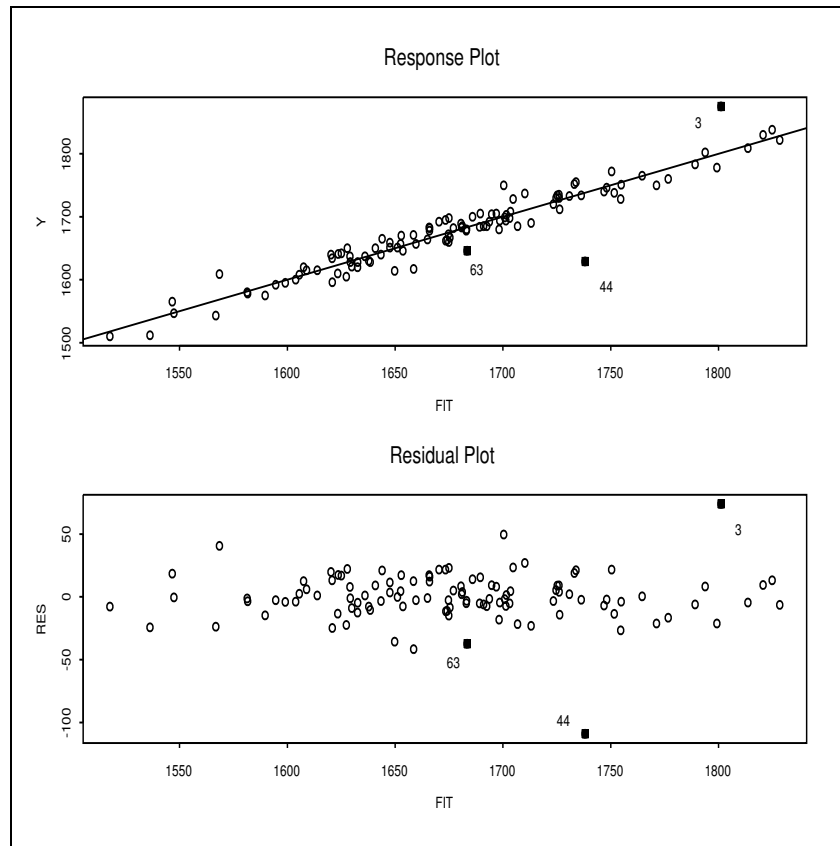
**Notation:** In this text, a plot of  $x$  versus  $Y$  will have  $x$  on the horizontal axis, and  $Y$  on the vertical axis. For the *additive error regression* model  $Y = m(\mathbf{x}) + e$ , the  $i$ th residual is  $r_i = Y_i - \hat{m}(\mathbf{x}_i) = Y_i - \hat{Y}_i$  where  $\hat{Y}_i = \hat{m}(\mathbf{x}_i)$  is the  $i$ th fitted value. The additive error regression model is a 1D regression model with sufficient predictor  $SP = h(\mathbf{x}) = m(\mathbf{x})$ .

For the additive error regression model, the response plot is a plot of  $\hat{Y}$  versus  $Y$  where the *identity line* with unit slope and zero intercept is added as a visual aid. The residual plot is a plot of  $\hat{Y}$  versus  $r$ . Assume the errors  $e_i$  are iid from a unimodal distribution that is not highly skewed. Then the plotted points should scatter about the identity line and the  $r = 0$  line (the horizontal axis) with no other pattern if the fitted model (that produces  $\hat{m}(\mathbf{x})$ ) is good.

## 10.2 Plots for MLR

**Theorem 10.1.** Suppose that the MLR estimator  $\mathbf{b}$  of  $\boldsymbol{\beta}$  is used to find the residuals  $r_i \equiv r_i(\mathbf{b})$  and the fitted values  $\hat{Y}_i \equiv \hat{Y}_i(\mathbf{b}) = \mathbf{x}_i^T \mathbf{b}$ . Then in the response plot of  $\hat{Y}_i$  versus  $Y_i$ , the vertical deviations from the identity line (that has unit slope and zero intercept) are the residuals  $r_i(\mathbf{b})$ .

**Proof.** The identity line in the response plot is  $Y = \mathbf{x}^T \mathbf{b}$ . Hence the vertical deviation is  $Y_i - \mathbf{x}_i^T \mathbf{b} = r_i(\mathbf{b})$ .  $\square$



**Fig. 10.1** Residual and Response Plots for the Tremearne Data

**Example 10.1.** Tremearne (1911) presents a data set of about 17 measurements on 115 people of Hausa nationality. We deleted 3 cases because of missing values and used *height* as the response variable  $Y$ . Along with a constant  $x_{i,1} \equiv 1$ , the five additional predictor variables used were *height*

when sitting, height when kneeling, head length, nasal breadth, and span (perhaps from left hand to right hand). Figure 6.1 presents the (ordinary) least squares (OLS) response and residual plots for this data set. These plots show that an MLR model  $Y = \mathbf{x}^T \boldsymbol{\beta} + e$  should be a useful model for the data since the plotted points in the response plot are linear and follow the identity line while the plotted points in the residual plot follow the  $r = 0$  line with no other pattern (except for a possible outlier marked 44). Note that many important acronyms, such as OLS and MLR, appear in Table 1.1.

To use the response plot to visualize the conditional distribution of  $Y|\mathbf{x}^T \boldsymbol{\beta}$ , use the fact that the fitted values  $\hat{Y} = \mathbf{x}^T \hat{\boldsymbol{\beta}}$ . For example, suppose the height given fit = 1700 is of interest. Mentally examine the plot about a narrow vertical strip about fit = 1700, perhaps from 1685 to 1715. The cases in the narrow strip have a mean close to 1700 since they fall close to the identity line. Similarly, when the fit =  $w$  for  $w$  between 1500 and 1850, the cases have heights near  $w$ , on average.

Cases 3, 44, and 63 are highlighted. The 3rd person was very tall while the 44th person was rather short. Beginners often label too many points as *outliers*: cases that lie far away from the bulk of the data. See Chapter 7. Mentally draw a box about the bulk of the data ignoring any outliers. Double the width of the box (about the identity line for the response plot and about the horizontal line for the residual plot). Cases outside of this imaginary doubled box are potential outliers. Alternatively, visually estimate the standard deviation of the residuals in both plots. In the residual plot look for residuals that are more than 5 standard deviations from the  $r = 0$  line. In Figure 6.1, the standard deviation of the residuals appears to be around 10. Hence cases 3 and 44 are certainly worth examining.

The identity line can also pass through or near an outlier or a cluster of outliers. Then the outliers will be in the upper right or lower left of the response plot, and there will be a large gap between the cluster of outliers and the bulk of the data. Figure 6.1 was made with the following *R* commands, using *lspack* function `MLRplot` and the *major.lsp* data set from the text's webpage.

```
major <- matrix(scan(), nrow=112, ncol=7, byrow=T)
#copy and paste the data set, then press enter
major <- major[,-1]
X<-major[,-6]
Y <- major[,6]
MLRplot(X,Y) #left click the 3 highlighted cases,
#then right click Stop for each of the two plots
```

### 10.2.1 Plots for Variable Selection

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

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

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

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

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

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

**Definition 10.2.**

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

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

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

**Theorem 10.2.** Suppose that a numerical variable selection method suggests several submodels with  $k$  predictors, including a constant, where  $2 \leq k \leq p$ .

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

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

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

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

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

**Remark 10.1.** Consider the model  $I_i$  that deletes the predictor  $x_i$ . Then the model has  $k = p - 1$  predictors including the constant, and the test statistic is  $t_i$  where

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

Using Definition 4.2 and  $C_p(I_{full}) = p$ , it can be shown that

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

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

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

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

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

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

Now  $k$  is the number of terms in the model  $I$  including a constant while  $p - k$  is the number of terms set to 0. As  $k \rightarrow 0$ , the partial  $F$  test will reject  $H_0: \boldsymbol{\beta}_O = \mathbf{0}$  (i.e. say that the full model should be used instead of the submodel  $I$ ) unless  $F_I$  is not much larger than 1. If  $p$  is very large and  $p - k$  is very small, then the partial  $F$  test will tend to suggest that there is a model  $I$  that is about as good as the full model even though model  $I$  deletes  $p - k$  predictors.

**Definition 10.3.** The “fit–fit” or *FF plot* is a plot of  $\hat{Y}_{I,i}$  versus  $\hat{Y}_i$  while a “residual–residual” or *RR plot* is a plot  $r_{I,i}$  versus  $r_i$ . A *response plot* is a plot of  $\hat{Y}_{I,i}$  versus  $Y_i$ . An *EE plot* is a plot of  $\text{ESP}(I)$  versus  $\text{ESP}$ . For MLR, the EE and FF plots are equivalent.

Six graphs will be used to compare the full model and the candidate submodel: the FF plot, RR plot, the response plots from the full and submodel, and the residual plots from the full and submodel. These six plots will contain a great deal of information about the candidate subset provided that Equation (4.1) holds and that a good estimator (such as OLS) for  $\hat{\beta}$  and  $\hat{\beta}_I$  is used.

**Application 10.1.** To visualize whether a candidate submodel using predictors  $\mathbf{x}_I$  is good, use the fitted values and residuals from the submodel and full model to make an RR plot of the  $r_{I,i}$  versus the  $r_i$  and an FF plot of  $\hat{Y}_{I,i}$  versus  $\hat{Y}_i$ . Add the OLS line to the RR plot and identity line to both plots as visual aids. The subset  $I$  is good if the plotted points cluster tightly about the identity line in *both plots*. In particular, the OLS line and the identity line should “nearly coincide” so that it is difficult to tell that the two lines intersect at the origin in the RR plot.

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

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

**Theorem 10.3.** Suppose that every submodel contains a constant and that  $\mathbf{X}$  is a full rank matrix.

**Response Plot:** i) If  $w = \hat{Y}_I$  and  $z = Y$  then the OLS line is the identity line.

ii) If  $w = Y$  and  $z = \hat{Y}_I$  then the OLS line has slope  $b = [\text{corr}(Y, \hat{Y}_I)]^2 = R^2(I)$  and intercept  $a = \bar{Y}(1 - R^2(I))$  where  $\bar{Y} = \sum_{i=1}^n Y_i/n$  and  $R^2(I)$  is the coefficient of multiple determination from the candidate model.

**FF or EE Plot:** iii) If  $w = \hat{Y}_I$  and  $z = \hat{Y}$  then the OLS line is the identity line. Note that  $ESP(I) = \hat{Y}_I$  and  $ESP = \hat{Y}$ .

iv) If  $w = \hat{Y}$  and  $z = \hat{Y}_I$  then the OLS line has slope  $b = [\text{corr}(\hat{Y}, \hat{Y}_I)]^2 = SSR(I)/SSR$  and intercept  $a = \bar{Y}[1 - (SSR(I)/SSR)]$  where  $SSR$  is the regression sum of squares.

**RR Plot:** v) If  $w = r$  and  $z = r_I$  then the OLS line is the identity line.  
vi) If  $w = r_I$  and  $z = r$  then  $a = 0$  and the OLS slope  $b = [\text{corr}(r, r_I)]^2$  and

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

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

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

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

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

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

ii) By (\*), the slope

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

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

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

iv) From iii),

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

Hence

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

and the slope

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

Also the slope

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

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

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

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

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

Hence

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

and the slope

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

Algebra shows that

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

**Remark 10.2.** Let  $I_{min}$  be the model that minimizes  $C_p(I)$  among the models  $I$  generated from the variable selection method such as forward selection. Assuming the full model  $I_p$  is one of the models generated, then  $C_p(I_{min}) \leq C_p(I_p) = p$ , and  $\text{corr}(r, r_{I_{min}}) \rightarrow 1$  as  $n \rightarrow \infty$  by Theorem 4.2 vi). Referring to Equation (4.1), if  $P(S \subseteq I_{min})$  does not go to 1 as  $n \rightarrow \infty$ , then the above correlation would not go to one. Hence  $P(S \subseteq I_{min}) \rightarrow 1$  as  $n \rightarrow \infty$ .



### 10.2.2 Plots for Response Transformations

## 10.3 Plots for GLMs and GAMs

## 10.4 Outlier Detection for the MLD Model

Now suppose the multivariate data has been collected into an  $n \times p$  matrix

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

where the  $i$ th row of  $\mathbf{W}$  is the  $i$ th case  $\mathbf{x}_i^T$  and the  $j$ th column  $\mathbf{v}_j$  of  $\mathbf{W}$  corresponds to  $n$  measurements of the  $j$ th random variable  $X_j$  for  $j = 1, \dots, p$ . Hence the  $n$  rows of the data matrix  $\mathbf{W}$  correspond to the  $n$  cases, while the  $p$  columns correspond to measurements on the  $p$  random variables  $X_1, \dots, X_p$ . For example, the data may consist of  $n$  visitors to a hospital where the  $p = 2$  variables *height* and *weight* of each individual were measured.

**Definition 10.36.** The *coordinatewise median*  $\text{MED}(\mathbf{W}) = (\text{MED}(X_1), \dots, \text{MED}(X_p))^T$  where  $\text{MED}(X_i)$  is the sample median of the data in column  $i$  corresponding to variable  $X_i$  and  $\mathbf{v}_i$ .

**Example 10.11.** Let the data for  $X_1$  be 1, 2, 3, 4, 5, 6, 7, 8, 9 while the data for  $X_2$  is 7, 17, 3, 8, 6, 13, 4, 2, 1. Then  $\text{MED}(\mathbf{W}) = (\text{MED}(X_1), \text{MED}(X_2))^T = (5, 6)^T$ .

**Definition 10.37: Rousseeuw and Van Driessen (1999).** The *DD plot* is a plot of the classical Mahalanobis distances  $\text{MD}_i$  versus robust Mahalanobis distances  $\text{RD}_i$ .

The DD plot is used as a diagnostic for multivariate normality, elliptical symmetry, and for outliers. Assume that the data set consists of iid vectors from an  $EC_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}, g)$  distribution with second moments. See Section 1.7 for notation. Then the classical sample mean and covariance matrix  $(T_M, \mathbf{C}_M) = (\bar{\mathbf{x}}, \mathbf{S})$  is a consistent estimator for  $(\boldsymbol{\mu}, c_{\mathbf{x}}\boldsymbol{\Sigma}) = (E(\mathbf{x}), \text{Cov}(\mathbf{x}))$ . Assume that an alternative algorithm estimator  $(T_A, \mathbf{C}_A)$  is a consistent estimator for  $(\boldsymbol{\mu}, a_A\boldsymbol{\Sigma})$  for some constant  $a_A > 0$ . By scaling the algorithm estimator, the DD plot can be constructed to follow the identity line with unit slope and zero intercept. Let  $(T_R, \mathbf{C}_R) = (T_A, \mathbf{C}_A/\tau^2)$  denote the scaled algorithm estimator where  $\tau > 0$  is a constant to be determined. Notice that  $(T_R, \mathbf{C}_R)$  is a valid estimator of location and dispersion. Hence the robust distances used in the DD plot are given by

$$\text{RD}_i = \text{RD}_i(T_R, C_R) = \sqrt{(\mathbf{x}_i - T_R(\mathbf{W}))^T [C_R(\mathbf{W})]^{-1} (\mathbf{x}_i - T_R(\mathbf{W}))}$$

$= \tau D_i(T_A, C_A)$  for  $i = 1, \dots, n$ .

The following theorem shows that if consistent estimators are used to construct the distances, then the DD plot will tend to cluster tightly about the line segment through  $(0, 0)$  and  $(\text{MD}_{n,\alpha}, \text{RD}_{n,\alpha})$  where  $0 < \alpha < 1$  and  $\text{MD}_{n,\alpha}$  is the  $100\alpha$ th sample percentile of the  $\text{MD}_i$ . Nevertheless, the variability in the DD plot may increase with the distances. Let  $K > 0$  be a constant, e.g. the 99th percentile of the  $\chi_p^2$  distribution.

**Theorem 10.32.** Assume that  $\mathbf{x}_1, \dots, \mathbf{x}_n$  are iid observations from a distribution with parameters  $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  where  $\boldsymbol{\Sigma}$  is a symmetric positive definite matrix. Let  $a_j > 0$  and assume that  $(\hat{\boldsymbol{\mu}}_{j,n}, \hat{\boldsymbol{\Sigma}}_{j,n})$  are consistent estimators of  $(\boldsymbol{\mu}, a_j \boldsymbol{\Sigma})$  for  $j = 1, 2$ .

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

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

$$D_{\mathbf{x}}^2(\hat{\boldsymbol{\mu}}_j, \hat{\boldsymbol{\Sigma}}_j) - \frac{1}{a_j} D_{\mathbf{x}}^2(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = O_p(n^{-\delta}).$$

c) Let  $D_{i,j} \equiv D_i(\hat{\boldsymbol{\mu}}_{j,n}, \hat{\boldsymbol{\Sigma}}_{j,n})$  be the  $i$ th Mahalanobis distance computed from  $(\hat{\boldsymbol{\mu}}_{j,n}, \hat{\boldsymbol{\Sigma}}_{j,n})$ . Consider the cases in the region  $R = \{i | 0 \leq D_{i,j} \leq K, j = 1, 2\}$ . Let  $r_n$  denote the correlation between  $D_{i,1}$  and  $D_{i,2}$  for the cases in  $R$  (thus  $r_n$  is the correlation of the distances in the “lower left corner” of the DD plot). Then  $r_n \rightarrow 1$  in probability as  $n \rightarrow \infty$ .

**Proof.** Let  $B_n$  denote the subset of the sample space on which both  $\hat{\boldsymbol{\Sigma}}_{1,n}$  and  $\hat{\boldsymbol{\Sigma}}_{2,n}$  have inverses. Then  $P(B_n) \rightarrow 1$  as  $n \rightarrow \infty$ .

a) and b):  $D_{\mathbf{x}}^2(\hat{\boldsymbol{\mu}}_j, \hat{\boldsymbol{\Sigma}}_j) = (\mathbf{x} - \hat{\boldsymbol{\mu}}_j)^T \hat{\boldsymbol{\Sigma}}_j^{-1} (\mathbf{x} - \hat{\boldsymbol{\mu}}_j) =$

$$\begin{aligned} & (\mathbf{x} - \hat{\boldsymbol{\mu}}_j)^T \left( \frac{\boldsymbol{\Sigma}^{-1}}{a_j} - \frac{\boldsymbol{\Sigma}^{-1}}{a_j} + \hat{\boldsymbol{\Sigma}}_j^{-1} \right) (\mathbf{x} - \hat{\boldsymbol{\mu}}_j) \\ &= (\mathbf{x} - \hat{\boldsymbol{\mu}}_j)^T \left( \frac{-\boldsymbol{\Sigma}^{-1}}{a_j} + \hat{\boldsymbol{\Sigma}}_j^{-1} \right) (\mathbf{x} - \hat{\boldsymbol{\mu}}_j) + (\mathbf{x} - \hat{\boldsymbol{\mu}}_j)^T \left( \frac{\boldsymbol{\Sigma}^{-1}}{a_j} \right) (\mathbf{x} - \hat{\boldsymbol{\mu}}_j) \\ &= \frac{1}{a_j} (\mathbf{x} - \hat{\boldsymbol{\mu}}_j)^T (-\boldsymbol{\Sigma}^{-1} + a_j \hat{\boldsymbol{\Sigma}}_j^{-1}) (\mathbf{x} - \hat{\boldsymbol{\mu}}_j) + \\ & (\mathbf{x} - \boldsymbol{\mu} + \boldsymbol{\mu} - \hat{\boldsymbol{\mu}}_j)^T \left( \frac{\boldsymbol{\Sigma}^{-1}}{a_j} \right) (\mathbf{x} - \boldsymbol{\mu} + \boldsymbol{\mu} - \hat{\boldsymbol{\mu}}_j) \\ &= \frac{1}{a_j} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \end{aligned}$$

$$\begin{aligned}
& + \frac{2}{a_j} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}_j) + \frac{1}{a_j} (\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}_j)^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}_j) \\
& + \frac{1}{a_j} (\mathbf{x} - \hat{\boldsymbol{\mu}}_j)^T [a_j \hat{\boldsymbol{\Sigma}}_j^{-1} - \boldsymbol{\Sigma}^{-1}] (\mathbf{x} - \hat{\boldsymbol{\mu}}_j) \quad (10.2)
\end{aligned}$$

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

c) Following the proof of a),  $D_j^2 \equiv D_{\hat{\boldsymbol{x}}}^2(\hat{\boldsymbol{\mu}}_j, \hat{\boldsymbol{\Sigma}}_j) \xrightarrow{P} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) / a_j$  for fixed  $\mathbf{x}$ , and the result follows.  $\square$

The above result implies that a plot of the  $\text{MD}_i$  versus the  $D_i(T_A, \mathbf{C}_A) \equiv D_i(A)$  will follow a line through the origin with some positive slope since if  $\mathbf{x} = \boldsymbol{\mu}$ , then both the classical and the algorithm distances should be close to zero. We want to find  $\tau$  such that  $\text{RD}_i = \tau D_i(T_A, \mathbf{C}_A)$  and the DD plot of  $\text{MD}_i$  versus  $\text{RD}_i$  follows the identity line. By Theorem 10.32, the plot of  $\text{MD}_i$  versus  $D_i(A)$  will follow the line segment defined by the origin  $(0, 0)$  and the point of observed median Mahalanobis distances,  $(\text{med}(\text{MD}_i), \text{med}(D_i(A)))$ . This line segment has slope

$$\text{med}(D_i(A)) / \text{med}(\text{MD}_i)$$

which is generally not one. By taking  $\tau = \text{med}(\text{MD}_i) / \text{med}(D_i(A))$ , the plot will follow the identity line if  $(\bar{\mathbf{x}}, \mathbf{S})$  is a consistent estimator of  $(\boldsymbol{\mu}, c_{\mathbf{x}} \boldsymbol{\Sigma})$  and if  $(T_A, \mathbf{C}_A)$  is a consistent estimator of  $(\boldsymbol{\mu}, a_A \boldsymbol{\Sigma})$ . (Using the notation from Theorem 10.32, let  $(a_1, a_2) = (c_{\mathbf{x}}, a_A)$ .) The classical estimator is consistent if the population has a nonsingular covariance matrix. The algorithm estimators  $(T_A, \mathbf{C}_A)$  from Theorem 8.29 are consistent on a large class of EC distributions that have a nonsingular covariance matrix, but tend to be biased for non-EC distributions. We recommend using RFCH or RMVN as the robust estimators in DD plots.

By replacing the observed median  $\text{med}(\text{MD}_i)$  of the classical Mahalanobis distances with the target population analog, say MED,  $\tau$  can be chosen so that the DD plot is *simultaneously* a diagnostic for elliptical symmetry and a diagnostic for the target EC distribution. That is, the plotted points follow the identity line if the data arise from a target EC distribution such as the multivariate normal distribution, but the points follow a line with non-unit slope if the data arise from an alternative EC distribution. In addition the DD plot can often detect departures from elliptical symmetry such as outliers, the presence of two groups, or the presence of a mixture distribution.

**Example 10.12.** We will use the multivariate normal  $N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  distribution as the target. If the data are indeed iid MVN vectors, then the  $(\text{MD}_i)^2$  are asymptotically  $\chi_p^2$  random variables, and  $\text{MED} = \sqrt{\chi_{p,0.5}^2}$  where  $\chi_{p,0.5}^2$  is the median of the  $\chi_p^2$  distribution. Since the target distribution is Gaussian, let

$$RD_i = \frac{\sqrt{\chi_{p,0.5}^2}}{\text{med}(D_i(A))} D_i(A) \quad \text{so that} \quad \tau = \frac{\sqrt{\chi_{p,0.5}^2}}{\text{med}(D_i(A))}. \quad (10.3)$$

Since every nonsingular estimator of multivariate location and dispersion defines a hyperellipsoid, the DD plot can be used to examine which points are in the robust hyperellipsoid

$$\{\mathbf{x} : (\mathbf{x} - T_R)^T C_R^{-1} (\mathbf{x} - T_R) \leq RD_{(h)}^2\} \quad (10.4)$$

where  $RD_{(h)}^2$  is the  $h$ th smallest squared robust Mahalanobis distance, and which points are in a classical hyperellipsoid

$$\{\mathbf{x} : (\mathbf{x} - \bar{\mathbf{x}})^T \mathbf{S}^{-1} (\mathbf{x} - \bar{\mathbf{x}}) \leq MD_{(h)}^2\}. \quad (10.5)$$

In the DD plot, points below  $RD_{(h)}$  correspond to cases that are in the hyperellipsoid given by Equation (10.19) while points to the left of  $MD_{(h)}$  are in a hyperellipsoid determined by Equation (10.20). In particular, we can use the DD plot to examine which points are in the nonparametric prediction region (4.11).

**Application 10.5.** Consider the DD plot with RFCH or RMVN. The DD plot can be used *simultaneously* as a diagnostic for whether the data arise from a multivariate normal distribution or from another EC distribution with nonsingular covariance matrix. EC data will cluster about a straight line through the origin; MVN data in particular will cluster about the identity line. Thus the DD plot can be used to assess the success of numerical transformations towards elliptical symmetry. The DD plot can be used to detect multivariate outliers. Use the DD plot to detect outliers and leverage groups if  $n \geq 10p$  for the predictor variables in regression.

**Fig. 10.2** 4 DD Plots

For this application, the RFCH and RMVN estimators may be best. For MVN data, the  $RD_i$  from the RFCH estimator tend to have a higher correlation with the  $MD_i$  from the classical estimator than the  $RD_i$  from the FCH estimator, and the `cov.mcd` estimator may be inconsistent.

Figure 10.12 shows the DD plots for 3 artificial data sets using `cov.mcd`. The DD plot for 200  $N_3(\mathbf{0}, \mathbf{I}_3)$  points shown in Figure 10.12a resembles the identity line. The DD plot for 200 points from the elliptically contoured distribution  $0.6N_3(\mathbf{0}, \mathbf{I}_3) + 0.4N_3(\mathbf{0}, 25\mathbf{I}_3)$  in Figure 10.12b clusters about a line through the origin with a slope close to 2.0.

A *weighted DD plot* magnifies the lower left corner of the DD plot by omitting the cases with  $RD_i \geq \sqrt{\chi_{p,.975}^2}$ . This technique can magnify features that are obscured when large  $RD_i$ 's are present. If the distribution of  $\mathbf{x}$  is EC with nonsingular  $\Sigma$ , Theorem 8.32 implies that the correlation of the points in the weighted DD plot will tend to one and that the points will cluster about a line passing through the origin. For example, the plotted points in the weighted DD plot (not shown) for the non-MVN EC data of Figure 10.12b are highly correlated and still follow a line through the origin with a slope close to 2.0.

Figures 10.12c and 10.12d illustrate how to use the weighted DD plot. The  $i$ th case in Figure 10.12c is  $(\exp(x_{i,1}), \exp(x_{i,2}), \exp(x_{i,3}))^T$  where  $\mathbf{x}_i$  is the  $i$ th case in Figure 10.12a; i.e. the marginals follow a lognormal distribution. The plot does not resemble the identity line, correctly suggesting that the distribution of the data is not MVN; however, the correlation of the plotted points is rather high. Figure 10.12d is the weighted DD plot where cases with  $RD_i \geq \sqrt{\chi_{3,.975}^2} \approx 3.06$  have been removed. Notice that the correlation of the plotted points is not close to one and that the best fitting line in Figure 10.12d may not pass through the origin. These results suggest that the distribution of  $\mathbf{x}$  is not EC.

**Fig. 10.3** DD Plots for the Buxton Data

**Example 10.13.** Buxton (1920, pp. 232-5) gave 20 measurements of 88 men. We will examine whether the multivariate normal distribution is a reasonable model for the measurements *head length*, *nasal height*, *bigonal breadth*, and *cephalic index* where one case has been deleted due to missing values. Figure 10.13a shows the DD plot. Five head lengths were recorded to be around 5 feet and are massive outliers. Figure 10.13b is the DD plot computed after deleting these points and suggests that the multivariate normal distribution is reasonable. (The recomputation of the DD plot means that the plot is not a weighted DD plot which would simply omit the outliers and then rescale the vertical axis.)

```
library(MASS)
x <- cbind(buxy,buxx)
ddplot(x,type=3) #Figure 7.13a), right click Stop

zx <- x[-c(61:65),]
ddplot(zx,type=3) #Figure 7.13b), right click Stop
```

**10.5 Summary**

**10.6 Complements**

**10.7 Problems**