# 7 Regression Diagnostics

This section discusses some graphical and analytical regression diagnostic techniques for detecting outliers and assessing whether the assumptions of our regression model are met.

## 7.1 Leverage values

Leverage values $h_{ii}$ indicate how much influence an observation $\boldsymbol{X}_i$ has on the regression fit. They are calculated as

$$h_{ii} = \boldsymbol{X}_i'(\boldsymbol{X}'\boldsymbol{X})^{-1}\boldsymbol{X}_i$$

and represent the diagonal entries of the hat-matrix

$$\boldsymbol{P} = \boldsymbol{X}(\boldsymbol{X}'\boldsymbol{X})^{-1}\boldsymbol{X}'.$$

A low leverage implies the presence of many regressor observations similar to $\boldsymbol{X}_i$ in the sample, while a high leverage indicates a lack of similar observations near $\boldsymbol{X}_i$.

An observation with a high leverage $h_{ii}$ but a response value $Y_i$ that is close to the true regression line $\boldsymbol{X}_i'\boldsymbol{\beta}$ (indicating a small error $u_i$) is considered a **good leverage point**. It positively influences the model, especially in data-sparse regions.

Conversely, a **bad leverage point** occurs when both $h_{ii}$ and the error $u_i$ are large, indicating both unusual regressor and response values. This can misleadingly impact the regression fit.

The actual error term is unknown, but standardized residuals can be used to differentiate between good and bad leverage points.

## 7.2 Standardized residuals

Many regression diagnostic tools rely on the residuals of the OLS estimation $\hat{u}_i$ because they provide insight into the properties of the unknown error terms $u_i$.

Under the homoskedastic linear regression model (A1)–(A5), the errors are independent and have the property

$$Var[u_i|\boldsymbol{X}] = \sigma^2.$$

Since $PX = X$ and, therefore,

$$\hat{u} = (I_n - P)Y = (I_n - P)(X\beta + u) = (I_n - P)u,$$

the residuals have a different property:

$$Var[\hat{u}|X] = \sigma^2(I_n - P).$$

The $i$-th residual satisfies

$$Var[\hat{u}_i|X] = \sigma^2(1 - h_{ii}),$$

where $h_{ii}$ is the $i$-th leverage value.

Under the assumption (A5), the variance of $\hat{u}_i$ depends on $\mathbf{X}$, while the variance of $u_i$ does not. Dividing by $\sqrt{1 - h_{ii}}$ removes the dependency:

$$Var\left[\frac{\hat{u}_i}{\sqrt{1 - h_{ii}}}\bigg|X\right] = \sigma^2$$

The **standardized residuals** are defined as follows:

$$r_i := \frac{\hat{u}_i}{\sqrt{s_{\hat{u}}^2(1 - h_{ii})}}.$$

Standardized residuals are available using the R command `rstandard()`.
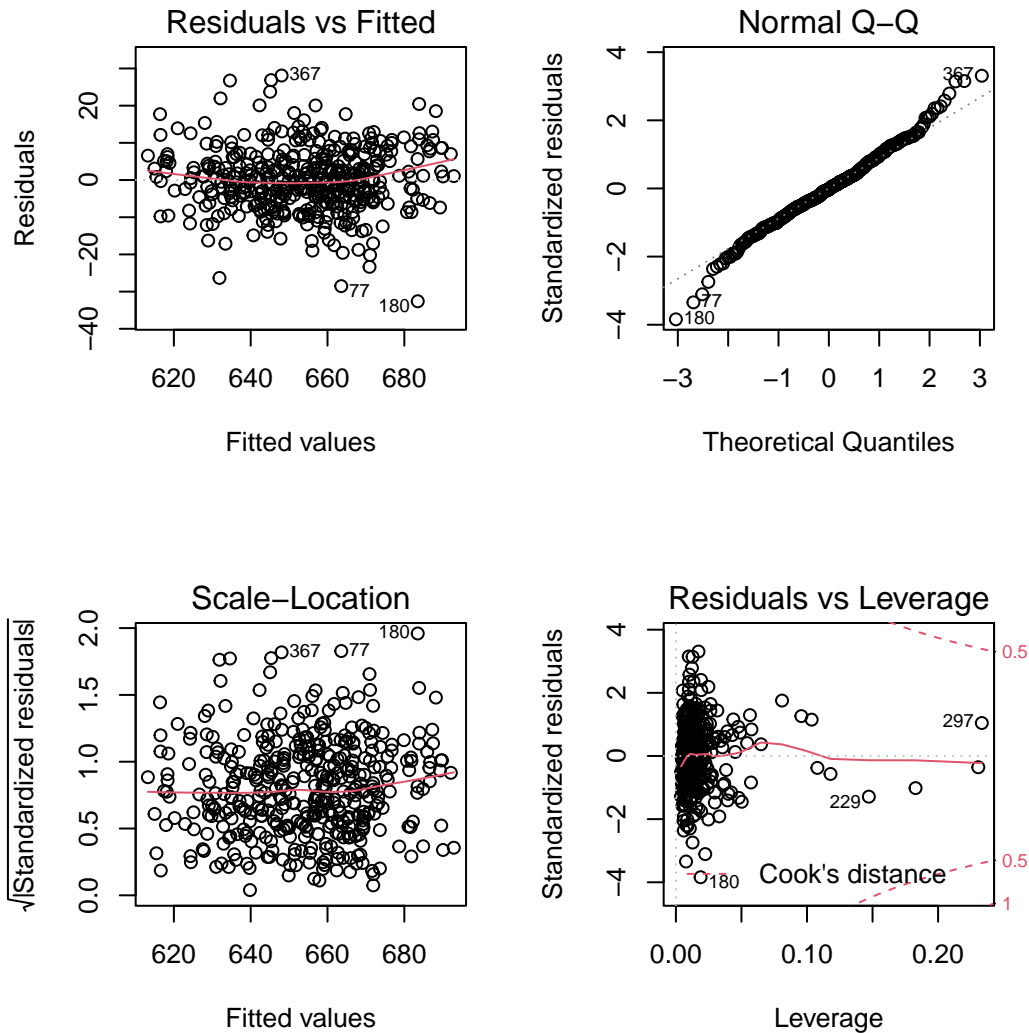
## 7.3 Diagnostics plots

Let's consider the `CASchools` dataset from the previous subsection:

```
library(AER)
data(CASchools)
CASchools$STR <- CASchools$students/CASchools$teachers
CASchools$score <- (CASchools$read + CASchools$math)/2
TS_mod7 <- lm(score ~ STR + I(STR^2) + I(STR^3)
              + english + lunch + log(income),
              data = CASchools)
```

The `plot()` function applied to an `lm` object returns four diagnostics plots:

```
par(mfrow=c(2,2))
plot(TS_mod7)
```

**Residuals vs Fitted**

Residuals

20 0 −20 −40

367

77
180

620 640 660 680

Fitted values

**Normal Q–Q**

Standardized residuals

4 2 0 −2 −4

367

77
180

−3 −1 0 1 2 3

Theoretical Quantiles

**Scale–Location**

√|Standardized residuals|

2.0 1.5 1.0 0.5 0.0

180
367 77

620 640 660 680

Fitted values

**Residuals vs Leverage**

Standardized residuals

4 2 0 −2 −4

0.5

297

229

180    Cook's distance

0.5

1

0.00 0.10 0.20

Leverage

These plots show different scatterplots of the fitted values $\widehat{Y}_i$, residuals $\hat{u}_i$, quantiles of the standard normal distribution, leverage values, and standardized residuals.

The red solid line indicates a local scatterplot smoother, which is a smooth locally weighted line through the points on the scatterplot to visualize the general pattern of the data.
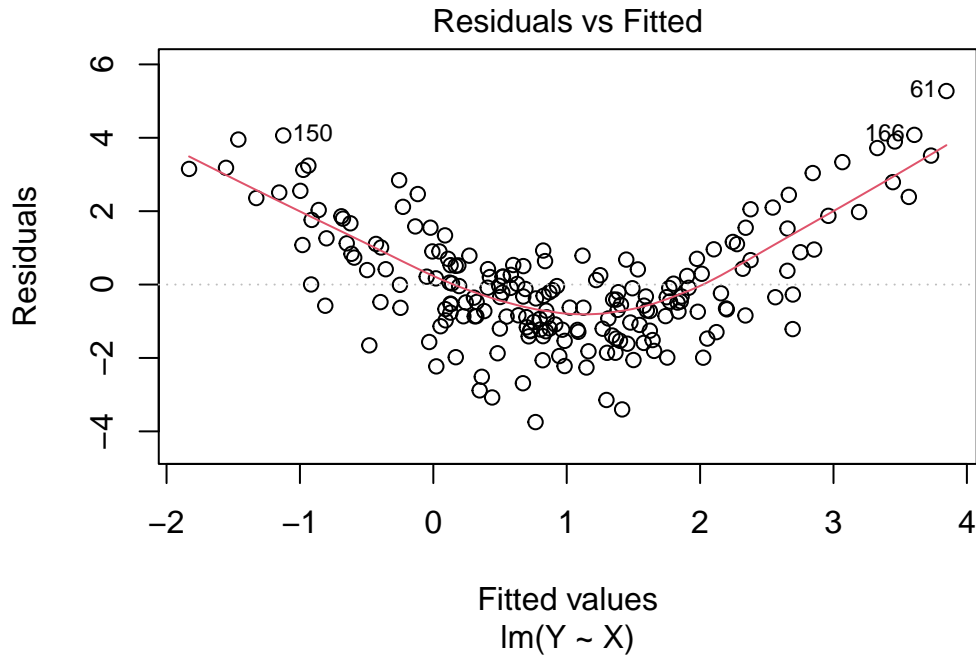
## Plot 1: Residuals vs Fitted

This plot indicates whether there are strong hidden nonlinear relationships between the response and the regressors that are not captured by the model. If a linear model is estimated but the relationship is nonlinear, then the assumption (A1) $E[u_i \mid \boldsymbol{X}_i] = 0$ is violated.

The residuals serve as a proxy for the unknown error terms. If you find equally spread residuals around a horizontal line without distinct patterns, that is a good indication you don't have non-linear relationships.

In the `CASchools` regression, there is only little indication for an omitted non-linear relationship. Here is an example of a strong omitted nonlinear pattern:

```
# Set seed for reproducability
set.seed(1)
# Simulate normally distributed regressors
X = rnorm(200)
# Simulate response nonlinearly
Y = X + X^2 + rnorm(200)
# Omit the nonlinearity in the regression
plot(lm(Y ~ X), which = 1)
```
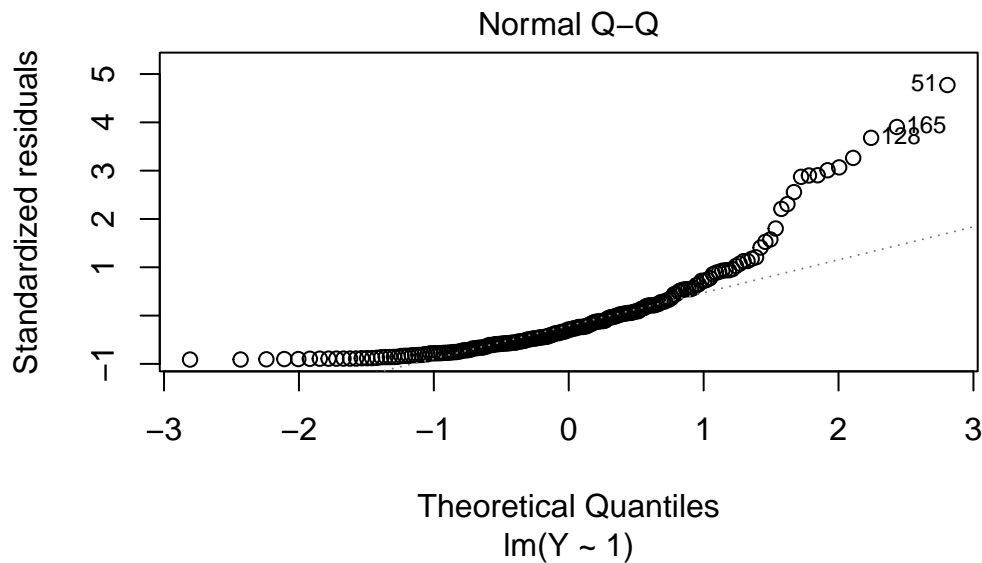


### Plot 2: Normal Q-Q

The QQ plot is a graphical tool to help us assess if the errors are conditionally normally distributed, i.e. whether assumption (A6) is satisfied.

Let $r_{(i)}$ be the order statistics of the standardized residuals (sorted standardized residuals). The QQ plot plots the ordered standardized residuals $u^*_{(i)}$ against the $((i-0.5)/n)$-quantiles of the standard normal distribution.

If the residuals are lined well on the straight dashed line, there is indication that the distribution of the residuals is close to a normal distribution.

In the `CASchools` regression, we see a slight deviation from normality in the tails. Here is an extrem example with a strong deviation from normality:

```r
# Exponentially distributed response variable
Y = rexp(200)
# Intercept only regression model
plot(lm(Y ~ 1), which = 2)
```
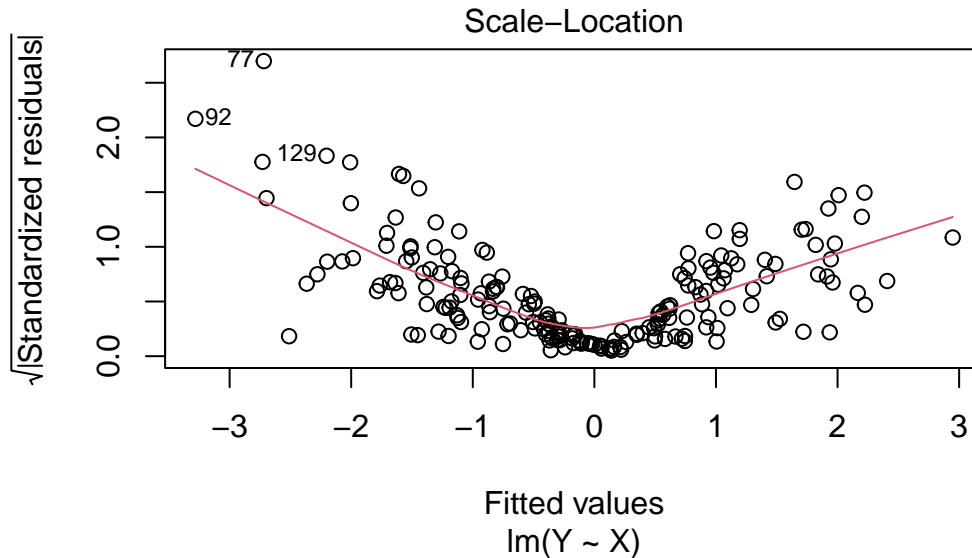


## Plot 3: Scale-Location

This plot shows if error terms are spread equally along the ranges of regressor values, which is how you can check the assumption of homoskedasticity (A5).

If you see a horizontal line with equally spread points, there is no indication for heteroskedasticity.

In the `CASchools` regression, we have some indication for weak heteroskedasticity. Here is an example with extreme heteroskedasticity:

```r
## simulate regressor values
X = rnorm(200)
## error variance varies with the regressor value
u = rnorm(200)*X^2
```

94

```
## response value
Y = X + u
plot(lm(Y ~ X), which = 3)
```



**Plot 4: Residuals vs Leverage**

Plotting standardized residuals against leverage values provides a graphical tool for detecting outliers. High leverage points have a strong influence on the regression fit. High leverage values with standardized residuals close to 0 are good leverage points, and high leverage values with large standardized residuals are bad leverage points.

The plot also shows Cook's distance thresholds. Cook's distance for observation $i$ is defined as

$$D_i = \frac{(\hat{\boldsymbol{\beta}}_{(-i)} - \hat{\boldsymbol{\beta}})' \boldsymbol{X}' \boldsymbol{X}(\hat{\boldsymbol{\beta}}_{(-i)} - \hat{\boldsymbol{\beta}})}{k s_{\hat{u}}^2},$$

where

$$\hat{\boldsymbol{\beta}}_{(-i)} = \hat{\boldsymbol{\beta}} - (\boldsymbol{X}'\boldsymbol{X})^{-1} \boldsymbol{X}_i \hat{u}_i (1 - h_{ii})^{-1}$$

is the $i$-th leave-one-out estimator (the OLS estimator when the $i$-th observation is left out).
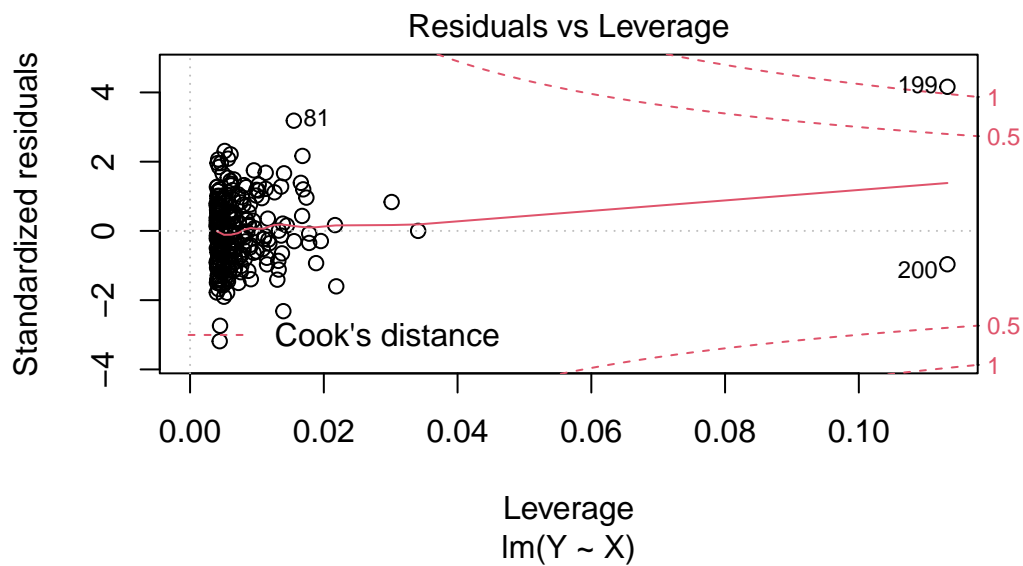
We should pay special attention to points outside Cook's distance thresholds of 0.5 and 1 and check for measurement errors or other anomalies.

Here is an example with two high leverage points. Observation $i = 200$ is a good leverage point and $i = 199$ is a bad leverage point:

```
## simulate regressors and errors
X = rnorm(250)
u = rnorm(250)
## set some unusual observations manually
X[199] = 6
X[200] = 6
u[199] = 5
u[200] = 0
## define dependent variable
Y = X + u
## residuals vs leverage plot
plot(lm(Y ~ X), which = 5)
```



## 7.4  Diagnostics tests

The asymptotic properties of the OLS estimator and inferential methods using HC-type standard errors do not depend on the validity of the homoskedasticity and normality assumptions (A5)–(A6).

However, if you are interested in exact inference, verifying the assumptions (A5)–(A6) becomes crucial, especially in small samples.

### 7.4.1 Breusch-Pagan Test (Koenker's version)

Under homoskedasticity, the variance of the error term does not depend on the values of the regressors.

To test for heteroskedasticity, we regress the squared residuals on the regressors.

$$\hat{u}_i^2 = \boldsymbol{X}_i'\boldsymbol{\gamma} + v_i, \quad i = 1, \dots, n. \tag{7.1}$$

Here, $\boldsymbol{\gamma}$ are the auxiliary coefficients and $v_i$ are the auxiliary error terms. Under homoskedasticity, the regressors should not be able to explain any variation in the residuals.

Let $R_{aux}^2$ be the r-squared coefficient of the auxiliary regression of Equation 7.1. The test statistic:

$$BP = nR_{aux}^2$$

Under the null hypothesis of homoskedasticity, we have

$$BP \xrightarrow{D} \chi_{k-1}^2$$

Test decision rule: Reject $H_0$ if $BP$ exceeds $\chi_{(1-\alpha,k-1)}^2$.

In R we can apply the `bptest()` function from the `lmtest` package to the `lm` object of our regression.

### 7.4.2 Jarque-Bera Test

A general property of any normally distributed random variable is that it has a skewness of 0 and a kurtosis of 3.

Under (A5)–(A6), we have $u_i \sim \mathcal{N}(0, \sigma^2)$, which implies $E[u_i^3] = 0$ and $E[u_i^4] = 3\sigma^4$.

Consider the sample skewness and the sample kurtosis of the residuals from your regression:

$$\widehat{skew}_{\hat{u}} = \frac{1}{n\hat{\sigma}_{\hat{u}}^3} \sum_{i=1}^n \hat{u}_i^3, \quad \widehat{kurt}_{\hat{u}} = \frac{1}{n\hat{\sigma}_{\hat{u}}^4} \sum_{i=1}^n \hat{u}_i^4$$

Jarque-Bera test statistic and null distribution if (A5)–(A6) hold:

$$JB = n\left(\frac{1}{6}(\widehat{skew}_{\hat{u}})^2 + \frac{1}{24}(\widehat{kurt}_{\hat{u}} - 3)^2\right) \xrightarrow{D} \chi_2^2.$$

Test decision rule: Reject the null hypothesis of normality if $JB$ exceeds $\chi_{(1-\alpha,2)}^2$.

The Jarque-Bera test is sensitive to outliers.

In R we apply use the `jarque.test()` function from the `moments` package to the residual vector from our regression.

## 7.5 R-codes

methods-sec07.R

# Part III

# C) Panel Data Methods