

# Linear model

MS-C2128 Prediction and Time Series Analysis

Fall term 2020

# Week 1: Linear model

- 1 Definition and assumptions
- 2 Matrix representation of linear models
- 3 Least squares estimation
- 4 Fitted values and residuals
- 5 Variance decomposition and coefficient of determination
- 6 Significance testing
- 7 Confidence intervals
- 8 Prediction interval
- 9 Random explanatory variables
- 10 On optimality of least squares estimators

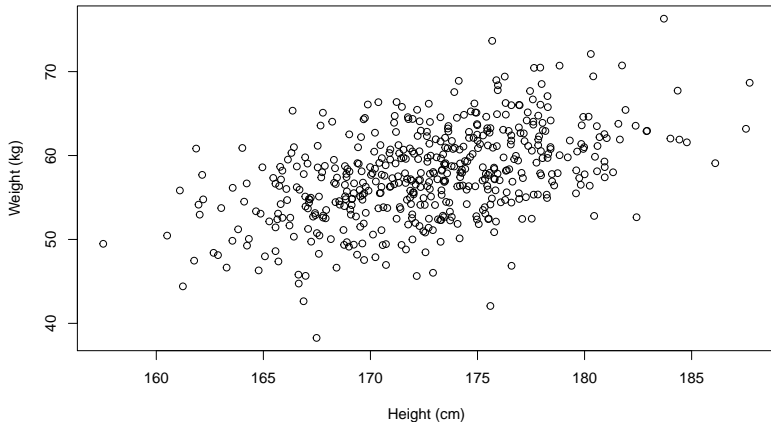
- 1 Definition and assumptions
- 2 Matrix representation of linear models
- 3 Least squares estimation
- 4 Fitted values and residuals
- 5 Variance decomposition and coefficient of determination
- 6 Significance testing
- 7 Confidence intervals
- 8 Prediction interval
- 9 Random explanatory variables
- 10 On optimality of least squares estimators

# Response variable and explanatory variables

- Variation in a *response variable*  $Y$  is explained by variation in *explanatory variables*  $X_1, \dots, X_k$ .
  - $y_1, \dots, y_n$  are the observed values of the response variable  $Y$ .
  - $x_{1j}, \dots, x_{nj}$  are the observed values of the explanatory variable  $X_j$ .

# Response variable and explanatory variable

Height and weight of adolescents



- Multiple linear model is given by

$$Y_i = \beta_0 + \beta_1 X_{i1} + \dots + \beta_k X_{ik} + \epsilon_i, \quad i = 1, \dots, n,$$

where

- The variables  $Y_i$  are random and the variables  $X_{i1}, \dots, X_{ik}$  are non-random.
- The residual (error term)  $\epsilon_i$  are random.
- The coefficients (regression parameters)  $\beta_0, \beta_1, \dots, \beta_k$  are constants.

## Remark

For simplicity, we assumed that the variables  $X_{i1}, \dots, X_{ik}$  are non-random. Then all the randomness of  $Y_i$  comes from the residual  $\epsilon_i$ .

## Remark

Conventional notation

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + \epsilon_i, \quad i = 1, \dots, n,$$

is used even though some of the variables are random and some are not.

# Standard assumptions

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + \epsilon_i, \quad i = 1, \dots, n,$$

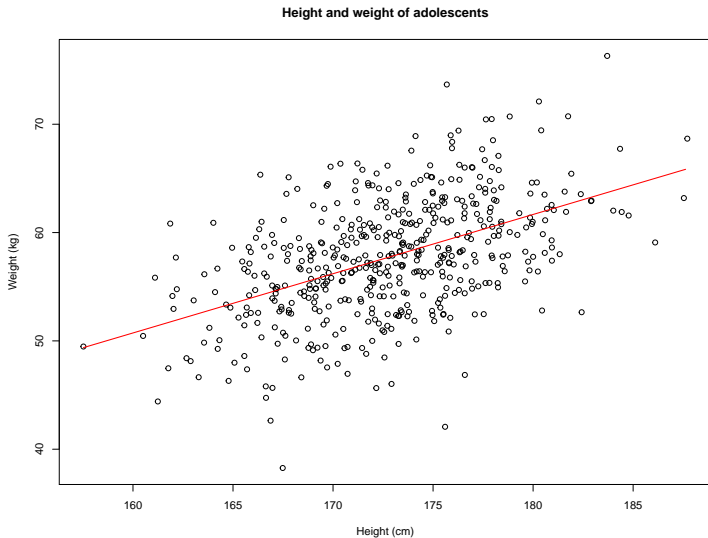
- (i) The explanatory variables are non-random.
- (ii) The explanatory variables are linearly independent.
- (iii)  $E[\epsilon_i] = 0$  for all  $i = 1, \dots, n$
- (iv)  $\text{var}(\epsilon_i) = \sigma^2$  for all  $i = 1, \dots, n$
- (v)  $\text{cor}(\epsilon_i, \epsilon_l) = 0$ , for  $i \neq l$

## Remark

In addition to the standard assumptions, it is customary to assume that the residuals are independent and identically distributed. Moreover, it is often assumed that the residuals are normally distributed, but this is not a necessary assumption.



# Response variable and explanatory variable



## (i) The variables $X_j$ are non-random.

- This is a strong assumption.
- However, linear regression analysis works perfectly well also in the case when the variables  $X_j$  are assumed to be random. In that case, the notations are a bit more complicated. (One has to consider conditional expected values instead of expected values.) Estimation procedures are exactly the same in the case when the  $X_j$  are assumed to be non-random and in the case when the  $X_j$  are assumed to be random.

# Standard assumptions (ii) and (iii):

## (ii) The explanatory variables are linearly independent.

- Multicollinearity makes it difficult to assess the effect of the different explanatory variables separately.
- If  $X_j$  can be given as a linear combination of the other explanatory variables, it can be removed from the model.
- Assumption (ii) guarantees that the so called least square estimators for the parameters  $\beta_0, \beta_1, \dots, \beta_k$  can be given in closed form.

## (iii) $E[\epsilon_i] = 0$ for all $i = 1, \dots, n$ .

- Assumption (ii) guarantees that there is no systematic bias. That is,

$$\begin{aligned} E[y_i] &= E[\beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + \epsilon_i] \\ &= \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + E[\epsilon_i] \\ &= \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik}. \end{aligned}$$

## Standard assumption (iv):

(iv)  $\text{var}(\epsilon_i) = \sigma^2$  for all  $i \neq l$ .

- The parameter  $\sigma^2$  is called the **residual variance**.
- If Assumption (iv) holds, the residuals  $\epsilon_i$  are said to be **homoscedastic**.
- If Assumption (iv) does not hold, the residuals  $\epsilon_i$  are said to be **heteroscedastic**.
- Heteroscedasticity makes standard least squares estimators unstable.
  - In the presence of heteroscedasticity one could apply the so called generalized least squares estimators. (This is beyond the scope of this course.)
- Homoscedasticity can be tested. (See Week 2, regression diagnostics.)

## Standard assumption (v):

**(v)  $\text{cor}(\epsilon_i, \epsilon_l) = 0$  for all  $i = 1, \dots, n$ .**

- The residuals are uncorrelated.
- If Assumption (v) does not hold, the residuals  $\epsilon_i$  are correlated.
- Correlatedness can make the regression parameter estimators inefficient and even biased.
  - In the presence of correlated residuals, one could apply dynamic regression models.
- Uncorrelatedness can be tested. (See Week 2, regression diagnostics.)

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + \epsilon_i, \quad i = 1, \dots, n,$$

If Assumptions (i)-(vi) hold, then

$$(iii)' \quad E[y_i] = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik}, \quad i = 1, \dots, n$$

$$(iv)' \quad \text{var}(y_i) = \sigma^2, \quad i = 1, \dots, n$$

$$(v)' \quad \text{cor}(y_i, y_j) = 0, \quad i \neq j$$

## Remark

If the residuals are assumed to be normally distributed, then

$$y_i \sim N(E[y_i], \sigma^2), \quad i = 1, \dots, n.$$

# Model part and random part

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + \epsilon_i, \quad i = 1, \dots, n,$$

- Linear model can be given as a sum of two parts:

$$y_i = E[y_i] + \epsilon_i, \quad i = 1, \dots, n$$

- The expected value  $E[y_i] = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik}$  is called the **systematic part** or the **model part**. This part depends on the  $x_j$ .
- The residual  $\epsilon_i$  forms the **random part**. This part is independent of  $x_j$ .

# Regression plane and regression parameters

- The systematic part  $E[y_i]$  defines the **regression plane**

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$$

in  $\mathbb{R}^{k+1}$ .

- The variance  $\sigma^2$  measures the variation of  $(x_{i1}, \dots, x_{ik}, y_i) \in \mathbb{R}^{k+1}$  around the regression plane.
- The regression parameters  $\beta_j$  can be interpreted as follows:
  - Assume that the value of the explanatory variable  $x_j$  grows by one ( $x_j \rightarrow x_j + 1$ ) and assume that the values of all the other explanatory variables remain unchanged. The parameter  $\beta_j$  models the change in the expected value of the response variable  $y$  as the value of  $x_j$  changes by one unit:

$$E[y] \rightarrow E[y] + \beta_j.$$



# Content

- 1 Definition and assumptions
- 2 Matrix representation of linear models**
- 3 Least squares estimation
- 4 Fitted values and residuals
- 5 Variance decomposition and coefficient of determination
- 6 Significance testing
- 7 Confidence intervals
- 8 Prediction interval
- 9 Random explanatory variables
- 10 On optimality of least squares estimators

# Matrix representation

A linear model can be represented as  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ , where

$$\mathbf{y} = (y_1, y_2, \dots, y_n)^\top = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad \boldsymbol{\beta} = (\beta_0, \dots, \beta_k)^\top = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix}$$
$$\mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1k} \\ 1 & x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{nk} \end{bmatrix} \quad \boldsymbol{\epsilon} = (\epsilon_1, \dots, \epsilon_n)^\top = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

# Standard assumptions: Matrix representation

A linear model can be represented as  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$

- (i) The elements of  $\mathbf{X}$  are non-random.
- (ii) The columns of  $\mathbf{X}$  are linearly independent.
- (iii)  $E[\boldsymbol{\epsilon}] = \mathbf{0}$ .
- (iv)-(v)  $\text{cov}(\boldsymbol{\epsilon}) = \sigma^2 \mathbf{I}$ .

## Remark

- Let  $\mathbf{z} = (z_1, z_2, \dots, z_p)^\top$  be a random vector, where the elements  $z_1, z_2, \dots, z_p$  are random.
- The expected value  $\boldsymbol{\mu} = E[\mathbf{z}]$  of a vector  $\mathbf{z}$  is taken componentwise

$$\boldsymbol{\mu} = E[\mathbf{z}] = (E[z_1], E[z_2], \dots, E[z_p])^\top \in \mathbb{R}^p$$

- The covariance  $\Sigma = \text{cov}(\mathbf{z})$  of  $\mathbf{z}$  refers to the matrix

$$\begin{aligned} \Sigma = \text{cov}(\mathbf{z}) &= E \left[ (\mathbf{z} - E[\mathbf{z}])(\mathbf{z} - E[\mathbf{z}])^\top \right] \\ &= \begin{bmatrix} \text{var}(z_1) & \text{cov}(z_1, z_2) & \text{cov}(z_1, z_3) & \cdots & \text{cov}(z_1, z_p) \\ \text{cov}(z_2, z_1) & \text{var}(z_2) & \text{cov}(z_2, z_3) & \cdots & \text{cov}(z_2, z_p) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \text{cov}(z_p, z_1) & \text{cov}(z_p, z_2) & \text{cov}(z_p, z_3) & \cdots & \text{var}(z_p) \end{bmatrix} \in \mathbb{R}^{p \times p} \end{aligned}$$

- 1 Definition and assumptions
- 2 Matrix representation of linear models
- 3 Least squares estimation**
- 4 Fitted values and residuals
- 5 Variance decomposition and coefficient of determination
- 6 Significance testing
- 7 Confidence intervals
- 8 Prediction interval
- 9 Random explanatory variables
- 10 On optimality of least squares estimators

# Least squares estimation

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + \epsilon_i, \quad i = 1, \dots, n,$$

The parameters  $\beta_0, \beta_1, \dots, \beta_k$  are usually estimated using the **least squares estimation** method:

- The sum of the squared residuals

$$\sum_{i=1}^n \epsilon_i^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_k x_{ik})^2$$

is minimized with respect to the parameters  $\beta_0, \beta_1, \dots, \beta_k$ :

- The partial derivatives with respect to  $\beta_0, \beta_1, \dots, \beta_k$  are calculated.
- $k + 1$  equations are obtained by setting the derivatives equal to zero.
- If the explanatory variable matrix  $\mathbf{X}$  is of full rank, there exist a unique solution.
- The **least squares estimators**  $b_j$  for  $\beta_j$  are obtained.

# Estimator vs estimate

- The *random* object  $b_j$  is called an **estimator**.
- If we calculate the numerical value for  $b_j$  from an observed sample, we obtain an **estimate**, that is a non-random realization of the corresponding estimator.

## Remark

An estimator and the corresponding estimate are often denoted by the same symbol. The interpretation depends on the context.

# Least squares estimator

- Let  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$  and let the matrix  $\mathbf{X}$  be of full rank. Then the **least squares estimator of  $\boldsymbol{\beta}$**  is

$$\mathbf{b} = (b_0, \dots, b_k)^\top = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}.$$

- If the standard assumptions (i)-(v) hold, then

$$E[\mathbf{b}] = \boldsymbol{\beta} \quad \text{and} \quad \text{cov}(\mathbf{b}) = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}.$$

- As the expected value  $E[\mathbf{b}] = \boldsymbol{\beta}$ , the least squares estimator  $\mathbf{b}$  is an **unbiased** estimator of  $\boldsymbol{\beta}$ .
- If the residuals are normally distributed, then

$$\mathbf{b} \sim N_{k+1} \left( \boldsymbol{\beta}, \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} \right)$$



# Content

- 1 Definition and assumptions
- 2 Matrix representation of linear models
- 3 Least squares estimation
- 4 Fitted values and residuals**
- 5 Variance decomposition and coefficient of determination
- 6 Significance testing
- 7 Confidence intervals
- 8 Prediction interval
- 9 Random explanatory variables
- 10 On optimality of least squares estimators

# Fitted values and residuals

- The **fitted values** of the estimated model are given as follows:

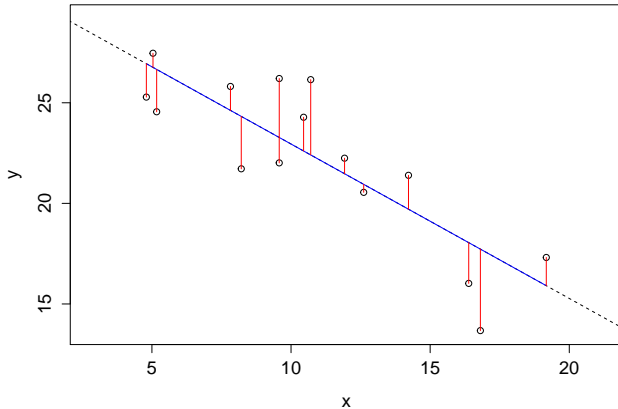
$$\hat{y}_i = b_0 + b_1 x_{i1} + \dots + b_k x_{ik}$$

- This gives the estimated values of the response variable  $y$  in points  $i$ .
- The estimated **residuals**:

$$e_i = y_i - \hat{y}_i = y_i - b_0 - b_1 x_{i1} - \dots - b_k x_{ik}$$

- The residuals are the differences between the observed values of the response variables and the estimated values of the response variables.
- The model explains the variation in  $y$  the better the smaller are the estimated residuals  $e_i$ .

# Fitted values and residuals



# Fitted values and residuals

Under the standard assumptions (i)-(v), we have that:

- For the fitted values:

$$\hat{\mathbf{y}} = \mathbf{X}\mathbf{b} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = \mathbf{P}\mathbf{y}$$

$$E[\hat{\mathbf{y}}] = \mathbf{X}\boldsymbol{\beta}$$

$$\text{cov}(\hat{\mathbf{y}}) = \sigma^2 \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top = \sigma^2 \mathbf{P}.$$

- For the residuals:

$$\mathbf{e} = \mathbf{y} - \hat{\mathbf{y}} = \left( \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \right) \mathbf{y} = (\mathbf{I} - \mathbf{P})\mathbf{y} = \mathbf{M}\mathbf{y}$$

$$E[\mathbf{e}] = \mathbf{0}$$

$$\text{cov}(\mathbf{e}) = \sigma^2 \left( \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \right) = \sigma^2 (\mathbf{I} - \mathbf{P}) = \sigma^2 \mathbf{M}.$$

# Fitted values and residuals

- The  $n \times n$ -matrices

$$\mathbf{P} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$$

$$\mathbf{M} = \mathbf{I} - \mathbf{P} = \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$$

are **symmetric** ja **idempotent**:

$$\mathbf{P}^\top = \mathbf{P}$$

$$\mathbf{P}^2 = \mathbf{P}$$

$$\mathbf{M}^\top = \mathbf{M}$$

$$\mathbf{M}^2 = \mathbf{M}$$

- Moreover:  $\mathbf{PM} = \mathbf{MP} = \mathbf{0}$
- The above mentioned properties of the matrices  $\mathbf{P}$  and  $\mathbf{M}$  play a crucial role when distributions related to estimating and testing the linear regression parameters are derived.

# Content

- 1 Definition and assumptions
- 2 Matrix representation of linear models
- 3 Least squares estimation
- 4 Fitted values and residuals
- 5 Variance decomposition and coefficient of determination**
- 6 Significance testing
- 7 Confidence intervals
- 8 Prediction interval
- 9 Random explanatory variables
- 10 On optimality of least squares estimators

# Residual variance

- If the standard assumptions (i)-(v) hold, then an **unbiased estimator** for  $\text{var}(\epsilon_i) = \sigma^2$  can be given as

$$s^2 = \frac{1}{n - k - 1} \sum_{i=1}^n e_i^2,$$

where  $e_i = y_i - \hat{y}_i$ ,  $k$  is the number of the explanatory variables and where  $n$  is the number of the observations.

- This is actually simply the sample variance of the  $e_i$  as there are  $k + 1$  estimated parameters and as

$$\sum_{i=1}^n e_i = 0 \implies \bar{e} = \frac{1}{n} \sum_{i=1}^n e_i = 0 \quad \text{and}$$
$$s^2 = \frac{1}{n - k - 1} \sum_{i=1}^n (e_i - \bar{e})^2 = \frac{1}{n - k - 1} \sum_{i=1}^n e_i^2$$

# Variance decomposition

The goal in regression analysis is to explain the variation in the response variable by the variation in the explanatory variables. Success in this can be assessed using the so called variance decomposition.



# Variance decomposition: Definition

- The total sum of squares  $SST = \sum_{i=1}^n (y_i - \bar{y})^2$ 
  - measures the variation of the observations  $y_i$
  - the sample variance of the  $y_i$ ,  $s_y^2 = SST/(n - 1)$
- The error sum of squares  $SSE = \sum_{i=1}^n e_i^2$ 
  - measures the variation of the estimated residuals  $e_i$
  - the sample variance of the  $e_i$ ,  $s^2 = SSE/(n - k - 1)$
- The model sum of squares  $SSM = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$ 
  - measures the part of the variation of the observations  $y_i$  that is explained by the estimated regression model.
- The variance decomposition:  $SST = SSM + SSE$

# Coefficient of determination: Definition

- The variance decomposition  $SST = SSM + SSE$  tells about the goodness of the estimated regression model.
  - The larger the proportion of the model sum of squares  $SSM$  (that is, the smaller the proportion of the error sum of squares  $SSE$ ) is of the total sum of squares  $SST$ , the better the estimated model explains the variation in the response variable.
- This motivates the use of the **coefficient of determination**

$$R^2 = 1 - \frac{SSE}{SST} = \frac{SSM}{SST} \in [0, 1]$$

as a measure of goodness of a regression model.

- The coefficient of determination measures the proportion of the variance in the response variable that is predictable from the explanatory variables.

# Coefficient of determination: Properties

$$R^2 = 1 - \frac{SSE}{SST} = \frac{SSM}{SST} \in [0, 1]$$

- The following conditions are equivalent:
  - 1  $R^2 = 1$ .
  - 2  $e_i = 0$  for all  $i = 1, 2, \dots, n$ .
  - 3 all the observed  $(x_{i1}, x_{i2}, \dots, x_{ik}, y_i)$  lie on the same plane in  $\mathbb{R}^{k+1}$ .
  - 4 "The model explains the variation in the response variable perfectly."
- The following conditions are equivalent:
  - 1  $R^2 = 0$ .
  - 2  $b_1 = b_2 = \dots = b_k = 0$ .
  - 3 "The model does not explain the variation in the response variable at all."
- The coefficient of determination is equal to the square of the sample Pearson correlation coefficient of the observed values  $y_i$  and the fitted values  $\hat{y}_i$ .

# Content

- 1 Definition and assumptions
- 2 Matrix representation of linear models
- 3 Least squares estimation
- 4 Fitted values and residuals
- 5 Variance decomposition and coefficient of determination
- 6 Significance testing**
- 7 Confidence intervals
- 8 Prediction interval
- 9 Random explanatory variables
- 10 On optimality of least squares estimators

# Significance testing

Let the standard assumptions (i-v) hold. Now, the expected value and the covariance matrix of the least squares estimator  $\mathbf{b}$  are

$$\begin{aligned}E[\mathbf{b}] &= \boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_k) \\D^2(\mathbf{b}) &= \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1}.\end{aligned}$$

An unbiased estimator for the covariance matrix  $D^2(\mathbf{b})$  is

$$\hat{D}^2(\mathbf{b}) = s^2(\mathbf{X}^\top \mathbf{X})^{-1}$$

where  $s^2$  is the unbiased estimator of the residual variance  $\sigma^2$ ,

$$s^2 = \frac{1}{n - k - 1} \sum_{i=1}^n e_i^2.$$

# Significance testing

- The coefficient of determination measures how big part of the variation in the response variable is explained by the estimated regression model. However, it does not tell whether the coefficient of determination differs significantly from zero or not.
- The significance of the coefficient of determination can be tested by examining how likely it is to obtain as large value as was obtained, under the assumption that the residuals explain all of the variation in the response variable.
- Note that if the residuals explain all of the variation in the response variable, then  $\beta_j = 0$  for all  $j = 1, \dots, k$ .

# Significance testing

The significance of the coefficient of determination can be tested by applying the **permutation test** as follows:

- 1 The null hypothesis is  $H_0 : \beta_j = 0$  for all  $j = 1, \dots, k$  and the alternative hypothesis is  $H_1 : \beta_j \neq 0$  for at least one  $j$ .
- 2 Calculate the value  $R^2$  from the original observations  $(y_i, \mathbf{x}_i)$ ,  $\mathbf{x}_i = (x_{i1}, \dots, x_{ik})$ .
- 3 Form  $n$  new pairs  $(y_i, \mathbf{x}_i)$ , such that each original explanatory variable vector  $\mathbf{x}_i$  (and each original response variables  $y_i$ ) is used exactly once in the new sample. Construct all the possible  $n!$  permuted samples.
- 4 For each sample  $p$ , calculate the corresponding coefficient of determination,  $R_p^2$ . You obtain  $n!$  values.
- 5 Order the values  $R_p^2$  from the smallest to the largest and calculate the empirical  $(1 - \alpha) \cdot 100$ th percentile from the sample. If the original  $R^2$  is larger than the calculated percentile, the coefficient of determination is considered significant on level  $\alpha$ .

# Significance of one single regression parameter

Permutation test can be applied also in testing the significance of one single regression parameter  $\beta_j$ .

- 1 The null hypothesis is  $H_0 : \beta_j = 0$  and the alternative hypothesis is  $H_1 : \beta_j \neq 0$ .
- 2 Calculate the value  $R^2$  from the original observations  $(y_i, \mathbf{x}_i)$ ,  $\mathbf{x}_i = (x_{i1}, \dots, x_{ik})$ .
- 3 Pair each  $y_i$  with the original  $\mathbf{x}_i = (x_{i1}, \dots, x_{ik})$ ,  $i = 1, \dots, n$ , but permute the components (and only the components)  $x_{ij}$  that correspond to the parameter  $\beta_j$ . You can construct  $n!$  permuted samples. Each sample size is  $n$ .
- 4 For each sample  $p$ , calculate the corresponding coefficient of determination,  $R_p^2$ . You obtain  $n!$  values.
- 5 Order the values  $R_p^2$  from the smallest to the largest and calculate the empirical  $(1 - \alpha) \cdot 100$ th percentile from the sample. If the original  $R^2$  is larger than the calculated percentile, the null hypothesis is rejected.



## Remark

If the sample size  $n \geq 10$  it is not reasonable to construct all the  $n!$  permuted samples. In that case, instead of considering all the possible permutations, one can take for example 1000 or 10000 randomly chosen permuted samples and base the test on using these.

Let the standard assumptions hold. Assume also that the residuals are normally distributed. Now, the significance of the regression model can be tested using the  $F$  test statistic.

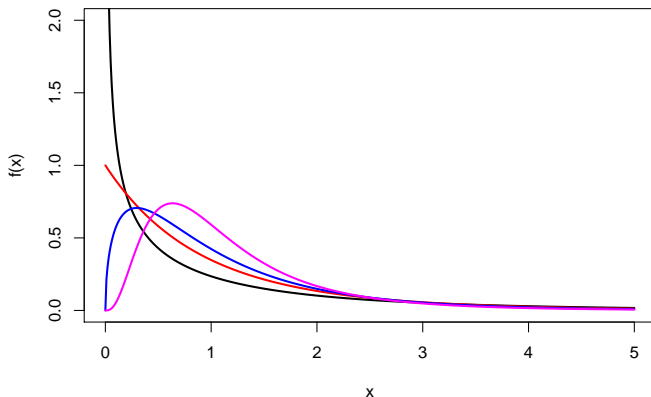
- 1 The null hypothesis is  $H_0 : \beta_j = 0$  for all  $j = 1, \dots, k$  and the alternative hypothesis is  $H_1 : \beta_j \neq 0$  for at least one  $j$ .
- 2 The  $F$  test statistic:

$$F = \frac{n - k - 1}{k} \frac{R^2}{1 - R^2} = \frac{n - k - 1}{k} \frac{SSM}{SSE}$$

follows, under the null hypothesis, the  $F(k, n - k - 1)$  distribution.

Note that this test statistic is reliable only when the residuals are normally distributed!

Probability density function of F-distribution



Probability density functions of the  $F(k, n - k - 1)$  distribution for  $n = 20$  and  $k = 1$  (black),  $k = 2$  (red),  $k = 3$  (blue) and  $k = 8$  (purple).

# Significance of one single regression parameter, $\epsilon \sim N(0, \sigma^2)$

Let the standard assumptions hold. Assume also that the residuals are normally distributed.

- Now

$$\mathbf{b} \sim N_{k+1}(\boldsymbol{\beta}, D^2(\mathbf{b})),$$

where  $D^2(\mathbf{b}) = \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1}$ .

- Thus

$$\frac{b_j - \beta_j}{s_{bj}} \sim t(n - k - 1),$$

where  $t(n - k - 1)$  denotes the  $t$  distribution with  $n - k - 1$  degrees of freedom and  $s_{bj}^2 = [\hat{D}^2(\mathbf{b})]_{jj}$  is the estimated variance of  $\beta_j$ . Here  $[\hat{D}^2(\mathbf{b})]_{jj}$  denotes the  $jj$  element of the estimated covariance matrix  $\hat{D}^2(\mathbf{b}) = \mathbf{s}^2(\mathbf{X}^\top \mathbf{X})^{-1}$ .

# Significance of one single regression parameter,

$$\epsilon \sim N(0, \sigma^2)$$

Let the standard assumptions hold. Assume also that the residuals are normally distributed. Now, the significance of the regression parameter  $\beta_j$  can be tested using the  $t$  test statistic.

- 1 The null hypothesis is  $H_0 : \beta_j = 0$  and the alternative hypothesis is  $H_1 : \beta_j \neq 0$ .
- 2 The  $t$  test statistic:

$$t = \frac{b_j}{s_{bj}}, \quad j = 0, 1, 2, \dots, k,$$

follows, under the null hypothesis, the  $t$  distribution with  $n - k - 1$  degrees of freedom.

Note that this test statistic is reliable only when the residuals are normally distributed!

# Content

- 1 Definition and assumptions
- 2 Matrix representation of linear models
- 3 Least squares estimation
- 4 Fitted values and residuals
- 5 Variance decomposition and coefficient of determination
- 6 Significance testing
- 7 Confidence intervals**
- 8 Prediction interval
- 9 Random explanatory variables
- 10 On optimality of least squares estimators

## Fact

A level  $(1 - \alpha)$  confidence interval for a parameter  $\theta$  is a random interval that contains the true (non-random) parameter value  $\theta$  with probability  $(1 - \alpha)$ .

If one calculates a  $(1 - \alpha)$  confidence interval for a parameter  $\theta$  from 500 independent samples of i.i.d. observations, then approximately  $500 \times (1 - \alpha)$  of the intervals do contain the true value  $\theta$ .

## Bootstrap confidence interval for $\beta_j$

A  $(1 - \alpha)$  bootstrap confidence interval for the regression parameter  $\beta_j, j = 0, 1, \dots, k$  can be obtained as follows.

- 1 Select  $n$  data points randomly with replacement from the original observations  $(y_i, \mathbf{x}_i), i = 1, \dots, n$ . Each data point  $(y_i, \mathbf{x}_i)$  can be selected once, multiple times, or not at all.
- 2 Calculate a new parameter vector estimate  $\mathbf{b}_b$  from the new sample formed in the previous step.
- 3 Repeat the previous steps  $m - 1$  times. (For example 999 times.)
- 4 Order all the obtained estimates  $b_j = (\mathbf{b}_b)_j$  from the smallest to the largest. Include also the original estimate.
- 5 Set the lower end of the bootstrap confidence interval to be smaller than or equal to the  $[\frac{\alpha}{2} \times m]$ th ordered estimate and set the upper end of the bootstrap confidence interval to be larger than or equal to the  $[(1 - \frac{\alpha}{2}) \times m]$ th ordered estimate.



Some remarks:

- In regression settings, one can calculate bootstrap confidence intervals for  $\beta_j$  also by bootstrapping the residuals.
- Bootstrap confidence intervals are distribution free.
- The estimate is the better the larger the original sample size is and the larger the number ( $m$ ) of the bootstrap samples is.

## Confidence interval for $\beta_j, \epsilon \sim N(0, \sigma^2)$

Let the standard assumptions hold. Assume also that the residuals are normally distributed.

- Now, as noted above,

$$\frac{b_j - \beta_j}{s_{bj}} \sim t(n - k - 1),$$

where  $t(n - k - 1)$  is the  $t$  distribution with  $n - k - 1$  degrees of freedom and  $s_{bj}^2$  is the estimated variance of  $\beta_j$ .

## Confidence interval for $\beta_j, \epsilon \sim N(0, \sigma^2)$

Let the standard assumptions hold. Assume also that the residuals are normally distributed. Under these assumptions, a level  $(1 - \alpha)$  confidence interval for  $\beta_j$  can be given as

$$(b_j - t_{1-\alpha/2} s_{bj}, b_j + t_{1-\alpha/2} s_{bj}),$$

where  $t_{1-\alpha/2}$  is the  $(1 - \alpha/2) \cdot 100$ th percentile of the  $t(n - k - 1)$  distribution.

# Content

- 1 Definition and assumptions
- 2 Matrix representation of linear models
- 3 Least squares estimation
- 4 Fitted values and residuals
- 5 Variance decomposition and coefficient of determination
- 6 Significance testing
- 7 Confidence intervals
- 8 Prediction interval**
- 9 Random explanatory variables
- 10 On optimality of least squares estimators

# Prediction interval

The predicted value of  $y$  for a fixed  $\tilde{\mathbf{x}} = (\tilde{x}_1, \dots, \tilde{x}_k)$  is

$$\hat{y} = b_0 + b_1\tilde{x}_1 + b_2\tilde{x}_2 + \dots + b_k\tilde{x}_k.$$

This prediction is done under uncertainties as

- 1 the regression parameters have been estimated and
- 2 the residuals explain part of the variation in  $y$ .

$\implies$  We wish to obtain an interval estimator that contains the true value with probability  $(1 - \alpha)$ .

# Prediction interval

One can apply bootstrapping in constructing interval estimates for the predicted value of  $y$ , for fixed  $\tilde{\mathbf{x}} = (\tilde{x}_1, \dots, \tilde{x}_k)$ , as follows.

- 1 Calculate one bootstrap estimate  $\mathbf{b}_b$  as described above.
- 2 Select a residual  $e_b$  randomly from the original estimated residuals  $e_i$  and calculate a bootstrap prediction  $y_b = \tilde{\mathbf{x}}_*^\top \mathbf{b}_b + e_b$ , where  $\tilde{\mathbf{x}}_* = (1, \tilde{x}_1, \dots, \tilde{x}_k)^\top$ .
- 3 Repeat the previous steps  $m - 1$  times. (For example 999 times.)
- 4 Order all the obtained predictions  $y_b$  from the smallest to the largest. Include also the original prediction.
- 5 Set the lower end of the bootstrap prediction interval to be smaller than or equal to the  $[\frac{\alpha}{2} \times m]$ th ordered bootstrap prediction and set the upper end of the bootstrap prediction interval to be larger than or equal to the  $[(1 - \frac{\alpha}{2}) \times m]$ th ordered bootstrap prediction.

Let the standard assumptions hold. Assume also that the residuals are normally distributed. Now, a prediction interval for  $y$  conditioned on  $\tilde{\mathbf{x}} = (\tilde{x}_1, \dots, \tilde{x}_k)$  can be given as

$$\tilde{\mathbf{x}}_*^\top \mathbf{b} \pm t_{1-\alpha/2} s \left[ 1 + \tilde{\mathbf{x}}_*^\top (\mathbf{X}^\top \mathbf{X})^{-1} \tilde{\mathbf{x}}_* \right]^{\frac{1}{2}},$$

where  $s^2$  is the estimated residual variance and  $t_{1-\alpha/2}$  is the  $(1 - \alpha/2) \cdot 100$ th percentile of the  $t(n - k - 1)$  distribution.

# Content

- 1 Definition and assumptions
- 2 Matrix representation of linear models
- 3 Least squares estimation
- 4 Fitted values and residuals
- 5 Variance decomposition and coefficient of determination
- 6 Significance testing
- 7 Confidence intervals
- 8 Prediction interval
- 9 Random explanatory variables**
- 10 On optimality of least squares estimators



# Random explanatory variables

Everything above can be applied directly also when the matrix  $\mathbf{X}$  is random.

# Random explanatory variables

Standard assumptions, when the explanatory variables are random are given as follows.

- (i) The explanatory variables are random.
- (ii) The explanatory variables are linearly independent.
- (iii)  $E[\epsilon | \mathbf{X}] = \mathbf{0}$  for all  $i = 1, \dots, n$
- (iv)-(v)  $\text{cov}(\epsilon | \mathbf{X}) = \sigma^2 \mathbf{I}$ .

# Content

- 1 Definition and assumptions
- 2 Matrix representation of linear models
- 3 Least squares estimation
- 4 Fitted values and residuals
- 5 Variance decomposition and coefficient of determination
- 6 Significance testing
- 7 Confidence intervals
- 8 Prediction interval
- 9 Random explanatory variables
- 10 On optimality of least squares estimators**

# Positive definite matrices

A  $d \times d$  matrix  $\mathbf{C}$  (with real valued or complex valued elements) is positive definite, if  $\mathbf{a}^* \mathbf{C} \mathbf{a} > 0$ , for all  $\mathbf{a} \in \mathbb{C}^d \setminus \{\mathbf{0}\}$ .

A  $d \times d$  matrix  $\mathbf{C}$  (with real valued or complex valued elements) is positive semidefinite, if  $\mathbf{a}^* \mathbf{C} \mathbf{a} \geq 0$ , for all  $\mathbf{a} \in \mathbb{C}^d \setminus \{\mathbf{0}\}$ .

# Positive definite matrices

It follows from the definition of positive definiteness that all positive definite matrices are Hermite symmetric. That is, if a square matrix  $\mathbf{C}$  is positive definite, then  $\mathbf{C} = \mathbf{C}^*$ . For real square matrices, positive definiteness can be defined as follows.

A  $d \times d$  matrix  $\mathbf{M}$  with real valued elements is positive definite if it is symmetric and if  $\mathbf{a}^\top \mathbf{M} \mathbf{a} > 0$  for all  $\mathbf{a} \in \mathbb{R}^d \setminus \{\mathbf{0}\}$ .

A  $d \times d$  matrix  $\mathbf{M}$  with real valued elements is positive semidefinite if it is symmetric and if  $\mathbf{a}^\top \mathbf{M} \mathbf{a} \geq 0$  for all  $\mathbf{a} \in \mathbb{R}^d \setminus \{\mathbf{0}\}$ .

If the standard assumptions hold, then the least squares estimator  $\mathbf{b}$  is

- 1 **unbiased** ( $E[\mathbf{b}] = \beta$ )
- 2 **efficient** in the sense that if  $\mathbf{b}^*$  is another unbiased linear estimator for  $\beta$ , then the matrix  $\mathbf{C} := \text{cov}(\mathbf{b}^*) - \text{cov}(\mathbf{b})$  is positive semidefinite.
- 3 **consistent**.

Standard least squares estimator is not always efficient.

- Standard least squares estimator is not necessarily efficient if the standard assumption (iv) homoscedasticity and/or the standard assumption (v) uncorrelatedness do not hold.
- If the standard assumptions hold, but the vector  $\beta$  has linear restrictions, then the corresponding restricted least squares estimator is efficient.

## References:

- 1 J. S. Milton, J.S., Arnold, J.C. (1995): Introduction to Probability and Statistics, McGraw-Hill Inc
- 2 Hogg, R.V., McKean, J.W., Craig, A.T. (2005): Introduction to Mathematical Statistics, Pearson Education.
- 3 Davison, A.C., Hinkley, D.V. (2009): Bootstrap Methods and their Applications, Cambridge University Press
- 4 Belsley, D.A., Kuh, E., Welsch, R.E. (2005): Regression Diagnostics: Identifying Influential Data and Sources of Collinearity, Wiley
- 5 Harrell, F. E. Jr. (2015): Regression Modeling Strategies: With Applications to Linear Models, Logistic and Ordinal Regression, and Survival Analysis, Springer



- Regression diagnostics
  - Regression graphics
  - Outlying observations
  - Constant/Non-constant regression parameters
  - Multicollinearity
  - Heteroscedasticity
  - Normality/Non-normality
  - Prediction capability
- Model selection
  - Model selection tests and strategies
  - Model selection criteria
  - Linearization