

Diagnostics

Motivation

- Estimation and inference of the regression model $Y = X\beta + \epsilon$ depend on several assumptions.
- These assumptions should be checked using regression diagnostics.
- The potential problems can be divided into three categories:
 - *Error*. We have assumed that $\epsilon \sim N(0, \sigma^2 I)$ or in words, that the errors are (1) independent, have (2) equal variance and are (3) normally distributed.
 - *Unusual observations*. Sometimes just a few observations do not fit the model. These few observations might change the choice and fit of the model.
 - *Model*. We have assumed that the structural part of the model, $E(Y | X) = X\beta$, is correct.

1. Checking Error Assumptions

- Goal: check error assumptions on (1) independence, (2) constant variance, and (3) normality.
- Errors are not observable \rightarrow examine the sample residuals $\hat{\epsilon} = Y - X\hat{\beta}$.
- Similarity: $E(\hat{\epsilon} | X) = 0$
- But there are differences:
 - Note $\hat{\epsilon} = (I - P)Y = (I - P)(X\beta + \epsilon) = (I - P)\epsilon$.
 - $\text{cov}(\epsilon | X) = \sigma^2 I_n$ whereas $\text{cov}(\hat{\epsilon} | X) = \sigma^2(I - P)$.
 - * Even if the errors have equal variance and are uncorrelated, the residuals do not.
 - * When this impact is small, diagnostics can reasonably be applied to residuals to check the assumptions on errors.

1.1 Constant variances

- homoscedasticity: constant symmetrical variation in the vertical $\hat{\epsilon}$ direction.

- heteroscedasticity: nonconstant variance.

Plot of residuals $\hat{\epsilon}$ against fitted values \hat{y}

- $\text{cov}(\hat{\epsilon}, \hat{Y}) = \text{cov}(PY, (I - P)Y) = \sigma^2 P(I - P) = 0$

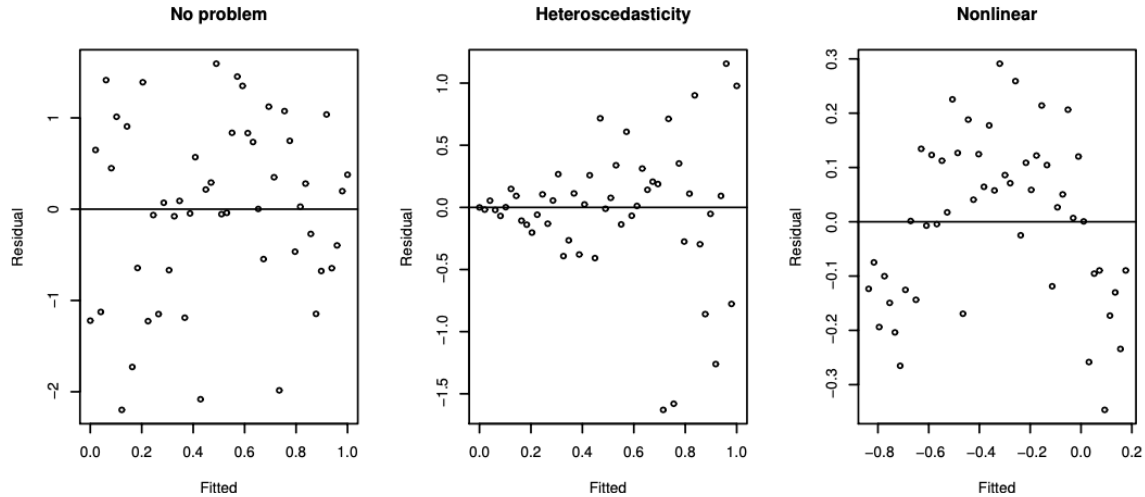
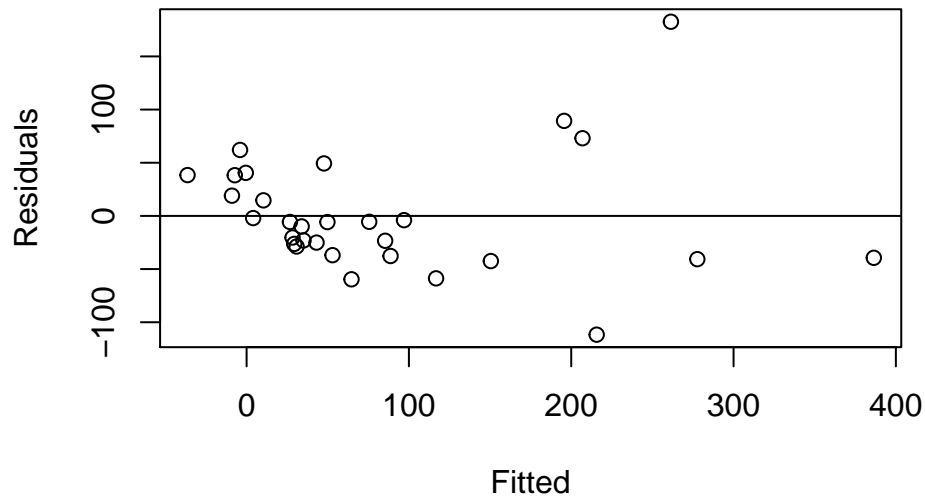


Figure 1: Residuals vs. fitted plots - the first suggests no change to the current model; the second shows nonconstant variance; the third indicates some nonlinearity, which should prompt some change in the structural form of the model.

- Other plots $\hat{\epsilon}$ versus x , $|\hat{\epsilon}|$ versus \hat{y} , or $\sqrt{|\hat{\epsilon}|}$ versus \hat{y} can all be used.
- It could be hard to judge residual plots without prior experience so it is helpful to generate some artificial plots where the true relationship is known. Sample codes provided below.

```
par(mfrow=c(3,3))
for(i in 1:9) plot(1:50,rnorm(50)) #constant variance
for(i in 1:9) plot(1:50,(1:50)*rnorm(50)) #strong non-constant variance
for(i in 1:9) plot(1:50,sqrt((1:50))*rnorm(50)) #mild non-constant variance
for(i in 1:9) plot(1:50,cos((1:50)*pi/25)+rnorm(50)) #nonlinearity
par(mfrow=c(1,1))

data(gala, package = "faraway")
lmod <- lm(Species ~ Area + Elevation + Scrub + Nearest + Adjacent, gala)
plot(fitted(lmod),residuals(lmod),xlab="Fitted",ylab="Residuals")
abline(h=0)
```

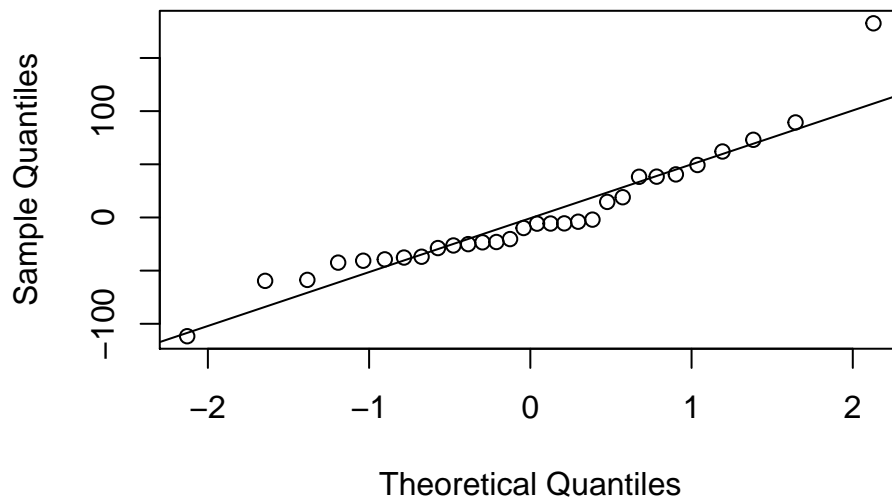


1.2 Normality

- Tests and confidence intervals are based on the assumption of normal errors.
- The normality of residuals can be assessed by a Q-Q plot. (Histogram is usually not preferred.)
 - Sort sample residuals $\hat{\epsilon}_{(1)} < \dots < \hat{\epsilon}_{(n)}$.
 - Plot them against $\Phi^{-1}(\frac{1}{n+1}) < \dots < \Phi^{-1}(\frac{n}{n+1})$ where $\Phi(\cdot)$ denotes the CDF (cumulative distribution function).
 - Normal residuals should follow the line approximately.

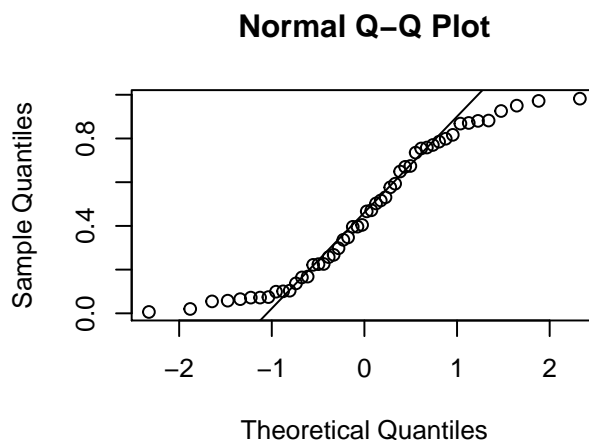
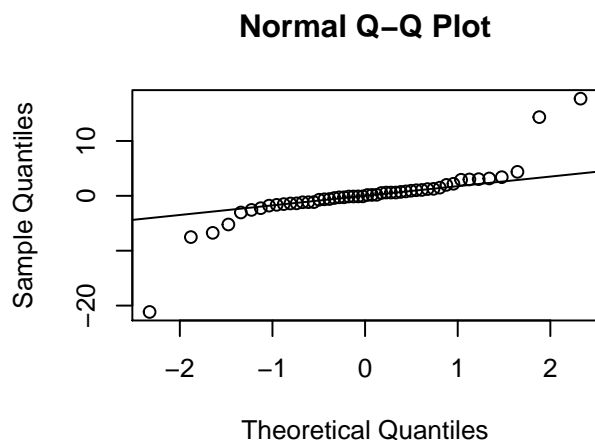
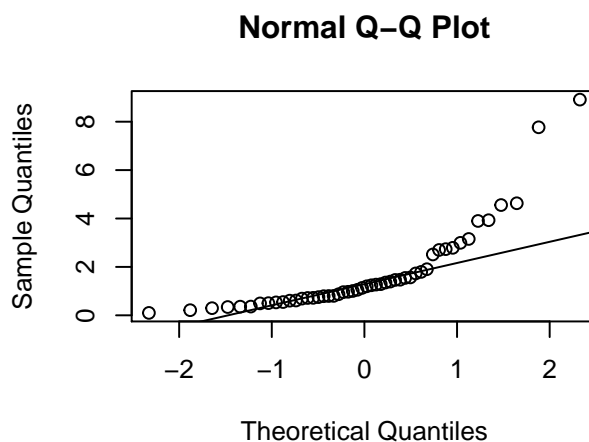
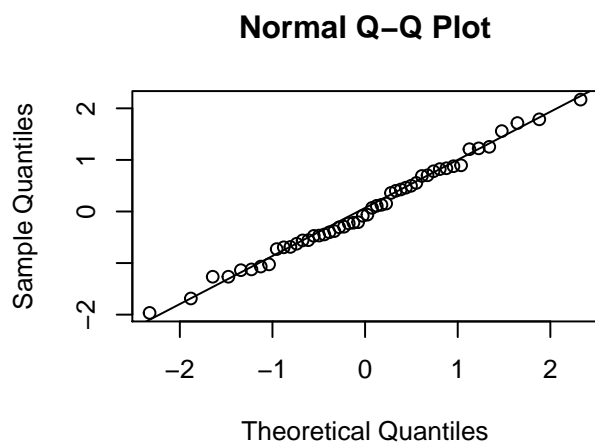
```
qqnorm(residuals(lmod))
qqline(residuals(lmod)) #passes through first and third quantiles
```

Normal Q-Q Plot



We can generate sample data from different distributions.

```
par(mfrow=c(2,2))
set.seed(123)
n = 50
#1. normal
x = rnorm(n); qqnorm(x); qqline(x)
#2. log normal: an example of a skewed distribution
x = exp(rnorm(n)); qqnorm(x); qqline(x)
#3. Cauchy: an example of a long-tailed distribution
x = rcauchy(n); qqnorm(x); qqline(x)
#4. Uniform: an example of a short-tailed distribution
x = runif(n); qqnorm(x); qqline(x)
```



```
par(mfrow=c(1,1))
```

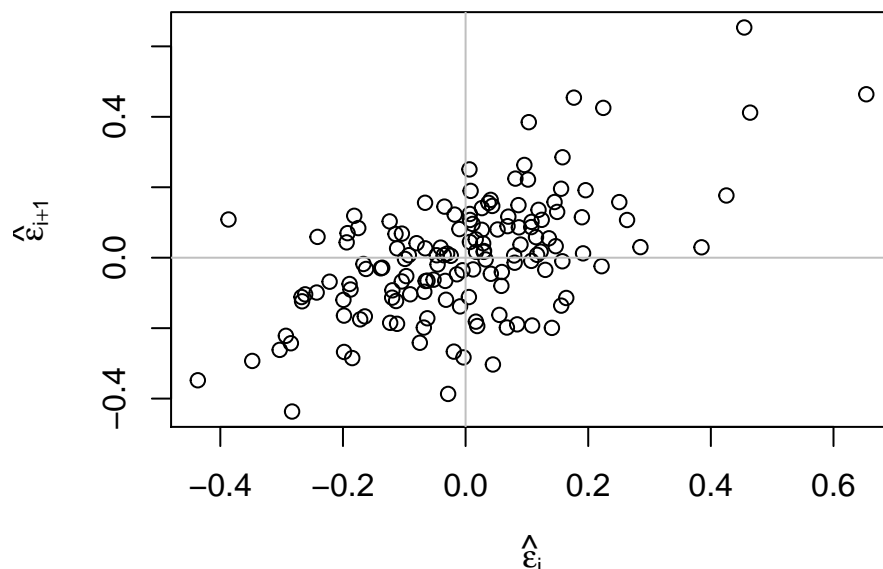
1.3 Correlated errors

- It is difficult to check for correlated errors in general because there are just too many possible patterns of correlation that may occur.
 - Data collected over time may have some correlation in successive errors.
 - Spatial data may have correlation in the errors of nearby measurements.
 - Data collected in blocks may show correlated errors within those blocks.

As an example, we consider a serial data on records of annual temperatures.

- The data contains temprature information 1856 through 2000 by Jones and Mann (2004) Climate over past millennia.
- We can build a linear model to predict temperature since 1856 and then subsequently use this to predict earlier temperatures based on proxy information.

```
data(globwarm,package="faraway")
lmod <- lm(nhtemp ~ wusa + jasper + westgreen + chesapeake +
  tornetrask + urals + mongolia + tasman, globwarm)
n <- length(residuals(lmod))
plot(tail(residuals(lmod),n-1) ~ head(residuals(lmod),n-1), xlab=
  expression(hat(epsilon)[i]),ylab=expression(hat(epsilon)[i+1]))
abline(h=0,v=0,col=grey(0.75))
```



- We can see a positive correlation again indicating positive serial correlation. If you have some doubt as to the significance of the correlation.
- We can model the observed correlation directly by linear regression:

```
summary(lm(tail(residuals(lmod),n-1) ~ head(residuals(lmod),n-1) -1))$coefficients

##                                Estimate Std. Error  t value      Pr(>|t|)
## head(residuals(lmod), n - 1) 0.5950759 0.06931205  8.585462 1.390651e-14
```

- Note that the the serial correlation is confirmed.

2. Finding unusual observations

1. Leverage: A leverage point is extreme in the predictor space.
2. Outliers: Some observations do not fit the model well.
3. Influential: Some observations change the fit of the model in a substantive manner.

2.1 Leverage

- $h_{ii} = H_{ii}$, i.e., i -th diagonal of hat matrix $H = X(X^T X)^{-1} X^T$ are called **leverages**.
- $\text{var}(\hat{\epsilon}_i) = \sigma^2(1 - h_{ii})$
 - a large leverage h_i , will make $\text{var}(\hat{\epsilon}_i)$ small.
 - The fit will be attracted toward y_i .
 - Large values of h_i are usually due to extreme values in the X -space.
- The value of h_{ii} depends only on X and not y .
- $\sum_{i=1}^n h_{ii} = \text{tr}(H) = p$ (See notes Sep 26)
 - Average $\sum_{i=1}^n h_{ii}/n = p/n$.
 - A rough rule is that leverages of more than $2p/n$ should be looked at more closely.

```
data(gala, package = "faraway")
lmod <- lm(Species ~ Area + Elevation + Scrub + Nearest + Adjacent, gala)
hatv <- hatvalues(lmod)
head(hatv)
```

```
##      Baltra      Bartolome      Caldwell      Champion      Coamano Daphne.Major
## 0.07871937 0.09135324 0.06231443 0.07237676 0.16878374 0.07163790
```

```
sum(hatv) #number of parametrs in the model
```

```
## [1] 6
```

- Standardized residuals:

$$r_i = \frac{\hat{\varepsilon}_i}{\hat{\sigma}\sqrt{1-h_i}}$$

- Standardization can only correct for the natural non-constant variance in residuals when the errors have constant variance.
- If there is some underlying heteroscedasticity in the errors, standardization cannot correct for it.
- When there are unusually large leverages, there could be differences between raw and standardized residuals in their plots.

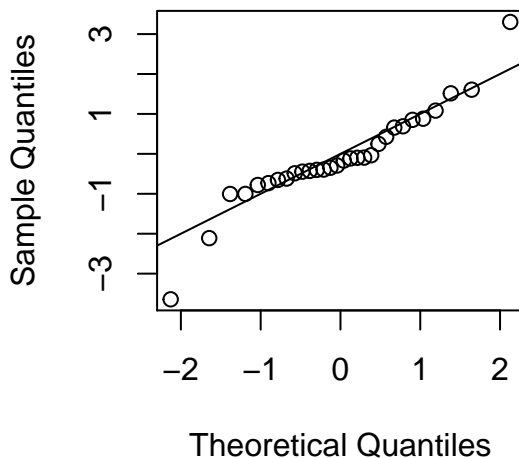
- Leave-one-out studentized residuals:

$$\tilde{r}_i = \frac{\hat{\varepsilon}_i}{\hat{\sigma}_{(i)}\sqrt{1-h_i}}$$

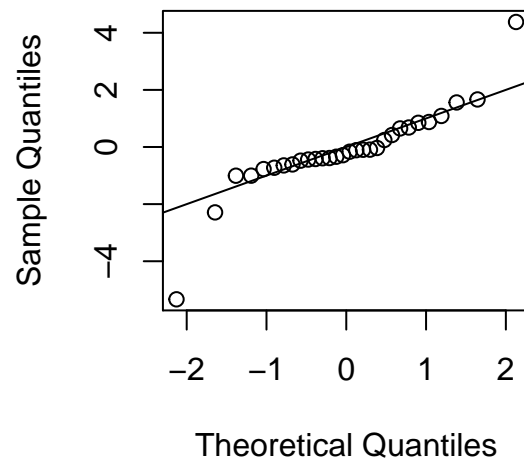
- $\hat{\sigma}_{(i)}$: estimate σ without i -th observation.

```
par(mfrow=c(1,2))  
qqnorm(rstandard(lmod))  
abline(0,1)  
  
qqnorm(rstudent(lmod))  
abline(0,1)
```

Normal Q–Q Plot



Normal Q–Q Plot



```
par(mfrow=c(1,1))
```

2.2 Outliers

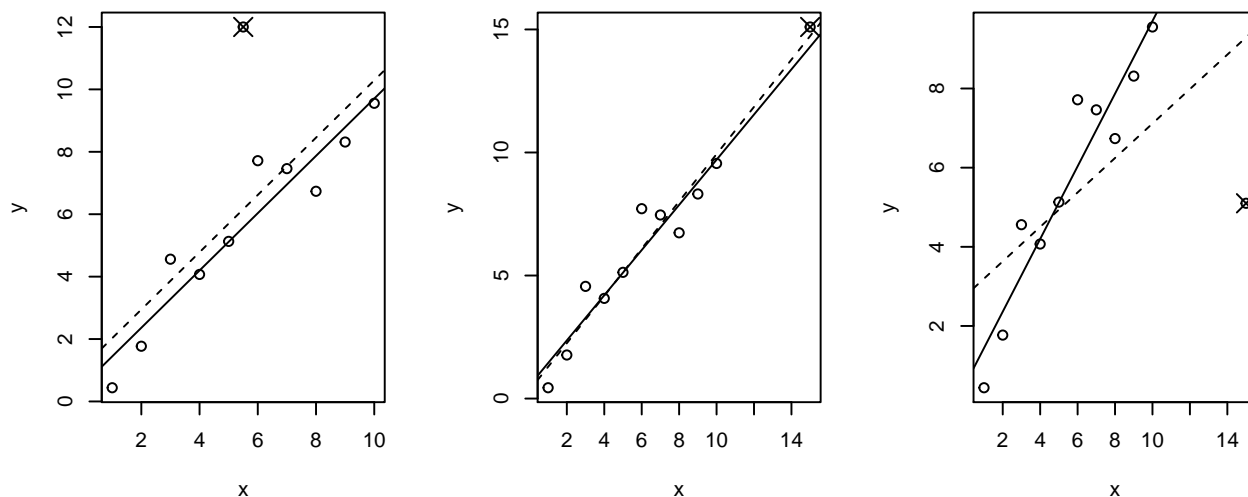
- An outlier is a point that does not fit the current model well.
- Outliers may or may not affect the fit substantially.

```
set.seed(123)
testdata <- data.frame(x=1:10,y=1:10+rnorm(10))
lmod0 <- lm(y ~ x, testdata)
```

```
par(mfrow=c(1,3))
p1 <- c(5.5,12)
lmod1 <- lm(y ~ x, rbind(testdata, p1))
plot(y ~ x, rbind(testdata, p1))
points(5.5,12,pch=4,cex=2)
abline(lmod0)
abline(lmod1, lty=2)
```

```
p2 <- c(15,15.1)
lmod2 <- lm(y ~ x, rbind(testdata, p2))
plot(y ~ x, rbind(testdata, p2))
points(15,15.1,pch=4,cex=2)
abline(lmod0)
abline(lmod2,lty=2)
```

```
p3 <- c(15,5.1)
lmod3 <- lm(y ~ x, rbind(testdata, p3))
plot(y ~ x, rbind(testdata, p3))
points(15,5.1,pch=4,cex=2)
abline(lmod0)
abline(lmod3,lty=2)
```

```
par(mfrow=c(1,1))
```

- A solid regression line shows the fit without the additional point marked with a cross.
- The dashed line shows the fit with the extra point.
- Left panel:
 - Added point is an outlier.
 - But it does not have large leverage or influential (on the fit).
- Middle panel:
 - Added point has large leverage. (well outside of the range of X .)
 - But is not an outlier and is not influential.
- Right panel:
 - Added point changes the fitted line substantially. (influential point.)
 - This is both an outlier and an influential point.

To detect outliers,

- we exclude the point i and recompute the estimates to get $\hat{\beta}_{(i)}$ and $\hat{\sigma}_{(i)}$.
- Let $\hat{y}_{(i)} = x_i^T \hat{\beta}_{(i)}$
- We have

$$\text{var}(y_i - \hat{y}_{(i)}) = \sigma^2 \left(1 + x_i^T (X_{(i)}^T X_{(i)})^{-1} x_i \right)$$

- Define

$$t_i = \frac{y_i - \hat{y}_{(i)}}{\hat{\sigma}_{(i)} \left(1 + x_i^T \left(X_{(i)}^T X_{(i)} \right)^{-1} x_i \right)^{1/2}},$$

where $X_{(i)}$ represents the design matrix deleting i -th observation.

- It can be proved $t_i = \frac{\hat{\epsilon}_i}{\hat{\sigma}_{(i)} \sqrt{1-h_{ii}}} = r_i \sqrt{\frac{(n-p-1)}{n-p-r_i^2}}$. (Easy-to-compute. Proof see Theorem 10.1 in Lee and Seber.)
- If i -th case is not outlier, model is correct, and $\epsilon \sim N(0, \sigma^2 I_n)$, $t_i \sim t_{(n-1)-p}$, where $n-1$ is the sample size.
- Test outliers
 - Practically, $|t_i| > 3$ can imply possible outliers.
 - If we want a level α test,
 - * $P(\text{all tests accept}) = 1 - P(\text{at least one rejects}) \geq 1 - \sum_i P(\text{test } i \text{ rejects}) = 1 - n\alpha$.
 - * Each test should use level α/n . (Bonferroni correction.)

2.3 Influential point

- An influential point is one whose removal from the dataset would cause a large change in the fit.
 - An influential point may or may not be an outlier,
 - and may or may not have large leverage,
 - but it will tend to have at least one of these two properties.
- Measure of influence: Cook's distance statistic Cook (1977)

$$D_i = \frac{(\hat{y} - \hat{y}_{(i)})^T (\hat{y} - \hat{y}_{(i)})}{p \hat{\sigma}^2}$$

- It can also be computed as $D_i = \frac{1}{p} r_i^2 \frac{h_{ii}}{1-h_{ii}}$, where r_i represents i -th standardized residual.

```
cook <- cooks.distance(lmod)
```

```
summary(lmod)$coefficients
```

```
##              Estimate Std. Error    t value    Pr(>|t|)
## (Intercept)  7.068220709 19.15419782  0.369016796 7.153508e-01
## Area        -0.023938338  0.02242235 -1.067610554 2.963180e-01
## Elevation    0.319464761  0.05366280  5.953187968 3.823409e-06
## Scruz        -0.240524230  0.21540225 -1.116628222 2.752082e-01
## Nearest      0.009143961  1.05413595  0.008674366 9.931506e-01
## Adjacent    -0.074804832  0.01770019 -4.226216850 2.970655e-04
```

```
summary(lmod)$r.squared
```

```
## [1] 0.7658469
```

```
lmodi <- lm(Species ~ Area + Elevation + Scruz + Nearest + Adjacent,
            gala, subset = (cook < max(cook)))
```

```
summary(lmodi)$coefficients
```

```
##              Estimate Std. Error    t value    Pr(>|t|)
## (Intercept) 22.58614473 13.40191356  1.6852925 1.054542e-01
## Area         0.29574351  0.06186188  4.7807068 8.042013e-05
## Elevation    0.14039023  0.04970484  2.8244782 9.613092e-03
## Scruz        -0.09010457  0.14979821 -0.6015063 5.533860e-01
## Nearest      -0.25518223  0.72167754 -0.3535959 7.268624e-01
## Adjacent     -0.06503051  0.01222732 -5.3184596 2.124483e-05
```

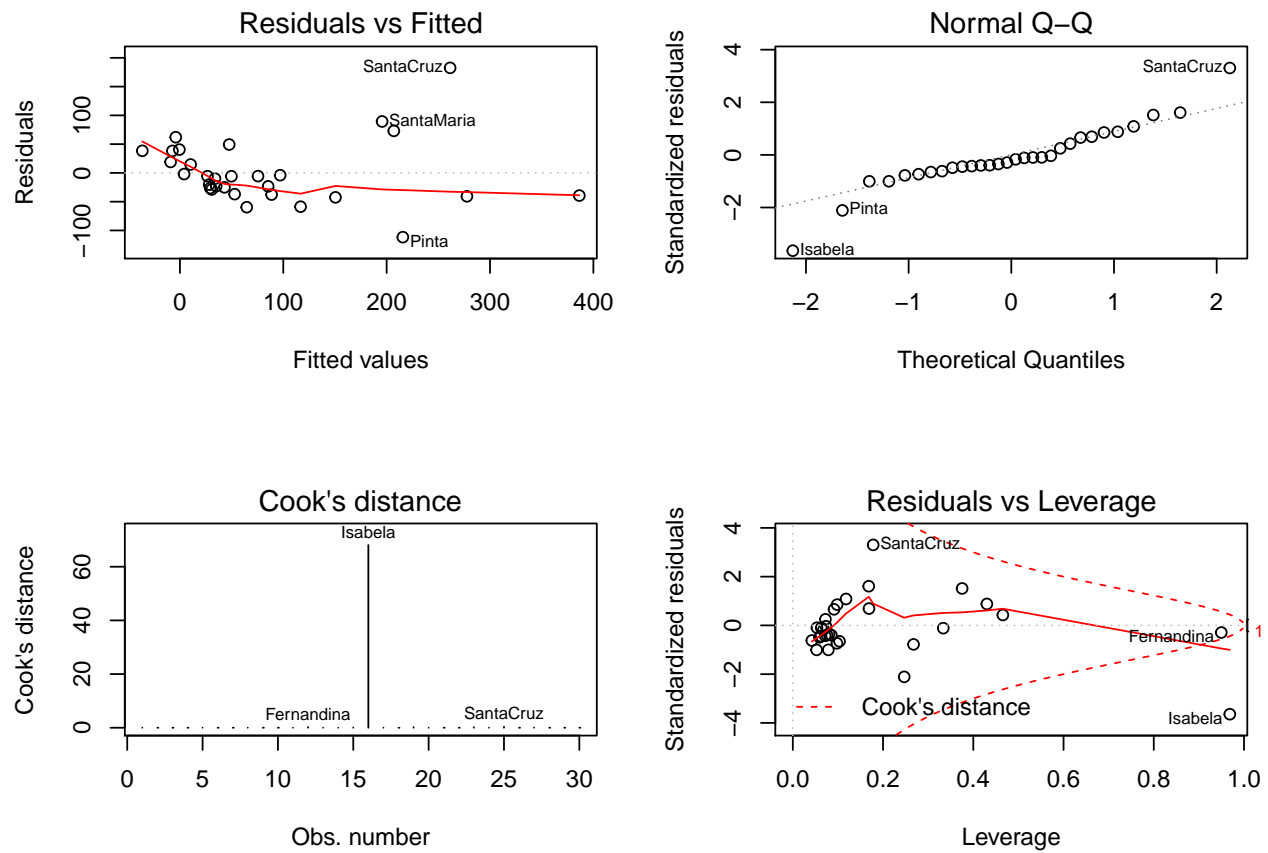
```
summary(lmodi)$r.squared
```

```
## [1] 0.8714011
```

- p-value of covariate Area changes significantly.
- We usually do not want estimates to be so sensitive to the presence/deletion of just one observation.

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

```
plot(lmod, which = c(1,2,4, 5), cook.levels = 1) #R codes for multiple diagnostic plots
```



```
par(mfrow=c(1,1))
```

- As the Cook statistics represent a function of standardized residuals and leverage, we can plot contours (the above plot shows contours with Cook's distance = 1).
- Any point that lies beyond these contours might well be influential and require closer attention.

3. Checking the Structure of the Model

- Check linear structure $E(Y | X) = X\beta$
- Residual plots can suggest transformations of the variables which might improve the structural form of the model.
- A formal lack of test may be conducted in some cases.

Diagnostics 2

I. Issues of predictors

Collinearity of predictors

- When some predictors are linear combinations of others, $X^\top X$ is singular.
 - No unique OLS solution. Theoretical derivations could be problematic.
 - May also be called exact collinearity.
 - A solution could be removing redundant predictors that can be represented by other predictors.
- Another challenging case is $X^\top X$ is nearly singular but not exactly so.
 - Known as **collinearity** or **multicollinearity**.
 - This could lead to imprecise estimates of β . Intuitively, $(X^\top X)^{-1}$ blows up, and variances are large.
- Detecting (multi)collinearity:
 - Examine correlation matrix of predictors: values close to -1 or $+1$ indicate large pairwise correlations/collinearities.
 - * $Y \sim X_1 + X_2 + X_3 + X_4$. If $X^\top X$ is nearly singular, $\text{corr}(X_i, X_j)$ for $i \neq j$ may not be large.
 - * We want to some “generalized correlations”.
 - Regressing each X_j on all other predictors X_{-j} for $j = 1, \dots, p$.
 - * Let $R_{X_j|X_{-j}}^2$ denote the R^2 obtained from $X_j \sim X_{-j}$.
 - * Values close to 1 indicate a problem, because it means one predictor can almost be predicted exactly by a linear combination of other predictors.
 - * We can obtain

$$\text{var}(\hat{\beta}_j) = \sigma^2 \frac{1}{1 - R_{X_j|X_{-j}}^2} \times \frac{1}{\sum_{i=1}^n (x_{ij} - \bar{x}_j)^2}$$

- Equivalently, people also use Variance Inflation Factor (VIF)

$$\text{VIF}(\hat{\beta}_j) = \frac{\text{var}(\hat{\beta}_j)_{\text{Full model } Y \sim X_1 + \dots + X_p}}{\text{var}(\hat{\beta}_j)_{\text{Univariate model } Y \sim X_j}} = \frac{1}{1 - R_{X_j|X_{-j}}^2}.$$

- If we have the ability to choose the X , we can minimize the variance of the regression coefficients by
 - $R_{X_j|X_{-j}}^2 = 0$ (Orthogonality design.)
 - Maximizing $\sum_{i=1}^n (x_{ij} - \bar{x}_j)^2 = \widehat{\text{var}}(X_j)$ by spreading X as much as possible.
 - Note this relies on linear model. In practice, instead of putting all X at maximum/minimum, we also put points in the middle.

Data Example:

- Car drivers like to adjust the seat position for their own comfort.
- Car designers would find it helpful to know where different drivers will position the seat depending on their size and age.
- Researchers at the HuMoSim laboratory at the University of Michigan collected data on 38 drivers.
- They measured age in years, weight in pounds, height with shoes and without shoes in centimeters, seated height arm length, thigh length, lower leg length and hipcenter, the horizontal distance of the midpoint of the hips from a fixed location in the car in millimeters.
- We fit a model with all the predictors:

```
data(seatpos, package="faraway")
lmod <- lm(hipcenter ~ ., seatpos)
summary(lmod)$coefficients
```

##	Estimate	Std. Error	t value	Pr(> t)
## (Intercept)	436.43212823	166.5716187	2.62008697	0.01384361
## Age	0.77571620	0.5703288	1.36012113	0.18427175
## Weight	0.02631308	0.3309704	0.07950283	0.93717877
## HtShoes	-2.69240774	9.7530351	-0.27605845	0.78446097
## Ht	0.60134458	10.1298739	0.05936348	0.95306980
## Seated	0.53375170	3.7618942	0.14188376	0.88815293
## Arm	-1.32806864	3.9001969	-0.34051323	0.73592450
## Thigh	-1.14311888	2.6600237	-0.42974011	0.67056106
## Leg	-6.43904627	4.7138601	-1.36598163	0.18244531

- This model shows the signs of collinearity.
 - S.E. is large. R^2 is not small, but none of the individual predictors is significant.
- We can take a look at the pairwise correlations:

```
round(cor(seatpos[, -9]), 2)
```

##	Age	Weight	HtShoes	Ht	Seated	Arm	Thigh	Leg
## Age	1.00	0.08	-0.08	-0.09	-0.17	0.36	0.09	-0.04


```
## Weight    0.08    1.00    0.83  0.83    0.78 0.70    0.57  0.78
## HtShoes -0.08    0.83    1.00  1.00    0.93 0.75    0.72  0.91
## Ht       -0.09    0.83    1.00  1.00    0.93 0.75    0.73  0.91
## Seated  -0.17    0.78    0.93  0.93    1.00 0.63    0.61  0.81
## Arm      0.36    0.70    0.75  0.75    0.63 1.00    0.67  0.75
## Thigh    0.09    0.57    0.72  0.73    0.61 0.67    1.00  0.65
## Leg     -0.04    0.78    0.91  0.91    0.81 0.75    0.65  1.00
```

- There are several large pairwise correlations between predictors.
- We can check the VIFs.

```
x <- model.matrix(lmod)[-1]
(rsq_1 <- summary(lm(x[,1] ~ x[-1]))$r.squared)
```

```
## [1] 0.4994823
```

```
#calculate VIF of first covariate by definition
1/(1-rsq_1)
```

```
## [1] 1.997931
```

```
#By R function
library(faraway)
vif(x)
```

```
##      Age      Weight    HtShoes      Ht      Seated      Arm      Thigh
## 1.997931  3.647030 307.429378 333.137832  8.951054  4.496368  2.762886
##      Leg
## 6.694291
```

II. Some remedies of error issues

- We have seen that assumptions of errors can be violated and we must then consider alternatives.
- When the errors are dependent, we can use generalized least squares (GLS).
- When the errors are independent, but not identically distributed, we can use weighted least squares (WLS), which is a special case of GLS.

II.1 Generalized Least Squares

- We have assumed $\epsilon = \sigma^2 I$.

If Σ is known:

- Suppose instead $\text{var}(\epsilon) = \sigma^2 \Sigma$.
 - σ^2 is unknown but Σ is known.
 - that is, we know the correlation and relative variance between the errors,
 - but we do not know the absolute scale of the variation.
- Write $\Sigma = SS^\top$. We can transform the regression model as

$$\begin{aligned} Y &= X\beta + \epsilon \\ S^{-1}Y &= S^{-1}X\beta + S^{-1}\epsilon \quad \Rightarrow \quad \tilde{Y} = \tilde{X}\beta + \tilde{\epsilon} \end{aligned}$$

- Then OLS can be conducted to the transformed variables \tilde{Y} and \tilde{X} .
- Then

$$\begin{aligned} \hat{\beta} &= (\tilde{X}^\top \tilde{X})^{-1} \tilde{X}^\top \tilde{Y} = (X^\top \Sigma^{-1} X)^{-1} X^\top \Sigma^{-1} Y \\ \text{var}(\hat{\beta}) &= \sigma^2 (\tilde{X}^\top \tilde{X})^{-1} = \sigma^2 (X^\top \Sigma^{-1} X)^{-1}. \end{aligned}$$

- Also diagnostics should be applied to the transformed residuals $S^{-1}\hat{\epsilon}$, which should be approximately i.i.d.

If Σ is unknown:

- We need to estimate Σ . Can be done through R function `gls`.
- Recall the temperature data that we investigated where serial correlation was observed.