

# Regression Analysis

## Chapter 8 and 9: Transformations and Regression Diagnostics

Shaobo Jin

Department of Mathematics

# Diagnostics

Regression diagnostics are used after model fitting to check if a fitted mean function and assumptions are consistent with observed data.

Recall that the residuals are

$$\hat{e} = \mathbf{y} - \hat{\mathbf{y}} = (\mathbf{I} - \mathbf{H}) \mathbf{y}.$$

- If OLS is used to estimate  $\beta$ , then

$$\mathbf{H} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T.$$

- If WLS is used to estimate  $\beta$ , then

$$\mathbf{H} = \mathbf{W}^{1/2} \mathbf{X} (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}^{1/2}.$$

# Properties of Hat Matrix

$$\text{Hat matrix } \mathbf{H} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T.$$

- The hat matrix is symmetric and idempotent. Hence it is an **orthogonal projection matrix** to the column space of  $\mathbf{X}$ .
  - Let  $\mathbf{X} = [\mathbf{X}_1 \quad \mathbf{X}_2]$  and  $\mathbf{H}_1 = \mathbf{X}_1 (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T$ . Then,

$$\mathbf{H} \mathbf{H}_1 \mathbf{y} = \mathbf{H}_1 \mathbf{y}.$$

- The hat matrix is not full rank:  $\text{rank}(\mathbf{H}) = \text{tr}(\mathbf{H}) = p$ , if  $\mathbf{X}$  is a  $n \times p$  matrix.
- The diagonal entries satisfy  $h_{ii} \leq 1$  for all  $i$ .
  - If the intercept is included, then we also have  $1/n \leq h_{ii} \leq 1$ .

# Leverage

Suppose that an intercept is included in the model. Then,

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

The predicted value is  $\hat{\mathbf{y}} = \mathbf{H}\mathbf{y}$ . If  $h_{ii}$  is close to 1, then

- ① the predicted value tends to be close to  $y_i$ , i.e.,

$$\hat{y}_i = \sum_{j=1}^n h_{ij}y_j = h_{ii}y_i + \sum_{j \neq i} h_{ij}y_j$$

- ② and  $\text{Var}(\hat{e}_i | \mathbf{X}) = \sigma^2(1 - h_{ii})$  is close to 0.

For this reason,  $h_{ii}$  is called the **leverage** of the  $i$ th observation.

# Leverage on $\hat{\beta}$

Let  $\mathbf{X} = [\mathbf{1} \quad \mathbf{X}_{-1}]$  and  $\bar{\mathbf{x}}_{-1} = \frac{1}{n} \mathbf{X}_{-1}^T \mathbf{1}$ . Let

$$\mathcal{X} = \begin{bmatrix} x_{11} - \bar{x}_1 & x_{1,p-1} - \bar{x}_{p-1} \\ \vdots & \vdots \\ x_{n1} - \bar{x}_1 & x_{n,p-1} - \bar{x}_{p-1} \end{bmatrix}$$

be the demeaned data matrix. It can be shown that

$$h_{ii} = \frac{1}{n} + (\mathbf{x}_i^* - \bar{\mathbf{x}}_{-1})^T (\mathcal{X}^T \mathcal{X})^{-1} (\mathbf{x}_i^* - \bar{\mathbf{x}}_{-1}),$$

where  $\mathbf{x}_i^T = [1 \quad (\mathbf{x}_i^*)^T]$ . Hence, if  $h_{ii}$  is close to 1, it means that  $\mathbf{x}_i^*$  is far from the average  $\bar{\mathbf{x}}_{-1}$  and the influence of the  $i$ th observation to the estimation of our regression model is large.

## Different Residuals

Suppose that the estimator is the minimizer of

$$\text{RSS}(\boldsymbol{\beta}) = \sum_{i=1}^n w_i (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2.$$

- ① The residual is  $\hat{e}_i = y_i - \mathbf{x}_i^T \hat{\boldsymbol{\beta}}$ .
- ② The **Pearson residual** is  $\sqrt{w_i} (y_i - \mathbf{x}_i^T \hat{\boldsymbol{\beta}})$ .
  - The sum of squared Pearson residuals is the same as  $\text{RSS}(\hat{\boldsymbol{\beta}})$ .
- ③ Under some assumptions, we can show that

$$\text{Var}(\hat{\mathbf{e}} \mid \mathbf{X}) = \sigma^2 \mathbf{W}^{-1/2} (\mathbf{I} - \mathbf{H}) \mathbf{W}^{-1/2}.$$

A **standardized residual** is

$$\frac{\hat{e}_i}{\hat{\sigma} \sqrt{(1 - h_{ii}) / w_{ii}}}.$$

Its variance is closer to 1 than  $\hat{e}_i$ .

# Properties of Residual

Suppose that the intercept is included in our model and OLS is used. Then, regardless of whether our model is correct or not,

- ① we always have  $\sum_i \hat{e}_i = 0$ .
- ② the sample correlation between residual and regressors is zero.
- ③ the sample correlation between residual and the fitted value is zero.

However, if the model is misspecified, we may still observe patterns. We can plot residuals to see whether something has gone wrong.

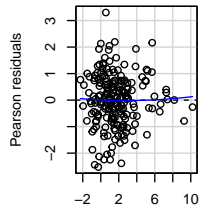
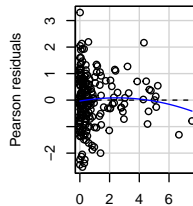
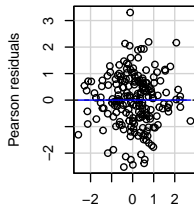
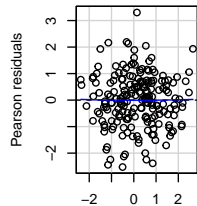
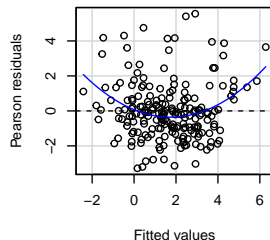
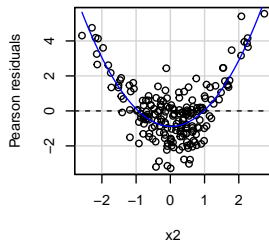
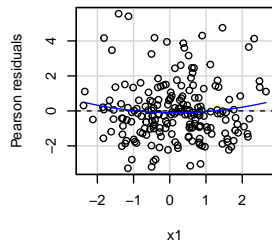
# Detect Lack of Fit Using Residuals (1)

The residuals should be plotted against each of the explanatory variables of the model, and against the fitted value.

- If the model fits the data well, no systematic patterns can be observed.
- A systematic pattern suggests another model form or additional terms.



# Detect Lack of Fit Using Residuals (1)



# Normality Assumption

We often assume  $\mathbf{e}$  is normally distributed.

- ① If  $\mathbf{e} \mid \mathbf{X}$  is normally distributed, then

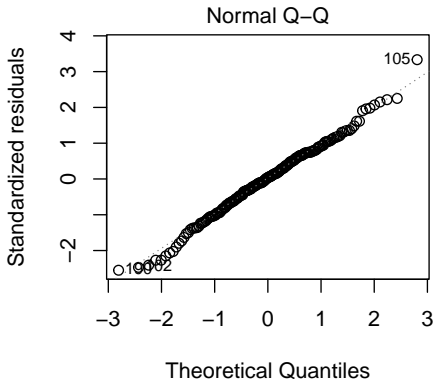
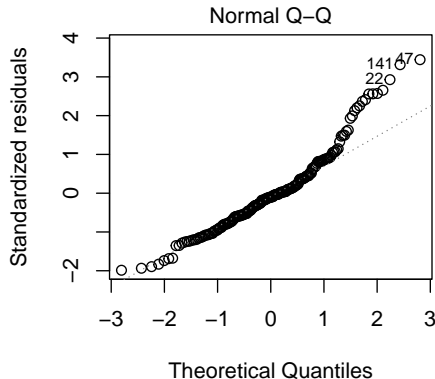
$$\hat{\mathbf{e}} \mid \mathbf{X} \sim N(\mathbf{0}, \sigma^2(\mathbf{I} - \mathbf{H})).$$

We can test whether the standardized residuals

## Detect Lack of Fit Using Residuals (2)

QQ plot: plot the sample quantiles of the **standardized** residuals against the expected quantiles from a assumed distribution.

- Departures from the straight line indicate departures from the assumed distribution.



# Studentized Residual

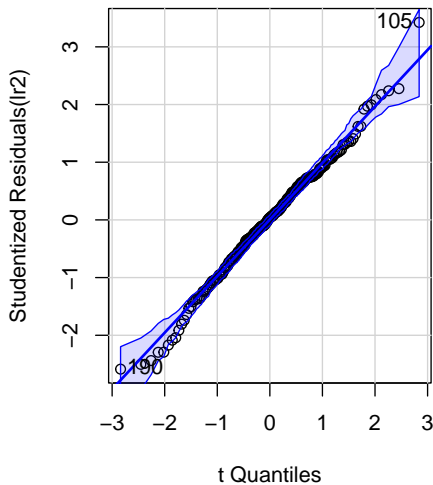
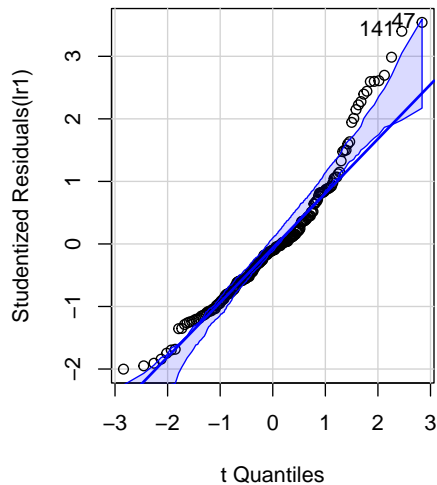
Let  $\hat{\boldsymbol{\beta}}$  be the WLS estimator based on the entire data set. A **studentized residual** is

$$\frac{\sqrt{w_i} \left( y_i - \mathbf{x}_i^T \hat{\boldsymbol{\beta}} \right)}{\hat{\sigma}_{(i)} \sqrt{1 - h_{ii}}},$$

where  $\sqrt{w_i} \left( y_i - \mathbf{x}_i^T \hat{\boldsymbol{\beta}} \right)$  is the Pearson residual, and  $h_{ii}$  is the  $(i, i)$ th entry of  $\mathbf{H}$ .

- Here,  $\hat{\sigma}_{(i)}^2$  is an estimator of  $\sigma^2$ , but not from the model based on the entire data set.
- Instead, we delete the  $i$ th observation and refit the model.  $\hat{\sigma}_{(i)}^2$  is the estimator of  $\sigma^2$  from the refitted model.

# More QQ Plot

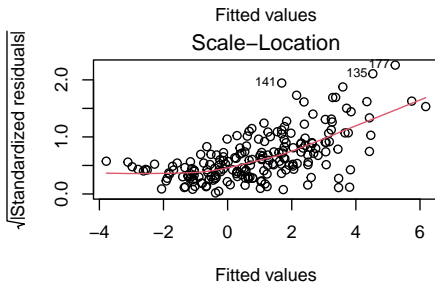
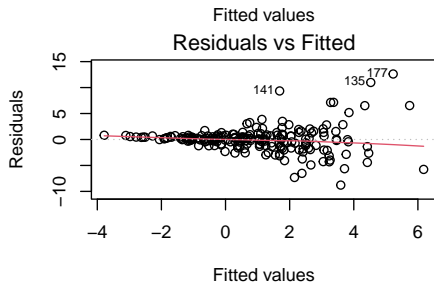
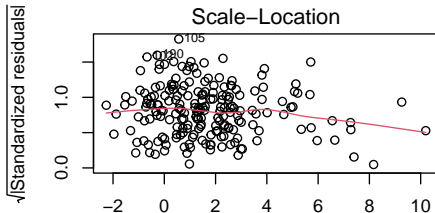
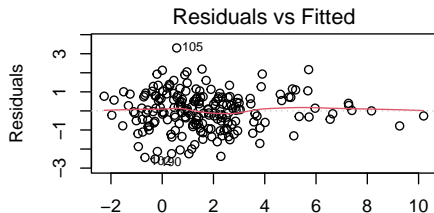


## Detect Lack of Fit Using Residuals (3)

The residuals should also be plotted against the fitted values to detect changes in variance.

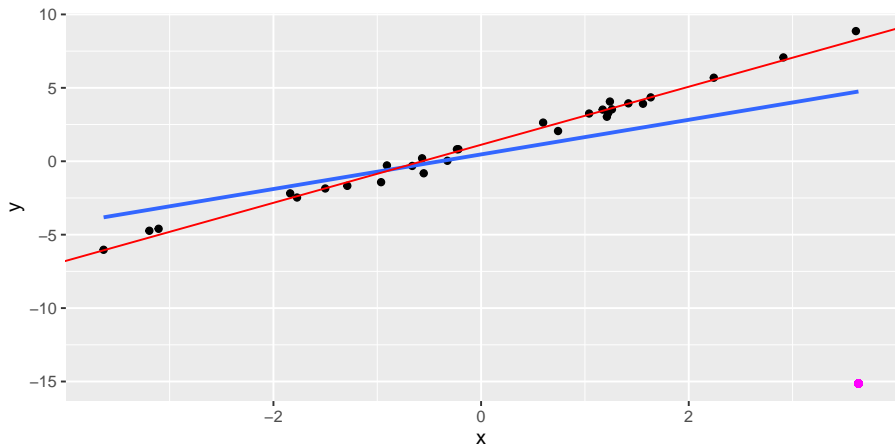
- If a systematic pattern is observed, the variance of standardized residuals is a function of fitted values, not homoskedasticity.
- Some plots do not show systematic deviation from zero but the variation depends on  $\hat{y}$ .

# Detect Lack of Fit Using Residuals (3)



# Influential Point or Outlier

An **influential point** or an **outlier** is a point that differs greatly from any other points.





# Influence Analysis

The general idea of **influence analysis** is to study changes in a specific part of the analysis when the data are perturbed.

- 1 The studentized residual is an example of influence analysis.
- 2 The **Cook's distance** is defined as

$$D_i = \frac{\left(\mathbf{X}\hat{\boldsymbol{\beta}}_{(i)} - \mathbf{X}\hat{\boldsymbol{\beta}}\right)^T \left(\mathbf{X}\hat{\boldsymbol{\beta}}_{(i)} - \mathbf{X}\hat{\boldsymbol{\beta}}\right)}{p\hat{\sigma}^2},$$

where  $\hat{\boldsymbol{\beta}}$  is the estimator with the whole data set and  $\hat{\boldsymbol{\beta}}_{(i)}$  is the estimator after deleting the  $i$ th observation.

- If the  $i$ th observation has a substantial influence, then we expect  $\hat{\boldsymbol{\beta}}$  differ much from  $\hat{\boldsymbol{\beta}}_{(i)}$ .
- A rule-of-thumb is that if  $D_i > 1$ , then it deserves some consideration.

# Cook's Distance and Leverage

In fact, we do not need to refit the model  $n$  times. It can be shown that

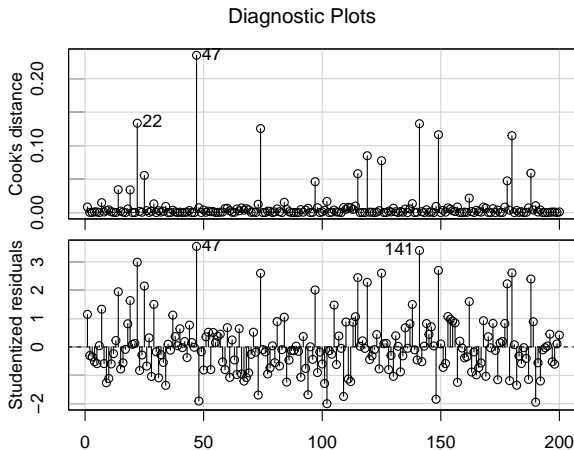
$$D_i = \frac{1}{p} r_i^2 \frac{h_{ii}}{1 - h_{ii}},$$

where  $r_i$  is the standardized residual, and  $h_{ii}$  is the leverage.

- If the leverage is close to 1, we probably have a large Cook's distance.

## Detect Lack of Fit Using Residuals (4)

Use the Cook's distance and other tools to detect influential observations/outliers, and other plots.



# DFBETAS

Another quantity to measure the influence of an observation is

$$\text{DFBETAS}_j = \frac{\hat{\beta}_j - \hat{\beta}_{j(i)}}{\hat{\sigma}_{(i)} \sqrt{\left[ (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \right]_{jj}}},$$

where  $\hat{\beta}_{j(i)}$  is the estimator of  $\beta_j$  without the  $i$ th observation, and  $\hat{\sigma}_{(i)}^2$  is an estimator of  $\sigma^2$  using the data set without  $i$ th observation.

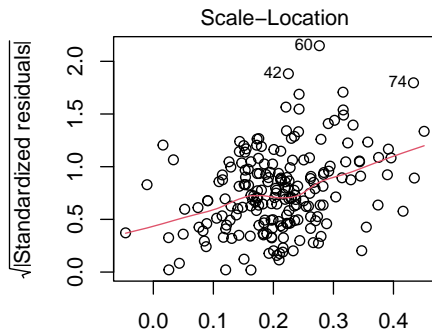
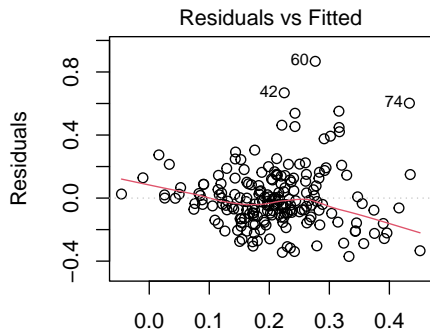
It measures how much the coefficient changes when the  $i$ th observation is deleted.

## OBS!

But keep in mind that the residual plots only suggest that something is wrong but never tell what is definitely wrong. We generate data from

$$E(Y | \mathbf{x}) = \frac{|x_1|}{2 + (1.5 + x_2)^2}$$

but fit a linear model  $E(Y | \mathbf{x}) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$ .



## Test Outlier: T-Test

Suppose that the  $i$ th observation is an outlier.

- The mean function for all other points are  $E(Y_j | \mathbf{x}_j) = \mathbf{x}_j^T \boldsymbol{\beta}$ ,  $j \neq i$ .
- The mean function for the outlier is  $E(Y_i | \mathbf{x}_i) = \mathbf{x}_i^T \boldsymbol{\beta} + \delta$  for some  $\delta$ .

Now we suspect the  $i$ th observation to be an outlier.

- ① Define a variable  $U$  that has 1 for the  $i$ th observation and 0 for all other elements.
- ② We can test  $H_0: \delta = 0$  versus  $H_1: \delta \neq 0$  using a t-test if we assume  $\text{Var}(Y | \mathbf{x}) = \sigma^2$  and data are normally distributed.

## Test Outlier: Leave-One-Out

Another approach works as follows.

- 1 Delete the  $i$ th observation from the data. Denote the remaining design matrix by  $\mathbf{X}_{(i)}$ .
- 2 Estimate  $\beta$  and  $\sigma^2$  using the remaining data. Denote the estimator by  $\hat{\beta}_{(i)}$  and  $\hat{\sigma}_{(i)}^2$ .
- 3 The predicted value of  $y_i$  is  $\hat{y}_{i(i)} = \mathbf{x}_i^T \hat{\beta}_{(i)}$ , which is independent of  $y_i$  if we assume observations are mutually independent.
- 4 Under the assumptions that  $\text{Var}(Y | \mathbf{x}) = \sigma^2$ , we can obtain

$$\text{Var}(y_i - \hat{y}_{i(i)} | \mathbf{X}) = \sigma^2 + \sigma^2 \mathbf{x}_i^T \left( \mathbf{X}_{(i)}^T \mathbf{X}_{(i)} \right)^{-1} \mathbf{x}_i.$$

- 5 Test  $E(y_i - \hat{y}_{i(i)} | \mathbf{X})$  using a t-test under the normality assumption.

It turns out that this procedure is equivalent to the t-test using regression.

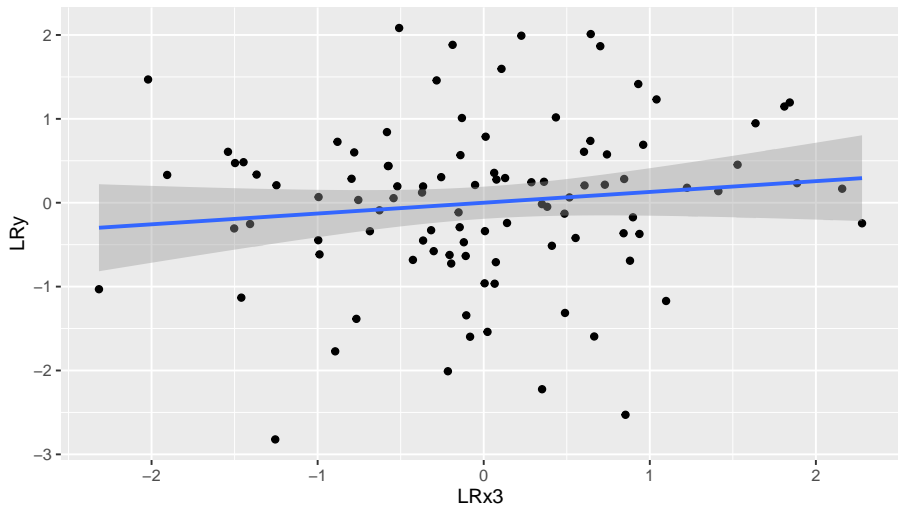
# Added Variable Plot

Suppose that we have fitted a model  $\hat{y} = \mathbf{x}^T \hat{\boldsymbol{\beta}}$  using the variables in  $\mathbf{x}$ , but we have an extra variable  $Z$ . An **added variable plot** or **partial regression plot** can reveal whether we can use  $Z$  to improve the model.

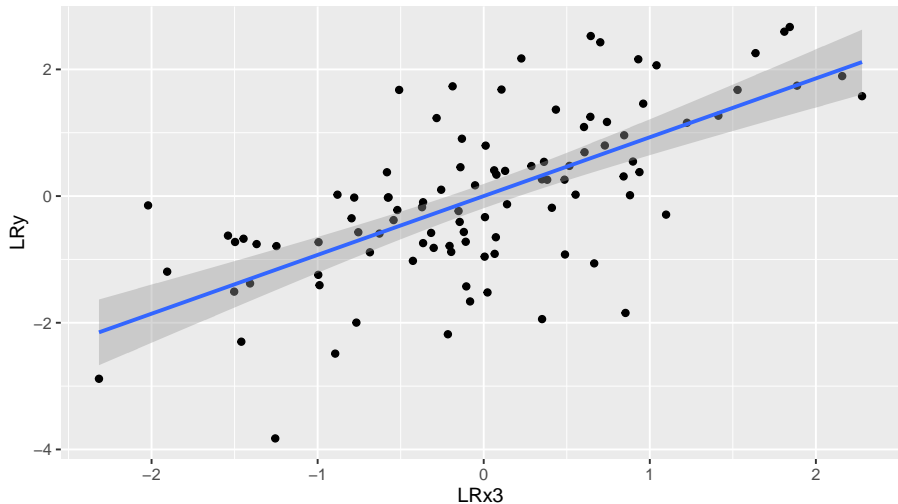
- Let  $\hat{\mathbf{e}}(Y | \mathbf{x})$  be the residual vector when regress  $Y$  on  $\mathbf{x}$  (part of  $Y$  not explained by  $\mathbf{x}$ ).
- Let  $\hat{\mathbf{e}}(Z | \mathbf{x})$  be the residual vector when regress  $Z$  on  $\mathbf{x}$  (part of  $Z$  not explained by  $\mathbf{x}$ ).
- We regress  $\hat{\mathbf{e}}(Y | \mathbf{x})$  on  $\hat{\mathbf{e}}(Z | \mathbf{x})$  and check whether part of  $Y$  not explained by  $\mathbf{x}$  can be explained by part of  $Z$  not explained by  $\mathbf{x}$ .



# Illustration 1

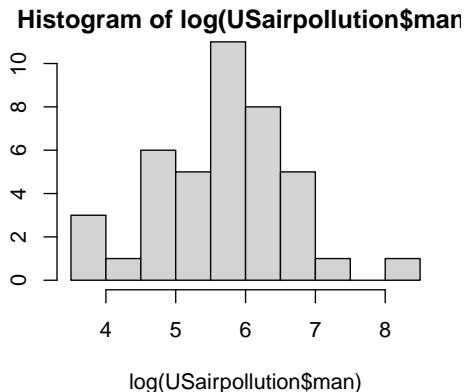
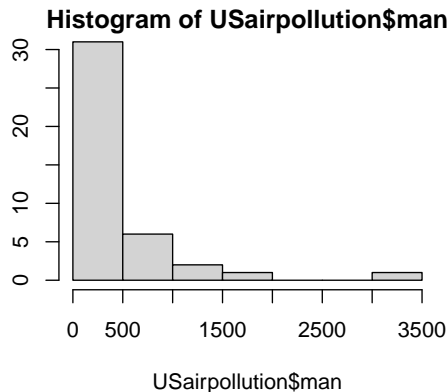


## Illustration 2



# Purpose of Transformation

The purpose of transformation (of the response variable and the regressors) is to make our assumptions more plausible.



# Box-Cox Transformation

The [Box-Cox transformation](#) is often used to transform the response variable so that the normal assumption is more plausible:

$$\psi(y, \lambda) = \begin{cases} \frac{y^\lambda - 1}{\lambda}, & \text{if } \lambda \neq 0, \\ \log(y), & \text{if } \lambda = 0. \end{cases}$$

In case  $y$  is not positive, we have [two-parameter Box-Cox transformation](#).

$$\psi(y, \lambda, \varepsilon) = \begin{cases} \frac{(y+\varepsilon)^\lambda - 1}{\lambda}, & \text{if } \lambda \neq 0, \\ \log(y + \varepsilon), & \text{if } \lambda = 0. \end{cases}$$

## Which $\lambda$ to Choose

The idea is to find the value of  $\lambda$  such that the residual from our regression model is most normal.

- The regression model is

$$\psi(y, \lambda) = \mathbf{x}^T \boldsymbol{\beta} + e.$$

- With normal error,  $\psi(y, \lambda) - \mathbf{x}^T \boldsymbol{\beta}$  should be normal.
- Find the  $\lambda$  value that maximizes the normal likelihood.

# Bootstrap

Another approach when the assumptions are violated is the **bootstrap** that can be used to handle

- ① residuals are not normally distributed,
- ② homoscedasticity is violated.

method = "case" eller method = "residual"?

Boot {car}

R Documentation

## Bootstrapping for regression models

### Description

This function provides a simple front-end to the `boot` function in the **boot** package that is tailored to bootstrapping based on regression models. Whereas `boot` is very general and therefore has many arguments, the `Boot` function has very few arguments.

### Usage

```
Boot(object, f=coef, labels=names(f(object)), R=999,  
      method=c("case", "residual"), ncores=1, ...)
```

method = "case" eller method = "residual"?

method = "residual" requires

- homoscedasticity
- $Y$  is random, but  $\mathbf{x}$  is not random
  - For example, we want to compare the caffeine content of Lindvalls Mörkrost och Lindvalls Brygg.

method = "case"

- does not require homoscedasticity
- requires both  $Y$  and  $\mathbf{x}$  are random.