

# Robust Multivariate Linear Regression

David J. Olive\*  
Southern Illinois University

May 3, 2013

## Abstract

A robust multivariate linear regression estimator can be obtained by replacing the least squares estimator with the robust **hbreg** estimator. Then the robust multivariate linear regression estimator is asymptotically equivalent to the classical multivariate linear regression estimator since the probability that the robust estimator is equal to the classical estimator goes to one in probability as the sample size  $n \rightarrow \infty$  for a large class of iid zero mean error distributions. This paper discusses the robust estimator and some tests and prediction regions using the robust estimator that are asymptotically equivalent to those using the classical estimator. A second robust estimator that is useful for outlier detection is derived.

**KEY WORDS:** outliers; prediction regions.

---

\*David J. Olive is Associate Professor, Department of Mathematics, Southern Illinois University, Mailcode 4408, Carbondale, IL 62901-4408, USA. E-mail address: dolive@siu.edu.

# 1 INTRODUCTION

Olive (2013b), using results from Su and Cook (2012) and Kakizawa (2009), derived a useful prediction region for the classical multivariate linear regression model, and gave  $F$  approximations to the widely used Wilk, Pillai, and Hotelling Lawley test statistics. This paper will show that these large sample tests and prediction regions also work for the robust multivariate linear regression estimator that replaces least squares with the **h**reg estimator. This section reviews the multivariate linear regression model and the results from Olive (2013b). Section 2 reviews the **h**reg estimator and derives the robust estimator and section 3 gives some examples and simulations. Section 4 derives a robust estimator that is useful for outlier detection.

## 1.1 The Multivariate Linear Regression Model

The *multivariate linear regression model*  $\mathbf{y}_i = \mathbf{B}^T \mathbf{x}_i + \boldsymbol{\epsilon}_i$  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  $(\mathbf{x}_i^T, \mathbf{y}_i^T) = (x_{i1}, x_{i2}, \dots, x_{ip}, Y_{i1}, \dots, Y_{im})$  where the constant  $x_{i1} = 1$  could be omitted from the case. 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  $E(\boldsymbol{\epsilon}_k) = \mathbf{0}$  and  $\text{Cov}(\boldsymbol{\epsilon}_k) = \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}} = (\sigma_{ij})$  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  $i, j = 1, \dots, m$  where  $\mathbf{I}_n$  is the  $n \times n$  identity matrix and  $\mathbf{e}_i$  is defined below. Then the  $p \times m$  coefficient matrix  $\mathbf{B} = [\beta_1 \ \beta_2 \ \dots \ \beta_m]$  and the  $m \times m$  covariance matrix  $\boldsymbol{\Sigma}_{\boldsymbol{\epsilon}}$  are to be estimated, and  $E(\mathbf{Z}) = \mathbf{X}\mathbf{B}$  while  $E(Y_{ij}) = \mathbf{x}_i^T \boldsymbol{\beta}_j$ .

The  $n \times m$  matrix of response variables

$$\mathbf{Z} = \begin{bmatrix} \mathbf{Y}_1 & \mathbf{Y}_2 & \dots & \mathbf{Y}_m \end{bmatrix} = \begin{bmatrix} \mathbf{y}_1^T \\ \vdots \\ \mathbf{y}_n^T \end{bmatrix}.$$

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

$$\mathbf{X} = \begin{bmatrix} \mathbf{v}_1 & \mathbf{v}_2 & \dots & \mathbf{v}_p \end{bmatrix} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix}$$

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

The  $n \times m$  matrix of errors

$$\mathbf{E} = \begin{bmatrix} \mathbf{e}_1 & \mathbf{e}_2 & \dots & \mathbf{e}_m \end{bmatrix} = \begin{bmatrix} \boldsymbol{\epsilon}_1^T \\ \vdots \\ \boldsymbol{\epsilon}_n^T \end{bmatrix}.$$

Least squares is the classical method for fitting the multivariate linear model. The *least squares estimators* are  $\hat{\mathbf{B}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Z} = [\hat{\beta}_1 \ \hat{\beta}_2 \ \dots \ \hat{\beta}_m]$ . The *predicted*

values or fitted values  $\hat{\mathbf{Z}} = \mathbf{X}\hat{\mathbf{B}} = [\hat{\mathbf{Y}}_1 \quad \hat{\mathbf{Y}}_2 \quad \dots \quad \hat{\mathbf{Y}}_m]$ . The residuals

$$\hat{\mathbf{E}} = \mathbf{Z} - \hat{\mathbf{Z}} = \mathbf{Z} - \mathbf{X}\hat{\mathbf{B}} = \begin{bmatrix} \hat{\boldsymbol{\epsilon}}_1^T \\ \hat{\boldsymbol{\epsilon}}_2^T \\ \vdots \\ \hat{\boldsymbol{\epsilon}}_n^T \end{bmatrix} = [\hat{\mathbf{r}}_1 \quad \hat{\mathbf{r}}_2 \quad \dots \quad \hat{\mathbf{r}}_m].$$

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

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

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  $\hat{\boldsymbol{\epsilon}}_i$  since the sample mean of the  $\hat{\boldsymbol{\epsilon}}_i$  is  $\mathbf{0}$ . Let  $\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon}} = \hat{\boldsymbol{\Sigma}}_{\boldsymbol{\epsilon},p}$  be the unbiased estimator of  $\boldsymbol{\Sigma}_{\boldsymbol{\epsilon}}$ .

The  $\boldsymbol{\epsilon}_i$  are assumed to be iid. Some important joint distributions for  $\boldsymbol{\epsilon}$  are completely specified by an  $m \times 1$  population *location* vector  $\boldsymbol{\mu}$  and an  $m \times m$  symmetric positive definite population *dispersion* matrix  $\boldsymbol{\Sigma}$ . An important model is the elliptically contoured  $EC_m(\boldsymbol{\mu}, \boldsymbol{\Sigma}, g)$  distribution with probability density function

$$f(\mathbf{z}) = k_m |\boldsymbol{\Sigma}|^{-1/2} g[(\mathbf{z} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{z} - \boldsymbol{\mu})]$$

where  $k_m > 0$  is some constant and  $g$  is some known function. The multivariate normal (MVN)  $N_m(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  distribution is a special case.

Some additional notation will be useful. Assume that  $\mathbf{x}_1, \dots, \mathbf{x}_n$  are iid from a multivariate distribution. The classical estimator  $(\bar{\mathbf{x}}, \mathbf{S})$  of multivariate location and dispersion is the sample mean and sample covariance matrix where

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

Let the  $p \times 1$  column vector  $T$  be a multivariate location estimator, and let the  $p \times p$  symmetric positive definite matrix  $\mathbf{C}$  be a dispersion estimator. Then the  $i$ th *squared sample Mahalanobis distance* is the scalar

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

for each observation  $\mathbf{x}_i$ . Notice that the Euclidean distance of  $\mathbf{x}_i$  from the estimate of center  $T$  is  $D_i(T, \mathbf{I}_p)$ . The classical Mahalanobis distance uses  $(T, \mathbf{C}) = (\bar{\mathbf{x}}, \mathbf{S})$ . Following Johnson (1987, pp. 107-108), the population squared Mahalanobis distance

$$U \equiv D^2(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}), \quad (3)$$

and for elliptically contoured distributions,  $U$  has probability density function (pdf)

$$h(u) = \frac{\pi^{p/2}}{\Gamma(p/2)} k_p u^{p/2-1} g(u). \quad (4)$$

## 1.2 A Prediction Region

Following Olive (2013b), 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 vector  $\mathbf{y}_f$ . Then a large sample  $(1 - \delta)100\%$  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.

If the  $\epsilon_i$  are iid from an  $EC_m(\mathbf{0}, \Sigma, g)$  distribution with continuous decreasing  $g$  and nonsingular covariance matrix  $\Sigma_\epsilon = c\Sigma$  for some constant  $c > 0$ , then the population asymptotically optimal prediction region is  $\{\mathbf{y} : D_{\mathbf{y}}(\mathbf{B}^T \mathbf{x}_f, \Sigma_\epsilon) \leq D_{1-\delta}\}$  where  $P(D_{\mathbf{y}}(\mathbf{B}^T \mathbf{x}_f, \Sigma_\epsilon) \leq D_{1-\delta}) = 1 - \delta$ . For example, if the iid  $\epsilon_i \sim N_m(\mathbf{0}, \Sigma_\epsilon)$ , 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.

*Theorem 1, see Olive (2013b).* Suppose  $\mathbf{y}_i = E(\mathbf{y}_i) + \epsilon_i = \hat{\mathbf{y}}_i + \hat{\epsilon}_i$  where  $\text{Cov}(\epsilon_i) = \Sigma_\epsilon > 0$ , and where the zero mean  $\epsilon_f$  and the  $\epsilon_i$  are iid for  $i = 1, \dots, n$ . Given  $\mathbf{x}_f$ , suppose the fitted model produces  $\hat{\mathbf{y}}_f$  and nonsingular  $\hat{\Sigma}_\epsilon$ . Let  $\hat{\mathbf{z}}_i = \hat{\mathbf{y}}_f + \hat{\epsilon}_i$  and

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

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

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

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  $D_i$ . Let the nominal  $100(1 - \delta)\%$  prediction region for  $\mathbf{y}_f$  be given by  $\{\mathbf{z} : (\mathbf{z} - \hat{\mathbf{y}}_f)^T \hat{\Sigma}_\epsilon^{-1} (\mathbf{z} - \hat{\mathbf{y}}_f) \leq D_{(U_n)}^2\} = \{\mathbf{z} : D_{\mathbf{z}}^2(\hat{\mathbf{y}}_f, \hat{\Sigma}_\epsilon) \leq D_{(U_n)}^2\} =$

$$\{\mathbf{z} : D_{\mathbf{z}}(\hat{\mathbf{y}}_f, \hat{\Sigma}_\epsilon) \leq D_{(U_n)}\}. \quad (5)$$

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  $(\hat{\mathbf{y}}_f, \hat{\Sigma}_\epsilon)$  is a consistent estimator of  $(E(\mathbf{y}_f), \Sigma_\epsilon)$ , then (5) is a large sample  $100(1 - \delta)\%$  prediction region for  $\mathbf{y}_f$ .

c) If  $(\hat{\mathbf{y}}_f, \hat{\Sigma}_\epsilon)$  is a consistent estimator of  $(E(\mathbf{y}_f), \Sigma_\epsilon)$ , and the  $\epsilon_i$  come from an elliptically contoured distribution such that the highest density region is  $\{\mathbf{z} : D_{\mathbf{z}}(\mathbf{0}, \Sigma_\epsilon) \leq D_{1-\delta}\}$ , then the prediction region (5) is asymptotically optimal.

Notice that for the data  $(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)$ , if  $\hat{\Sigma}_\epsilon^{-1}$  exists, then  $100q_n\%$  of the  $n$  cases are in their corresponding prediction region, and  $q_n \rightarrow 1 - \delta$  even if  $(\hat{\mathbf{y}}_i, \hat{\Sigma}_\epsilon)$  is not a good estimator. Hence the coverage  $q_n$  of the training data is robust to model assumptions. Of course the volume of the prediction region could be large if a poor estimator  $(\hat{\mathbf{y}}_i, \hat{\Sigma}_\epsilon)$  is used or if the  $\epsilon_i$  do not come from an elliptically contoured distribution. The nonparametric region uses  $(\hat{\mathbf{y}}_f, \hat{\Sigma}_\epsilon) = (\hat{\mathbf{y}}_f, \mathbf{S}_r)$  from the classical estimator in (5).

### 1.3 Testing

Following Olive (2013b), next 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. 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 the error or residual sum of squares and cross products matrix be

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

Then  $\mathbf{W}_e / (n - p) = \hat{\Sigma}_\epsilon$ . 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 Wilk's  $\Lambda$  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  $V(\mathbf{L}) = tr[(\mathbf{H} + \mathbf{W}_e)^{-1} \mathbf{H}] = \sum_{i=1}^m \frac{\lambda_i}{1 + \lambda_i}$ .

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

*Theorem 2, Olive (2013b).* The Hotelling-Lawley trace statistic

$$U(\mathbf{L}) = \frac{1}{n - p} [vec(\mathbf{L} \hat{\mathbf{B}})]^T [\hat{\Sigma}_\epsilon^{-1} \otimes (\mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T)^{-1}] [vec(\mathbf{L} \hat{\mathbf{B}})]. \quad (6)$$

Some notation is useful to show (6) 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

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

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  $\max_{1 \leq i \leq n} h_i \rightarrow 0$  as  $n \rightarrow \infty$ , assume that the zero mean iid errors have finite fourth moments, and assume that  $\frac{1}{n} \mathbf{X}^T \mathbf{X} \xrightarrow{P} \mathbf{W}^{-1}$ .

Then for the least squares estimator, Su and Cook (2012) show that if assumption D1 holds, then  $\hat{\Sigma}_\epsilon$  is  $\sqrt{n}$  consistent and  $\sqrt{n} \ vec(\hat{\mathbf{B}} - \mathbf{B}) \xrightarrow{D} N_{pm}(\mathbf{0}, \Sigma_\epsilon \otimes \mathbf{W})$ .

*Theorem 3, Olive (2013b).* If assumption D1 holds and if  $H_0$  is true, then  $(n - p)U(\mathbf{L}) \xrightarrow{D} \chi_{rm}^2$ .

Kakizawa (2009) shows, under stronger assumptions than Theorem 3, 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 noncentral  $\chi_{rm}^2(\omega^2)$  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 chi-square approximation is very good if  $n \geq 3(m^2 + p^2)$  for multivariate normal errors.

Theorems 2 and 3 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  $\hat{\Sigma}_e^{-1} = 1/\hat{\sigma}^2$ . Hence the scaled Hotelling Lawley test statistic is the partial  $F$  test statistic extended to  $m > 1$  predictor variables by Theorem 2.

By Theorem 3, for example,  $rF_R \xrightarrow{D} \chi_r^2$  for a large class of nonnormal error distribution. 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

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

Following Khattree and Naik (1999, p. 67) for the Roy's largest root test, if  $h =$

$\max(r, m)$ , use

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

Simulations in Olive (2013b) suggest that this approximation is good for  $r = 1$  but poor for  $r > 1$ . Anderson (1984, p. 333) states that Roy's largest root test has the greatest power if  $r = 1$  but is an inferior test for  $r > 1$ .

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. 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.

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 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.

The MANOVA partial F test is used to test whether a reduced model is good where the reduced model deletes  $r$  of the variables from the full model. For this test, 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.

## 2 Robust Estimators

### 2.1 Resistant Regression Estimators for Multiple Linear Regression

Consider the multiple linear regression model, written in matrix form as  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$ . This model is a special case of the multivariate linear regression model with  $m = 1$ .

Resistant estimators are useful for detecting certain types of outliers. Resistant estimators are often created by computing several trial fits  $\mathbf{b}_i$  that are estimators of  $\boldsymbol{\beta}$ . Then a criterion is used to select the trial fit to be used in the resistant estimator. Suppose  $c \approx n/2$ . The LMS( $c$ ) criterion is  $Q_{LMS}(\mathbf{b}) = r_{(c)}^2(\mathbf{b})$  where  $r_{(1)}^2 \leq \dots \leq r_{(n)}^2$  are the ordered squared residuals, and the LTS( $c$ ) criterion is  $Q_{LTS}(\mathbf{b}) = \sum_{i=1}^c r_{(i)}^2(\mathbf{b})$ . The LTA( $c$ ) criterion is  $Q_{LTA}(\mathbf{b}) = \sum_{i=1}^c |r(\mathbf{b})|_{(i)}$  where  $|r(\mathbf{b})|_{(i)}$  is the  $i$ th ordered absolute residual. Three impractical high breakdown robust estimators are the Hampel (1975) least median of squares (LMS) estimator, the Rousseeuw (1984) least trimmed sum of squares (LTS) estimator, and the Hössjer (1991) least trimmed sum of absolute deviations (LTA) estimator. These estimators correspond to the  $\hat{\boldsymbol{\beta}}_L \in \mathcal{R}^p$  that minimizes the corresponding criterion.

A good resistant estimator is the Olive (2005) *median ball algorithm* (MBA or mbareg). The Euclidean distance of the  $i$ th vector of predictors  $\mathbf{x}_i$  from the  $j$ th vector of predictors  $\mathbf{x}_j$  is

$$D_i(\mathbf{x}_j) = \sqrt{(\mathbf{x}_i - \mathbf{x}_j)^T (\mathbf{x}_i - \mathbf{x}_j)}.$$

For a fixed  $\mathbf{x}_j$  consider the ordered distances  $D_{(1)}(\mathbf{x}_j), \dots, D_{(n)}(\mathbf{x}_j)$ . Next, let  $\hat{\beta}_j(\alpha)$  denote the OLS fit to the  $\min(p + 3 + \lfloor \alpha n / 100 \rfloor, n)$  cases with the smallest distances where the approximate percentage of cases used is  $\alpha \in \{1, 2.5, 5, 10, 20, 33, 50\}$ . (Here  $\lfloor x \rfloor$  is the greatest integer function so  $\lfloor 7.7 \rfloor = 7$ . The extra  $p + 3$  cases are added so that OLS can be computed for small  $n$  and  $\alpha$ .) This yields seven OLS fits corresponding to the cases with predictors closest to  $\mathbf{x}_j$ . A fixed number  $K$  of cases are selected at random without replacement to use as the  $\mathbf{x}_j$ . Hence  $7K$  OLS fits are generated. We use  $K = 7$  as the default. A robust criterion  $Q$  is used to evaluate the  $7K$  fits and the OLS fit to all of the data. Hence  $7K + 1$  OLS fits are generated and the MBA estimator is the fit that minimizes the criterion. The median squared residual is a good choice for  $Q$ .

Three ideas motivate this estimator. First,  $\mathbf{x}$ -outliers, which are outliers in the predictor space, tend to be much more destructive than  $Y$ -outliers which are outliers in the response variable. Suppose that the proportion of outliers is  $\gamma$  and that  $\gamma < 0.5$ . We would like the algorithm to have at least one “center”  $\mathbf{x}_j$  that is not an outlier. The probability of drawing a center that is not an outlier is approximately  $1 - \gamma^K > 0.99$  for  $K \geq 7$  and this result is free of  $p$ . Secondly, by using the different percentages of coverages, for many data sets there will be a center and a coverage that contains no outliers. Thirdly, the MBA estimator is a  $\sqrt{n}$  consistent estimator.

The Olive and Hawkins (2011) **hbreg** estimator is a robust estimator that is asymptotically equivalent to the least squares estimator for many error distributions. Assume that the multiple linear regression model contains an intercept and that the median absolute deviation (MAD) of the  $Y$  values from their median is finite. Make an OLS fit to the  $c_n$  cases whose  $Y$  values are closest to the median  $Y$ , and use this fit as the start for concentration: find the  $c_n$  cases with the smallest squared residuals and fit OLS to these cases. Use 10 concentration steps and let the attractor be the final estimator, denoted by  $\hat{\beta}_B$ . It can be shown that  $\hat{\beta}_B$  is a high breakdown estimator.

With these preliminaries, we now define our high breakdown procedure. This is made up of three components.

- 1) The OLS estimator  $\hat{\beta}_C$  that is consistent for clean data.
- 2) The practical  $\sqrt{n}$  consistent **mbareg** estimator  $\hat{\beta}_A$  that is effective for outlier identification.
- 3) The practical high-breakdown estimator  $\hat{\beta}_B$ .

By selecting one of these three estimators according to the features each of them uncovers in the data, we may inherit the good properties of each of them.

The **hbreg** estimator  $\hat{\beta}_H$  is defined as follows. Pick a constant  $a > 1$  and set  $\hat{\beta}_H = \hat{\beta}_C$ . If  $aQ_L(\hat{\beta}_A) < Q_L(\hat{\beta}_C)$ , set  $\hat{\beta}_H = \hat{\beta}_A$ . If  $aQ_L(\hat{\beta}_B) < \min[Q_L(\hat{\beta}_C), aQ_L(\hat{\beta}_A)]$ , set  $\hat{\beta}_H = \hat{\beta}_B$ .

That is, find the smallest of the three scaled criterion values  $Q_L(\hat{\beta}_C)$ ,  $aQ_L(\hat{\beta}_A)$ ,  $aQ_L(\hat{\beta}_B)$ . According to which of the three estimators attains this minimum, set  $\hat{\beta}_H$  to  $\hat{\beta}_C$ ,  $\hat{\beta}_A$  or  $\hat{\beta}_B$  respectively.

Large sample theory for **hbreg** is simple and given in the following theorem. Let  $\hat{\beta}_L$  be the LMS, LTS or LTA estimator that minimizes the criterion  $Q_L$ . Note that the impractical estimator  $\hat{\beta}_L$  is never computed. The following theorem shows that  $\hat{\beta}_H$  is asymptotically equivalent to  $\hat{\beta}_C = \hat{\beta}_{OLS}$ . The clean data are in general position if any  $p$

clean cases give a unique estimate of  $\hat{\beta}$ . The LTA criterion will be used in the simulations.

*Theorem 4, Olive and Hawkins (2011).* Assume the clean data are in general position, and suppose that both  $\hat{\beta}_L$  and  $\hat{\beta}_C$  are consistent estimators of  $\beta$  where the regression model contains a constant. Then the **hbreg** estimator  $\hat{\beta}_H$  is high breakdown and is asymptotically equivalent to  $\hat{\beta}_C$  since the probability that  $\hat{\beta}_H = \hat{\beta}_C$  goes to one as  $n \rightarrow \infty$ .

## 2.2 Robust Multivariate Linear Regression

The classical multivariate linear regression estimator is found from  $m$  least squares multiple linear regressions of  $Y_j$  on the predictors. The robust multivariate linear regression estimator is found from  $m$  **hbreg** multiple linear regressions of  $Y_j$  on the predictors. By Theorem 4, the probability that the robust estimator is equal to the classical estimator goes to one as  $n \rightarrow \infty$  for a large class of error distributions.

Hence the large sample nonparametric prediction region and the large sample Wilk's test, Pillai's test and Hotelling Lawley test using the robust estimator are asymptotically equivalent to their analogs using the classical estimator for a large class of error distributions. The next section investigates whether reasonable sample sizes result in good results for the robust estimator.

## 3 Plots, Examples and Simulations

### 3.1 Plots

A *response plot* for the  $j$ th response variable is a plot of the fitted values  $\hat{Y}_{ij}$  versus the response  $Y_{ij}$  where  $i = 1, \dots, n$ . 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  $r_{ij}$ .

Make the  $m$  response and residual plots for the multivariate linear regression model. In a response plot, the vertical deviations from the identity line are the residuals  $r_{ij} = Y_{ij} - \hat{Y}_{ij}$ . The plotted points in the response plot 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 changed or corrected. The response and residual plots are used just as for multiple linear regression where  $m = 1$ . See Olive and Hawkins (2005) and Cook and Weisberg (1999, p. 432).

The Rousseeuw and Van Driessen (1999) DD plot is a plot of classical Mahalanobis distances versus robust Mahalanobis distances. Results from Olive (2002) suggest the plotted points in the DD plot will cluster about the identity line if the  $\epsilon_i$  are iid from a multivariate normal  $N_m(\mathbf{0}, \Sigma_\epsilon)$  distribution and about some line through the origin with slope greater than one for a large class of elliptically contoured distributions. Make a DD plot of the residuals  $\hat{\epsilon}_i$  to check the error distribution. Make a DD plot of the continuous predictor variables to check for  $\mathbf{x}$ -outliers.

The Olive and Hawkins (2010) RMVN estimator  $(T_{RMVN}, \mathbf{C}_{RMVN})$  is an easily computed  $\sqrt{n}$  consistent estimator of  $(\boldsymbol{\mu}, c\boldsymbol{\Sigma})$  under regularity conditions (E1) that include a large class of elliptically contoured distributions, and  $c = 1$  for the  $N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  distribution. The RMVN estimator also gives useful estimates of  $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  for  $N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  data even when certain types of outliers are present, and will be the robust estimator used in the DD plots. Also see Zhang, Olive and Ye (2012).

Consider the DD plot applied to the  $\hat{\mathbf{z}}_i$  based on the robust estimator. The non-parametric region based on the robust estimator uses the sample mean and sample covariance matrix applied to the  $\hat{\mathbf{z}}_i$ . The DD plot will have a vertical line at the cutoff  $D_{(U_n)}$ . Hence points to the left of the line correspond to cases that are in the non-parametric region. The RMVN estimator can be applied to the  $\hat{\mathbf{z}}_i$ . The region that uses  $D_i(T_{RMVN}, \mathbf{C}_{RMVN}) \leq D_{(U_n)}(T_{RMVN}, \mathbf{C}_{RMVN})$  will be called the semiparametric region, while the parametric MVN region uses  $D_i(T_{RMVN}, \mathbf{C}_{RMVN}) \leq \sqrt{\chi_{p, q_n}^2}$  where  $P(W \leq \chi_{p, q_n}^2) = q_n$  if  $W \sim \chi_p^2$ . These two regions are only conjectured to be large sample prediction regions, but are added to the DD plot as visual aids. Cases below the horizontal line that crosses the identity line correspond to the semiparametric region while cases below the horizontal line that ends at the identity line correspond to the parametric MVN region. A vertical line dropped down from this point of intersection does correspond to a large sample prediction region for multivariate normal data. Note that  $\hat{\mathbf{z}}_i = \hat{\mathbf{y}}_f + \hat{\boldsymbol{\epsilon}}_i$ , 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.

### 3.2 Examples and Simulations

**Example 1.** Cook and Weisberg (1999, p. 351, 433, 447) give 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. Figures 1 and 2 give the response and residual plots for  $Y_1$  and  $Y_2$ . For  $Y_1$ , case 79 sticks out while for  $Y_2$ , cases 8, 25 and 48 are not fit well. Figure 3 shows the DD plot of the residual vectors. 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. For this data set, the classical and robust estimators were identical, and hence the Cook (1977) distances can be computed. Highlighted cases had Cook's distance  $> \min(0.5, 2p/n)$ . The response, residual and DD plots are effective for finding influential cases, for checking linearity and 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 that are not in their prediction region. These are the cases corresponding to residual vectors with large Mahalanobis distances. Also adding a constant does not change the distance, so the DD plot for the residuals is the same as the DD plot for the  $\hat{\mathbf{z}}_i$ .

**Example 2.** Buxton (1920) gives various measurements of 88 men. *Head length* and person's *height* were the response variables while an intercept, *nasal height*, *bigonal*

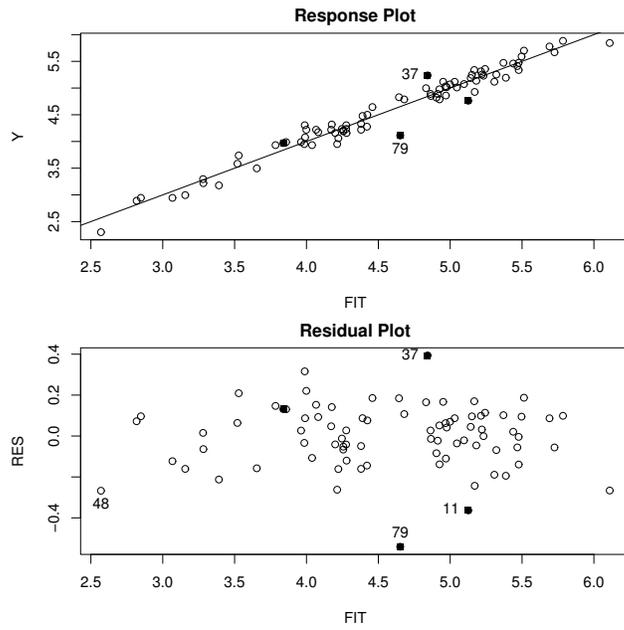


Figure 1: Plots for  $Y_1 = \log(S)$ .

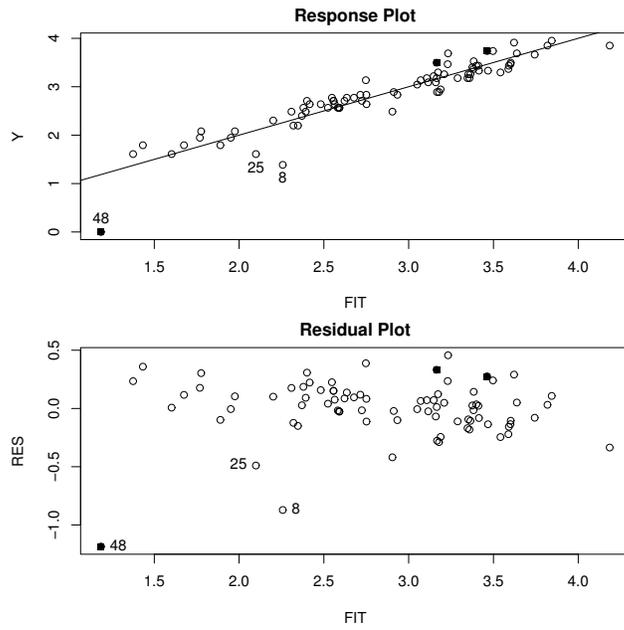


Figure 2: Plots for  $Y_2 = \log(M)$ .

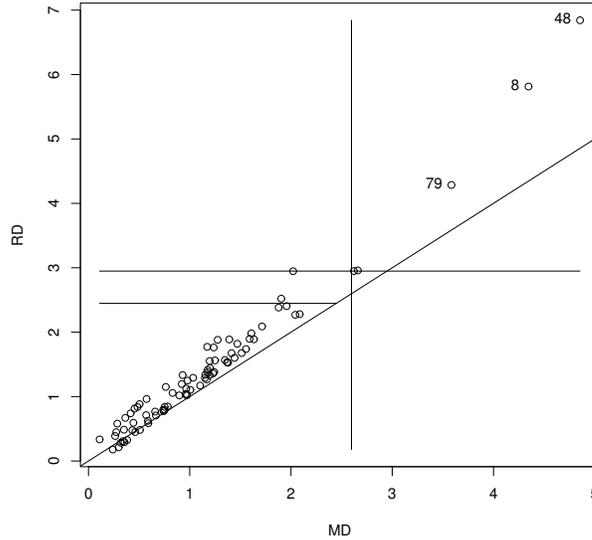


Figure 3: DD Plot of the Residual Vectors.

*breadth*, and *cephalic index* were used as predictors in the multiple linear regression model. Observation 9 was deleted since it had missing values. Five individuals, numbers 62–66, were reported to be about 0.75 inches tall with head lengths well over five feet! Figure 4 shows the response and residual plots corresponding to  $Y_1$  for the robust estimator. The response plot for the classical estimator, not shown, has the identity line tilted slightly above most of the plotted points in the lower part of the plot, while the plotted points in the lower part of the residual plot follow a line with negative slope instead of the  $r = 0$  line. Figure 5 shows the response and residual plots corresponding to  $Y_2$  for the robust estimator. The response plot for the classical estimator, not shown, has the identity line tilted slightly below most of the plotted points in the upper part of the plot, while the plotted points in the upper part of the residual plot follow a line with negative slope instead of the  $r = 0$  line. Figure 6 shows the DD plot. The 90% semiparametric and nonparametric regions use the 95th percentile which is a linear combination of an outlying case with a nonoutlying case. The parametric MVN region contains cases below the  $RD = 2.448$  line, which is obscured by the identity line. The tests of hypotheses for the robust estimator are not robust to outliers because all  $n = 87$  residual vectors are used to make  $\hat{\Sigma}\epsilon$ . As is typically the case, outliers can be detected with the plots using the classical or robust estimator.

A small simulation was used to study the prediction region and the Wilk's 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. These test statistics were computed with the robust estimator  $\hat{\mathbf{B}}$  instead of the classical estimator. 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

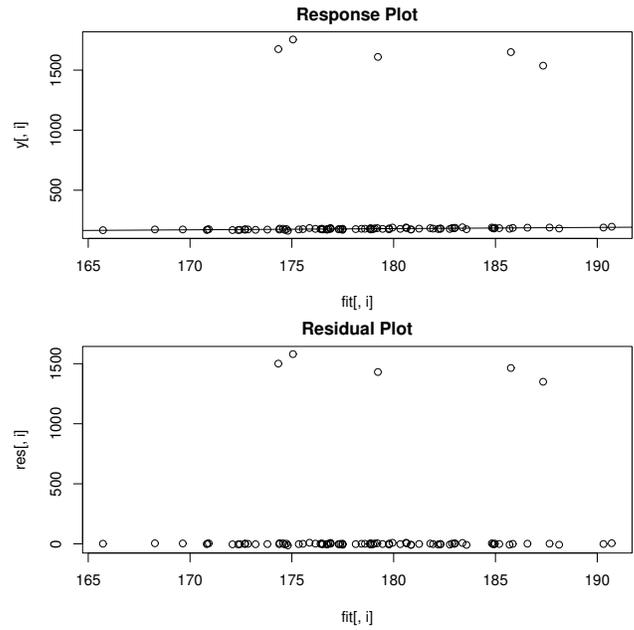


Figure 4: Plots for  $Y_1 = \text{head length}$ .

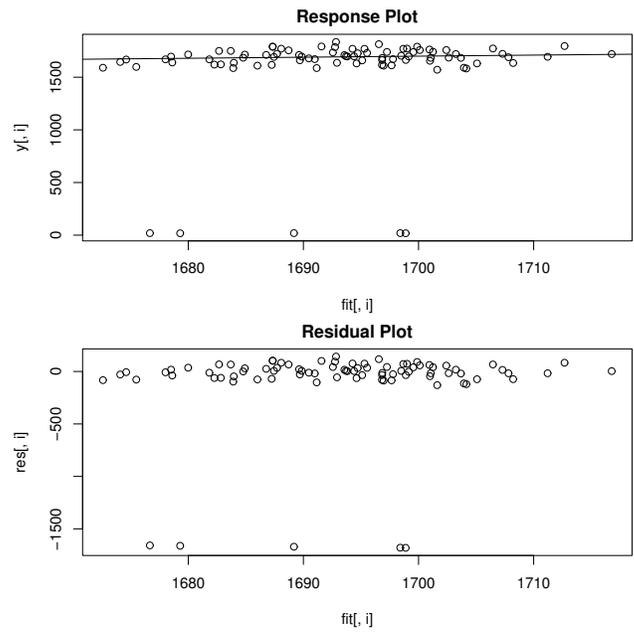


Figure 5: Plots for  $Y_2 = \text{height}$ .

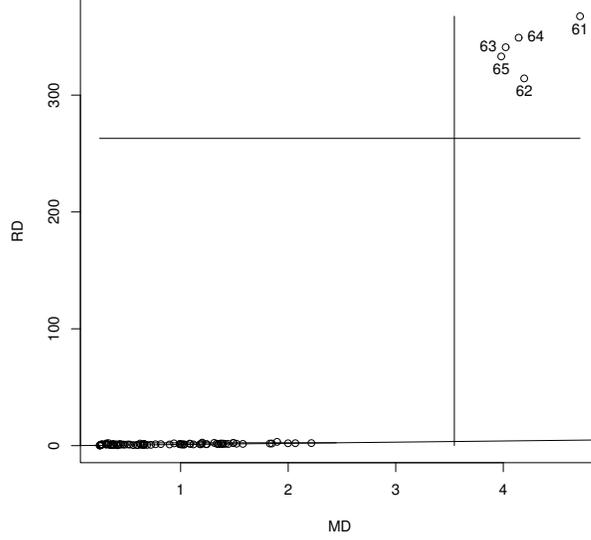


Figure 6: DD Plot of the Residual Vectors for the Buxton Data.

is false, then the second to last row of  $\mathbf{B}$  is  $(1, 0, \dots, 0)$ , the third to last row is  $(1, 1, 0, \dots, 0)$  etcetera as long as the first row is not changed from  $\mathbf{1}^T$ . First  $m$  iid errors  $\mathbf{w}_i$  are generated such that the  $m$  errors 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  $\hat{\boldsymbol{\Sigma}}\boldsymbol{\epsilon} = \sigma^2\mathbf{A}\mathbf{A}^T = (\sigma_{ij})$  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  $\psi = 0.10$ . Hence the correlations  $(2\psi + (m-2)\psi^2)/(1 + (m-1)\psi^2)$ . As  $\psi$  gets close to 1, the data clusters about the line in the direction of  $(1, \dots, 1)^T$ . 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 -  $E(\text{lognormal})$ : 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.

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. 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.

The **hbreg** estimator is asymptotically equivalent to OLS provided that OLS is consistent and the probability that **hbreg** selects  $\hat{\beta}_B$  goes to 0 as  $n \rightarrow \infty$ . For the simulated data with symmetric error distributions,  $\hat{\beta}_B$  appeared to give biased estimates of the slopes. However, for the simulated data with right skewed error distributions,  $\hat{\beta}_B$  appeared to give good estimates of the slopes but not the constant, and the probability that the **hbreg** estimator selected  $\hat{\beta}_B$  appeared to go to one. Removing  $\hat{\beta}_B$  from the **hbreg** estimator results in a  $\sqrt{n}$  consistent estimator when OLS is  $\sqrt{n}$  consistent, but massive sample sizes were still needed to get good estimates of the constants for highly skewed error distributions. Although the **mbareg** estimator is a  $\sqrt{n}$  consistent estimator of  $\beta$ , if  $m = 1$  and OLS needed  $n = 1000$  to estimate the constant well, **mbareg** might need  $n >$  one million. Getting all  $m$  **mbareg** estimators to estimate the constant well needs even larger sample sizes.

This paragraph will explain why **mbareg** needs large samples to give a good estimate of the constant for highly skewed error distributions when  $m = 1$ . Note that the LMS, LTA and LMS criteria use half sets. For simplicity, consider the LMS criterion that minimizes the median squared residual. Heuristically, for highly right skewed data, let the “left tail half set” shift the constant of the OLS hyperplane down so that the half set of cases closest to the plane are the half set with the smallest OLS residual values. These cases will have negative residuals and residuals close to zero, which are roughly the cases corresponding to the half set of errors in the left tail of the error distribution. Let the “OLS half set” correspond to the half set of cases with the smallest absolute OLS residuals, so the cases closest to the OLS hyperplane. Since the distribution is highly right skewed, the “OLS half set” has much more variability than the “left tail half set.” (For the location model, OLS is the sample mean which is greater than the sample median for right skewed data. The “left tail half set” shifts the mean down to the midpoint  $c$  of the minimum value and the median value, and often  $c \approx \text{median}/2$  if the support of the highly right skewed distribution is  $(0, \infty)$ .) A trial fit that uses the same OLS slope estimates but which shifts the intercept down to use the “left tail half set” will have a smaller median squared residual than the median squared residual using OLS. The trial fits for **mbareg** are  $\sqrt{n}$  consistent, so estimate the OLS intercept eventually. However, the trial fit that uses 1% of the data has less than 1% efficiency since the  $\mathbf{x}$  values are close in distance rather than spread out. Unless the sample size is large, the **mbareg** estimator tends to produce some trial fits that shift the intercept down, and one of these trial fits is selected to be the final **mbareg** estimator since it has a smaller median squared residual than the other trial fits.

In the simulations, **hbreg** estimated the slopes well for the highly skewed lognormal data, but not the constant. This results in incorrect residual vectors and test statistics. The **hbreg** tests can be used as diagnostics if the plotted points in the DD plot cluster tightly about a line through the origin. Also compare  $\hat{\mathbf{B}}$  for the robust and classical

Table 1: Test Coverages: MANOVA  $F$   $H_0$  is True.

$w$ dist	$n$	test	$F_1$	$F_2$	$F_{p-1}$	$F_p$	$F_M$
MVN	50	W	1	0.051	0.051	0.054	0.029
MVN	50	P	1	0.036	0.033	0.038	0.006
MVN	50	HL	1	0.104	0.110	0.119	0.130
MVN	50	R	1	0.098	0.100	0.110	0.727
MVN	200	W	1	0.044	0.046	0.044	0.044
MVN	200	P	1	0.042	0.042	0.042	0.035
MVN	200	HL	1	0.056	0.056	0.054	0.060
MVN	200	R	1	0.050	0.051	0.050	0.520
MIX	200	W	1	0.043	0.043	0.044	0.034
MIX	200	P	1	0.041	0.040	0.040	0.027
MIX	200	HL	1	0.054	0.053	0.053	0.048
MIX	200	R	1	0.050	0.049	0.048	0.518
MVT(7)	200	W	1	0.040	0.042	0.043	0.036
MVT(7)	200	P	1	0.038	0.040	0.040	0.028
MVT(7)	200	HL	1	0.049	0.049	0.053	0.051
MVT(7)	200	R	1	0.046	0.046	0.049	0.524

estimator. If the slopes are similar but not the intercepts, there may be highly skewed data or  $y$ -outliers.

Hence the simulation for tests of hypotheses used the symmetric elliptically contoured distributions. In Olive (2013b) for the classical estimator when the null hypothesis  $H_0$  was true, simulated values started to get close to nominal levels for  $n \geq 0.75(m + p)^2$ , 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  $r > 0$ . Roy's test was very good for the  $F_j$  tests but very poor for the MANOVA  $F$  test.

The robust estimator needed larger values of  $n$ , and results are shown in Table 1 for  $m = p = 5$ . Want values for the  $F_1$  test to be close to 1 since  $H_0$  is false for the  $F_1$  test and want values close to 0.05, otherwise. Results from Berndt and Savin (1977) suggest that Pillai's test will reject  $H_0$  less often than Wilk's test which will reject less often than the Hotelling Lawley test. For Table 2 want values near 1 except for the  $F_p$  column.

The same type of data and 5000 runs were used to simulate the prediction regions for  $\mathbf{y}_f$  given  $\mathbf{x}_f$  for multivariate regression. 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. Following Olive (2013a), consider the prediction region  $\{\mathbf{z} : (\mathbf{z} - T)^T \mathbf{C}^{-1}(\mathbf{z} - T) \leq h^2\} = \{\mathbf{z} : D_{\mathbf{z}}^2 \leq h^2\} = \{\mathbf{z} : D_{\mathbf{z}} \leq h\}$ . Then the ratio of the prediction region volumes

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

was recorded where  $i = 1$  was the nonparametric region,  $i = 2$  was the semiparametric

Table 2: Test Coverages: MANOVA  $F$   $H_0$  is False.

$n$	$\mathbf{w}$ dist	$m = p$	test	$F_1$	$F_2$	$F_{p-1}$	$F_p$	$F_M$
30	MVN	5	W	0.807	0.522	0.152	0.000	0.043
30	MVN	5	P	0.000	0.000	0.000	0.000	0.000
30	MVN	5	HL	0.995	0.889	0.570	0.009	0.878
30	MVN	5	R	1	0.956	0.756	0.076	0.999
50	MVN	5	W	1	0.999	0.921	0.057	1
50	MVN	5	P	1	0.998	0.911	0.042	0.997
50	MVN	5	HL	1	0.999	0.943	0.116	1
50	MVN	5	R	1	0.999	0.941	0.107	1
200	MVN	5	W	1	1	1	0.051	1
200	MVN	5	P	1	1	1	0.051	1
200	MVN	5	HL	1	1	1	0.060	1
200	MVN	5	R	1	1	1	0.056	1
50	MIX	5	W	0.803	0.569	0.278	0.010	0.502
50	MIX	5	P	0.747	0.518	0.236	0.006	0.230
50	MIX	5	HL	0.897	0.689	0.392	0.025	0.723
50	MIX	5	R	0.888	0.675	0.377	0.022	0.970
50	MVT	5	W	1	0.997	0.892	0.040	0.998
50	MVT	5	P	1	0.996	0.871	0.025	0.989
50	MVT	5	HL	1	0.999	0.925	0.087	1
50	MVT	5	R	1	0.998	0.922	0.080	1
450	MVN	20	W	1	1	1	0.015	1
450	MVN	20	P	1	1	1	0.013	1
450	MVN	20	HL	1	1	1	0.030	1
450	MVN	20	R	1	1	1	0.053	1

Table 3: Coverages for 90% Prediction Regions.

$\mathbf{w}$ dist	$n$	$m = p$	ncov	scov	mcov	nvol	mvol
MVN	300	5	0.903	0.899	0.902	1.006	1.017
MIX	300	5	0.897	0.907	0.688	0.885	0.001
MVT(7)	300	5	0.901	0.910	0.775	0.913	0.291
LN	300	5	0.915	0.916	0.592	0.688	0.008

region, and  $i = 3$  was the parametric MVN region. Here  $h_1$  and  $h_2$  were the cutoff  $D_{(U_n)}(T_i, \mathbf{C}_i)$  for  $i = 1, 2$ , and  $h_3 = \sqrt{\chi_{m,q_n}^2}$ .

If, as conjectured, the RMVN estimator is a consistent estimator when applied to the residual vectors instead of iid data, then the volume ratios converge in probability to 1 if the iid zero mean errors  $\sim N_m(\mathbf{0}, \Sigma_\epsilon)$ , and the volume ratio converges to 1 for  $i = 1$  for a large class of elliptically contoured distributions. These volume ratios were denoted by voln and volm for the nonparametric and parametric MVN regions. The coverage was the proportion of times the prediction region contained  $\mathbf{y}_f$  where ncov, scov and mcov are for the nonparametric, semiparametric and parametric MVN regions.

In the simulations, took  $n = 3(m + p)^2 = 300$  and  $m = p = 5$ . Table 3 shows that the coverage of the nonparametric region was close to 0.9 in all cases. The volume ratio voln was fairly close to 1 for the three elliptically contoured distributions. Since the volume of the prediction region is proportional to  $h^m$ , the volume can be very small if  $h$  is too small and  $m$  is large. Parametric prediction regions usually give poor estimates of  $h$  when the parametric distribution is misspecified. Hence the parametric MVN region only performed well for multivariate normal data. The results using the robust estimator were nearly the same as those using the classical estimator. See Table 3 in Olive (2013b).

## 4 Another Robust Estimator

First we will review some results for multiple linear regression. Let  $\mathbf{x} = (1, \mathbf{w}^T)^T$  and let

$$\text{Cov}(\mathbf{w}) = E[(\mathbf{w} - E(\mathbf{w}))(\mathbf{w} - E(\mathbf{w}))^T] = \Sigma_{\mathbf{w}}$$

and  $\text{Cov}(\mathbf{w}, Y) = E[(\mathbf{w} - E(\mathbf{w}))(Y - E(Y))] = \Sigma_{\mathbf{w}Y}$ . Let  $\boldsymbol{\beta} = (\alpha, \boldsymbol{\eta}^T)^T$  be the population OLS coefficients from the regression of  $Y$  on  $\mathbf{x}$  ( $\mathbf{w}$  and a constant), where  $\alpha$  is the constant and  $\boldsymbol{\eta}$  is the vector of slopes. Let the OLS estimator be  $\hat{\boldsymbol{\beta}} = (\hat{\alpha}, \hat{\boldsymbol{\eta}}^T)^T$ . Then the population coefficients from an OLS regression of  $Y$  on  $\mathbf{x}$  are

$$\alpha = E(Y) - \boldsymbol{\eta}^T E(\mathbf{w}) \quad \text{and} \quad \boldsymbol{\eta} = \Sigma_{\mathbf{w}}^{-1} \Sigma_{\mathbf{w}Y}. \quad (7)$$

Then the OLS estimator  $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$ . The sample covariance matrix of  $\mathbf{w}$  is

$$\hat{\Sigma}_{\mathbf{w}} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{w}_i - \bar{\mathbf{w}})(\mathbf{w}_i - \bar{\mathbf{w}})^T \quad \text{where the sample mean} \quad \bar{\mathbf{w}} = \frac{1}{n} \sum_{i=1}^n \mathbf{w}_i.$$

Similarly, define the sample covariance vector of  $\mathbf{w}$  and  $Y$  to be

$$\hat{\Sigma}_{\mathbf{w}Y} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{w}_i - \bar{\mathbf{w}})(Y_i - \bar{Y}).$$

Suppose that  $(Y_i, \mathbf{w}_i^T)^T$  are iid random vectors such that  $\Sigma_{\mathbf{w}}^{-1}$  and  $\Sigma_{\mathbf{w}Y}$  exist. Then

$$\hat{\alpha} = \bar{Y} - \hat{\boldsymbol{\eta}}^T \bar{\mathbf{w}} \xrightarrow{D} \alpha$$

and

$$\hat{\boldsymbol{\eta}} = \hat{\Sigma}_{\mathbf{w}}^{-1} \hat{\Sigma}_{\mathbf{w}Y} \xrightarrow{D} \boldsymbol{\eta} \text{ as } n \rightarrow \infty.$$

Now for multivariate linear regression,  $\hat{\boldsymbol{\beta}}_j = (\hat{\alpha}_j, \hat{\boldsymbol{\eta}}_j^T)^T$  where  $\hat{\alpha}_j = \bar{Y}_j - \hat{\boldsymbol{\eta}}_j^T \bar{\mathbf{w}}$  and  $\hat{\boldsymbol{\eta}}_j = \hat{\Sigma}_{\mathbf{w}}^{-1} \hat{\Sigma}_{\mathbf{w}Y_j}$ . Let  $\hat{\Sigma}_{\mathbf{w}\mathbf{y}} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{w}_i - \bar{\mathbf{w}})(\mathbf{y}_i - \bar{\mathbf{y}})^T$  which has  $j$ th column  $\hat{\Sigma}_{\mathbf{w}Y_j}$  for  $j = 1, \dots, m$ . Let

$$\mathbf{u} = \begin{pmatrix} \mathbf{w} \\ \mathbf{y} \end{pmatrix}, \quad E(\mathbf{u}) = \boldsymbol{\mu}_{\mathbf{u}} = \begin{pmatrix} E(\mathbf{w}) \\ E(\mathbf{y}) \end{pmatrix} = \begin{pmatrix} \boldsymbol{\mu}_{\mathbf{w}} \\ \boldsymbol{\mu}_{\mathbf{y}} \end{pmatrix}, \quad \text{and} \quad \text{Cov}(\mathbf{u}) = \Sigma_{\mathbf{u}} = \begin{pmatrix} \Sigma_{\mathbf{w}\mathbf{w}} & \Sigma_{\mathbf{w}\mathbf{y}} \\ \Sigma_{\mathbf{y}\mathbf{w}} & \Sigma_{\mathbf{y}\mathbf{y}} \end{pmatrix}.$$

Let the vector of constants be  $\boldsymbol{\alpha}^T = (\alpha_1, \dots, \alpha_m)$  and the matrix of slope vectors  $\mathbf{B}_S = \begin{bmatrix} \boldsymbol{\eta}_1 & \boldsymbol{\eta}_2 & \dots & \boldsymbol{\eta}_m \end{bmatrix}$ . Then the population least squares coefficient matrix is

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

where  $\boldsymbol{\alpha} = \boldsymbol{\mu}_{\mathbf{y}} - \mathbf{B}_S^T \boldsymbol{\mu}_{\mathbf{w}}$  and  $\mathbf{B}_S = \Sigma_{\mathbf{w}}^{-1} \Sigma_{\mathbf{w}\mathbf{y}}$  where  $\Sigma_{\mathbf{w}} = \Sigma_{\mathbf{w}\mathbf{w}}$ .

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

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

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

Let  $(T, \mathbf{C}) = (\tilde{\boldsymbol{\mu}}_{\mathbf{u}}, \tilde{\Sigma}_{\mathbf{u}})$  be a robust estimator of multivariate location and dispersion. If  $\tilde{\boldsymbol{\mu}}_{\mathbf{u}}$  is a consistent estimator of  $\boldsymbol{\mu}_{\mathbf{u}}$  and  $\tilde{\Sigma}_{\mathbf{u}}$  is a consistent estimator of  $c \Sigma_{\mathbf{u}}$  for some constant  $c > 0$ , then a robust estimator of multivariate linear regression is the plug in estimator  $\tilde{\boldsymbol{\alpha}} = \tilde{\boldsymbol{\mu}}_{\mathbf{y}} - \tilde{\mathbf{B}}_S^T \tilde{\boldsymbol{\mu}}_{\mathbf{w}}$  and  $\tilde{\mathbf{B}}_S = \tilde{\Sigma}_{\mathbf{w}}^{-1} \tilde{\Sigma}_{\mathbf{w}\mathbf{y}}$ .

If  $(T, \mathbf{C})$  is the RMVN estimator applied to the  $\mathbf{u}_i$ , then  $(T, \mathbf{C})$  is a  $\sqrt{n}$  consistent estimator of  $(\boldsymbol{\mu}_{\mathbf{u}}, c \Sigma_{\mathbf{u}})$  if  $\mathbf{u}$  is from a large class of  $EC_d(\boldsymbol{\mu}_{\mathbf{u}}, \Sigma_{\mathbf{u}}, g)$  distributions where  $d = m + p - 1$ . Thus the classical and robust estimators of multivariate linear regression are both  $\sqrt{n}$  consistent estimator of  $\mathbf{B}$  if the  $\mathbf{u}_i$  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  $\mathbf{u}$  is not elliptically contoured, it is possible that the robust estimator is bad and very different from the good classical least squares estimator.

Now the RMVN estimator computes the classical estimator of multivariate location and dispersion on the RMVN set of cases in highly concentrated ellipsoidal region, and then multiplies the dispersion estimator by a constant  $c$ . Hence the plug in robust multivariate linear regression estimator using RMVN is equivalent to the least squares multivariate linear regression estimator applied to the cases in the RMVN set. Call this estimator the `rmreg2` estimator.

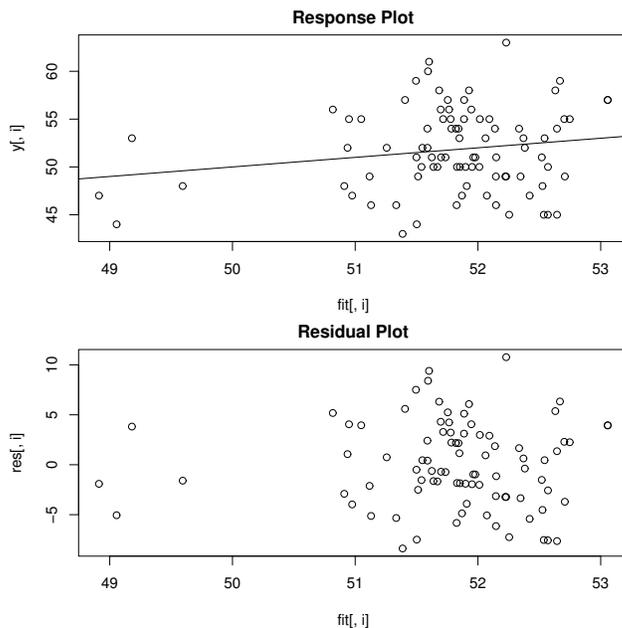


Figure 7: Plots for  $Y_1 = \text{nasal height}$  using `hblog`.

**Example 2, continued.** The plots for the `rmreg2` estimator were very similar to Figures 5 and 6. Now let  $Y_1 = \text{nasal height}$  and  $Y_2 = \text{height}$  with  $x_2 = \text{head length}$ ,  $x_3 = \text{bigonal breadth}$  and  $x_4 = \text{cephalic index}$ . Then  $Y_1$  and  $x_2$  have massive outliers. Then the response and residual plots for the classical estimator and the robust estimator using `hblog` were nearly identical. Figures 7 and 8 show that the fit using `hblog` went right through the outliers. Figures 9 and 10 show that the response and residual plots corresponding to `rmreg2` do not have fits that pass through the outliers.

## 5 Conclusions

Multivariate linear regression is a semiparametric method that is nearly as easy to use as multiple linear regression if  $m$  is small. The  $m$  response and residual plots should be made as well as the DD plot. For the classical estimator, response and residual plots can look good for  $n \geq 10p$ , but for testing and prediction regions, may need  $n \geq k(m + p)^2$

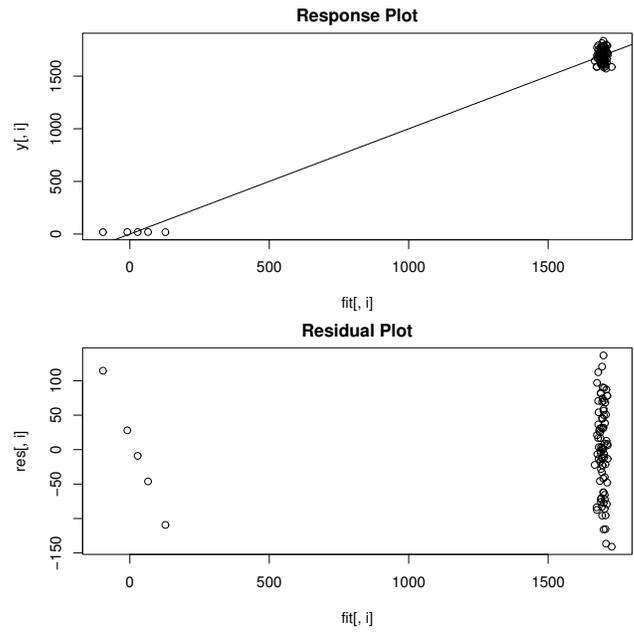


Figure 8: Plots for  $Y_2 = \text{height}$  using `hbrreg`.

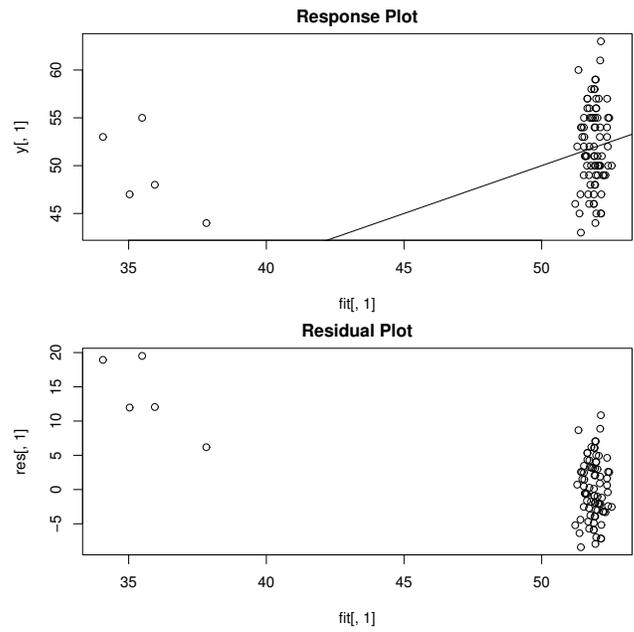


Figure 9: Plots for  $Y_1 = \text{nasal height}$  using `rmreg2`.

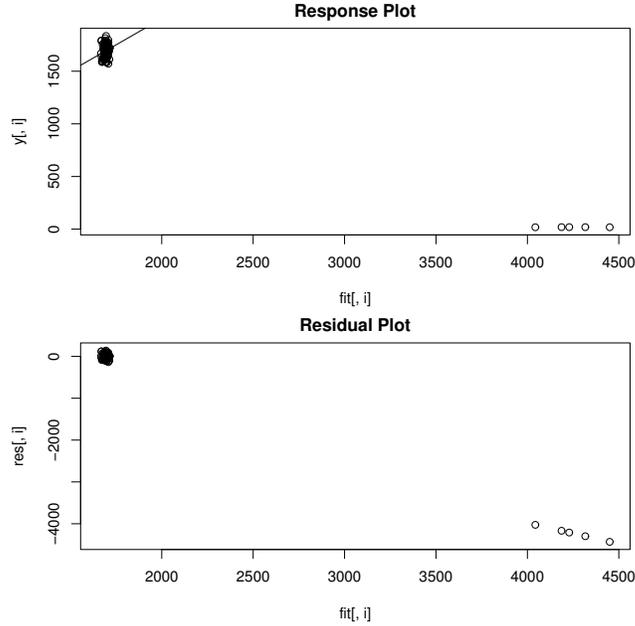


Figure 10: Plots for  $Y_2 = \text{height}$  using `rmreg2`.

where  $0.5 \leq k \leq 3$  even for well behave elliptically contoured error distributions. For the robust estimator, larger sample sizes are needed, and for highly skewed data, the robust tests may fail. If the plotted points in the DD plot cluster tightly about a line through the origin, then an elliptically contoured error distribution may be reasonable, and then the first row of  $\hat{\mathbf{B}}$  corresponding to the intercepts should be similar for both the robust and classical estimators.

The *R* software was used to make plots and software. See R Development Core Team (2011). The programs in the collection of functions *mpack.txt* are available at ([www.math.siu.edu/olive/mpack.txt](http://www.math.siu.edu/olive/mpack.txt)). The function `rmpredsim` was used to simulate the prediction regions, `rmregsim` was used to simulate the tests of hypotheses, and `rmregddsim` simulated the DD plots for various distributions. The function `rmltreg` makes the response and residual plots and computes the  $F_j$ , MANOVA  $F$  and MANOVA partial  $F$  test pvalues while the function `ddplot4` makes the DD plots. Similar functions for the classical estimator delete the initial “r.” The highly outlier resistant multivariate linear regression estimator based on RMVN can be computed using the function `rmreg2`. When  $m = 1$ , this estimator is a highly outlier resistant multiple linear regression estimator.

The two methods for robust multivariate regression were to plug in a robust multiple linear regression estimator like `hbreg` in place of OLS, and to use a robust estimator of multivariate location and dispersion as a plug in estimator. Neither idea is new, but using practical highly outlier resistant estimators backed by theory, like `hbreg` and RMVN, is new. Rousseeuw, Van Aelst, Van Driessen and Agulló (2004) use FLTS and FMCD. Agulló, Croux and Van Aelst (2008) use the F-MLTS estimator. Kudraszow and Maronna (2011) use F-MM-type estimators. Also see Maronna and Morgenthaler (1986) for using a robust estimator of multivariate location and dispersion estimator to make a

robust regression estimator when  $m = 1$ .

## References

- Agulló, J., Croux, C. and Van Aelst, S. (2008), "The Multivariate Least-Trimmed Squares Estimator," *Journal of Multivariate Analysis*, 99, 311-338.
- Anderson, T.W. (1984), *An Introduction to Multivariate Statistical Analysis*, 2nd ed., Wiley, New York, NY.
- Berndt, E.R. and Savin, N.E. (1977), "Conflict Among Criteria for Testing Hypotheses in the Multivariate Linear Regression Model," *Econometrica*, 45, 1263-1277.
- Buxton, L.H.D. (1920), "The Anthropology of Cyprus," *The Journal of the Royal Anthropological Institute of Great Britain and Ireland*, 50, 183-235.
- Cook, R.D., and Weisberg, S. (1999), *Applied Regression Including Computing and Graphics*, Wiley, New York, NY.
- Hampel, F.R. (1975), "Beyond Location Parameters: Robust Concepts and Methods," *Bulletin of the International Statistical Institute*, 46, 375-382.
- Henderson, H.V., and Searle, S.R. (1977), "Vec and Vech Operators for Matrices, with Some Uses in Jacobians and Multivariate Statistics," *The Canadian Journal of Statistics*, 7, 65-81.
- Hössjer, O. (1991), *Rank-Based Estimates in the Linear Model with High Breakdown Point*, Ph.D. Thesis, Report 1991:5, Department of Mathematics, Uppsala University, Uppsala, Sweden.
- Johnson, M.E. (1987), *Multivariate Statistical Simulation*, Wiley, New York, NY.
- Kakizawa, Y. (2009), "Third-order Power Comparisons for a Class of Tests for Multivariate Linear Hypothesis Under General Distributions," *Journal of Multivariate Analysis*, 100, 473-496.
- Khattree, R., and Naik, D.N. (1999), *Applied Multivariate Statistics with SAS Software*, 2nd ed., SAS Institute, Cary, NC.
- Kudraszow, N.L., and Maronna, R.A. (2011), "Estimates of MM Type for the Multivariate Linear Model," *Journal of Multivariate Analysis*, 102, 1280-1292.
- Kshirsagar, A.M. (1972), *Multivariate Analysis*, Marcel Dekker, New York, NY.
- Maronna, R.A., and Morgenthaler, S. (1986), "Robust Regression Through Robust Covariances," *Communications in Statistics: Theory and Methods*, 15, 1347-1365.
- Olive, D.J. (2002), "Applications of Robust Distances for Regression," *Technometrics*, 44, 64-71.
- Olive, D.J. (2005), "Two Simple Resistant Regression Estimators," *Computational Statistics and Data Analysis*, 49, 809-819.
- Olive, D.J. (2013a), "Asymptotically Optimal Regression Prediction Intervals and Prediction Regions for Multivariate Data," *International Journal of Statistics and Probability*, 2, 90-100.
- Olive, D.J. (2013b), "Plots, Prediction and Testing for the Multivariate Linear Model," preprint, see ([www.math.siu.edu/olive/ppmultreg.pdf](http://www.math.siu.edu/olive/ppmultreg.pdf)).
- Olive, D.J., and Hawkins, D.M. (2005), "Variable Selection for 1D Regression Models," *Technometrics*, 47, 43-50.
- Olive, D.J., and Hawkins, D.M. (2010), "Robust Multivariate Location and Dispersion," Preprint, see ([www.math.siu.edu/olive/preprints.htm](http://www.math.siu.edu/olive/preprints.htm)).

- Olive, D.J., and Hawkins, D.M. (2011), “Practical High Breakdown Regression,” Preprint, see ([www.math.siu.edu/olive/preprints.htm](http://www.math.siu.edu/olive/preprints.htm)).
- R Development Core Team (2011), “R: a Language and Environment for Statistical Computing,” R Foundation for Statistical Computing, Vienna, Austria, ([www.R-project.org](http://www.R-project.org)).
- Rousseeuw, P.J. (1984), “Least Median of Squares Regression,” *Journal of the American Statistical Association*, 79, 871-880.
- Rousseeuw, P.J., Van Aelst, S., Van Driessen, K., and Agulló, J. (2004), “Robust Multivariate Regression,” *Technometrics*, 46, 293-305.
- Rousseeuw, P.J., and Van Driessen, K. (1999), “A Fast Algorithm for the Minimum Covariance Determinant Estimator,” *Technometrics*, 41, 212-223.
- Seber, G.A.F., and Lee, A.J. (2003), *Linear Regression Analysis*, 2nd ed., Wiley, New York, NY.
- Su, Z., and Cook, R.D. (2012), “Inner Envelopes: Efficient Estimation in Multivariate Linear Regression,” *Biometrika*, 99, 687-702.
- Zhang, J., Olive, D.J., and Ye, P. (2012), “Robust Covariance Matrix Estimation With Canonical Correlation Analysis,” *International Journal of Statistics and Probability*, 1, 119-136.