

Analysis of residuals and influence

The methods for obtaining estimates, tests and other summaries developed so far tell only half the story of regression analysis. All of these methods are computed as if the model and assumptions are correct. But, in any practical problem, some assumptions used in regression analysis are in doubt. A second phase of analysis designed to check assumptions and to build a model is usually required. In this chapter we will study methods for detecting outliers and influential observations.

1 Residuals

The residuals provide information regarding assumptions about error terms and the appropriateness of the model. Any complete data analysis requires examination of the residuals. Here we will present the outline of analysis of residuals. We will look at various residuals such as:

1. Raw residuals: $\hat{\epsilon}_i = Y_i - \hat{Y}_i$.
2. Semi-Studentized residuals: $\hat{\epsilon}_i^* = \hat{\epsilon}_i / s$.
3. PRESS residuals: $\hat{\epsilon}_{(i)} = Y_i - \hat{Y}_{(i)}$.
4. Standardized residuals: $r_i = \hat{\epsilon}_i / (s\sqrt{1 - h_{ii}})$.

(textbook: Studentized residuals, Internally Studentized residuals).

5. Studentized residuals (Jackknifed residuals): $r_{(i)} = \hat{\epsilon}_i / (s_{(i)} \sqrt{1 - h_{ii}})$.

(textbook: Studentized deleted residuals, Externally Studentized residuals).

1.1 Raw residuals

Usual regression model is $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ with $\boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$, so the *true residuals*, $\boldsymbol{\epsilon} = (\epsilon_1, \dots, \epsilon_n)'$, are *iid* $N(0, \sigma^2)$. But the true residuals can not be obtained because we do not know the true value of $\boldsymbol{\beta}$. We can only estimate $\boldsymbol{\epsilon}$ by

$$\hat{\boldsymbol{\epsilon}} = \mathbf{Y} - \hat{\mathbf{Y}} = \mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}},$$

which I call the (estimated) *raw residuals*. The raw residuals $\hat{\boldsymbol{\epsilon}} = (\hat{\epsilon}_1, \dots, \hat{\epsilon}_n)'$ are consistent since $\hat{\boldsymbol{\beta}} \xrightarrow{P} \boldsymbol{\beta}$ as $n \rightarrow \infty$. Thus, if n is very large relative to the number of parameters p , the raw residuals are essentially *iid* $N(0, \sigma^2)$. In small to moderate samples, however, they are not. Recall $\hat{\boldsymbol{\epsilon}} = \mathbf{Y} - \hat{\mathbf{Y}} = (\mathbf{I} - \mathbf{H})\mathbf{Y}$, which implies

$$\hat{\boldsymbol{\epsilon}} \sim N(\mathbf{0}, \sigma^2(\mathbf{I} - \mathbf{H})).$$

This tells us the followings.

1. The raw residuals $\hat{\epsilon}_1, \dots, \hat{\epsilon}_n$ are not independent because the off-diagonal elements of $(\mathbf{I} - \mathbf{H})$ are not zero.
2. The raw residuals are not identically distributed because their variances

$$\text{Var}(\hat{\epsilon}_i) = \sigma^2(1 - h_{ii})$$

are not equal. Here, h_{ij} is the (i, j) th element of the hat matrix \mathbf{H} , *i.e.*, $\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' = [h_{ij}]$. Notice that the raw residuals have a mean of zero.

3. The variance $\text{Var}(\hat{\epsilon}_i)$ decreases as x_i moves away from the center of X -range. For example, in a simple linear regression with intercept,

$$h_{ij} = \frac{1}{n} + \frac{(x_i - \bar{x})(x_j - \bar{x})}{S_{xx}}. \quad (10.1)$$

Especially when $i = j$, the h_{ii} is called the *leverage* value, or the *potential* value. As x_i moves away from \bar{x} , the leverage value of h_{ii} increases. Thus, $\text{Var}(\hat{\epsilon}_i)$ decreases as x_i moves away from the center in X -range. The leverage is “a measure of how far from the center in X -range.” This generalizes to multiple linear regression. Points with high leverage tend to decrease residuals.

4. For models with an intercept, we have $\mathbf{H}\mathbf{1} = \mathbf{1}$, or in scalar form:

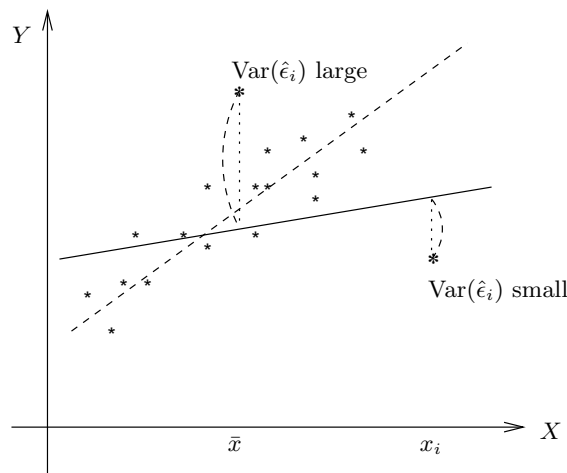
$$\sum_{i=1}^n h_{ij} = \sum_{j=1}^n h_{ij} = 1.$$

As x_i moves far away from \bar{x} , the term $(x_i - \bar{x})^2/S_{xx}$ gets close to one. Thus, as can be seen from $\text{Var}(\hat{\epsilon}_i) = \sigma^2(1 - h_{ii})$, with large values of h_{ii} (*i.e.*, x_i moves far away from \bar{x}), $\text{Var}(\hat{\epsilon}_i)$ will have a small value. We can also point this out using a scalar form of $\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{H}\mathbf{Y}$:

$$\hat{Y}_i = \sum_{j=1}^n h_{ij}Y_j = h_{ii}Y_i + \sum_{j \neq i}^n h_{ij}Y_j.$$

In combination with $h_{ij} \approx 0$ for large n (see Eq. (10.1)) and $\sum_{j=1}^n h_{ij} = 1$, this shows that as the leverage h_{ii} approaches 1 (*i.e.*, $\mathbf{H} \rightarrow \mathbf{I}$), the fitted value \hat{Y}_i approaches Y_i (*i.e.*, $\hat{\epsilon}_i = Y_i - \hat{Y}_i \rightarrow 0$).

Raw residuals $\hat{\epsilon}_i$ with large and small leverages



5. For very large n , all $h_{ij} \approx 0$. Thus, if $n \gg p$, then we can usually ignore the dependencies of $\hat{\epsilon}_i$.

Example 10.1. Illustration of Hat Matrix. See Table 10.2 on Page 393 of the textbook.

Minitab

Read Data

```

1 MTB > READ c1 c2 c11
2 DATA> 14 25 301
3 DATA> 19 32 327
4 DATA> 12 22 246
5 DATA> 11 15 187
6 DATA> END
7 4 rows read.
8 MTB > name c1 'X1'
9 MTB > name c2 'X2'
10 MTB > name c11 'Y'
```

$$\text{Model: } Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon$$

```

1 MTB > REGR c11 2 c1 c2 ;
2 SUBC> fits c12 ;
3 SUBC> residuals c21 ;
4 SUBC> mse k1 ;
5 SUBC> hi c24 .
6
7 Regression Analysis: Y versus X1, X2
8
9 The regression equation is
10 Y = 80.9 - 5.8 X1 + 11.3 X2
11 Predictor      Coef  SE Coef      T      P
12 Constant      80.93   57.94    1.40  0.396
13 X1             -5.84   11.74   -0.50  0.706
14 X2            11.325   5.931    1.91  0.307
15
16 S = 23.9768   R-Sq = 95.0%   R-Sq(adj) = 85.1%
17
18 Analysis of Variance
19 Source          DF      SS      MS      F      P
20 Regression       2  11009.9  5504.9   9.58  0.223
21 Residual Error   1    574.9   574.9
22 Total           3  11584.7
23
24 Source  DF  Seq SS
25 X1       1  8913.8
26 X2       1  2096.1
27
28 Unusual Observations
29 Obs    X1      Y    Fit  SE Fit  Residual  St Resid
30   4  11.0  187.0  186.5   24.0      0.5      1.00 X
31 X denotes an observation whose X value gives it large leverage.
32
33 Residual Plots for Y
34 MTB > Let c33 = k1 * (1-c24)
35 MTB > print c1 c2 c11 c12 c21 c24 c33
36
37 Data Display
38 Row  X1  X2  Y      C12      C21      C24      C33
39   1  14  25  301  282.238  18.7621  0.387681  352.016
40   2  19  32  327  332.292  -5.2919  0.951288  28.004
41   3  12  22  246  259.951 -13.9513  0.661433  194.638
42   4  11  15  187  186.519   0.4811  0.999597   0.231

```

R

④ Read Data

```

1 > X1 = c(14, 19, 12, 11)
2 > X2 = c(25, 32, 22, 15)
3 > Y = c(301, 327, 246, 187)

```

④ Model: $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon$

```

1 > source("https://raw.githubusercontent.com/AppliedStat/LM/master/Diagnostics.R")
2 > LM = lm (Y ~ X1 + X2)
3 > summary(LM)
4
5 Call:
6 lm(formula = Y ~ X1 + X2)
7
8 Coefficients:
9             Estimate Std. Error t value Pr(>|t|)
10 (Intercept)   80.930     57.944   1.397   0.396
11 X1            -5.845     11.745  -0.498   0.706
12 X2            11.325     5.931   1.909   0.307
13
14 Residual standard error: 23.98 on 1 degrees of freedom

```

```

15 Multiple R-Squared: 0.9504, Adjusted R-squared: 0.8511
16 F-statistic: 9.576 on 2 and 1 DF, p-value: 0.2228
17
18 > X = model.matrix(LM)
19 > X
20 (Intercept) X1 X2
21 1          1 14 25
22 2          1 19 32
23 3          1 12 22
24 4          1 11 15
25 attr(,"assign")
26 [1] 0 1 2
27
28 > Y.fit = fitted(LM)
29
30 > e      = resid(LM)
31
32 > hii     = hatvalues(LM)
33
34 > mse = MSE(LM)
35 > mse
36 [1] 574.8893
37
38 > s2 = mse * (1-hii)
39
40 > cbind(X1, X2, Y, Y.fit, e, hii, s2)
41   X1 X2   Y   Y.fit      e    hii      s2
42 1 14 25 301 282.2379 18.7620773 0.3876812 352.0155444
43 2 19 32 327 332.2919 -5.2918680 0.9512882 28.0038665
44 3 12 22 246 259.9513 -13.9512882 0.6614332 194.6384437
45 4 11 15 187 186.5189  0.4810789 0.9995974  0.2314369
46 >
47 > H = X %%% solve( t(X) %%% X ) %%% t(X)
48
49 > Cov.e = mse * ( diag(1, length(Y)) - H )
50
51 > Y.hat = H %%% Y
52
53 > cbind(Y.fit, Y.hat)
54   Y.fit
55 1 282.2379 282.2379
56 2 332.2919 332.2919
57 3 259.9513 259.9513
58 4 186.5189 186.5189
59
60 # Note: trace(Hat matrix) = p
61 > sum( hii )
62 [1] 3

```

||

1.2 Internally Studentized residuals

The detection of outlying or extreme Y observations based on an examination of the residuals has been considered in earlier chapters. We used the *raw residuals* given by

$$\hat{\epsilon}_i = Y_i - \hat{Y}_i, \quad i = 1, \dots, n$$

or the *semi-Studentized residuals* given by

$$\hat{\epsilon}_i^* = \frac{\hat{\epsilon}_i}{\sqrt{\text{MSE}}} = \frac{\hat{\epsilon}_i}{s},$$

where $s = \sqrt{\text{MSE}}$. The raw residuals and semi-Studentized residuals have some of the difficulty in detecting outliers when the leverages are high. The variance of $\hat{\epsilon}_i$ is $\text{Var}(\hat{\epsilon}_i) = \sigma^2(1 - h_{ii})$ and this can be estimated by $s^2(1 - h_{ii})$. Thus, it is better to use

$$r_i = \frac{\hat{\epsilon}_i}{s\sqrt{1 - h_{ii}}}$$

which have the mean 0 and the variance 1 approximately. We call these r_i the *internally Studentized residuals*. Some textbooks, Minitab and R language call them *standardized residuals*.”

Remark 10.1.

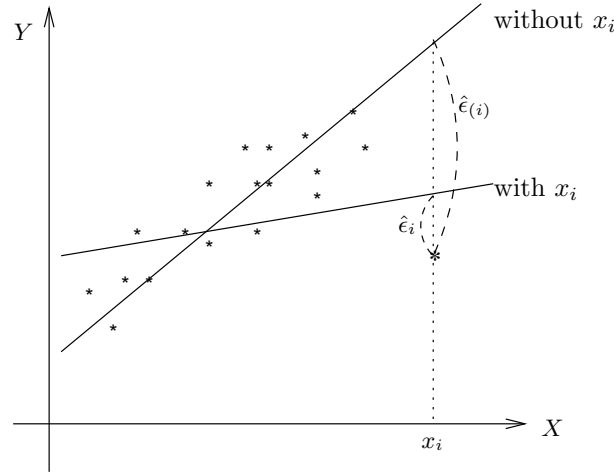
1. The internally Studentized residuals r_i are *identically distributed*, but still not independent.
2. The distribution of r_i is something like t -distribution with $\text{df} = n - p$ because σ^2 is replaced by s^2 . But it is not exactly the t -distribution because the numerator and denominator are not independent. Note that we used the p normal equations to estimate the parameters $\beta_0, \dots, \beta_{p-1}$. Hence if $n \gg p$, then we can usually ignore the dependencies of r_i .
3. For very large $n \gg p$, all $h_{ij} \approx 0$ (of course, the leverage $h_{ii} \approx 0$ also) and the r_i 's are nearly proportional to $\hat{\epsilon}_i$'s. Thus, for large samples, plots of r_1, \dots, r_n look nearly the same as plots of $\hat{\epsilon}_1^*, \dots, \hat{\epsilon}_n^*$.

△

1.3 PRESS residuals

The PRESS residual is defined as $\hat{\epsilon}_{(i)} = Y_i - \hat{Y}_{(i)}$ where $\hat{Y}_{(i)}$ is calculated with x_i omitted from the regression fit. These measure the prediction errors. Large values indicate that points are far from what the model predicts. Using the following theorem, we can calculate the PRESS residuals from the ordinary residuals.

Raw residuals $\hat{\epsilon}_i$ and PRESS residuals $\hat{\epsilon}_{(i)}$



Lemma 10.1 (Sherman-Morrison-Woodbury). Let \mathbf{U} and \mathbf{V} be $m \times k$ matrices, and \mathbf{A} be an $m \times m$ square matrix. Then we have

$$(\mathbf{A} + \mathbf{U}\mathbf{V}')^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{U}(\mathbf{I}_k + \mathbf{V}'\mathbf{A}^{-1}\mathbf{U})^{-1}\mathbf{V}'\mathbf{A}^{-1},$$

where \mathbf{I}_k is an identity matrix.

Proof. See § 2.1.3 of Golub and Van Loan (1996). □

Theorem 10.2. The PRESS residual is obtained as

$$\hat{\epsilon}_{(i)} = Y_i - \hat{Y}_{(i)} = \frac{\hat{\epsilon}_i}{1 - h_{ii}},$$

where h_{ii} is the i th diagonal element of the hat matrix, \mathbf{H} .

Proof. Let \mathbf{x}'_i be the i th row of the data matrix \mathbf{X} and $\mathbf{X}_{(i)}$ be the data matrix without the use of the i th observed data. Similarly, let $\mathbf{Y}_{(i)}$ be the column vector with Y_i omitted. Then we can easily show that

$$\mathbf{X}'_{(i)}\mathbf{X}_{(i)} = \mathbf{X}'\mathbf{X} - \mathbf{x}_i\mathbf{x}'_i \quad (10.2)$$

$$\mathbf{X}'_{(i)}\mathbf{Y}_{(i)} = \mathbf{X}'\mathbf{Y} - \mathbf{x}_iY_i. \quad (10.3)$$

It is immediate upon using Lemma 10.1 that the inverse of (10.2) is

$$\begin{aligned}
(\mathbf{X}'_{(i)}\mathbf{X}_{(i)})^{-1} &= (\mathbf{X}'\mathbf{X})^{-1} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}_i\left(1 - \mathbf{x}'_i(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}_i\right)^{-1}\mathbf{x}'_i(\mathbf{X}'\mathbf{X})^{-1} \\
&= (\mathbf{X}'\mathbf{X})^{-1} + \frac{(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}_i\mathbf{x}'_i(\mathbf{X}'\mathbf{X})^{-1}}{1 - \mathbf{x}'_i(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}_i} \\
&= (\mathbf{X}'\mathbf{X})^{-1} + \frac{(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}_i\mathbf{x}'_i(\mathbf{X}'\mathbf{X})^{-1}}{1 - h_{ii}}.
\end{aligned} \tag{10.4}$$

Let $\hat{\boldsymbol{\beta}}_{(i)}$ be the vector of the estimated regression coefficients with the i th observed data omitted. Then we have

$$\hat{\boldsymbol{\beta}}_{(i)} = (\mathbf{X}'_{(i)}\mathbf{X}_{(i)})^{-1}\mathbf{X}'_{(i)}\mathbf{Y}_{(i)}.$$

Then we have

$$\hat{Y}_{(i)} = \mathbf{x}'_i\hat{\boldsymbol{\beta}}_{(i)} = \mathbf{x}'_i(\mathbf{X}'_{(i)}\mathbf{X}_{(i)})^{-1}\mathbf{X}'_{(i)}\mathbf{Y}_{(i)}. \tag{10.5}$$

Substituting (10.4) into (10.5) with (10.3) gives

$$\begin{aligned}
\hat{Y}_{(i)} &= \mathbf{x}'_i \left[(\mathbf{X}'\mathbf{X})^{-1} + \frac{(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}_i\mathbf{x}'_i(\mathbf{X}'\mathbf{X})^{-1}}{1 - h_{ii}} \right] \mathbf{X}'_{(i)}\mathbf{Y}_{(i)} \\
&= \left[\mathbf{x}'_i(\mathbf{X}'\mathbf{X})^{-1} + \frac{h_{ii}\mathbf{x}'_i(\mathbf{X}'\mathbf{X})^{-1}}{1 - h_{ii}} \right] (\mathbf{X}'\mathbf{Y} - \mathbf{x}_iY_i) \\
&= \frac{1}{1 - h_{ii}} \mathbf{x}'_i(\mathbf{X}'\mathbf{X})^{-1}(\mathbf{X}'\mathbf{Y} - \mathbf{x}_iY_i) \\
&= \frac{\hat{Y}_i - h_{ii}Y_i}{1 - h_{ii}}.
\end{aligned}$$

Thus, the PRESS residual is given by

$$\hat{\epsilon}_{(i)} = Y_i - \hat{Y}_{(i)} = Y_i - \frac{\hat{Y}_i - h_{ii}Y_i}{1 - h_{ii}} = \frac{Y_i - \hat{Y}_i}{1 - h_{ii}} = \frac{\hat{\epsilon}_i}{1 - h_{ii}}.$$

□

Remark 10.2.

1. For very large $n \gg p$, all $h_{ii} \approx 0$ and $\hat{\epsilon}_{(i)} \approx \hat{\epsilon}_i$.
2. $\hat{\epsilon}_{(i)}$ far from $\hat{\epsilon}_i$ indicates influential point.

3. The PRESS residuals have unequal variances:

$$\text{Var}(\hat{\epsilon}_{(i)}) = \frac{1}{(1 - h_{ii})^2} \text{Var}(\hat{\epsilon}_i) = \frac{1}{(1 - h_{ii})^2} \sigma^2 (1 - h_{ii}) = \frac{\sigma^2}{1 - h_{ii}}.$$

Dividing $\hat{\epsilon}_{(i)}$ by the estimated standard deviation $s/\sqrt{1 - h_{ii}}$ gives

$$r_i = \frac{\hat{\epsilon}_i}{s\sqrt{1 - h_{ii}}},$$

which is the internally Studentized residual. Thus, the standardized (internally Studentized) residuals r_1, \dots, r_n can also be thought of as standardized PRESS residuals.

△

1.4 Externally Studentized residuals

Recall that the internally Studentized residual

$$r_i = \frac{\hat{\epsilon}_i}{s\sqrt{1 - h_{ii}}}, \quad i = 1, \dots, n$$

is not exactly t -distributed because the numerator and denominator are dependent. But if we replace $s = \sqrt{\text{MSE}}$ by $s_{(i)} = \sqrt{\text{MSE}_{(i)}}$, where $\text{MSE}_{(i)}$ is the MSE from the model fit without the i th observation. It follows that

$$r_{(i)} = \frac{\hat{\epsilon}_i}{s_{(i)}\sqrt{1 - h_{ii}}}, \quad i = 1, \dots, n$$

which has a t -distribution with $(n - 1 - p)$ degrees of freedom because $\hat{\epsilon}_i$ and $s_{(i)}$ are independent. We will call these the *externally Studentized* residuals. The textbook calls these the Studentized deleted residuals and use t_i notation instead of $r_{(i)}$. **Minitab** and **R** language call them “Studentized residuals.”

Remark 10.3.

1. Each $r_{(i)}$ is distributed as t_{n-1-p} under the model. But $r_{(1)}, \dots, r_{(n)}$ are *not* independent.
2. The externally Studentized residuals are traditionally used for *outlier detection* with respect to Y values.
3. Under the model, we can test the hypothesis that a single observation deviates from the model

by comparing $r_{(i)}$ to t -distribution:

$$\begin{aligned} p\text{-value} &= 2 \times \text{Prob}[t_{n-1-p} \geq |r_{(i)}|] \\ &= 2 \times \left\{ 1 - \text{Prob}[t_{n-1-p} \leq |r_{(i)}|] \right\}, \end{aligned}$$

where t_{n-1-p} is the random variable having a t -distribution with $\text{df} = n - 1 - p$. Note that even if the model holds for every observation (*i.e.*, there are no outliers), one expects about 5% of the observations to have p -values less than 0.05 when the significance level $\alpha = 5\%$ is used. So, we should not automatically call all the observations with p -values below 0.05 outliers, especially when n is large. We can conduct a formal test by means of the Bonferroni test procedure. That is, if $|r_{(i)}| > t(1 - \alpha/(2n); n - 1 - p)$, then we conclude that the i th observation is an outlier.

4. Minitab subcommands:

Residuals	Subcommand
Raw	RESIDUALS C21; $\hat{e}_i = \text{C21}$
Internally Studentized	SRESIDUALS C22; $r_i = \text{C22}$
Externally Studentized	TRESIDUALS C23; $r_{(i)} = \text{C23}$
Leverage	HI C24; $h_{ii} = \text{C24}$

5. R functions:

Residuals	R Functions	Package
Raw	<code>resid()</code>	intrinsic
Semi-Studentized	<code>semiresid()</code>	Class Web
Internally Studentized	<code>rstandard()</code>	intrinsic
	<code>stdres()</code>	MASS
Externally Studentized	<code>rstudent()</code>	intrinsic
	<code>studres()</code>	MASS
leverage	<code>hatvalues()</code>	intrinsic

Class Web: <https://github.com/AppliedStat/LM/blob/master/Diagnostics.R>

△

Example 10.2. Residuals, Diagonal of Hat Matrix, Studentized Deleted Residuals: Body Fat Data in Table 7.1 with Two Predictors (X_1 and X_2).

Minitab

Read Data

```
1 MTB > read c1 c2 c3 c11 ;
2 SUBC>      file "S:\LM\CH07TA01.txt" .
3 Entering data from file: S:\LM\CH07TA01.TXT
4 20 rows read.
5 MTB > name c1 'X1'
6 MTB > name c2 'X2'
7 MTB > name c11 'Y'
```

Model: $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon$

```
1 MTB > REGR c11 2 c1 c2 ;
2 SUBC> residuals c21 ;
3 SUBC> sresiduals c22 ;
4 SUBC> tresiduals c23 ;
5 SUBC> hi          c24 .
6
7 Regression Analysis: Y versus X1, X2
8 The regression equation is
9 Y = - 19.2 + 0.222 X1 + 0.659 X2
10
11 Predictor      Coef    SE Coef      T      P
12 Constant     -19.174    8.361   -2.29   0.035
13 X1             0.2224    0.3034    0.73   0.474
14 X2             0.6594    0.2912    2.26   0.037
15
16 S = 2.54317    R-Sq = 77.8%    R-Sq(adj) = 75.2%
17
18 Analysis of Variance
19 Source          DF      SS      MS      F      P
20 Regression        2   385.44   192.72   29.80   0.000
21 Residual Error   17   109.95    6.47
22 Total            19   495.39
23
24 Source   DF   Seq SS
25 X1        1   352.27
26 X2        1    33.17
27
28 Residual Plots for Y
29
30 MTB > print c21 c24 c23 c22
31 Data Display
32 Row      C21      C24      C23      C22
33 1    -1.68271  0.201013  -0.72999  -0.74023
34 2     3.64293  0.058895   1.53425   1.47658
35 3    -3.17597  0.371933  -1.65433  -1.57579
36 4    -3.15847  0.110940  -1.34848  -1.31715
37 5    -0.00029  0.248010  -0.00013  -0.00013
38 6    -0.36082  0.128616  -0.14755  -0.15199
39 7     0.71620  0.155517   0.29813   0.30645
40 8     4.01473  0.096288   1.76009   1.66061
41 9     2.65511  0.114636   1.11765   1.10955
42 10    -2.47481  0.110244  -1.03373  -1.03165
43 11     0.33581  0.120337   0.13666   0.14078
44 12     2.22551  0.109266   0.92318   0.92722
45 13    -3.94686  0.178382  -1.82590  -1.71215
46 14     3.44746  0.148007   1.52476   1.46861
47 15     0.57059  0.333212   0.26715   0.27476
48 16     0.64230  0.095277   0.25813   0.26552
49 17    -0.85095  0.105595  -0.34451  -0.35380
50 18    -0.78292  0.196793  -0.33441  -0.34350
51 19    -2.85729  0.066954  -1.17617  -1.16313
52 20     1.04045  0.050085   0.40936   0.41976
```

R

Read Data

```

1 > mydata =
  read.table("https://raw.githubusercontent.com/AppliedStat/LM/master/CH07TA01.txt")
2 > x1 = mydata[,1]
3 > x2 = mydata[,2]
4 > y = mydata[,4]

```

Ⓜ Model: $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon$

```

1 > source("https://raw.githubusercontent.com/AppliedStat/LM/master/Diagnostics.R")
2 > LM = lm ( y ~ x1 + x2 )
3
4 > e.raw = resid(LM)
5 > e.semi.Student = semiresid(LM)
6 > e.int.Student = rstandard(LM)
7 > e.ext.Student = rstudent(LM)
8 > hat.diagonal = hatvalues(LM)
9
10 > round( cbind( e.raw, hat.diagonal, e.ext.Student, e.int.Student, e.semi.Student ),
11          3 )

```

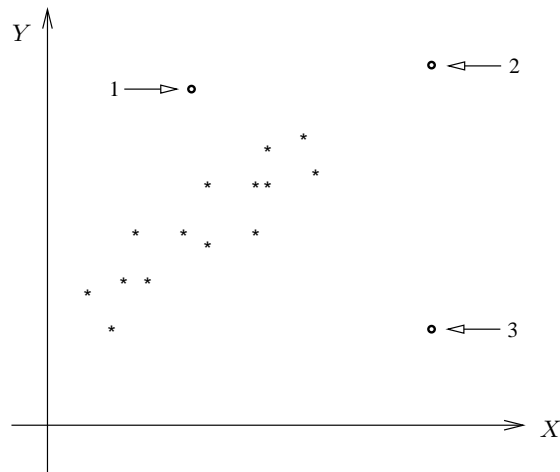
	e.raw	hat.diagonal	e.ext.Student	e.int.Student	e.semi.Student
1	-1.683	0.201	-0.730	-0.740	-0.662
2	3.643	0.059	1.534	1.477	1.432
3	-3.176	0.372	-1.654	-1.576	-1.249
4	-3.158	0.111	-1.348	-1.317	-1.242
5	0.000	0.248	0.000	0.000	0.000
6	-0.361	0.129	-0.148	-0.152	-0.142
7	0.716	0.156	0.298	0.306	0.282
8	4.015	0.096	1.760	1.661	1.579
9	2.655	0.115	1.118	1.110	1.044
10	-2.475	0.110	-1.034	-1.032	-0.973
11	0.336	0.120	0.137	0.141	0.132
12	2.226	0.109	0.923	0.927	0.875
13	-3.947	0.178	-1.826	-1.712	-1.552
14	3.447	0.148	1.525	1.469	1.356
15	0.571	0.333	0.267	0.275	0.224
16	0.642	0.095	0.258	0.266	0.253
17	-0.851	0.106	-0.345	-0.354	-0.335
18	-0.783	0.197	-0.334	-0.344	-0.308
19	-2.857	0.067	-1.176	-1.163	-1.124
20	1.040	0.050	0.409	0.420	0.409

||

2 Measures of Influence

“Influence” refers to the impact of a particular observation on the model. If a single suspect observation changes our conclusions, then our conclusions are not trustworthy. We shall consider an observation to be influential if its exclusion causes major changes in the fitted regression.

Scatter plot to illustrate outlier, leverage and influence



Observation	Outlier (Y-direction)	Leverage (Outlier in X -dir)	Influential
1	✓		
2		✓	
3	✓	✓	✓

2.1 Diagonals of hat matrix

The hat matrix $\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ has the following properties:

1. Symmetric: $\mathbf{H}' = \mathbf{H}$.
2. Idempotent: $\mathbf{H}\mathbf{H} = \mathbf{H}$.
3. $0 \leq h_{ii} \leq 1$ for every i .

With the intercept, $1/n \leq h_{ii} \leq 1$.

4. $\text{tr}(\mathbf{H}) = \sum_{i=1}^n h_{ii} = \text{rank}(\mathbf{X}) = p$.

In the special case of simple linear regression, we have

$$h_{ii} = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{S_{xx}}$$

$$\sum_{i=1}^n h_{ii} = \sum_{i=1}^n \frac{1}{n} + \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{S_{xx}} = 1 + 1 = 2.$$

As observations move away from the center of X -range, the leverages h_{ii} go up. Large h_{ii} indicates that an observation is *potentially* influential. The average of h_{ii} is p/n , so values of h_{ii} exceeding $2p/n$ are considered to be high leverages.

2.2 DFFITS

This statistic measures how much the fitted value for the i th observation changes when all n observations are used in fitting the regression function and when the i th observation is omitted. Denote \hat{Y}_i as the fitted value for the i th observation using all n observations and $\hat{Y}_{(i)}$ as the fitted value for the i th observation with i th observation omitted. DFFITS stands for the difference between the fitted values. The DFFITS_i is defined as

$$\text{DFFITS}_i = \frac{\hat{Y}_i - \hat{Y}_{(i)}}{s_{(i)}\sqrt{h_{ii}}},$$

where $s_{(i)} = \sqrt{\text{MSE}_{(i)}}$. Using $\hat{\epsilon}_{(i)} = Y_i - \hat{Y}_{(i)} = \frac{\hat{\epsilon}_i}{1-h_{ii}}$, we have $\hat{Y}_{(i)} = Y_i - \frac{\hat{\epsilon}_i}{1-h_{ii}}$. It follows that

$$\text{DFFITS}_i = \frac{\hat{Y}_i - Y_i + \frac{\hat{\epsilon}_i}{1-h_{ii}}}{s_{(i)}\sqrt{h_{ii}}} = \frac{\hat{\epsilon}_i \frac{h_{ii}}{1-h_{ii}}}{s_{(i)}\sqrt{h_{ii}}} = \frac{\hat{\epsilon}_i}{s_{(i)}\sqrt{1-h_{ii}}} \sqrt{\frac{h_{ii}}{1-h_{ii}}} = r_{(i)} \sqrt{\frac{h_{ii}}{1-h_{ii}}},$$

where $r_{(i)}$ is the externally Studentized residual. DFFITS_i is thus a residual, inflated or shrunk by leverage.

As a guideline for identifying influential cases, we suggest considering an observation influential if the $|\text{DFFITS}_i|$ exceeds 1 for small to medium data sets and $2\sqrt{p/n}$ for large data sets (say, $n \geq 30$).

DFFITS_i combines leverage h_{ii} and externally Studentized residual $r_{(i)}$ into one overall measure of how unusual an observation is.

2.3 Cook's distance

In contrast to the DFFITS_i which considers the influence of the i th observation on the fitted value \hat{Y}_i , Cook's distance considers the influence of the i th observation on all n fitted values. Cook's distance is defined as

$$D_i = \frac{\sum_{j=1}^n (\hat{Y}_j - \hat{Y}_{j(i)})^2}{p \cdot \text{MSE}},$$

where $\hat{Y}_{j(i)}$ is the fitted value for the j th observation with the i th observation omitted.

Using matrix notation, it can be expressed as

$$D_i = \frac{(\hat{\mathbf{Y}} - \hat{\mathbf{Y}}_{(i)})'(\hat{\mathbf{Y}} - \hat{\mathbf{Y}}_{(i)})}{p \cdot \text{MSE}} = \frac{(\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}_{(i)})'(\mathbf{X}'\mathbf{X})(\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}_{(i)})}{p \cdot \text{MSE}},$$

where $\hat{\mathbf{Y}}_{(i)}$ is the vector of the fitted values and $\hat{\boldsymbol{\beta}}_{(i)}$ is the vector of the estimated regression coefficients with the i th observation omitted. It has been found useful to relate D_i to the $F(p, n - p)$ distribution.

It has been suggested that observations with D_i values greater than the 50% percentile point of the F -distribution with p and $n - p$ degrees of freedom are classified as influential points. Because for most F -distributions, the 50% percentile point is near 1, the practical operational rule is to classify observations with $D_i > 1$ as being influential.

Fortunately, Cook's distance D_i can be calculated without fitting a new regression function each time a different observation is deleted. An algebraically equivalent expression is

$$D_i = \frac{\hat{\epsilon}_i^2}{p \cdot \text{MSE}} \cdot \frac{h_{ii}}{(1 - h_{ii})^2} = \frac{r_i^2}{p} \cdot \frac{h_{ii}}{1 - h_{ii}},$$

where r_i is the internally Studentized residual.

Cook's distance D_i combines leverage h_{ii} and internally Studentized residual r_i into one overall measure of how unusual an observation is.

2.4 DFBETAS

These statistics measure how much the values of the parameter estimates $(\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_{p-1})$ change when all n observations are used in estimating the regression parameters $(\beta_0, \beta_1, \dots, \beta_{p-1})$ and when the i th observation is omitted. We also standardize these statistics by dividing them by corresponding sample standard deviations. These measures, denoted by DFBETAS, are then defined by

$$\text{DFBETAS}_{k(i)} = \frac{\hat{\beta}_k - \hat{\beta}_{k(i)}}{\sqrt{\text{MSE}_{(i)} \cdot c_{kk}}},$$

where c_{kk} is the k th diagonal element of $(\mathbf{X}'\mathbf{X})^{-1}$, $k = 0, 1, 2, \dots, p-1$, and $i = 1, 2, \dots, n$.

The positive/negative sign of $\text{DFBETAS}_{k(i)}$ indicates that the i th observation leads to an increase/decrease in the k th parameter estimate and its absolute value indicates the amount of impact of the i th observation on the k th parameter estimate.

As a guideline for identifying influential cases, we suggest considering an observation influential if the $|\text{DFBETAS}_{k(i)}|$ exceeds 1 for small to medium data sets and $2/\sqrt{n}$ for large data sets (say, $n \geq 30$).

Example 10.3. DFFITS, Cook's distances, DFBETAS – Body Fat Data with two predictors. See Table 10.4 on Page 402.

Minitab

Read Data

```
1 MTB > read c1 c2 c3 c11 ;
2 SUBC>      file "S:\LM\CH07TA01.txt" .
3 Entering data from file: S:\LM\CH07TA01.TXT
4 20 rows read.
5 MTB > name c1 'X1'
6 MTB > name c2 'X2'
7 MTB > name c11 'Y'
```

Model: $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon$

```
1 MTB > REGR c11 2 c1 c2 ;
2 SUBC>   DFITS c31 ;
3 SUBC>   COOKD c32 .
4
5 Regression Analysis: Y versus X1, X2
6 The regression equation is
7 Y = - 19.2 + 0.222 X1 + 0.659 X2
8
9 Predictor      Coef  SE Coef      T      P
```

```

10 Constant    -19.174      8.361    -2.29    0.035
11 X1           0.2224     0.3034     0.73    0.474
12 X2           0.6594     0.2912     2.26    0.037
13
14 S = 2.54317    R-Sq = 77.8%    R-Sq(adj) = 75.2%
15
16 Analysis of Variance
17 Source      DF      SS      MS      F      P
18 Regression    2    385.44   192.72   29.80   0.000
19 Residual Error 17    109.95    6.47
20 Total         19    495.39
21
22 Source  DF  Seq SS
23 X1      1   352.27
24 X2      1    33.17
25
26 Residual Plots for Y
27
28 MTB > print c31 c32
29 Data Display
30 Row      C31      C32
31 1   -0.36615  0.045951
32 2    0.38381  0.045481
33 3   -1.27307  0.490157
34 4   -0.47635  0.072162
35 5   -0.00007  0.000000
36 6   -0.05669  0.001137
37 7    0.12794  0.005765
38 8    0.57452  0.097939
39 9    0.40216  0.053134
40 10   -0.36387  0.043957
41 11    0.05055  0.000904
42 12    0.32334  0.035154
43 13   -0.85078  0.212150
44 14    0.63551  0.124893
45 15    0.18885  0.012575
46 16    0.08377  0.002475
47 17   -0.11837  0.004926
48 18   -0.16553  0.009636
49 19   -0.31507  0.032360
50 20    0.09400  0.003097

```

R

④ Read Data

```

1 > mydata =
    read.table("https://raw.githubusercontent.com/AppliedStat/LM/master/CH07TA01.txt")
2 > x1 = mydata[,1]
3 > x2 = mydata[,2]
4 > y = mydata[,4]

```

④ Model: $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon$

```

1 > LM = lm ( y ~ x1 + x2 )
2
3 > DFF = dffits (LM)
4 > COOK = cooks.distance (LM)
5 > BETA = dfbetas (LM)
6 >
7 > round ( cbind(DFF, COOK, BETA), 3)
8      DFF  COOK (Intercept)      x1      x2
9 1  -0.366  0.046    -0.305 -0.131  0.232
10 2   0.384  0.045     0.173  0.115 -0.143
11 3  -1.273  0.490    -0.847 -1.183  1.067
12 4  -0.476  0.072    -0.102 -0.294  0.196
13 5   0.000  0.000     0.000  0.000  0.000
14 6  -0.057  0.001     0.040  0.040 -0.044

```

15	7	0.128	0.006	-0.078	-0.016	0.054
16	8	0.575	0.098	0.261	0.391	-0.332
17	9	0.402	0.053	-0.151	-0.295	0.247
18	10	-0.364	0.044	0.238	0.245	-0.269
19	11	0.051	0.001	-0.009	0.017	-0.002
20	12	0.323	0.035	-0.130	0.022	0.070
21	13	-0.851	0.212	0.119	0.592	-0.389
22	14	0.636	0.125	0.452	0.113	-0.298
23	15	0.189	0.013	-0.003	-0.125	0.069
24	16	0.084	0.002	0.009	0.043	-0.025
25	17	-0.118	0.005	0.080	0.055	-0.076
26	18	-0.166	0.010	0.132	0.075	-0.116
27	19	-0.315	0.032	-0.130	-0.004	0.064
28	20	0.094	0.003	0.010	0.002	-0.003

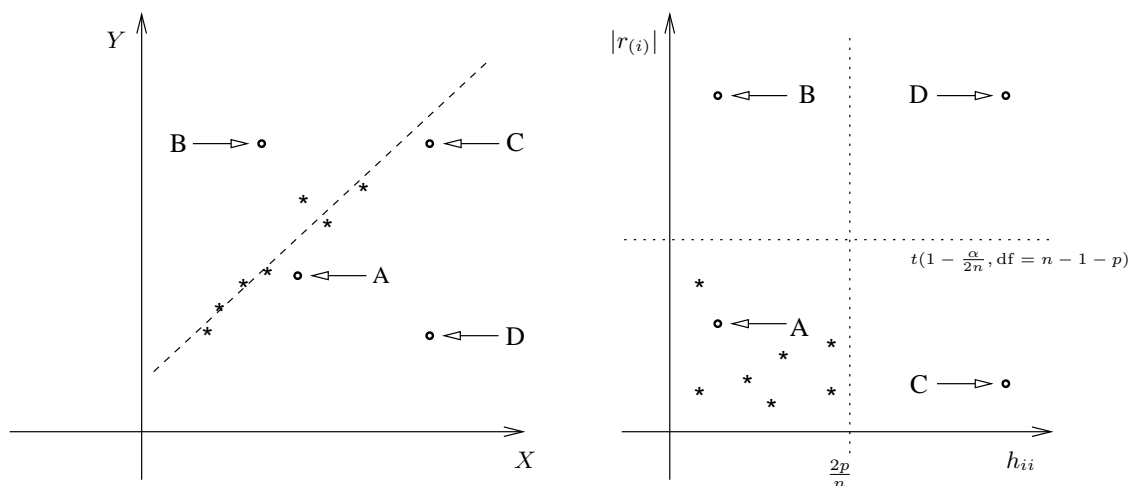
||

2.5 Strategy to find influential observations

Influential observations can be detected by finding observations which have *high leverage* values and are *outlying* with respect to Y . Cook's distance D_i and DFFITS $_i$ combines leverage h_{ii} and outlying measure into one. They mix together deviation in X -direction with deviation in Y -direction.

My personal suggestion is to look jointly at leverages h_{ii} and externally Studentized residuals $r_{(i)}$ in a plot of $r_{(i)}$ versus h_{ii} , or a plot of $|r_{(i)}|$ versus h_{ii} .

Scatter plot and $|r_{(i)}|$ vs. h_{ii} plot



Minitab Commands

Subcommands	Residual		Note
RESIDUALS C21;	raw	$\hat{\epsilon}_i = \text{C21}$	
SRESIDUALS C22;	internally Studentized	$r_i = \text{C22}$	
TRESIDUALS C23;	externally Studentized	$r_{(i)} = \text{C23}$	Y deviation
HI C24;	leverage	$h_{ii} = \text{C24}$	X deviation
COOKD C25;	Cook's distance	$D_i = \text{C25}$	X, Y mixed
DFITS C26;	DFFITS	$\text{DFFITS}_i = \text{C26}$	X, Y mixed

R functions

Subcommands	Residual		Note
<code>resid</code>	raw	$\hat{\epsilon}_i$	
<code>semiresid</code>	semi-Studentized	$\hat{\epsilon}_i^*$	
<code>rstandard, stdres</code>	internally Studentized	r_i	
<code>rstudent, studres</code>	externally Studentized	$r_{(i)}$	Y deviation
<code>hatvalues</code>	leverage	h_{ii}	X deviation
<code>cooks.distance</code>	Cook's distance	D_i	X, Y mixed
<code>dffits</code>	DFFITS	DFFITS_i	X, Y mixed
<code>dfbetas</code>	DFBETAS	$\text{DFBETAS}_{k(i)}$	X, Y mixed

References

Golub, G. H. and Van Loan, C. F. (1996). *Matrix Computations*. Johns Hopkins University Press, Baltimore and London, 3rd edition.