# Statistical Methods for Economists

## Lecture 4

Multiple Linear Regression

**SILESIAN UNIVERSITY**
SCHOOL OF BUSINESS
ADMINISTRATION IN KARVINA

**David Bartl**
Statistical Methods for Economists
INM/BASTE

# Outline of the lecture

- Introduction:  Simple Linear Regression & Least Squares Method

- Multiple Linear Regression:  Introduction

- Multiple Linear Regression:  Summary & Background

- The Classical Assumptions

- The Coefficient of Determination ($R^2$)

- Further Theorems, Tests of Hypotheses and Confidence Intervals

- Two-sample $t$-test for the difference of the population means // $\sigma_X = \sigma_Y$

- Simple linear regression without the intercept term

# Introduction

- Simple Linear Regression

- Motivation

- Example

- Least Squares Method

- Generalization

- Multiple Linear Regression:  Introduction

- Multiple Linear Regression:  Notation

# Simple Linear Regression: Motivation

## Motivation:

Assume a dataset $(y_i, x_{i1})_{i=1}^{n}$ of $n$ statistical units, i.e. we are given $n$ pairs $(y_1, x_{11})$, $(y_2, x_{21})$, ..., $(y_n, x_{n1})$ of quantitative variables $(x_{i1}, y_i \in \mathbb{R})$, such as

- $x_{i1}$ = investments        and   $y_i$ = the resulting revenues

- $x_{i1}$ = particular times     and   $y_i$ = the price of a stock at the given time

- $x_{i1}$ = the quantity of some goods supplied to a market

                                     and   $y_i$ = the resulting unit price for the goods

- etc.

# Simple Linear Regression:  Motivation

Given the $n$ pairs $(y_1, x_{11})$, $(y_2, x_{21})$, ..., $(y_n, x_{n1})$ of the measurements, we assume that there is a <u>simple linear relationship</u> between the values of $X_1$ and $Y$ of the form

$$Y \approx \beta_0 + \beta_1 X_1 \qquad \text{for some} \quad \beta_0, \beta_1 \in \mathbb{R}$$

or rather

$$Y = \beta_0 + \beta_1 X_1 + \varepsilon \qquad \text{for some} \quad \beta_0, \beta_1 \in \mathbb{R}$$

where $\varepsilon$ is a random deviation.

We do not know the parameters $\beta_0$ and $\beta_1$, however...

Based on the $n$ pairs $(y_1, x_{11})$, $(y_2, x_{21})$, ..., $(y_n, x_{n1})$ of the measurements, it is our purpose to find

of
$$\text{the estimates} \qquad b_0 \qquad \text{and} \qquad b_1$$
$$\text{the unknown} \qquad \beta_0 \qquad \text{and} \qquad \beta_1$$

The estimates $b_0$ and $b_1$ are also denoted by $\hat{\beta}_0$ and $\hat{\beta}_1$, respectively, sometimes, i.e. the estimates are

$$b_0 = \hat{\beta}_0 \qquad \text{and} \qquad b_1 = \hat{\beta}_1$$
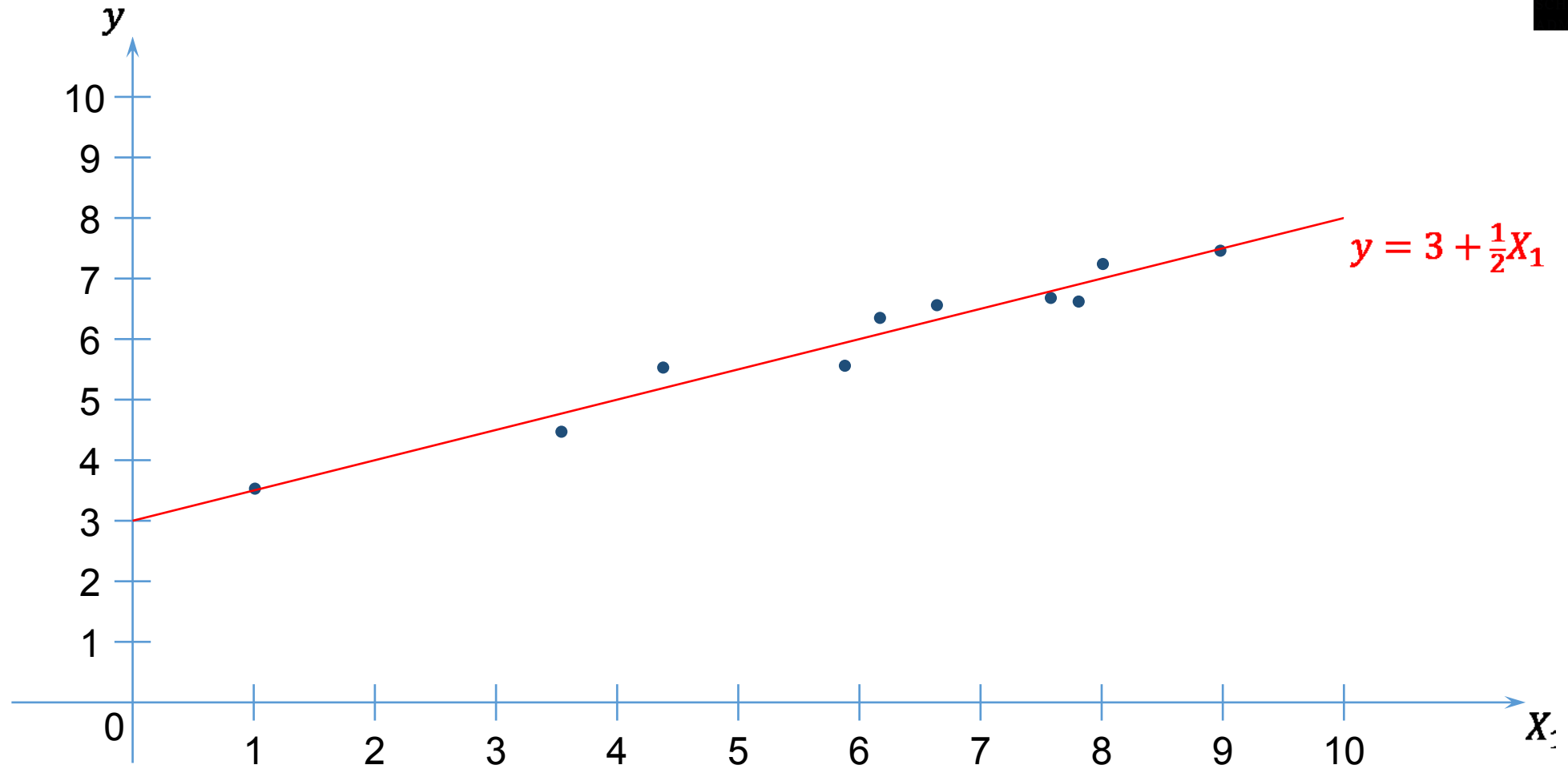
# Simple Linear Regression: Example

We have got a sample of $n = 10$ observations:

| $i$ | $x_{i1}$ | $y_i$ |
|---|---|---|
| ☐ 1 | 8.01 | 7.24 |
| ☐ 2 | 7.81 | 6.62 |
| ☐ 3 | 4.38 | 5.53 |
| ☐ 4 | 3.54 | 4.47 |
| ☐ 5 | 6.17 | 6.35 |
| ☐ 6 | 6.64 | 6.56 |
| ☐ 7 | 7.58 | 6.68 |
| ☐ 8 | 8.98 | 7.46 |
| ☐ 9 | 1.01 | 3.53 |
| 10 | 5.88 | 5.56 |

E.g.:   $x_{i1}$ = temperature  &  $y_i$ = the length of a metal rod

# Simple Linear Regression:  Example



$$y = 3 + \tfrac{1}{2}X_1$$

# Simple Linear Regression: Least Squares Method

We have got the $n$ pairs $(y_1, x_{11})$, $(y_2, x_{21})$, ..., $(y_n, x_{n1})$ of the observations.

For any $b_0, b_1 \in \mathbb{R}$, **the $i$-th estimated value is**

$$\hat{y}_i = b_0 + b_1 x_{i1} \qquad \text{for} \quad i = 1, 2, \ldots, n$$

**The $i$-th residual is the difference**

$$\hat{\varepsilon}_i = e_i = y_i - \hat{y}_i \qquad \text{for} \quad i = 1, 2, \ldots, n$$

The **residual sum of squares** is

$$\text{RSS} = \sum_{i=1}^{n} e_i^2 = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 = \sum_{i=1}^{n} (y_i - b_0 - b_1 x_{i1})^2$$

Given the $n$ pairs $(y_1, x_{11})$, $(y_2, x_{21})$, ..., $(y_n, x_{n1})$ of the observations,

find $b_0, b_1 \in \mathbb{R}$ so that the residual sum of squares

$$\text{RSS} = \sum_{i=1}^{n} (b_0 + b_1 x_{i1} - y_i)^2 \quad \longrightarrow \quad \min$$

is minimized.

The first-order optimality conditions are

$$\frac{\partial \text{RSS}}{\partial b_0} = 0 \qquad \text{and} \qquad \frac{\partial \text{RSS}}{\partial b_1} = 0$$

# Simple Linear Regression: Least Squares Method

Given $\text{RSS} = \sum_{i=1}^{n}(b_0 + b_1 x_{i1} - y_i)^2$, we obtain the system of two equations of two unknowns:

$$\frac{\partial \text{RSS}}{\partial b_0} = \sum_{i=1}^{n} 2(b_0 + b_1 x_{i1} - y_i) = 0 \quad \text{and} \quad \frac{\partial \text{RSS}}{\partial b_1} = \sum_{i=1}^{n} 2(b_0 + b_1 x_{i1} - y_i)x_i = 0$$

or

$$n\, b_0 + \sum_{i=1}^{n} x_{i1}\, b_1 = \sum_{i=1}^{n} y_i$$

$$\sum_{i=1}^{n} x_{i1}\, b_0 + \sum_{i=1}^{n} x_{i1}^2\, b_1 = \sum_{i=1}^{n} x_{i1} y_i$$

the normal equation

# Simple Linear Regression: Least Squares Method

Hence,
given the observations $(y_1, x_{11})$, $(y_2, x_{21})$, ..., $(y_n, x_{n1})$, the estimates are:

$$\hat{\beta}_0 = b_0 = \frac{1}{n}\left(\sum_{i=1}^{n} y_i - \sum_{i=1}^{n} x_{i1} b_1\right) =$$

$$= \frac{\sum_{i=1}^{n} x_{i1}x_{i1} \sum_{j=1}^{n} y_j - \sum_{i=1}^{n} x_{i1} \sum_{j=1}^{n} x_{j1}y_j}{n \sum_{i=1}^{n} x_{i1}x_{i1} - \sum_{i=1}^{n} x_{i1} \sum_{j=1}^{n} x_{j1}}$$

and

$$\hat{\beta}_1 = b_1 = \frac{n \sum_{i=1}^{n} x_{i1}y_i - \sum_{i=1}^{n} x_{i1} \sum_{j=1}^{n} y_j}{n \sum_{i=1}^{n} x_{i1}x_{i1} - \sum_{i=1}^{n} x_{i1} \sum_{j=1}^{n} x_{j1}}$$

Given the $n$ pairs $(y_1, x_{11})$, $(y_2, x_{21})$, ..., $(y_n, x_{n1})$ of the measurements, we have assumed the <u>simple linear relationship</u> of the form

$$Y \approx \beta_0 + \beta_1 X_1 \qquad \text{for some} \quad \beta_0, \beta_1 \in \mathbb{R}$$

The simple linear relationship can be <u>generalized to the form</u>

$$Y \approx \beta_0 X_0 + \beta_1 X_1 \qquad \text{with} \quad X_0 = 1$$

$$\text{for some} \quad \beta_0, \beta_1 \in \mathbb{R}$$

In general, we can have any $n$ triples $(y_1, x_{10}, x_{11})$, $(y_2, x_{20}, x_{21})$, ..., $(y_n, x_{n0}, x_{n1})$

We shall now study

Multiple Linear Regression

That is, we are given a dataset $(y_i, x_{i0}, x_{i1}, x_{i2}, \ldots, x_{ik})_{i=1}^{n}$ of $n$ statistical units $((k+2)$-tuples):

$$( \; y_1, \quad x_{10}, \; x_{11}, \; x_{12}, \; \ldots, \; x_{1k} \; )$$

$$( \; y_2, \quad x_{20}, \; x_{21}, \; x_{22}, \; \ldots, \; x_{2k} \; )$$

$$\ldots$$

$$( \; y_n, \quad x_{n0}, \; x_{n1}, \; x_{n2}, \; \ldots, \; x_{nk} \; )$$

where $\quad y_i, x_{i0}, x_{i1}, x_{i2}, \ldots, x_{ik} \in \mathbb{R} \quad$ for every $\quad i = 1, 2, \ldots, n.$

Given the $n$ $(k+2)$-tuples $(y_i, x_{i0}, x_{i1}, \ldots, x_{ik})_{i=1}^{n}$, such as measurements, we assume the <u>multiple linear relationship</u> between the values of $X_0, X_1, \ldots, X_k$ and $Y$ of the form

$$Y \approx \beta_0 X_0 + \beta_1 X_1 + \cdots + \beta_k X_k \qquad \text{for some} \quad \beta_0, \beta_1, \ldots, \beta_k \in \mathbb{R}$$

or rather

$$Y = \beta_0 X_0 + \beta_1 X_1 + \cdots + \beta_k X_k + \varepsilon \qquad \text{for some} \quad \beta_0, \beta_1, \ldots, \beta_k \in \mathbb{R}$$

where $\varepsilon$ is a random deviation.

We do not know the parameters $\beta_0, \beta_1, \ldots, \beta_k$, however...

We have the dataset of the $n$ $(k+2)$-tuples $(y_i, x_{i0}, x_{i1}, \ldots, x_{ik})_{i=1}^n$.

The values $(y_i)_{i=1}^n$ constitute an $n$-component column vector $y$, which is an $n \times 1$ matrix, and the $(k+1)$-tuples $(x_{i0}, x_{i1}, \ldots, x_{ik})_{i=1}^n$ constitute an $n \times (1+k)$ matrix $X$:

$$y = (y_i)_{i=1}^n = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} \qquad X = (x_{i0}, x_{i1}, \ldots, x_{ik})_{i=1}^n = \begin{pmatrix} x_{10} & x_{11} & \cdots & x_{1k} \\ x_{20} & x_{21} & \cdots & x_{2k} \\ \vdots & \vdots & \vdots & \vdots \\ x_{n0} & x_{n1} & \cdots & x_{nk} \end{pmatrix}$$

We often have $x_{i0} = 1$ for every $i = 1, 2, \ldots, n$,

# Multiple Linear Regression: Notation

Assuming the <u>multiple linear relationship</u> between the values of $X_0, X_1, \ldots, X_k$ and $Y$ of the form

$$Y \approx X_0 \beta_0 + X_1 \beta_1 + \cdots + X_k \beta_k$$

<u>the unknown parameters</u> $\beta_0, \beta_1, \ldots, \beta_k \in \mathbb{R}$ constitute a $(k+1)$-component column vector $\boldsymbol{\beta}$, which is a $(k+1) \times 1$ matrix:

$$\boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_k \end{pmatrix}$$

All in all, we have the $n$ equations

$$y_1 = x_{10}\beta_0 + x_{11}\beta_1 + x_{12}\beta_2 + \cdots + x_{1k}\beta_k + \varepsilon_1$$

$$y_2 = x_{20}\beta_0 + x_{21}\beta_1 + x_{22}\beta_2 + \cdots + x_{2k}\beta_k + \varepsilon_2$$

$$\vdots$$

$$y_n = x_{n0}\beta_0 + x_{n1}\beta_1 + x_{n2}\beta_2 + \cdots + x_{nk}\beta_k + \varepsilon_n$$

where

- the dataset $(y_i, x_{i0}, x_{i1}, \ldots, x_{ik})_{i=1}^n$ is given,

- the parameters $\beta_0, \beta_1, \ldots, \beta_k \in \mathbb{R}$ are unknown (to be estimated), and

- the values $\varepsilon_1, \varepsilon_2, \ldots, \varepsilon_n \in \mathbb{R}$ are random deviations (random errors).

The (unknown) <u>random deviations</u> $\varepsilon_1, \varepsilon_2, \ldots, \varepsilon_n \in \mathbb{R}$ constitute an $n$-component column vector $\boldsymbol{\varepsilon}$, which is an $n \times 1$ matrix:

$$\boldsymbol{\varepsilon} = (\varepsilon_l)_{l=1}^n = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

Moreover, the values $(x_{i0}, x_{i1}, \ldots, x_{ik})$ are seen as a $(1+k)$-component row vector $x_i$, which is a $1 \times (1+k)$ matrix:

$$x_l = \begin{pmatrix} x_{i0} & x_{i1} & \ldots & x_{ik} \end{pmatrix} \qquad \text{for} \quad i = 1, 2, \ldots, n$$

To sum up, assuming $n \geq 2$ and $k \geq 0$, we have:

$$y = (y_i)_{i=1}^n = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} \qquad X = (x_{i0}, x_{i1}, \ldots, x_{ik})_{i=1}^n = \begin{pmatrix} x_{10} & x_{11} & \ldots & x_{1k} \\ x_{20} & x_{21} & \ldots & x_{2k} \\ \vdots & \vdots & \vdots & \vdots \\ x_{n0} & x_{n1} & \ldots & x_{nk} \end{pmatrix} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$$

$$\beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_k \end{pmatrix} \qquad \varepsilon = (\varepsilon_i)_{i=1}^n = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

# Multiple Linear Regression:  Notation

The $n$ equations

$$y_1 = x_{10}\beta_0 + x_{11}\beta_1 + x_{12}\beta_2 + \cdots + x_{1k}\beta_k + \varepsilon_1 = x_1\boldsymbol{\beta} + \varepsilon_1$$

$$y_2 = x_{20}\beta_0 + x_{21}\beta_1 + x_{22}\beta_2 + \cdots + x_{2k}\beta_k + \varepsilon_2 = x_2\boldsymbol{\beta} + \varepsilon_2$$

$$\vdots$$

$$y_n = x_{n0}\beta_0 + x_{n1}\beta_1 + x_{n2}\beta_2 + \cdots + x_{nk}\beta_k + \varepsilon_n = x_n\boldsymbol{\beta} + \varepsilon_n$$

can then be written briefly as

$$\boldsymbol{y} = \boldsymbol{X\beta} + \boldsymbol{\varepsilon}$$

# Random vectors

- Random variable

- Random vector

- Mean value

- Variance-covariance matrix

- Uncorrelated random variables

- Independent random variables

Let $(\Omega, \mathcal{F}, P)$ be a **probability space**. That is,

- $\Omega$ — the <u>sample space</u> (a non-empty set),

- $\mathcal{F}$ — the <u>event space</u> (a σ-algebra on the sample space $\Omega$)

- $P$ — the <u>probability</u> measure on $(\Omega, \mathcal{F})$.

Recall that a **random variable** is a function

$$X : \Omega \longrightarrow \mathbb{R}$$

which is measurable, i.e. the preimage of any open interval is an event $(X^{-1}((a,b)) = \{\omega \in \Omega : X(\omega) \in (a,b)\} \in \mathcal{F}$ for every $a, b \in \mathbb{R}$ such that $a < b$).

# Random vector

Let $(\Omega, \mathcal{F}, P)$ be the probability space as above, and let $n$ random variables $X_1, X_2, \ldots, X_n$ be given. We can then stack the random variables into an $n$-dimensional **random vector**

$$X = \begin{pmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{pmatrix}$$

which is a (measurable) mapping

$$X: \Omega \to \mathbb{R}^n$$

<u>Remark:</u> The fact that the mapping $X: \Omega \to \mathbb{R}^n$ is measurable means that the preimage of any open set is an event:

$$X^{-1}(G) = \{\omega \in \Omega : X(\omega) \in G\} \in \mathcal{F} \qquad \text{for every} \quad \text{open} \ G \subseteq \mathbb{R}^n$$

**We assume for simplicity** that $\Omega = \mathbb{R}^n$ and that the mapping is the identity:

$$X: \omega \mapsto \omega$$

(the event space $\mathcal{F}$ is the collection of all Lebesgue measurable subsets of $\mathbb{R}^n$)

<u>Remark:</u> The mapping $X: \Omega \to \mathbb{R}^n$ is measurable if and only if

# Random vector:  Expected value

Given the random variables $X_1, X_2, \ldots, X_n$,

the **expected value** of the random vector $X$ is:

$$E[X] = \begin{pmatrix} E[X_1] \\ E[X_2] \\ \vdots \\ E[X_n] \end{pmatrix}$$

Given the probability space $(\Omega, \mathcal{F}, P)$, let $X: \Omega \to \mathbb{R}$ and $Y: \Omega \to \mathbb{R}$ be some two random variables.

The **variance** of the random variable $X$ is:

$$\mathrm{Var}(X) = \mathrm{E}[(X - \mathrm{E}[X])^2]$$

The **covariance** of the random variables $X$ and $Y$ is:

$$\mathrm{cov}(X, Y) = \mathrm{E}[(X - \mathrm{E}[X])(Y - \mathrm{E}[Y])]$$

Observation:

Given the random variables $X_1, X_2, \ldots, X_n$,

the **variance-covariance matrix** of the random vector $X$ is:

$$Var(X) = \begin{pmatrix} Var(X_1) & cov(X_1, X_2) & cov(X_1, X_3) & \ldots & cov(X_1, X_n) \\ cov(X_2, X_1) & Var(X_2) & cov(X_2, X_3) & \ldots & cov(X_2, X_n) \\ cov(X_3, X_1) & cov(X_3, X_2) & Var(X_3) & \ldots & cov(X_3, X_n) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ cov(X_n, X_1) & cov(X_n, X_2) & cov(X_n, X_3) & \ldots & Var(X_n) \end{pmatrix}$$

where

$$cov(X_i, X_j) = E\big[(X_i - E[X_i])(X_j - E[X_j])\big] \quad \text{and} \quad Var(X_i) = cov(X_i, X_i)$$

# Random vector:  Uncorrelated random variables

The random variables $X_1, X_2, \ldots, X_n$ are (pairwise) **uncorrelated** if and only if

$$\text{cov}(X_i, X_j) = 0 \quad \text{if} \quad i \neq j \quad \text{for all} \quad i, j = 1, 2, \ldots, n$$

Let $(\Omega, \mathcal{F}, P)$ be a probability space and let $A_1, A_2, \ldots, A_n \in \mathcal{F}$ be events.

Recall that the events $A_1, A_2, \ldots, A_n$ are **mutually independent** if and only if

$$P(A_1 \cap A_2 \cap \cdots \cap A_n) = P(A_1) \times P(A_2) \times \cdots \times P(A_n)$$

Let $(\Omega, \mathcal{F}, P)$ be the underlying probability space and let $X_1, X_2, \ldots, X_n : \Omega \to \mathbb{R}$

be random variables.

The random variables $X_1, X_2, \ldots, X_n$ are **mutually independent** if and only if

$$P\begin{pmatrix} \{\omega \in \Omega : a_1 < X_1(\omega) < b_1\} \cap \\ \cap \{\omega \in \Omega : a_2 < X_2(\omega) < b_2\} \cap \\ \ldots \\ \cap \{\omega \in \Omega : a_n < X_n(\omega) < b_n\} \end{pmatrix} = \begin{matrix} P\{\omega \in \Omega : a_1 < X_1(\omega) < b_1\} \times \\ \times P\{\omega \in \Omega : a_2 < X_2(\omega) < b_2\} \times \\ \ldots \\ \times P\{\omega \in \Omega : a_n < X_n(\omega) < b_n\} \end{matrix}$$

for every $a_1, b_1, a_2, b_2, \ldots, a_n, b_n \in \mathbb{R} \cup \{-\infty, +\infty\}$ such that $a_i < b_i$

**Theorem:** If the random variables

$X_1, X_2, \ldots, X_n$ are **mutually independent**, then they are **pairwise uncorrelated**.

**Remark:** ¡ The converse does not hold true in general !

**Remark:** The proof of the theorem is easy if the sample space is finite

$(\Omega = \{1, 2, \ldots, N\})$ or countable $(\Omega = \{1, 2, 3, \ldots\})$. The proof is somewhat involved in the general case (requires some knowledge of the theory of the Lebesgue integral, uses limiting steps – Levi's Theorem).

# Multivariate normal distribution

Consider a probability space $(\Omega, \mathcal{F}, P)$ where the sample space $\Omega = \mathbb{R}$, the

event space $\mathcal{F}$ is the collection of all Lebesgue measurable subsets of $\mathbb{R}$,

and the probability $P$ is given by its probability density function

$$\varphi(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \, e^{-\frac{(x-\mu)^2}{2\sigma^2}} \qquad \text{for} \quad x \in \mathbb{R}$$

for some $\mu \in \mathbb{R}$ and $\sigma^2 \in \mathbb{R}^+$, so that $\sigma^2 > 0$.

That is, the probability is

$$P(A) = \int_A \varphi(x) \, dx \qquad \text{for any} \quad A \in \mathcal{F}$$

Given the above probability space $(\Omega, \mathcal{F}, P)$, the probability $P$ being given by its density

$$\varphi(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{\frac{(x-\mu)^2}{2\sigma^2}} \qquad \text{for} \quad x \in \mathbb{R}$$

then the identity random variable $X \colon \mathbb{R} \to \mathbb{R}$

$$X(x) = x \qquad \text{for} \quad x \in \mathbb{R}$$

follows the Gaussian normal distribution.

We then say that $X$ is a **Gaussian normal random variable** and write

$$X \sim \mathcal{N}(\mu, \sigma^2)$$

**Theorem:** Let $(\Omega, \mathcal{F}, P)$ be a probability space. (Consider $\Omega = \mathbb{R}^n$ for simplicity.)

If the random variables $X_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$, $X_2 \sim \mathcal{N}(\mu_2, \sigma_2^2)$, ..., $X_n \sim \mathcal{N}(\mu_n, \sigma_n^2)$

are **mutually independent** and normally distributed, then

$$X_1 + X_2 + \cdots + X_n \sim \mathcal{N}(\mu_1 + \mu_2 + \cdots + \mu_n, \ \sigma_1^2 + \sigma_2^2 + \cdots + \sigma_n^2)$$

that is, their sum is also normally distributed.

# Variance-covariance matrix

**Theorem:** Let $(\Omega, \mathcal{F}, P)$ be a probability space (with $\Omega = \mathbb{R}^n$ for simplicity) and let $X_1, X_2, \ldots, X_n : \Omega \to \mathbb{R}$ be any random variables, which are stacked into a random vector $X$. Then its <u>variance-covariance matrix</u>

$$\Sigma = \text{Var}(X)$$

is <u>symmetric</u> and <u>positively semi-definite</u>.

That is, it holds

$$\Sigma^{\mathrm{T}} = \Sigma \qquad \text{and} \qquad u^{\mathrm{T}} \Sigma u \geq 0 \qquad \text{for every} \quad u \in \mathbb{R}^n$$

**Theorem:**

Let $\Sigma \in \mathbb{R}^{n \times n}$ be any underlined symmetric and positively semi-definite matrix, and let

$$k = \operatorname{rank}(\Sigma)$$

Then there exists a matrix $A \in \mathbb{R}^{n \times k}$ such that

$$\Sigma = AA^{\mathrm{T}}$$

**Remark:** The matrix $A$ can be obtained • either from the spectral decomposition / eigendecomposition of the matrix $\Sigma$: $\Sigma = Q\Lambda Q^{\mathrm{T}}$ where $\Lambda$ is diagonal and $Q$ is orthonormal ($QQ^{\mathrm{T}} = I$); • or from the Cholesky decomposition: $\Sigma = LL^{\mathrm{T}}$ where

Consider a probability space $(\Omega, \mathcal{F}, P)$ where the sample space $\Omega = \mathbb{R}^k$, the event space $\mathcal{F}$ is the collection of all Lebesgue measurable subsets of $\mathbb{R}^k$, and the probability $P$ is given by the standardized normal density function

$$\varphi_k(x) = \frac{1}{\sqrt{(2\pi)^k}} e^{-\frac{x^T x}{2}} \qquad \text{for} \quad x \in \mathbb{R}^k$$

That is, the probability is

$$P(A) = \int_A \varphi_k(x)\, dx \qquad \text{for any} \quad A \in \mathcal{F}$$

Given the above probability space $(\Omega, \mathcal{F}, P)$, the probability $P$ being given by its density

$$\varphi_k(x) = \frac{1}{\sqrt{(2\pi)^k}} e^{-\frac{x^T x}{2}} \qquad \text{for} \quad x \in \mathbb{R}^k$$

then the identity random vector $Z: \mathbb{R}^k \to \mathbb{R}^k$

$$Z(x) = x \qquad \text{for} \quad x \in \mathbb{R}^k$$

follows the standard Gaussian multivariate normal distribution.

We then say that $Z$ is a **standard multivariate normal random vector** and write

$$Z \sim \mathcal{N}(0, I)$$

# Multivariate normal distribution

Let a vector $\mu \in \mathbb{R}^n$ (mean values) and a symmetric positively semi-definite matrix $\Sigma \in \mathbb{R}^{n \times n}$ (variance-covariance matrix) of rank $k$ be given. Moreover, let $A \in \mathbb{R}^{n \times k}$ be a matrix such that $\Sigma = AA^T$. Finally, consider the probability space $(\Omega, \mathcal{F}, P)$ with the sample space $\Omega = \mathbb{R}^k$ and the standard multivariate normal random variable $Z \sim \mathcal{N}(0, I)$.

Then the random vector $X: \mathbb{R}^k \to \mathbb{R}^n$ defined so that

$$X(x) = AZ(x) \qquad \text{for} \quad x \in \mathbb{R}^k$$

follows the standard Gaussian multivariate normal distribution.

We then say that $X$ is a **multivariate normal random vector** and write

# Multivariate normal distribution: Density

If the variance-covariance matrix $\Sigma \in \mathbb{R}^{n \times n}$ is underline{non-singular},

that is $\operatorname{rank}(\Sigma) = k = n$, then the **probability density function**

of the multivariate normal probability distribution

$$\mathcal{N}(\mu, \Sigma)$$

is

$$f(x) = \frac{1}{\sqrt{(2\pi)^k \det(\Sigma)}} e^{\frac{(x-\mu)^\mathsf{T} \Sigma^{-1} (x-\mu)}{2}} \qquad \text{for} \quad x \in \mathbb{R}^n$$

If the variance-covariance matrix $\Sigma \in \mathbb{R}^{n \times n}$ is underline{singular}, that is $\operatorname{rank}(\Sigma) = k < n$,

then the **probability density function** of the multivariate normal probability

# Multivariate normal distribution: Another definition

Let $(\Omega, \mathcal{F}, P)$ be a probability space and let $X: \Omega \to \mathbb{R}^n$ be a random vector.

Then $X$ <u>follows a multivariate normal distribution</u>, that is $X \sim \mathcal{N}(\mu, \Sigma)$

for some $\mu \in \mathbb{R}^n$ and for some symmetric and positively semi-definite $\Sigma \in \mathbb{R}^{n \times n}$

if and only if,

for every $a \in \mathbb{R}^n$, the random variable $a^{\mathsf{T}} X$ is normally distributed;

that is, there exist a $\mu_a \in \mathbb{R}$ and a non-negative $\sigma_a^2 \in \mathbb{R}$ such that

$$a_1 X_1 + a_2 X_2 + \cdots + a_n X_n \sim \mathcal{N}(\mu_a, \sigma_a^2) \qquad \text{for every} \quad a \in \mathbb{R}^n$$

# Multivariate normal distribution: Linear transformation

**Theorem:** Let $(\Omega, \mathcal{F}, P)$ be a probability space and let $X : \Omega \to \mathbb{R}^n$ be a <u>multivariate normally distributed random vector</u>, that is $X \sim \mathcal{N}(\mu, \Sigma)$ for some $\mu \in \mathbb{R}^n$ and for some symmetric and positively semi-definite $\Sigma \in \mathbb{R}^{n \times n}$.

Then

$$AX \sim \mathcal{N}\left(A\mu, A\Sigma A^{\mathrm{T}}\right) \qquad \text{for any matrix} \quad A \in \mathbb{R}^{m \times n}$$

**Theorem:** Let random variables $X_1, X_2, \ldots, X_n$ be stacked into a multivariate

normally distributed random vector $X$, that is $X \sim \mathcal{N}(\mu, \Sigma)$ for some $\mu \in \mathbb{R}^n$

and for some symmetric and positively semi-definite $\Sigma \in \mathbb{R}^{n \times n}$.  Then

the random variables $X_1, X_2, \ldots, X_n$ are **mutually independent**

<u>if and only if</u>

the random variables $X_1, X_2, \ldots, X_n$ are **pairwise uncorrelated**

(that is, the variance-covariance matrix $\Sigma$ is diagonal).

**Remark:**

$\Longrightarrow$ holds true in general, see above

# Multiple Linear Regression: Summary & Background

- Summary

- Terminology

- Assumptions

- Random vectors

- The classical assumptions

- Notation

We have got the sample of the $n$ $(k+2)$-tuples

$$( y_i,\ x_i ) = ( y_i,\ x_{i0},\ x_{i1},\ x_{i2},\ \dots,\ x_{ik} ) \qquad \text{for} \quad i = 1, 2, \dots, n$$

of the observations, where $y_i \in \mathbb{R}$ and $x_i \in \mathbb{R}^{1 \times (1+k)}$ or $x_{i0}, x_{i1}, \dots, x_{ik} \in \mathbb{R}$ for $i = 1, 2, \dots, n$.

The sample could have been obtained in either of the following two ways:

(see the next two slides)

**First:**

— A sample of $n$ statistical units was selected from a larger population.

— Each of the statistical units was measured and we have obtained

the pairs $(y_i, x_i)$ for $i = 1, 2, \ldots, n$ thus.

— ¡¡¡ The values $x_i \in \mathbb{R}^{1 \times (1+k)}$ were measured / are known exactly !!!

(That is, the values $x_i$ are non-random.)

— We assume $y_i \approx x_i \beta$ and we have $y_i = x_i \beta + \varepsilon_i$,

where $\varepsilon_i$ is a random deviation (error).

— The random deviation is caused by the intrinsic properties of the statistical unit

**Second:**

— We prepared the values $x_1, x_2, \ldots, x_n \in \mathbb{R}^{1\times(1+k)}$ at the beginning.

— ¡¡¡ These values $x_1, x_2, \ldots, x_n$ are known exactly therefore !!!

— When making the $i$-th measurement,

we set up the system (adjust the system's setting to $x_i$ <u>exactly</u>) first and

we measure the value $y_i$ of the dependent variable then.

— The random deviation $\varepsilon_i$ here is caused

<u>either</u> by the intrinsic properties of the system (further unknown / "random" /

unconsidered factors),

# Multiple Linear Regression:  Summary

<u>Remarks:</u>

- In practice, the data may be obtained in either way (first or second).

- In either case (first or second), the independent values $x_1, x_2, ..., x_n$

  are assumed to be known exactly, i.e. without any measurement errors.

- Assuming $y_i \approx x_i\boldsymbol{\beta}$, even the dependent values $y_i$ may be measured exactly,

  i.e. without any measurement error, the random deviation $\varepsilon_i = y_i - x_i\boldsymbol{\beta}$ being

  caused by the intrinsic properties (other unknown / "random" / unconsidered

  factors).

- For the purpose of the mathematical analysis, we assume the second case only.

# Multiple Linear Regression:  Terminology

$$Y = X_0\beta_0 + X_1\beta_1 + \cdots + X_k\beta_k + \varepsilon$$

**Parameters**

Regression coefficients

**Regressand**

Predicand

Explained variable

Dependent variable

Endogenous variable

Controlled variable

Response

Outcome

Predicted variable

Measured variable

**Regressors**

Predictors

Explanatory variables

Independent variables

Exogenous variables

Control variables

Stimuli

Covariates

**Deviation**

Error term

Disturbance

Noise

# Multiple Linear Regression:  Terminology

If $X_0 = 1$:

The **intercept term**

**Parameters**

Regression coefficients

$$Y = \beta_0 + X_1\beta_1 + \cdots + X_k\beta_k + \varepsilon$$

| **Regressand** | **Regressors** | **Deviation** |
|---|---|---|
| Predicand | Predictors | Error term |
| Explained variable | Explanatory variables | Disturbance |
| Dependent variable | Independent variables | Noise |
| Endogenous variable | Exogenous variables | |
| Controlled variable | Control variables | |
| Response | Stimuli | |
| Outcome | Covariates | |
| Predicted variable | | |
| Measured variable | | |

- The $n$ row vectors $x_1, x_2, \ldots, x_n \in \mathbb{R}^{1 \times (1+k)}$ are known exactly, fixed, given before the measurements.

- We have $n$ random variables $Y_1, Y_2, \ldots, Y_n$

  and $n$ random variables $\varepsilon_1, \varepsilon_2, \ldots, \varepsilon_n$.

- We assume that the random variables $Y_1, Y_2, \ldots, Y_n$ are independent

  and the random variables $\varepsilon_1, \varepsilon_2, \ldots, \varepsilon_n$ are independent.

- <u>Remark:</u>

  It is enough to assume that the random variables $Y_1, Y_2, \ldots, Y_n$ are uncorrelated

  and the random variables $\varepsilon_1, \varepsilon_2, \ldots, \varepsilon_n$ are uncorrelated.

# Multiple Linear Regression: Assumptions

- Let $(\Omega, \mathcal{F}, P)$ be the underlying probability space.

- We stack the random variables $Y_1, Y_2, \ldots, Y_n$ and $\varepsilon_1, \varepsilon_2, \ldots, \varepsilon_n$ into random vectors:

$$Y = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix} \quad \text{and} \quad \varepsilon = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

which are (measurable) mappings

$$Y : \Omega \to \mathbb{R}^n \quad \text{and} \quad \varepsilon : \Omega \to \mathbb{R}^n$$

# Multiple Linear Regression: Random vectors

- We assume for simplicity that $\Omega = \mathbb{R}^n$ and that the mappings are identities:

$$Y: \omega \mapsto \omega \qquad \text{and} \qquad \varepsilon: \omega \mapsto \omega$$

(the event space $\mathcal{F}$ is the collection of all Lebesgue measurable subsets of $\mathbb{R}^n$)

- The **expected values** of the random vectors $Y$ and $\varepsilon$ are:

$$E[Y] = \begin{pmatrix} E[Y_1] \\ E[Y_2] \\ \vdots \\ E[Y_n] \end{pmatrix} \qquad \text{and} \qquad E[\varepsilon] = \begin{pmatrix} E[\varepsilon_1] \\ E[\varepsilon_2] \\ \vdots \\ E[\varepsilon_n] \end{pmatrix}$$

このセグメントはheaderナビゲーションではなくタイトル

# Multiple Linear Regression: Random vectors

- The **variance-covariance matrix** of the random vector $Y$ is:

$$\text{Var}(Y) = \begin{pmatrix} \text{Var}(Y_1) & \text{cov}(Y_1, Y_2) & \text{cov}(Y_1, Y_3) & \dots & \text{cov}(Y_1, Y_n) \\ \text{cov}(Y_2, Y_1) & \text{Var}(Y_2) & \text{cov}(Y_2, Y_3) & \dots & \text{cov}(Y_2, Y_n) \\ \text{cov}(Y_3, Y_1) & \text{cov}(Y_3, Y_2) & \text{Var}(Y_3) & \dots & \text{cov}(Y_3, Y_n) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \text{cov}(Y_n, Y_1) & \text{cov}(Y_n, Y_2) & \text{cov}(Y_n, Y_3) & \dots & \text{Var}(Y_n) \end{pmatrix}$$

Recall that

$$\text{cov}(Y_i, Y_j) = \text{E}\big[(Y_i - \text{E}[Y_i])(Y_j - \text{E}[Y_j])\big] \quad \text{and} \quad \text{Var}(Y_i) = \text{cov}(Y_i, Y_i)$$

- The **variance-covariance matrix** of the random vector $\varepsilon$ is:

$$\text{Var}(\varepsilon) = \begin{pmatrix} \text{Var}(\varepsilon_1) & \text{cov}(\varepsilon_1, \varepsilon_2) & \text{cov}(\varepsilon_1, \varepsilon_3) & \dots & \text{cov}(\varepsilon_1, \varepsilon_n) \\ \text{cov}(\varepsilon_2, \varepsilon_1) & \text{Var}(\varepsilon_2) & \text{cov}(\varepsilon_2, \varepsilon_3) & \dots & \text{cov}(\varepsilon_2, \varepsilon_n) \\ \text{cov}(\varepsilon_3, \varepsilon_1) & \text{cov}(\varepsilon_3, \varepsilon_2) & \text{Var}(\varepsilon_3) & \dots & \text{cov}(\varepsilon_3, \varepsilon_n) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \text{cov}(\varepsilon_n, \varepsilon_1) & \text{cov}(\varepsilon_n, \varepsilon_2) & \text{cov}(\varepsilon_n, \varepsilon_3) & \dots & \text{Var}(\varepsilon_n) \end{pmatrix}$$

Recall that

$$\text{cov}(\varepsilon_i, \varepsilon_j) = \text{E}\big[(\varepsilon_i - \text{E}[\varepsilon_i])(\varepsilon_j - \text{E}[\varepsilon_j])\big] \qquad \text{and} \qquad \text{Var}(\varepsilon_i) = \text{cov}(\varepsilon_i, \varepsilon_i)$$

- We have the underlying probability space $(\Omega, \mathcal{F}, P)$, with $\Omega = \mathbb{R}^n$ for simplicity.

- Let $\omega \in \Omega$ be the outcome of the random experiment.

- Recalling that $X$ is the $n \times (1 + k)$ design matrix, we have

$$y = Y(\omega) = X\beta + \varepsilon(\omega)$$

In other words:

- The measured values $y_1, y_2, \ldots, y_n$ are the numerical outcomes $Y_1(\omega), Y_2(\omega), \ldots, Y_n(\omega)$ of the random experiment.

- The numerical outcomes $Y_1(\omega), Y_2(\omega), \ldots, Y_n(\omega)$ are obtained so that the numerical outcomes $\varepsilon_1(\omega), \varepsilon_2(\omega), \ldots, \varepsilon_n(\omega)$ of the random experiment

# Multiple Linear Regression: The Classical Assumptions

<u>Recall:</u>

¡¡¡ The values of the regressors $x_1, x_2, \ldots, x_n$ are **<u>non-random</u> and <u>known</u>** !!!

¡¡¡ The values of the parameters $\beta_0, \beta_1, \ldots, \beta_k$ are <u>non-random but **unknown**</u> !!!

We assume that

$$Y \sim \mathcal{N}(X\beta, \sigma^2 I) \qquad \text{and} \qquad \varepsilon \sim \mathcal{N}(0, \sigma^2 I) \qquad \text{for some } \sigma^2 \in \mathbb{R}_0^+$$

where $I$ denotes the $n \times n$ identity matrix

and $0$ denotes the $n \times 1$ zero vector.

The classical assumptions $Y \sim \mathcal{N}(X\boldsymbol{\beta}, \sigma^2 I)$ and $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma^2 I)$ mean that

$$E[Y] = X\boldsymbol{\beta} \qquad \text{and} \qquad E[\boldsymbol{\varepsilon}] = 0$$

that is

$$E[Y] = \begin{pmatrix} E[Y_1] \\ E[Y_2] \\ \vdots \\ E[Y_n] \end{pmatrix} = \begin{pmatrix} x_1\boldsymbol{\beta} \\ x_2\boldsymbol{\beta} \\ \vdots \\ x_n\boldsymbol{\beta} \end{pmatrix} \qquad \text{and} \qquad E[\boldsymbol{\varepsilon}] = \begin{pmatrix} E[\varepsilon_1] \\ E[\varepsilon_2] \\ \vdots \\ E[\varepsilon_n] \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

That is

$$E[Y_i] = x_i\boldsymbol{\beta} \qquad \text{and} \qquad E[\varepsilon_i] = 0 \qquad \text{for} \quad i = 1, 2, \dots, n$$

The classical assumptions $Y \sim \mathcal{N}(X\beta, \sigma^2 I)$ and $\varepsilon \sim \mathcal{N}(0, \sigma^2 I)$ also mean that

$$\text{Var}(Y) = \text{Var}(\varepsilon) = \sigma^2 I = \begin{pmatrix} \sigma^2 & 0 & 0 & \cdots & 0 \\ 0 & \sigma^2 & 0 & \cdots & 0 \\ 0 & 0 & \sigma^2 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & \sigma^2 \end{pmatrix}$$

That is,

- $\text{Var}(Y_i) = \text{Var}(\varepsilon_i) = \sigma^2$ for $i = 1, 2, \ldots, n$ for some $\sigma^2 \in \mathbb{R}_0^+$ $\longleftarrow$

- and the random variables $Y_1, Y_2, \ldots, Y_n$ or $\varepsilon_1, \varepsilon_2, \ldots, \varepsilon_n$

are (pairwise) **uncorrelated**.

**homoscedasticity**, i.e. the variance is the same

# Multiple Linear Regression:  The Classical Assumptions

The classical assumptions $Y \sim \mathcal{N}(X\boldsymbol{\beta}, \sigma^2 I)$ and $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma^2 I)$ finally mean that

- $Y_i \sim \mathcal{N}(x_i\boldsymbol{\beta}, \sigma^2)$ and $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$ for $i = 1, 2, \ldots, n$

  where $\mathcal{N}(\mu, \sigma^2)$ denotes the normal (Gaussian) probability distribution

  with mean $\mu$ and variance $\sigma^2$ and that

- the random variables $Y_1, Y_2, \ldots, Y_n$ or $\varepsilon_1, \varepsilon_2, \ldots, \varepsilon_n$ are (pairwise) <u>uncorrelated</u>.

<u>Remark:</u> <u>It always holds:</u> If the random variables $Y_1, Y_2, \ldots, Y_n$ or $\varepsilon_1, \varepsilon_2, \ldots, \varepsilon_n$ are mutually independent, then they are (pairwise) uncorrelated.

<u>It also holds:</u> If $Y \sim \mathcal{N}(X\boldsymbol{\beta}, \sigma^2 I)$ and $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma^2 I)$, then the random variables

The classical assumption $Y \sim \mathcal{N}(X\beta, \sigma^2 I)$ implies that

linearity

$$E[Y] = X\beta$$

that is

$$E[Y_i] = x_i\beta = x_{i0}\beta_0 + x_{i1}\beta_1 + \cdots + x_{ik}\beta_k \qquad \text{for} \quad i = 1, 2, \ldots, n$$

# Multiple Linear Regression:  Notation

- The <u>unknown quantities</u>

    — unknown parameters $\beta_0, \beta_1, \ldots, \beta_k$

    — unknown (random) deviations $\varepsilon_1, \varepsilon_2, \ldots, \varepsilon_n$

  are denoted by <u>Greek letters</u>

- The <u>estimates</u> of the unknown parameters $\beta_0, \beta_1, \ldots, \beta_k$

  are denoted by the respective <u>Latin letters</u> $b_0, b_1, \ldots, b_k$

  or by the <u>hat</u> "^" $\hat{\beta}_0, \hat{\beta}_1, \ldots, \hat{\beta}_k,$

  so that $b_0 = \hat{\beta}_0, \; b_1 = \hat{\beta}_1, \; \ldots, \; b_k = \hat{\beta}_k$

# Multiple Linear Regression:  Notation

- The **predicted values** of the dependent variable are denoted by the  hat  "^":

$$\hat{y}_i = x_i b = x_{i0} b_0 + x_{i1} b_1 + \cdots + x_{ik} b_k \qquad \text{for} \quad i = 1, 2, \ldots, n$$

- The unknown (random) deviations $\varepsilon_1, \varepsilon_2, \ldots, \varepsilon_n$ are denoted by Greek letters. The **residuals** are denoted by the respective <u>Latin letters</u> $e_1, e_2, \ldots, e_n$ or by the  hat  "^"  $\hat{\varepsilon}_1, \hat{\varepsilon}_2, \ldots, \hat{\varepsilon}_n,$ so that

$$e_1 = \hat{\varepsilon}_1 = y_1 - \hat{y}_1 \qquad e_2 = \hat{\varepsilon}_2 = y_2 - \hat{y}_2 \qquad \ldots \qquad e_n = \hat{\varepsilon}_n = y_n - \hat{y}_n$$

# Basic Results (Theorems)

- The normal equation

- The predicted values

- Orthogonal projections

- Theorem 1: $\hat{\boldsymbol{y}}$ and $\boldsymbol{e}$ are independent

- Theorem 2: $\hat{\boldsymbol{y}} \sim \mathcal{N}(\boldsymbol{X\beta}, \sigma^2 \boldsymbol{H})$

- Theorem 3: $\boldsymbol{e} \sim \mathcal{N}(\boldsymbol{0}, \sigma^2 \boldsymbol{M})$

- Theorem 4: $\boldsymbol{b} \sim \mathcal{N}(\boldsymbol{\beta}, \sigma^2 \boldsymbol{C})$   if $\text{rank}(\boldsymbol{X}) = k + 1$

# Multiple Linear Regression:  The Normal Equation

Given the $n$ pairs $(y_1, x_1)$, $(y_2, x_2)$, ..., $(y_n, x_n)$ of the observations,

the **Residual Sum of Squares** for the estimates $b_0, b_1, ..., b_k \in \mathbb{R}$ is

$$\text{RSS} = \sum_{i=1}^{n} e_i^2 = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 = \sum_{i=1}^{n} (y_i - x_{i0}b_0 - x_{i1}b_1 - \cdots - x_{ik}b_k)^2 \longrightarrow \min$$

It is our purpose to find the estimates $b_0, b_1, ..., b_k \in \mathbb{R}$ so that the Residual Sum

of Squares  RSS  is <u>minimized</u>.  To this end, we let

$$\frac{\partial \text{RSS}}{\partial b_j} = \sum_{i=1}^{n} -2x_{ij}(y_i - x_{i0}b_0 - x_{i1}b_1 - \cdots - x_{ik}b_k) = 0 \qquad \text{for} \quad j = 0, 1, ..., k$$

# Multiple Linear Regression: The Normal Equation

From $\partial \text{RSS} / \partial b_j = \sum_{i=1}^{n} -2x_{ij}(y_i - x_{i0}b_0 - x_{i1}b_1 - \cdots - x_{ik}b_k) = 0$, we obtain

the Normal Equation:

$$\sum_{i=1}^{n} x_{ij}(x_{i0}b_0 + x_{i1}b_1 + \cdots + x_{ik}b_k) = \sum_{i=1}^{n} x_{ij}y_i \qquad \text{for} \quad j = 0, 1, \ldots, k$$

Recall the notation:

$$X = (x_{i0}, x_{i1}, \ldots, x_{ik})_{i=1}^{n} = \begin{pmatrix} x_{10} & x_{11} & \cdots & x_{1k} \\ x_{20} & x_{21} & \cdots & x_{2k} \\ \vdots & \vdots & \vdots & \vdots \\ x_{n0} & x_{n1} & \cdots & x_{nk} \end{pmatrix} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$$

and

$$y = (y_i)_{i=1}^{n} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

So the normal equation

$$\sum_{i=1}^{n} x_{ij}(x_{i0}b_0 + x_{i1}b_1 + \cdots + x_{ik}b_k) = \sum_{i=1}^{n} x_{ij}y_i \qquad \text{for} \quad j = 0, 1, \ldots, k$$

$$\sum_{i=1}^{n} x_{ij}(x_i b) = \sum_{i=1}^{n} x_{ij}y_i \qquad \text{for} \quad j = 0, 1, \ldots, k$$

can be written as

$$X^T X b = X^T y$$

Having the normal equation $(X^TXb = X^Ty)$, where $X$ is an $n \times (1+k)$ matrix,

let

$$p = \text{rank}(X)$$

**Assume for simplicity that the matrix $X$ is of full rank, that is,**

$$p = \text{rank}(X) = k + 1 \leq n$$

| **NO** |
| (perfect) |
| **multicollinearity** |

The matrix $X^TX$ is then non-singular; let:

$$C = (X^TX)^{-1}$$

We have

$$C = (X^{\mathrm{T}}X)^{-1} = \begin{pmatrix} c_{00} & c_{01} & c_{02} & \cdots & c_{0k} \\ c_{10} & c_{11} & c_{12} & \cdots & c_{1k} \\ c_{20} & c_{21} & c_{22} & \cdots & c_{2k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ c_{k0} & c_{k1} & c_{k2} & \cdots & c_{kk} \end{pmatrix}$$

The solution to the normal equation

$$X^{\mathrm{T}}Xb = X^{\mathrm{T}}y$$

is then

$$b = CX^{\mathrm{T}}y$$

Recall the $n$ equations

$$y_1 = x_{10}\beta_0 + x_{11}\beta_1 + x_{12}\beta_2 + \cdots + x_{1k}\beta_k + \varepsilon_1$$

$$y_2 = x_{20}\beta_0 + x_{21}\beta_1 + x_{22}\beta_2 + \cdots + x_{2k}\beta_k + \varepsilon_2$$

$$\vdots$$

$$y_n = x_{n0}\beta_0 + x_{n1}\beta_1 + x_{n2}\beta_2 + \cdots + x_{nk}\beta_k + \varepsilon_n$$

where

- the dataset $(y_i, x_{i0}, x_{i1}, \ldots, x_{ik})_{i=1}^n$ is given,

- the parameters $\beta_0, \beta_1, \ldots, \beta_k \in \mathbb{R}$ are unknown (<u>to be estimated</u>), and

- the values $\varepsilon_1, \varepsilon_2, \ldots, \varepsilon_n \in \mathbb{R}$ are <u>random deviations</u> (random errors).

Now the **predicted values** are:

$$\hat{y}_1 = x_{10}b_0 + x_{11}b_1 + x_{12}b_2 + \cdots + x_{1k}b_k$$

$$\hat{y}_2 = x_{20}b_0 + x_{21}b_1 + x_{22}b_2 + \cdots + x_{2k}b_k$$

$$\vdots$$

$$\hat{y}_n = x_{n0}b_0 + x_{n1}b_1 + x_{n2}b_2 + \cdots + x_{nk}b_k$$

where

- $b_0 = \hat{\beta}_0, \quad b_1 = \hat{\beta}_1, \quad \ldots, \quad b_k = \hat{\beta}_k$   are the estimates

  of the unknown parameters  $\beta_0, \beta_1, \ldots, \beta_k \in \mathbb{R}$,

- $\hat{y}_1, \hat{y}_2, \ldots, \hat{y}_n$  are the <u>predicted values</u>.

Shortly:

— the solution to the normal equation is

$$b = CX^Ty = (X^TX)^{-1}X^Ty$$

— the predicted values are

$$\hat{y} = Xb = XCX^Ty$$

**Introduce the notation:**

$$H = XCX^T = X(X^TX)^{-1}X^T$$

The letter "$H$" stands for "hat":

$$\hat{y} = Hy$$

By the construction (by the Least Squares Method: the vector $\hat{y}$ lies in the linear hull of the columns of the matrix $X$ and is as close to $y$ as possible in the Euclidean distance), the matrix

$$H = XCX^{\mathrm{T}} = X(X^{\mathrm{T}}X)^{-1}X^{\mathrm{T}}$$

is the **matrix of the orthogonal projection onto the linear subspace**

$$\{X\boldsymbol{\beta} : \boldsymbol{\beta} \in \mathbb{R}^{1+k}\}$$

(the linear hull of the columns of the matrix $X$ )

moreover

$$(I - H)$$

is the matrix of the orthogonal projection onto the orthogonal complement

The matrix $H = XCX^T = X(X^TX)^{-1}X^T$  therefore is:

— idempotent:

$$H = HH$$

— symmetric:

$$H = H^T$$

— and:

$$HX = X$$

The residuals are:

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

Therefore:

$$\hat{y} \perp e$$

The orthogonal decomposition
of the vector  $y \in \mathbb{R}^n$ :

$$y = \hat{y} + e \qquad \text{and} \qquad e \perp \hat{y}$$

$\{ X\beta : \beta \in \mathbb{R}^{1+k} \}^{\perp}$

(the orthogonal
complement =
= the space of
the residuals)

vector of the numerical outcomes
of the  $n$  random experiments

$(I - H)y = e$

$M$

$y$

By the Pythagoras Theorem:

$$\|y\|^2 = \|\hat{y}\|^2 + \|e\|^2$$

$$y^{\mathrm{T}}y = \hat{y}^{\mathrm{T}}\hat{y} + e^{\mathrm{T}}e$$

$0$

$\hat{y} = Hy$

$\{ X\beta : \beta \in \mathbb{R}^{1+k} \}$

(the linear hull of the columns of  $X$ )

**Residual Sum of Squares:**   $\mathrm{RSS} = \sum_{i=1}^{n}(y_i - \hat{y}_i)^2 = \sum_{i=1}^{n} e_i^2 = e^{\mathrm{T}}e$

# Multiple Linear Regression

Recalling that the regressors $x_1, x_2, \ldots, x_n \in \mathbb{R}^{1 \times (1+k)}$ are given, that we assume

$$Y_i = x_i \beta + \varepsilon_i \qquad \text{for} \quad i = 1, 2, \ldots, n$$

with

$$\mathrm{E}[Y_i] = x_i \beta \qquad \text{and} \qquad \mathrm{Var}(Y_i) = \sigma^2 \qquad \text{for} \quad i = 1, 2, \ldots, n$$

or

$$\mathrm{E}[\varepsilon_i] = 0 \qquad \text{and} \qquad \mathrm{Var}(\varepsilon_i) = \sigma^2 \qquad \text{for} \quad i = 1, 2, \ldots, n$$

where the random variables $Y_1, Y_2, \ldots, Y_n$, or $\varepsilon_1, \varepsilon_2, \ldots, \varepsilon_n$, respectively, are independent (or <u>uncorrelated</u>), and that $y_1, y_2, \ldots, y_n$ are some observations of the random variables $Y_1, Y_2, \ldots, Y_n$, **it follows that all the estimates**

$$b_0, b_1, \ldots, b_k \qquad (\hat{\beta}_0, \hat{\beta}_1, \ldots, \hat{\beta}_k) \qquad \hat{y}_1, \hat{y}_2, \ldots, \hat{y}_n \qquad \text{RSS} \qquad \text{etc.}$$

**Theorem 1:** The <u>random vectors</u> $\hat{y}$ and $e$ are <u>independent</u>.

That is,

$$P\left(\begin{array}{l} \{\omega \in \Omega : \hat{y}(\omega) \in G_{\hat{y}}\} \cap \\ \cap \{\omega \in \Omega : e(\omega) \in G_e\} \end{array}\right) = \begin{array}{l} P\{\omega \in \Omega : \hat{y}(\omega) \in G_{\hat{y}}\} \times \\ \times P\{\omega \in \Omega : e(\omega) \in G_e\} \end{array}$$

for every open set $G_{\hat{y}} \subseteq \mathbb{R}^n$ and for every open set $G_e \subseteq \mathbb{R}^n$.

**Corollary:** The <u>random vector</u> $\hat{y}$ and the <u>random variable</u> RSS are <u>independent</u>.

Recall the **Residual Sum of Squares** is $\text{RSS} = \sum_{i=1}^{n}(y_i - \hat{y}_i)^2 = \sum_{i=1}^{n} e_i^2 = e^T e$

That is,

$$P\left(\begin{array}{c} \{\omega \in \Omega : \hat{y}(\omega) \in G_{\hat{y}}\} \cap \\ \cap \{\omega \in \Omega : \text{RSS}(\omega) \in G_{\text{RSS}}\} \end{array}\right) = \begin{array}{c} P\{\omega \in \Omega : \hat{y}(\omega) \in G_{\hat{y}}\} \times \\ \times P\{\omega \in \Omega : \text{RSS}(\omega) \in G_{\text{RSS}}\} \end{array}$$

for every open set $G_{\hat{y}} \subseteq \mathbb{R}^n$ and for every open set $G_{\text{RSS}} \subseteq \mathbb{R}$

**Theorem 2:** It holds

$$\hat{y} \sim \mathcal{N}(X\beta, \sigma^2 H)$$

It holds in particular hence that

$$\mathrm{E}[\hat{y}_i] = \mathrm{E}[x_i b] = x_i \beta = \mathrm{E}[Y_i] \qquad \text{for} \quad i = 1, 2, \dots, n$$

and

$$\mathrm{Var}(\hat{y}_i) = \sigma^2 (XCX^\mathrm{T})_{ii} \qquad \text{for} \quad i = 1, 2, \dots, n$$

$$\mathrm{cov}(\hat{y}_i, \hat{y}_j) = \sigma^2 (XCX^\mathrm{T})_{ij} \qquad \text{for} \quad i, j = 1, 2, \dots, n$$

# Multiple Linear Regression: Theorem 3

**Theorem 3:** It holds

$$e \sim \mathcal{N}(0, \sigma^2 M)$$

where

$$M = I - H = I - XCX^{\mathrm{T}} = I - X(X^{\mathrm{T}}X)^{-1}X^{\mathrm{T}}$$

It holds in particular hence that

$$E[e_i] = 0 \qquad \text{for} \quad i = 1, 2, \dots, n$$

and

$$\mathrm{Var}(e_i) = \sigma^2 m_{ii} = \sigma^2 - \mathrm{Var}(\hat{y}_i) \qquad \text{for} \quad i = 1, 2, \dots, n$$

$$\mathrm{cov}(e_i, e_j) = \sigma^2 m_{ij} \qquad \text{for} \quad i, j = 1, 2, \dots, n$$

**Theorem 4:** If $\mathrm{rank}(X) = k + 1$, then

$$b \sim \mathcal{N}(\boldsymbol{\beta}, \sigma^2 C)$$

It holds in particular hence that

$$\mathrm{E}[b_j] = \beta_j \qquad \text{for} \quad j = 0, 1, \dots, k$$

and

$$\mathrm{Var}(b_j) = \sigma^2 c_{jj} \qquad \text{for} \quad j = 0, 1, \dots, k$$

$$\mathrm{cov}(b_j, b_i) = \sigma^2 c_{ji} \qquad \text{for} \quad j, i = 0, 1, \dots, k$$

# Residual Sum of Squares, $\chi^2$-test for the variance $\sigma^2$, and confidence intervals

- Residual Sum of Squares (RSS)
- Theorem 5: $\mathrm{RSS}/\sigma^2 \sim \chi^2_{n-p}$
- $\chi^2$-test for the variance $\sigma^2$
- Confidence intervals

Recall the **Residual Sum of Squares** is

$$\text{RSS} = \sum_{i=1}^{n} e_i^2 = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 = \sum_{i=1}^{n} (y_i - x_i b)^2 = (y - Xb)^{\text{T}}(y - Xb)$$

We know that the matrix $H$ is symmetric $(H^{\text{T}} = H)$ and idempotent $(HH = H)$.

Moreover, we know <u>by the Pythagoras Theorem</u> (see above) that

$$e^{\text{T}}e = y^{\text{T}}y - \hat{y}^{\text{T}}\hat{y} = y^{\text{T}}y - y^{\text{T}}H^{\text{T}}Hy =$$
$$= y^{\text{T}}y - y^{\text{T}}Hy = y^{\text{T}}y - y^{\text{T}}\hat{y} = y^{\text{T}}(y - Xb)$$

# Multiple Linear Regression: Mean Square Error

Put together, the **Residual Sum of Squares** is

$$\text{RSS} = \sum_{i=1}^{n} e_i^2 = e^{\mathrm{T}}e = (y - Xb)^{\mathrm{T}}(y - Xb) = y^{\mathrm{T}}(y - Xb)$$

Define the **residual variance** or the **Mean Square Error** as

$$s^2 = \frac{\text{RSS}}{n - p} = \frac{\sum_{i=1}^{n}(y_i - x_i b)^2}{n - p}$$

where

$$p = \text{rank}(X)$$

**Theorem 5:**  It holds

$$\frac{\text{RSS}}{\sigma^2} \sim \chi^2_{n-p}$$

where

$$p = \text{rank}(X) \qquad \text{and} \qquad \text{RSS} = e^{\mathrm{T}}e = y^{\mathrm{T}}y - y^{\mathrm{T}}Xb$$

Recall that, if $X \sim \chi^2_{n-p}$, then $\mathrm{E}[X] = n - p$.

Therefore:

$$\mathrm{E}[\text{RSS}] = \sigma^2(n-p)$$

$$\mathrm{E}[s^2] = \mathrm{E}\left[\frac{\text{RSS}}{n-p}\right] = \sigma^2$$

# Multiple Linear Regression:  Theorem 5

**Remark:**  Use Theorem 5  $(RSS/\sigma^2 \sim \chi^2_{n-p})$

— to obtain an unbiased estimate of the variance:

$$E[s^2] = \sigma^2 \qquad \text{that is} \qquad s^2 \approx \sigma^2$$

— for a $\chi^2$-test about the variance:

$$\frac{(n-p)s^2}{\sigma^2} \sim \chi^2_{n-p}$$

— or to establish the confidence intervals for the variance  $\sigma^2$

**Remark:**

**Theorem 5** $(\text{RSS}/\sigma^2 \sim \chi^2_{n-p})$ can be used to conduct the $\chi^2$-test for the variance.

- Let $\sigma_0^2 \in \mathbb{R}_0^+$ be a prescribed number.

- Formulate the null hypothesis:

$$H_0: \quad \sigma^2 = \sigma_0^2$$

- Formulate the alternative hypothesis

— two-sided: $\quad H_1: \quad \sigma^2 \neq \sigma_0^2$

— one-sided: $\quad H_1: \quad \sigma^2 < \sigma_0^2$

— one-sided: $\quad H_1: \quad \sigma^2 > \sigma_0^2$

<u>Notation:</u> Let

$$\chi^2_{n-p}(q)$$

denote the **quantile function of Pearson's $\chi^2$-distribution** with $n-p$ d.f., where $p = \mathrm{rank}(X)$.

The quantile function $\chi^2_{n-p}(q)$ is the function inverse to the cumulative distribution function $F(x)$ of **Pearson's $\chi^2$-distribution** with $n-p$ degrees of freedom, i.e.

$$\chi^2_{n-p}(q) = F^{-1}(q) \qquad \text{for} \quad q \in (0,1)$$

# $\chi^2$-test for the variance $\sigma^2$

<u>Notation:</u> Let

$$\chi^2_{n-p}(q)$$

denote the **quantile function of Pearson's $\chi^2$-distribution** with $n-p$ d.f.,

where $p = \text{rank}(X)$.

In other words, if $0 < q < 1$, then $x = \chi^2_{n-p}(q)$ is the unique value such that

$$\int_{-\infty}^{\chi^2_{n-p}(q)} f(t)\,dt = \int_{-\infty}^{x} f(t)\,dt = q$$

where $f(t)$ is the density of Pearson's $\chi^2$-distribution with $n-p$ d.f.

# $\chi^2$-test for the variance $\sigma^2$

Having chosen the value $\sigma_0^2 \in \mathbb{R}_0^+$ and assuming the null hypothesis $H_0: \sigma^2 = \sigma_0^2$ is true, calculate the statistic

$$X^2 = \frac{RSS}{\sigma^2} = \frac{RSS}{\sigma_0^2} = \frac{\sum_{i=1}^{n}(y_i - \hat{y}_i)^2}{\sigma_0^2}$$

# $\chi^2$-test for the variance $\sigma^2$

The $\chi^2$-test for $\sigma^2$ with two-sided alternative hypothesis ($\sigma^2 \neq \sigma_0^2$):

- choose **the level of significance**, a small number $\alpha > 0$, such as $\alpha = 5\%$, other popular values are $\alpha = 10\%$ or $\alpha = 1\%$ or $\alpha = 0.1\%$ etc.

- the **critical values** are $c = \chi^2_{n-p}\left(\frac{\alpha}{2}\right)$ and $d = \chi^2_{n-p}\left(1 - \frac{\alpha}{2}\right)$

- if $X^2 \in [0, c] \cup [d, +\infty)$, **the critical region**, then <u>**reject**</u> the null hypothesis

- if $X^2 \in (c, d)$, then <u>**do not reject**</u> (or <u>fail to reject</u>) the null hypothesis

# $\chi^2$-test for the variance $\sigma^2$

The $\chi^2$-test for $\sigma^2$ with one-sided alternative hypothesis ($\sigma^2 < \sigma_0^2$):

- choose **the level of significance**, a small number $\alpha > 0$, such as $\alpha = 5\,\%$, other popular values are $\alpha = 10\,\%$ or $\alpha = 1\,\%$ or $\alpha = 0.1\,\%$ etc.

- the **critical value** is $c = \chi^2_{n-p}(\alpha)$

- if $X^2 \in [0, c]$, **the critical region**, then <u>**reject**</u> the null hypothesis

- if $X^2 \in (c, +\infty)$, then <u>**do not reject**</u> (or <u>fail to reject</u>) the null hypothesis

# $\chi^2$-test for the variance $\sigma^2$

The $\chi^2$-test for $\sigma^2$ with one-sided alternative hypothesis $(\sigma^2 > \sigma_0^2)$:

- choose **the level of significance**, a small number $\alpha > 0$, such as $\alpha = 5\,\%$, other popular values are $\alpha = 10\,\%$ or $\alpha = 1\,\%$ or $\alpha = 0.1\,\%$ etc.

- the **critical value** is $d = \chi_{n-p}^2(1-\alpha)$

- if $X^2 \in [d, +\infty)$, **the critical region**, then <u>**reject**</u> the null hypothesis

- if $X^2 \in [0, d)$, then <u>**do not reject**</u> (or <u>fail to reject</u>) the null hypothesis

Let $x, y \in \mathbb{R}^+$ be any numbers such that $x < y$ and let $F(x)$ be the cumulative distribution function of Pearson's $\chi^2$-distribution with $n - p$ degrees of freedom. Then, by the definition of the cumulative distribution function and by Theorem 5, the probability

$$P\left(x < \frac{\text{RSS}}{\sigma^2} \leq y\right) = F(y) - F(x)$$

Therefore

$$P\left(\frac{\text{RSS}}{y} \leq \sigma^2 < \frac{\text{RSS}}{x}\right) = F(y) - F(x)$$

# Confidence interval for the variance $\sigma^2$

We have:

$$P\left(\frac{\text{RSS}}{y} \leq \sigma^2 < \frac{\text{RSS}}{x}\right) = F(y) - F(x)$$

Choose the level of significance, a small number $\alpha > 0$, such as $\alpha = 5\%$.

Let $y = \chi^2_{n-p}\left(1 - \frac{\alpha}{2}\right)$ and let $x = \chi^2_{n-p}\left(\frac{\alpha}{2}\right)$. Recall that $\chi^2_{n-p}(q) = F^{-1}(q)$.

Then, by the continuity of the cumulative distribution function, **the probability** that

$$\text{the unknown } \sigma^2 \in \left[\frac{\text{RSS}}{\chi^2_{n-p}\left(1 - \frac{\alpha}{2}\right)}, \frac{\text{RSS}}{\chi^2_{n-p}\left(1 - \frac{\alpha}{2}\right)}\right]$$

**is about** $1 - \alpha = 95\%$.

We have:

$$P\left(\frac{\text{RSS}}{y} \le \sigma^2 < \frac{\text{RSS}}{x}\right) = F(y) - F(x)$$

Choose the level of significance, a small number $\alpha > 0$, such as $\alpha = 5\,\%$.

Let $y = \chi^2_{n-p}(1-\alpha)$ and let $x \searrow 0$. Recall that $\chi^2_{n-p}(q) = F^{-1}(q)$.

Then, by the continuity of the cumulative distribution function, **the probability** that

$$\text{the unknown } \sigma^2 \in \left[\frac{\text{RSS}}{\chi^2_{n-p}(1-\alpha)}, \ +\infty\right)$$

**is about** $1-\alpha = 95\,\%$.

# Confidence interval for the variance $\sigma^2$

We have:

$$P\left(\frac{\text{RSS}}{y} \leq \sigma^2 < \frac{\text{RSS}}{x}\right) = F(y) - F(x)$$

Choose the level of significance, a small number $\alpha > 0$, such as $\alpha = 5\,\%$.

Let $y = +\infty$ and let $x = \chi^2_{n-p}(\alpha)$. Recall that $\chi^2_{n-p}(q) = F^{-1}(q)$.

Then, by the continuity of the cumulative distribution function, **the probability** that

$$\text{the unknown } \sigma^2 \in \left[0, \; \frac{\text{RSS}}{\chi^2_{n-p}(\alpha)}\right]$$

**is about** $1-\alpha = 95\,\%$.

# *t*-test for a single linear combination of the parameters $\beta_0, \beta_1, \ldots, \beta_k$ — e.g. an individual parameter $\beta_j$ — and confidence interval

- Theorem 6: $\boldsymbol{p}^{\mathrm{T}}\boldsymbol{b} \sim \mathcal{N}\left(\boldsymbol{p}^{\mathrm{T}}\boldsymbol{\beta},\ \sigma^2 \boldsymbol{p}^{\mathrm{T}}\boldsymbol{C}\boldsymbol{p}\right)$

  and $\quad \dfrac{\boldsymbol{p}^{\mathrm{T}}\boldsymbol{b} - \boldsymbol{p}^{\mathrm{T}}\boldsymbol{\beta}}{\sqrt{s^2}\sqrt{\boldsymbol{p}^{\mathrm{T}}\boldsymbol{C}\boldsymbol{p}}} \sim t_{n-(k+1)} \quad$ if $\operatorname{rank}(\boldsymbol{X}) = k+1$

- *t*-test for an individual parameter $\beta_j$

- Confidence interval for the $\beta_j$

**Theorem 6:**  Assume for simplicity that $\mathrm{rank}(X) = k + 1$

and let $p^{\mathrm{T}} \in \mathbb{R}^{1 \times (1+k)}$ be a non-zero row vector $(p^{\mathrm{T}} \neq 0^{\mathrm{T}})$.

Then

$$p^{\mathrm{T}}b \sim \mathcal{N}(p^{\mathrm{T}}\beta, \sigma^2 p^{\mathrm{T}}Cp)$$

and

$$\frac{p^{\mathrm{T}}b - p^{\mathrm{T}}\beta}{\sqrt{s^2}\sqrt{p^{\mathrm{T}}Cp}} \sim t_{n-(k+1)}$$

**Remark:**  The matrix $X^{\mathrm{T}}X$ is positively definite.

# Multiple Linear Regression: Prediction (Extrapolation)

**Remark:** Given a new row vector $x = (x_0, x_1, x_2 \ldots, x_k) \in \mathbb{R}^{1 \times (1+k)}$,

which is not included in the matrix $X$, we may wish to predict (extrapolate)

the value of the random variable $Y$ for this new statistical unit ($x$).

Assuming that the model is true, we should have

$$Y_x \approx x\beta$$

or

$$Y_x = x\beta + \varepsilon$$

where $\varepsilon$ is the random error.

Not knowing the parameters $\boldsymbol{\beta}$, we have to use their estimates $\boldsymbol{b}$ instead.

Then **the point estimate** of the value $Y_x$ is:

$$\tilde{Y}_x = \boldsymbol{xb}$$

**Remark:** If that $\text{rank}(\boldsymbol{X}) = k + 1$, then we can consider $\boldsymbol{p}^{\mathrm{T}} = \boldsymbol{x}$ and Theorem 6

$$\frac{\boldsymbol{p}^{\mathrm{T}}\boldsymbol{b} - \boldsymbol{p}^{\mathrm{T}}\boldsymbol{\beta}}{\sqrt{s^2}\sqrt{\boldsymbol{p}^{\mathrm{T}}\boldsymbol{Cp}}} \sim t_{n-(k+1)}$$

to obtain a confidence interval for the true value $\boldsymbol{x\beta}$.

Choose the level of significance, a small number $\alpha > 0$, such as $\alpha = 5\,\%$.

Let $c > 0$ be the value such that

$$\int_{-c}^{+c} f(t)\,\mathrm{d}t = 1 - \alpha$$

where $f(t)$ is the density of Student's $t$-distribution with $n - (k+1)$ d.f.

Then, by Theorem 6,

$$P\left(-c \leq \frac{p^{\mathrm{T}}b - p^{\mathrm{T}}\beta}{\sqrt{s^2}\sqrt{p^{\mathrm{T}}Cp}} \leq +c\right) = P\left(-c \leq \frac{xb - x\beta}{\sqrt{s^2}\sqrt{p^{\mathrm{T}}Cp}} \leq +c\right) = 1 - \alpha$$

By Theorem 6:

$$P\left(-c \le \frac{xb - x\beta}{\sqrt{s^2}\sqrt{p^{\mathrm{T}}Cp}} \le +c\right) = 1 - \alpha$$

$$P\left(-c\sqrt{s^2}\sqrt{p^{\mathrm{T}}Cp} \le xb - x\beta \le +c\sqrt{s^2}\sqrt{p^{\mathrm{T}}Cp}\right) = 1 - \alpha$$

$$P\left(xb - c\sqrt{s^2}\sqrt{p^{\mathrm{T}}Cp} \le x\beta \le xb + c\sqrt{s^2}\sqrt{p^{\mathrm{T}}Cp}\right) = 1 - \alpha$$

Having obtained

$$P\left(xb - c\sqrt{s^2}\sqrt{p^\mathsf{T} Cp} \leq x\beta \leq xb + c\sqrt{s^2}\sqrt{p^\mathsf{T} Cp}\right) = 1 - \alpha$$

the **probability** that the unknown

$$x\beta \in \left[xb - t_{n-(k+1)}\left(1 - \frac{\alpha}{2}\right)\sqrt{s^2}\sqrt{p^\mathsf{T} Cp}, \ xb + t_{n-(k+1)}\left(1 - \frac{\alpha}{2}\right)\sqrt{s^2}\sqrt{p^\mathsf{T} Cp}\right]$$

**is about** $1 - \alpha = 95\,\%$,

where $t_{n-(k+1)}(q)$ denotes the **quantile function of Student's _t_-distribution**

**Corollary:**  By considering

$$p^T = (0 \quad \ldots \quad 0 \quad 1 \quad 0 \quad \ldots \quad 0) \qquad \text{with the 1 at the } j-\text{th position}$$

we obtain:

$$\frac{b_j - \beta_j}{\sqrt{s^2}\sqrt{c_{jj}}} \sim t_{n-(k+1)} \qquad \text{for} \quad j = 0, 1, \ldots, k$$

**Remark:**  Use the Corollary

— for $t$-tests about the parameters $\beta_0, \beta_1, \ldots, \beta_k$  of the model,

— to establish the confidence intervals for the parameters $\beta_0, \beta_1, \ldots, \beta_k$,

- Choose any non-zero $p^T \in \mathbb{R}^{1\times(1+k)}$ and let $a \in \mathbb{R}$ be a prescribed number.

- We can then use Theorem 5 to test the null hypothesis $H_0$ that $p^T\beta = a$.

- By taking a particular choice of the non-zero $p^T \in \mathbb{R}^{1\times(1+k)}$, we can use

the Corollary $\left(\frac{b_j - \beta_j}{\sqrt{s^2}\sqrt{c_{jj}}} \sim t_{n-(k+1)}\right)$ to test the null hypothesis $H_0$ that $\beta_j = a$

or $\beta_j = 0$ (if we put $a = 0$ in particular).

<u>Notation:</u> Let

$$t_{n-(k+1)}(q)$$

denote the **quantile function of Student's *t*-distribution** with $n-(k+1)$ d.f.

The quantile function $t_{n-(k+1)}(q)$ is the function inverse to the cumulative distribution function $F(x)$ of **Student's *t*-distribution** with $n-(k+1)$ degrees of freedom, i.e.

$$t_{n-(k+1)}(q) = F^{-1}(q) \qquad \text{for} \quad q \in (0,1)$$

<u>Notation:</u> Let

$$t_{n-(k+1)}(q)$$

denote the **quantile function of Student's *t*-distribution** with $n-(k+1)$ d.f.

In other words, if $0 < q < 1$, then $x = t_{n-(k+1)}(q)$ is the unique value such that

$$\int_{-\infty}^{t_{n-(k+1)}(q)} f(t)\, dt = \int_{-\infty}^{x} f(t)\, dt = q$$

where $f(t)$ is the density of Student's *t*-distribution with $n-(k+1)$ d.f.

Choosing the index $j \in \{0, 1, \ldots, k\}$ and a value $b_{j0} \in \mathbb{R}$,

formulate the **null hypothesis**

$$H_0: \quad \beta_j = b_{j0}$$

formulate the **alternative hypothesis**

- two-sided: $\qquad H_1: \quad \beta_j \neq b_{j0}$

- one-sided: $\qquad H_1: \quad \beta_j < b_{j0}$

- one-sided: $\qquad H_1: \quad \beta_j > b_{j0}$

and use the aforementioned Corollary to conduct the test.

Having chosen the value $b_{j0} \in \mathbb{R}$, such as $b_{j0} = 0$, and assuming

the null hypothesis $H_0$: $\beta_j = b_{j0}$ is true, calculate the statistic

$$T = \frac{b_j - \beta_j}{\sqrt{s^2}\sqrt{c_{jj}}} = \frac{b_j - b_{j0}}{\sqrt{s^2}\sqrt{c_{jj}}} = \frac{b_j}{\sqrt{s^2}\sqrt{c_{jj}}}$$

The *t*-test for $\beta_j$ with two-sided alternative hypothesis ($\beta_j \neq b_{j0}$):

- choose **the level of significance**, a small number $\alpha > 0$, such as $\alpha = 5\,\%$, other popular values are $\alpha = 10\,\%$ or $\alpha = 1\,\%$ or $\alpha = 0.1\,\%$ etc.

- the **critical value** is $c = t_{n-(k+1)}\left(1 - \dfrac{\alpha}{2}\right)$

- if $T \in (-\infty, -c] \cup [+c, +\infty)$, **the critical region**, then **reject** the null hypothesis

- if $T \in (-c, +c)$, then **do not reject** (or fail to reject) the null hypothesis

The *t*-test for $\beta_j$ with one-sided alternative hypothesis ($\beta_j < b_{j0}$):

- choose **the level of significance**, a small number $\alpha > 0$, such as $\alpha = 5\,\%$, other popular values are $\alpha = 10\,\%$ or $\alpha = 1\,\%$ or $\alpha = 0.1\,\%$ etc.

- the **critical value** is $c = t_{n-(k+1)}(1 - \alpha)$

- if $T \in (-\infty, -c]$, **the critical region**, then **reject** the null hypothesis

- if $T \in (-c, +\infty)$, then **do not reject** (or fail to reject) the null hypothesis

The *t*-test for  $\beta_j$  with one-sided alternative hypothesis ($\beta_j > b_{j0}$):

- choose **the level of significance**, a small number  $\alpha > 0$,  such as  $\alpha = 5\,\%$, other popular values are  $\alpha = 10\,\%$  or  $\alpha = 1\,\%$  or  $\alpha = 0.1\,\%$  etc.

- the **critical value** is  $c = t_{n-(k+1)}(1-\alpha)$

- if  $T \in [+c, +\infty)$,  **the critical region**, then **<u>reject</u>** the null hypothesis

- if  $T \in (-\infty, +c)$,  then **<u>do not reject</u>** (or <u>fail to reject</u>) the null hypothesis

¡¡¡**WARNING!!!** **It usually makes no sense to test the null hypothesis** $\beta_0 = 0$ **if** $X_0 = 1$, **that is, if** $\beta_0$ **is the intercept term. Do not use the aforementioned test unless you know what and why you are doing.**

It can make sense to test the null hypothesis $\beta_0 = 0$ if the independent values $x_1, x_2, \ldots, x_n$ are from a neighbourhood of zero.

Otherwise (if the cluster of $x_1, x_2, \ldots, x_n$ is far from zero) it hardly makes sense to test the null hypothesis $\beta_0 = 0$ because the intercept term $\beta_0$ is just a constant

Let $x, y \in \mathbb{R}$ be any numbers such that $x < y$ and let $F(x)$ be the cumulative distribution function of Student's $t$-distribution with $n - (k + 1)$ degrees of freedom. Then, by the definition of the cumulative distribution function and by the Corollary, the probability

$$P\left( x < \frac{b_j - \beta_j}{\sqrt{s^2}\sqrt{c_{jj}}} \leq y \right) = F(y) - F(x)$$

Therefore

$$P\left( x\sqrt{s^2}\sqrt{c_{jj}} < b_j - \beta_j \leq y\sqrt{s^2}\sqrt{c_{jj}} \right) = F(y) - F(x)$$

$$P\left( b_j - y\sqrt{s^2}\sqrt{c_{jj}} \leq \quad \beta_j \quad < b_j - x\sqrt{s^2}\sqrt{c_{jj}} \right) = F(y) - F(x)$$

We have:

$$P\left(b_j - y\sqrt{s^2}\sqrt{c_{jj}} \leq \beta_j < b_j - x\sqrt{s^2}\sqrt{c_{jj}}\right) = F(y) - F(x)$$

Choose the level of significance, a small number $\alpha > 0$, such as $\alpha = 5\,\%$.

Let $y = t_{n-(k+1)}\left(1 - \frac{\alpha}{2}\right)$ and let $x = -y = -t_{n-(k+1)}\left(1 - \frac{\alpha}{2}\right) = t_{n-(k+1)}\left(\frac{\alpha}{2}\right)$.

Recall that $t_{n-(k+1)}(q) = F^{-1}(q)$.

Then, by the continuity of the cumulative distribution function, **the probability** that

the unknown $\beta_j \in \left[b_j - t_{n-(k+1)}\left(1 - \frac{\alpha}{2}\right)\sqrt{s^2}\sqrt{c_{jj}},\ b_j + t_{n-(k+1)}\left(1 - \frac{\alpha}{2}\right)\sqrt{s^2}\sqrt{c_{jj}}\right]$

We have:

$$P\left(b_j - y\sqrt{s^2}\sqrt{c_{jj}} \leq \beta_j < b_j - x\sqrt{s^2}\sqrt{c_{jj}}\right) = F(y) - F(x)$$

Choose the level of significance, a small number $\alpha > 0$, such as $\alpha = 5\,\%$.

Let $y = t_{n-(k+1)}(1-\alpha)$ and let $x = -\infty$. Recall that $t_{n-(k+1)}(q) = F^{-1}(q)$.

Then, by the continuity of the cumulative distribution function, **the probability** that

$$\text{the unknown } \beta_j \in \left[b_j - t_{n-(k+1)}(1-\alpha)\sqrt{s^2}\sqrt{c_{jj}},\ +\infty\right)$$

**is about** $1-\alpha = 95\,\%$.

We have:

$$P\left(b_j - y\sqrt{s^2}\sqrt{c_{jj}} \le \beta_j < b_j - x\sqrt{s^2}\sqrt{c_{jj}}\right) = F(y) - F(x)$$

Choose the level of significance, a small number $\alpha > 0$, such as $\alpha = 5\,\%$.

Let $y = +\infty$ and let $x = t_{n-(k+1)}(\alpha)$. Recall that $t_{n-(k+1)}(q) = F^{-1}(q)$.

Then, by the continuity of the cumulative distribution function, **the probability** that

$$\text{the unknown } \beta_j \in \left(-\infty,\ b_j + t_{n-(k+1)}(\alpha)\sqrt{s^2}\sqrt{c_{jj}}\right]$$

**is about** $1-\alpha = 95\,\%$.

## ¡¡¡ WARNING !!!

- Never use the above *t*-test for the parameters $\beta_0, \beta_1, \ldots, \beta_k$ consecutively!

- Never use the above construction of the confidence intervals consecutively!

- Use the following result (Theorem 7) instead !

## $F$-test for the significance of the model and confidence region & $F$-test for a system of linear combinations of the parameters $\beta_0, \beta_1, \ldots, \beta_k$

- Theorem 7:

$$\frac{(\boldsymbol{b}-\boldsymbol{\beta})^{\mathrm{T}}(\boldsymbol{A}^{\mathrm{T}}\boldsymbol{X})(\boldsymbol{b}-\boldsymbol{\beta})}{\mathrm{RSS}} \bigg/ \frac{k+1}{n-(k+1)} \sim F_{k+1, n-(k+1)}$$

- $F$-test for the significance of the model

- Confidence region

- Theorem 8:

$$\frac{(\boldsymbol{A}\boldsymbol{b}-\boldsymbol{a})^{\mathrm{T}}(\boldsymbol{A}\boldsymbol{C}\boldsymbol{A}^{\mathrm{T}})^{-1}(\boldsymbol{A}\boldsymbol{b}-\boldsymbol{a})}{\mathrm{RSS}} \bigg/ \frac{r}{n-(k+1)} \sim F_{r, n-(k+1)}$$

**Theorem 7:**  Assume for simplicity that $\text{rank}(X) = k + 1$.

It holds

$$\frac{(b - \beta)^\mathrm{T}(X^\mathrm{T}X)(b - \beta)}{\text{RSS}} \Bigg/ \frac{k + 1}{n - (k + 1)} \sim F_{k+1,\, n-(k+1)}$$

**Theorem 7\*:**   Assume for simplicity that  $\mathrm{rank}(X) = k + 1$.

Let  $a \in \mathbb{R}^{1+k}$  be a vector.

If

$$\beta = a$$

then

$$\frac{(b - a)^{\mathrm{T}}(X^{\mathrm{T}}X)(b - a)}{\mathrm{RSS}} \Bigg/ \frac{k + 1}{n - (k + 1)} \; \sim \; F_{k+1, n-(k+1)}$$

**Corollary:**  By considering

$$a = 0$$

that is the zero vector, we are testing the null hypothesis that

$$H_0: \qquad \beta = 0$$

that is

$$H_0: \qquad \beta_0 = \beta_1 = \cdots = \beta_k = 0$$

that is **we are testing the <u>overall significance of the model</u>**.

**Corollary:**  By considering

$$a = 0$$

we obtain:

$$\frac{b^T X^T X b}{\text{RSS}} \Big/ \frac{k+1}{n-(k+1)} = \frac{\hat{y}^T \hat{y}}{\text{RSS}} \Big/ \frac{k+1}{n-(k+1)} \sim F_{k+1,\,n-(k+1)}$$

**Remark:**  Use this Corollary

— for *F*-test about the significance of the model,

— to establish the confidence region.

The orthogonal decomposition
of the vector $y \in \mathbb{R}^n$ :

$$y = \hat{y} + e \qquad \text{and} \qquad e \perp \hat{y}$$

$\{X\boldsymbol{\beta} : \boldsymbol{\beta} \in \mathbb{R}^{1+k}\}^{\perp}$

(the orthogonal
complement =
= the space of
the residuals)

vector of the numerical outcomes
of the $n$ random experiments

$$(\cotan \varphi)^2 / \frac{k+1}{n-(k+1)} \sim F_{k+1,\, n-(k+1)}$$

$(I - H)y = e$

$M$

$y$

By the Corollary

By the Pythagoras Theorem:

$$\|y\|^2 = \|\hat{y}\|^2 + \|e\|^2$$

$$y^{\mathrm{T}}y = \hat{y}^{\mathrm{T}}\hat{y} + e^{\mathrm{T}}e$$

$$\cotan^2 \varphi = \frac{\hat{y}^{\mathrm{T}}\hat{y}}{\mathrm{RSS}}$$

$\varphi$

subspace
of dimension
$n - (k + 1)$

$0$

$\hat{y} = Hy$

$\{X\boldsymbol{\beta} : \boldsymbol{\beta} \in \mathbb{R}^{1+k}\}$

(the linear hull of the columns of $X$ )

subspace of dimension
$k + 1$

**Residual Sum of Squares:**   $\mathrm{RSS} = \sum_{i=1}^{n}(y_i - \hat{y}_i)^2 = \sum_{i=1}^{n} e_i^2 = e^{\mathrm{T}}e$

<u>Notation:</u> Let

$$F_{k+1,\,n-(k+1)}(q)$$

denote the **quantile function of Fisher's *F*-distribution** with $k+1$ and $n-(k+1)$ degrees of freedom.

The quantile function $F_{k+1,\,n-(k+1)}(q)$ is the function inverse to the cumulative distribution function $F(x)$ of **Fisher's *F*-distribution** with $k+1$ and $n-(k+1)$ degrees of freedom, i.e.

$$F_{k+1,\,n-(k+1)}(q) = F^{-1}(q) \qquad \text{for} \quad q \in (0,1)$$

Notation: Let

$$F_{k+1,\,n-(k+1)}(q)$$

denote the **quantile function of Fisher's *F*-distribution** with $k+1$ and $n-(k+1)$ degrees of freedom.

In other words, if $0 < q < 1$, then $x = F_{k+1,\,n-(k+1)}(q)$ is the unique value such that

$$\int_{-\infty}^{F_{k+1,\,n-(k+1)}(q)} f(t)\,dt = \int_{-\infty}^{x} f(t)\,dt = q$$

where $f(t)$ is the density of Fisher's *F*-distribution with $k+1$ and $n-(k+1)$ d.f.

Formulate the **null hypothesis**

$$H_0: \qquad \beta_0 = \beta_1 = \cdots = \beta_k = 0$$

- **Be cautious** because it usually makes no sense to test the value

  of the **intercept term** $\beta_0$ (see above).

  See also the Coefficient of Determination ($R^2$) below.

- The alternative hypothesis is simply $H_1 \equiv \neg H_0$, the logical negation of $H_0$,

  that is $\beta_j \neq 0$ for at least one $j \in \{0, 1, \ldots, k\}$.

- Calculate the statistic

$$F = \frac{\hat{y}^{\mathrm{T}}\hat{y}}{\mathrm{RSS}} \Big/ \frac{k+1}{n-(k+1)}$$

- Choose **the level of significance**, a small number $\alpha > 0$, such as $\alpha = 5\%$.

- The **critical value** is $c = F_{k+1, n-(k+1)}(1-\alpha)$, that is $\int_{c}^{+\infty} f(x)\, dx = \alpha$,

  where $f(x)$ is the density of Fisher's *F*-distribution with $k+1$ and $n-(k+1)$

  degrees of freedom,

- If $F \in [c, +\infty)$, the **critical region**, then **reject** the null hypothesis.

- If $F \in (0, c)$, then **do not reject** (fail to reject) the null hypothesis.

Consider $\bar{\boldsymbol{\beta}} = 0$, let $x \in \mathbb{R}$ be any real number, and let $F(x)$ be the cumulative distribution function of Fisher's $F$-distribution with $k + 1$ and $n - (k + 1)$ degrees of freedom. Then, by the definition of the cumulative distribution function and by the Corollary, the probability

$$P\left(\frac{(\boldsymbol{b} - \boldsymbol{\beta})^{\mathrm{T}}\boldsymbol{X}^{\mathrm{T}}\boldsymbol{X}(\boldsymbol{b} - \boldsymbol{\beta})}{\mathrm{RSS}} \bigg/ \frac{k + 1}{n - (k + 1)} \leq x\right) = F(x) \qquad \text{for any} \quad \boldsymbol{\beta} \in \mathbb{R}^{1+k}$$

Choose the level of significance, a small number $\alpha > 0$, such as $\alpha = 5\%$.

Then the probability that

the unknown $\boldsymbol{\beta} \in \left\{ \boldsymbol{\beta} \in \mathbb{R}^{1+k} : \dfrac{(\boldsymbol{b} - \boldsymbol{\beta})^{\mathrm{T}} \boldsymbol{X}^{\mathrm{T}} \boldsymbol{X} (\boldsymbol{b} - \boldsymbol{\beta})}{\mathrm{RSS}} \bigg/ \dfrac{k+1}{n - (k+1)} \leq F_{k+1,\, n-(k+1)}(1 - \alpha) \right\}$

**is about** $1 - \alpha = 95\%$.

Remark: This confidence region is an ellipsoid centred at $\boldsymbol{b}$.

The nominator $((\boldsymbol{b} - \boldsymbol{\beta})^{\mathrm{T}} \boldsymbol{X}^{\mathrm{T}} \boldsymbol{X} (\boldsymbol{b} - \boldsymbol{\beta}))$ is a quadratic expression in $\boldsymbol{\beta}$.

To gain a geometrical insight, calculate the spectral / eigendecomposition

**Theorem 8:** Assume for simplicity that $\mathrm{rank}(X) = k + 1$. Let $a \in \mathbb{R}^r$ be a vector and let $A \in \mathbb{R}^{r \times (1+k)}$ be an $r \times (1 + k)$ matrix of full-rank where $r \leq 1 + k$, that is

$$r = \mathrm{rank}(A) \leq \mathrm{rank}(X) = k + 1$$

If

$$A\beta = a$$

then

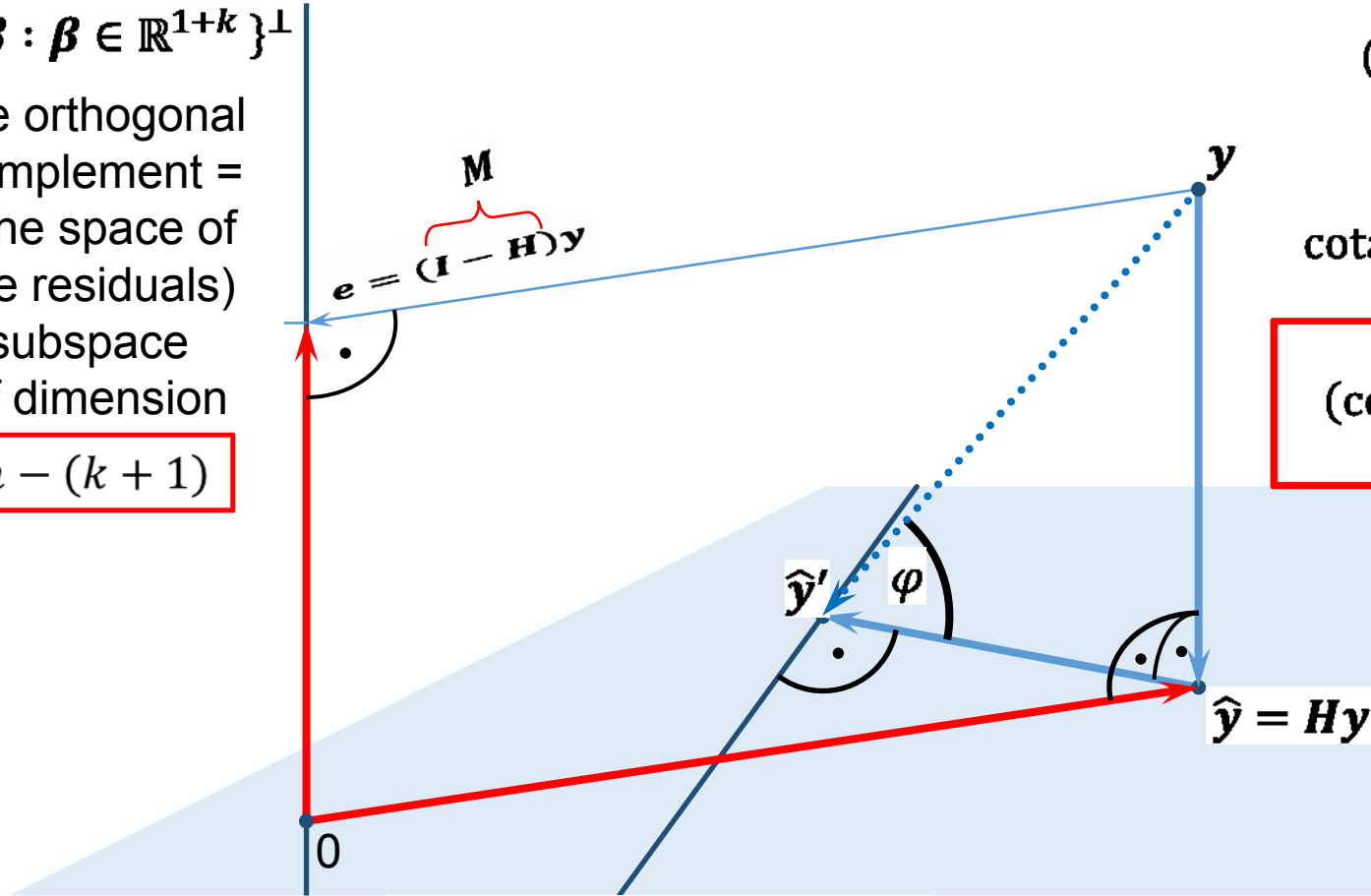$$\frac{(Ab - a)^{\mathrm{T}}(ACA^{\mathrm{T}})^{-1}(Ab - a)}{\mathrm{RSS}} \Big/ \frac{r}{n - (k + 1)} \sim F_{r, n-(k+1)}$$

$\{X\beta : \beta \in \mathbb{R}^{1+k}\}^{\perp}$

(the orthogonal complement = = the space of the residuals) subspace of dimension

$n - (k+1)$

$M$

$e = (I - H)y$

It holds:

$$(Ab - a)^{\mathrm{T}}(ACA^{\mathrm{T}})^{-1}(Ab - a) = \|\hat{y} - \hat{y}'\|^2 =$$
$$= (\hat{y} - \hat{y}')^{\mathrm{T}}(\hat{y} - \hat{y})$$

$$\mathrm{cotan}^2\, \varphi = \frac{(\hat{y} - \hat{y}')^{\mathrm{T}}(\hat{y} - \hat{y}')}{\mathrm{RSS}}$$

$$(\mathrm{cotan}\, \varphi)^2 \Big/ \frac{r}{n - (k+1)} \sim F_{r,\, n-(k+1)}$$

$y$

$\hat{y}'$

$\varphi$

$\hat{y} = Hy$

$0$

this is $\longrightarrow$ $\{X\beta : A\beta = a,\ \beta \in \mathbb{R}^{1+k}\}$

an affine subspace of dimension $k + 1 - r$

the dimension of its complement within the subspace of dimension $k + 1$ is $r$

$\{X\beta : \beta \in \mathbb{R}^{1+k}\}$

(the linear hull of the columns of $X$) subspace of dimension

$k + 1$

# Theorem 8

$$A\boldsymbol{\beta} = \boldsymbol{a} \qquad \Longrightarrow \qquad \frac{(A\boldsymbol{b}-\boldsymbol{a})^{\mathrm{T}}(AC A^{\mathrm{T}})^{-1}(A\boldsymbol{b}-\boldsymbol{a})}{\mathrm{RSS}} \Bigg/ \frac{r}{n-(k+1)} \sim F_{r,\,n-(k+1)}$$

is at the heart of the ANOVA method

and other results.

# The Coefficient of Determination ($R^2$)

- Assumption: $\mathbf{1} \in \{ \boldsymbol{X\beta} : \boldsymbol{\beta} \in \mathbb{R}^{1+k} \}$

- Motivation

- Some facts

- Theorem 8: Corollary

- $F$-test for the null hypothesis

  $H_0$: $\beta_1 = \ldots = \beta_k = 0$

**Assume** throughout this section that

$$1 \in \{X\beta : \beta \in \mathbb{R}^{1+k}\}$$

where **1** is the vector of $n$ ones:

$$1 = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}$$

<u>For example</u>, assume that

$$X_0 = 1$$

that is $\beta_0$ is the intercept term.

If $X_0 = 1$, that is $\beta_0$ is the intercept term, then it may be desirable to test the null hypothesis

$$H_0: \qquad \beta_1 = \cdots = \beta_k = 0$$

that is without the test for the parameter $\beta_0$.

To this end, apply <u>Theorem 8</u> with the $k \times (k+1)$ matrix and the $k$-vector

$$A = \begin{pmatrix} 0 & 1 & 0 & 0 & \cdots & 0 \\ 0 & 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & 1 \end{pmatrix} \qquad \text{and} \qquad a = \begin{pmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

Recall our assumption that

$$\mathbf{1} \in \{X\beta : \beta \in \mathbb{R}^{1+k}\}$$

Then the line

$$\{\mathbf{1}\lambda : \lambda \in \mathbb{R}\} \subset \{X\beta : \beta \in \mathbb{R}^{1+k}\}$$

In particular, if $X_0 = 1$, that is

$$X = \begin{pmatrix} 1 & x_{11} & \cdots & x_{1k} \\ 1 & x_{21} & \cdots & x_{2k} \\ 1 & x_{31} & \cdots & x_{3k} \\ \vdots & \vdots & \cdots & \vdots \\ 1 & x_{n1} & \cdots & x_{nk} \end{pmatrix} \qquad A = \begin{pmatrix} 0 & 1 & 0 & 0 & \cdots & 0 \\ 0 & 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & 1 \end{pmatrix} \quad \text{and} \quad a = \begin{pmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

then

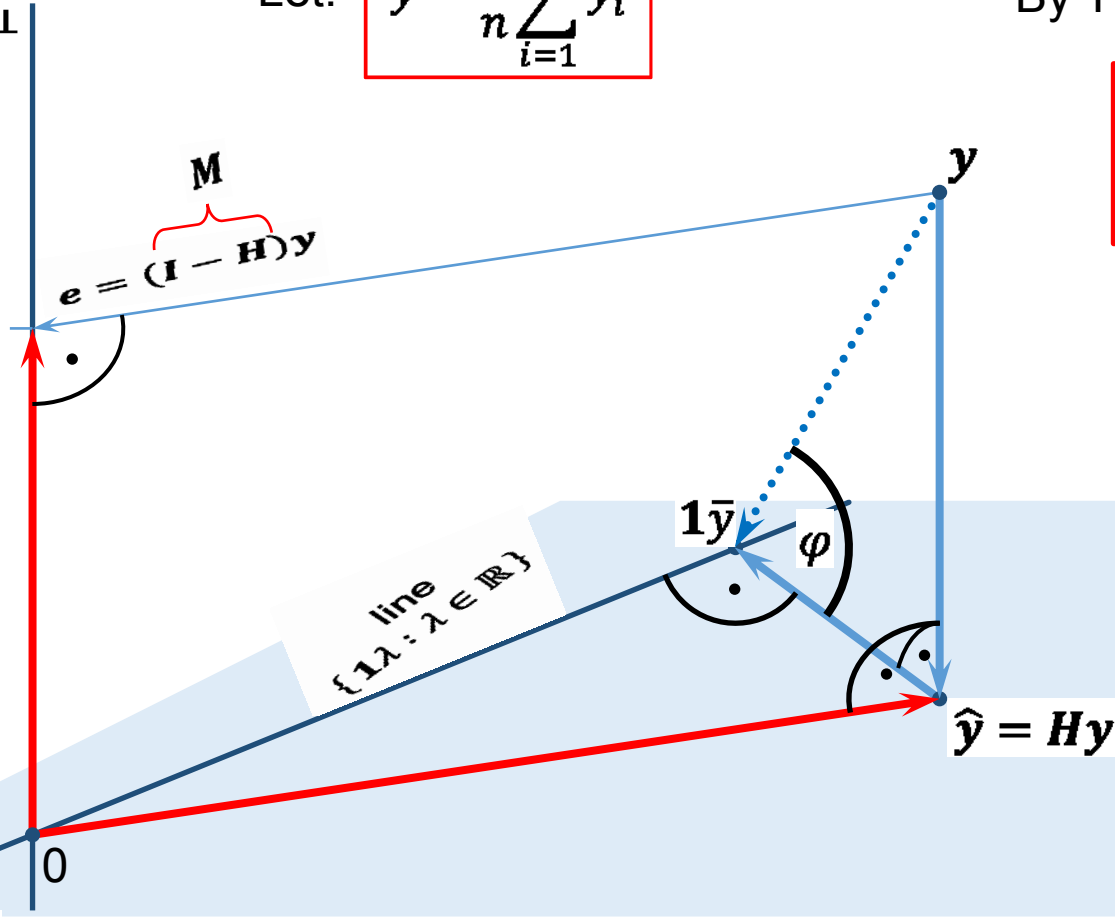# The Coefficient of Determination ($R^2$): Th. 8: Corollary

Let: $$\bar{y} = \frac{1}{n}\sum_{i=1}^{n} y_i$$

By Theorem 8:

$$(\cotan \varphi)^2 \Big/ \frac{k}{n-(k+1)} \sim F_{k,\, n-(k+1)}$$

$$\cotan^2 \varphi = \frac{(\hat{y} - \mathbf{1}\bar{y})^{\mathrm{T}}(\hat{y} - \mathbf{1}\bar{y})}{\text{RSS}} = \frac{R^2}{1 - R^2}$$

$\{X\beta : \beta \in \mathbb{R}^{1+k}\}^{\perp}$

(the orthogonal complement =
= the space of
the residuals)
subspace
of dimension

$n - (k+1)$

$\mathbf{M}$

$e = (I - H)y$

$y$

$\mathbf{1}\bar{y}$

$\varphi$

line $\{\mathbf{1}\lambda : \lambda \in \mathbb{R}\}$

$\hat{y} = Hy$

$0$

$\{X\beta : \beta \in \mathbb{R}^{1+k}\}$

(the linear hull of the columns of $X$)
subspace of dimension

$k + 1$

$\{\mathbf{1}\lambda : \lambda \in \mathbb{R}\}$

the line is a subspace
of dimension $1$

the dimension of its complement within
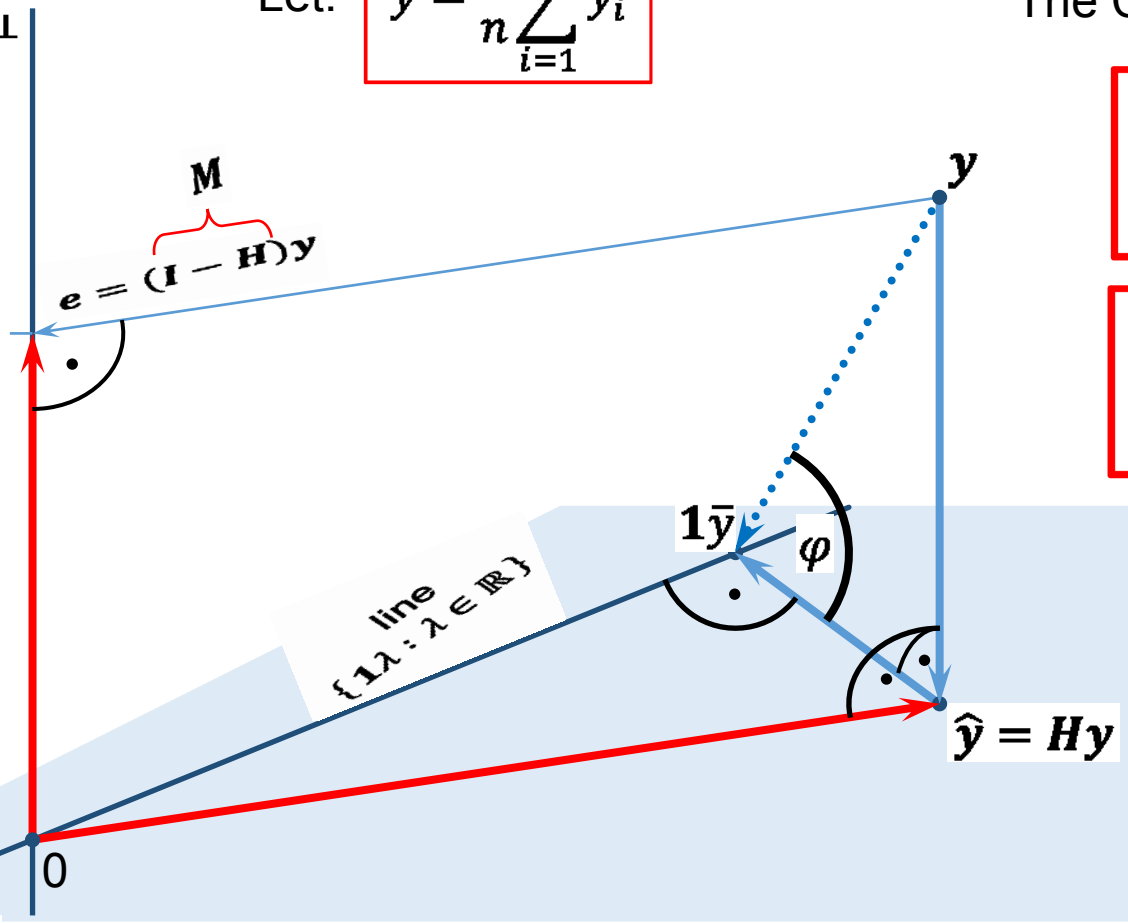the subspace of dimension $k + 1$ is $k$

Let: $\bar{y} = \dfrac{1}{n}\sum_{i=1}^{n} y_i$

The Coefficient of Determination:

$$R^2 = \cos^2 \varphi = \frac{(\hat{y} - \mathbf{1}\bar{y})^{\mathrm{T}}(\hat{y} - \mathbf{1}\bar{y})}{(y - \mathbf{1}\bar{y})^{\mathrm{T}}(y - \mathbf{1}\bar{y})}$$

$$\frac{R^2}{1 - R^2} = \frac{\cos^2 \varphi}{\sin^2 \varphi} = \cotan^2 \varphi$$

$\{X\beta : \beta \in \mathbb{R}^{1+k}\}^{\perp}$

(the orthogonal complement =
= the space of
the residuals)
subspace
of dimension

$n - (k + 1)$

$M$

$e = (I - H)y$

$y$

$\mathbf{1}\bar{y}$

$\varphi$

line $\{\mathbf{1}\lambda : \lambda \in \mathbb{R}\}$

$\hat{y} = Hy$

$0$

$\{X\beta : \beta \in \mathbb{R}^{1+k}\}$

(the linear hull of the columns of $X$)
subspace of dimension

$k + 1$

$\{\mathbf{1}\lambda : \lambda \in \mathbb{R}\}$

the line is a subspace
of dimension $\boxed{1}$

the dimension of its complement within
the subspace of dimension $k + 1$ is $k$

# The Coefficient of Determination ($R^2$): TSS=RSS+RegSS

Introduce the **Total Sum of Squares:**

$$\text{TSS} = (y - 1\bar{y})^{\text{T}}(y - 1\bar{y}) = \sum_{i=1}^{n}(y_i - \bar{y})^2$$

Introduce the **Regression Sum of Squares:**

$$\text{RegSS} = (\hat{y} - 1\bar{y})^{\text{T}}(\hat{y} - 1\bar{y}) = \sum_{i=1}^{n}(\hat{y}_i - \bar{y})^2$$

Recall the **Residual Sum of Squares:**

$$\text{RSS} = e^{\text{T}}e = \sum_{i=1}^{n}e_i^2 = (y - \hat{y})^{\text{T}}(y - \hat{y}) = \sum_{i=1}^{n}(y_i - \hat{y}_i)^2$$

# The Coefficient of Determination ($R^2$): TSS=RSS+RegSS

Let: $\bar{y} = \dfrac{1}{n}\sum\limits_{i=1}^{n} y_i$

$$\text{TSS} = (\boldsymbol{y} - \mathbf{1}\bar{y})^{\mathrm{T}}(\boldsymbol{y} - \mathbf{1}\bar{y})$$

$$\text{RegSS} = (\hat{\boldsymbol{y}} - \mathbf{1}\bar{y})^{\mathrm{T}}(\hat{\boldsymbol{y}} - \mathbf{1}\bar{y})$$

$$\text{RSS} = (\boldsymbol{y} - \hat{\boldsymbol{y}})^{\mathrm{T}}(\boldsymbol{y} - \hat{\boldsymbol{y}}) = \boldsymbol{e}^{\mathrm{T}}\boldsymbol{e}$$

By the Pythagoras Theorem:

$$\text{TSS} = \text{RSS} + \text{RegSS}$$

$\{\boldsymbol{X\beta} : \boldsymbol{\beta} \in \mathbb{R}^{1+k}\}^{\perp}$

(the orthogonal complement =
= the space of the residuals) subspace of dimension

$$n - (k + 1)$$

$M$

$\boldsymbol{e} = (\boldsymbol{I} - \boldsymbol{H})\boldsymbol{y}$

$\boldsymbol{y}$

$\mathbf{1}\bar{y}$

$\varphi$

line $\{\mathbf{1}\lambda : \lambda \in \mathbb{R}\}$

$\hat{\boldsymbol{y}} = \boldsymbol{Hy}$

$0$

$\{\mathbf{1}\lambda : \lambda \in \mathbb{R}\}$

the line is a subspace of dimension $1$

the dimension of its complement within the subspace of dimension $k+1$ is $k$

$\{\boldsymbol{X\beta} : \boldsymbol{\beta} \in \mathbb{R}^{1+k}\}$

(the linear hull of the columns of $\boldsymbol{X}$)
subspace of dimension

$k + 1$

# The Coefficient of Determination $(R^2)$: Some facts

**Proposition:** Under the assumption $\mathbf{1} \in \{X\boldsymbol{\beta} : \boldsymbol{\beta} \in \mathbb{R}^{1+k}\}$, it holds

$$\mathbf{1}^{\mathrm{T}}\mathbf{1}\bar{y} = \mathbf{1}^{\mathrm{T}}y = \mathbf{1}^{\mathrm{T}}\hat{y}$$

**In words:**

All three points $\mathbf{1}\bar{y}, \quad y, \quad \hat{y}$ lie in the hyperplane

$$\{\boldsymbol{\beta} \in \mathbb{R}^{1+k} : \mathbf{1}^{\mathrm{T}}\boldsymbol{\beta} = \mathbf{1}^{\mathrm{T}}\mathbf{1}\bar{y}\}$$

which is perpendicular to the line $\{\mathbf{1}\lambda : \lambda \in \mathbb{R}\}$.

# The Coefficient of Determination $(R^2)$: Some facts

**Proposition:** Under the assumption $1 \in \{X\beta : \beta \in \mathbb{R}^{1+k}\}$, it holds

$$1^T 1 \bar{y} = 1^T y = 1^T \hat{y}$$

**Corollary:**

$$1^T(y - \hat{y}) = 1^T e = \sum_{i=1}^{n} e_i = 0$$

# The Coefficient of Determination $(R^2)$: Some facts

**Proposition:** Under the assumption $1 \in \{X\beta : \beta \in \mathbb{R}^{1+k}\}$, it holds

$$1^T 1\bar{y} = 1^T y = 1^T \hat{y}$$

Proof:

The assumption equivalently says $H1 = 1$, where $H$ is the matrix

of the orthogonal projection onto the subspace $\{X\beta : \beta \in \mathbb{R}^{1+k}\}$.

Recall the matrix $H$ is symmetric $(H^T = H)$, therefore $1^T H^T = 1^T H = 1^T$.

Therefore $1^T y = 1^T H y = 1^T \hat{y}$.

The first equality is obvious:

$$1^T 1\bar{y} = \sum_{i=1}^{m} 1 \times 1 \times \sum_{i=1}^{n} y_i / n = n \times \sum_{i=1}^{n} y_i / n = \sum_{i=1}^{n} y_i = 1^T y.$$

**Proposition:** Under the assumption $\mathbf{1} \in \{X\boldsymbol{\beta} : \boldsymbol{\beta} \in \mathbb{R}^{1+k}\}$, it holds

$$TSS = RSS + RegSS$$

$$TSS = (\boldsymbol{y} - \mathbf{1}\bar{y})^T(\boldsymbol{y} - \mathbf{1}\bar{y})$$

$$RegSS = (\hat{\boldsymbol{y}} - \mathbf{1}\bar{y})^T(\hat{\boldsymbol{y}} - \mathbf{1}\bar{y})$$

$$RSS = (\boldsymbol{y} - \hat{\boldsymbol{y}})^T(\boldsymbol{y} - \hat{\boldsymbol{y}}) = \boldsymbol{e}^T\boldsymbol{e}$$

**Proof:**

The point $\hat{\boldsymbol{y}}$ is the orthogonal projection of the point $\boldsymbol{y}$ onto $\{X\boldsymbol{\beta} : \boldsymbol{\beta} \in \mathbb{R}^{1+k}\}$, therefore $(\boldsymbol{y} - \hat{\boldsymbol{y}}) \perp \{X\boldsymbol{\beta} : \boldsymbol{\beta} \in \mathbb{R}^{1+k}\}$.

We have $\hat{\boldsymbol{y}} \in \{X\boldsymbol{\beta} : \boldsymbol{\beta} \in \mathbb{R}^{1+k}\}$, and we assume $\mathbf{1} \in \{X\boldsymbol{\beta} : \boldsymbol{\beta} \in \mathbb{R}^{1+k}\}$, whence $\mathbf{1}\bar{y} \in \{X\boldsymbol{\beta} : \boldsymbol{\beta} \in \mathbb{R}^{1+k}\}$ follows, therefore $\hat{\boldsymbol{y}} - \mathbf{1}\bar{y} \in \{X\boldsymbol{\beta} : \boldsymbol{\beta} \in \mathbb{R}^{1+k}\}$ and $(\boldsymbol{y} - \hat{\boldsymbol{y}}) \perp (\hat{\boldsymbol{y}} - \mathbf{1}\bar{y})$. By using the Pythagoras Theorem,

# The Coefficient of Determination ($R^2$): Some facts

**Proposition:** Under the assumption $\mathbf{1} \in \{X\boldsymbol{\beta} : \boldsymbol{\beta} \in \mathbb{R}^{1+k}\}$,
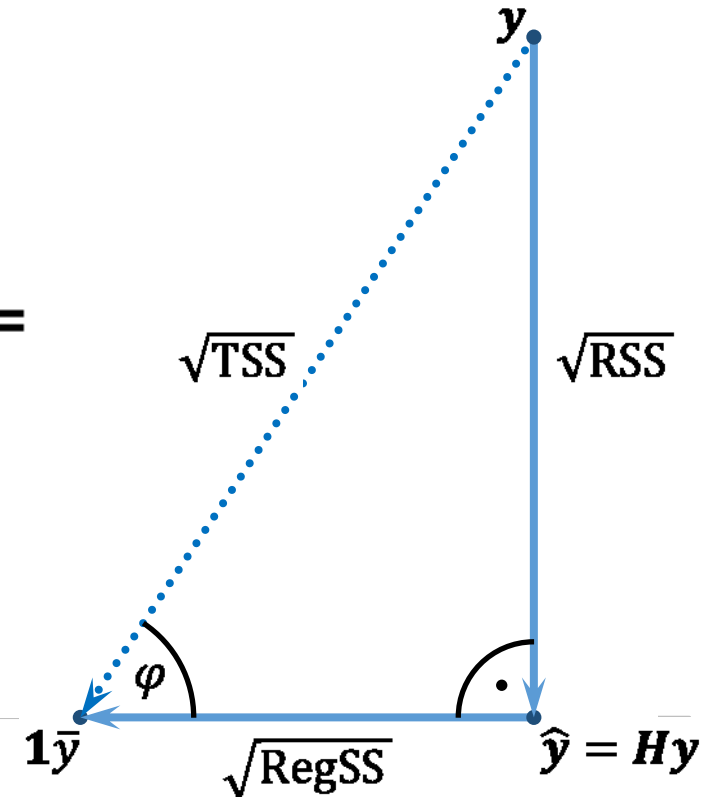
it holds

$$(\mathbf{y} - \mathbf{1}\bar{y})^{\mathrm{T}}(\hat{\mathbf{y}} - \mathbf{1}\bar{y}) = (\hat{\mathbf{y}} - \mathbf{1}\bar{y})^{\mathrm{T}}(\hat{\mathbf{y}} - \mathbf{1}\bar{y})$$

$$\mathrm{TSS} = (\mathbf{y} - \mathbf{1}\bar{y})^{\mathrm{T}}(\mathbf{y} - \mathbf{1}\bar{y})$$

$$\mathrm{RegSS} = (\hat{\mathbf{y}} - \mathbf{1}\bar{y})^{\mathrm{T}}(\hat{\mathbf{y}} - \mathbf{1}\bar{y})$$

$$\mathrm{RSS} = (\mathbf{y} - \hat{\mathbf{y}})^{\mathrm{T}}(\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{e}^{\mathrm{T}}\mathbf{e}$$

**Proof:**

$$(\mathbf{y} - \mathbf{1}\bar{y})^{\mathrm{T}}(\hat{\mathbf{y}} - \mathbf{1}\bar{y}) = \big((\mathbf{y} - \hat{\mathbf{y}}) + (\hat{\mathbf{y}} - \mathbf{1}\bar{y})\big)^{\mathrm{T}}(\hat{\mathbf{y}} - \mathbf{1}\bar{y}) =$$

$$= (\mathbf{y} - \hat{\mathbf{y}})^{\mathrm{T}}(\hat{\mathbf{y}} - \mathbf{1}\bar{y}) + (\hat{\mathbf{y}} - \mathbf{1}\bar{y})^{\mathrm{T}}(\hat{\mathbf{y}} - \mathbf{1}\bar{y}) =$$

$$= 0 + (\hat{\mathbf{y}} - \mathbf{1}\bar{y})^{\mathrm{T}}(\hat{\mathbf{y}} - \mathbf{1}\bar{y}) =$$

$$= (\hat{\mathbf{y}} - \mathbf{1}\bar{y})^{\mathrm{T}}(\hat{\mathbf{y}} - \mathbf{1}\bar{y})$$

q.e.d.

Assuming $1 \in \{X\beta : \beta \in \mathbb{R}^{1+k}\}$, define the

**Coefficient of Determination:**

$$TSS = (y - 1\bar{y})^T(y - 1\bar{y})$$
$$RegSS = (\hat{y} - 1\bar{y})^T(\hat{y} - 1\bar{y})$$
$$RSS = (y - \hat{y})^T(y - \hat{y}) = e^T e$$

$$R^2 = \frac{[(y - 1\bar{y})^T(\hat{y} - 1\bar{y})]^2}{(y - 1\bar{y})^T(y - 1\bar{y}) \times (\hat{y} - 1\bar{y})^T(\hat{y} - 1\bar{y})} =$$

$$= \frac{[(\hat{y} - 1\bar{y})^T(\hat{y} - 1\bar{y})]^2}{(y - 1\bar{y})^T(y - 1\bar{y}) \times (\hat{y} - 1\bar{y})^T(\hat{y} - 1\bar{y})} =$$

$$= \frac{(\hat{y} - 1\bar{y})^T(\hat{y} - 1\bar{y})}{(y - 1\bar{y})^T(y - 1\bar{y})} = \cos^2 \varphi = \frac{RegSS}{TSS}$$

# The Coefficient of Determination $(R^2)$

Assuming $1 \in \{X\boldsymbol{\beta} : \boldsymbol{\beta} \in \mathbb{R}^{1+k}\}$, define the

**Coefficient of Determination:**

$$TSS = (\boldsymbol{y} - \mathbf{1}\bar{y})^{\mathrm{T}}(\boldsymbol{y} - \mathbf{1}\bar{y})$$

$$RegSS = (\hat{\boldsymbol{y}} - \mathbf{1}\bar{y})^{\mathrm{T}}(\hat{\boldsymbol{y}} - \mathbf{1}\bar{y})$$

$$RSS = (\boldsymbol{y} - \hat{\boldsymbol{y}})^{\mathrm{T}}(\boldsymbol{y} - \hat{\boldsymbol{y}}) = \boldsymbol{e}^{\mathrm{T}}\boldsymbol{e}$$
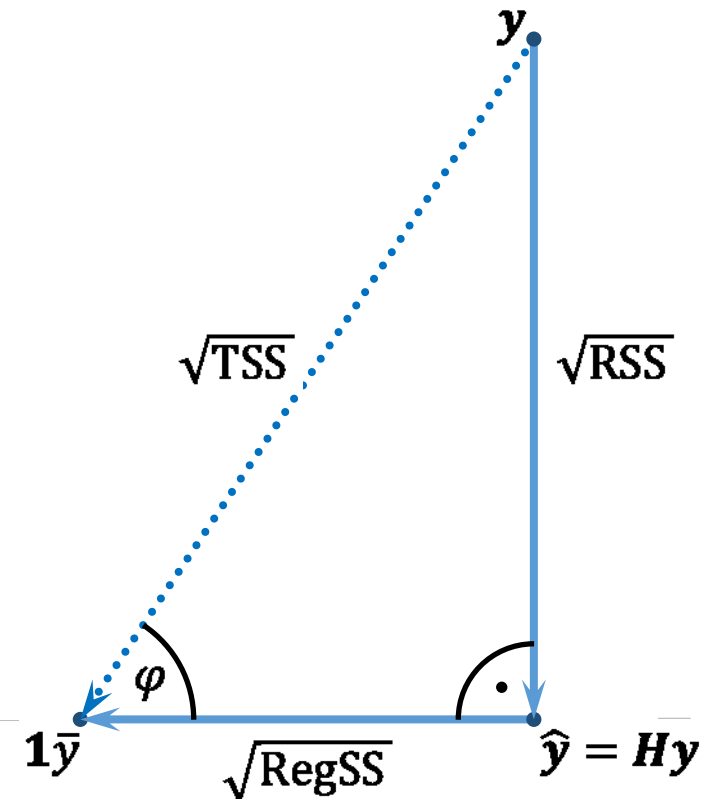
$$R^2 = \frac{RegSS}{TSS} = \frac{TSS - RSS}{TSS}$$

$$R^2 = \cos^2 \varphi = \frac{RegSS}{TSS} = 1 - \frac{RSS}{TSS}$$

$$\cotan^2 \varphi = \frac{\cos^2 \varphi}{\sin^2 \varphi} = \frac{R^2}{1 - R^2} = \frac{RegSS}{RSS}$$

**Theorem 8: Corollary:** Assume for simplicity that $\mathrm{rank}(X) = k + 1$

and assume that $1 \in \{X\boldsymbol{\beta} : \boldsymbol{\beta} \in \mathbb{R}^{1+k}\}$. Under the hypothesis that

$$\beta_1 = \cdots = \beta_k = 0$$

it holds

$$(\mathrm{cotan}\,\varphi)^2 \Big/ \frac{k}{n - (k+1)} = \frac{R^2}{1 - R^2} \Big/ \frac{k}{n - (k+1)} \sim F_{k,\,n-(k+1)}$$

$$= \frac{\mathrm{RegSS}}{\mathrm{RSS}} \Big/ \frac{k}{n - (k+1)} \sim F_{k,\,n-(k+1)}$$

# *F*-test for the null hypothesis $H_0$: $\beta_1 = \ldots = \beta_k = 0$

- Choose **the level of significance**, a small number $\alpha > 0$, such as $\alpha = 5\%$.

- Find the **critical value** $c > 0$ so that $\int_c^{+\infty} f(x)\,dx = \alpha$, where $f$ is the density of the *F*-distribution with $k$ and $n - (k+1)$ degrees of freedom.

- Calculate the statistic

$$F = \frac{R^2}{1 - R^2} \bigg/ \frac{k}{n - (k+1)} = \frac{\text{RegSS}}{\text{RSS}} \bigg/ \frac{k}{n - (k+1)} = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \hat{y}_i)^2} \bigg/ \frac{k}{n - (k+1)}$$

- If $F \in [c, +\infty)$, **the critical region**, then <u>**reject**</u> the null hypothesis.

- If $F \in [0, c)$, then <u>**do not reject**</u> (or <u>fail to reject</u>) the null hypothesis.

# The Coefficient of Determination ($R^2$)

Remark: The above $F$-test is one-factor ANOVA in fact.

The coefficient of determination

$$R^2 = \cos^2 \varphi = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = 1 - \frac{\text{RSS}}{\text{TSS}} = \frac{\text{RegSS}}{\text{TSS}}$$

is a "measure" (?) "how well the regression hyperplane $Y = b_0 + b_1 X_1 + \cdots + b_k X_k$
fits the observed data $(y_1, x_1), (y_2, x_2), \ldots, (y_n, x_n)$".

It holds

$$0 \le R^2 \le 1$$

# The Coefficient of Determination $(R^2)$

$$R^2 = \cos^2 \varphi$$

If $R^2 \nearrow 1$

— then

$$F = \frac{R^2}{1 - R^2} \Big/ \frac{k}{n - (k+1)} \quad \nearrow +\infty$$

— then
   — reject the null hypothesis that $(\beta_1 = \cdots = \beta_k = 0)$

— then
   — say "the fit is good"

# The Coefficient of Determination $(R^2)$

$$R^2 = \cos^2 \varphi$$

If $\quad R^2 \searrow 0$

— then

$$F = \frac{R^2}{1 - R^2} \Big/ \frac{k}{n - (k+1)} \quad \searrow 0$$

— then
- — fail to reject the null hypothesis that $(\beta_1 = \cdots = \beta_k = 0)$
- — it may be the case that

$$\mathrm{E}[y_i] = \beta_0 \qquad \text{for all} \quad i = 1, 2, \ldots, n \qquad \text{(cf. ANOVA)}$$

- — the sample $(y_1, x_1),\ (y_2, x_2)\ \ldots,\ (y_n, x_n)$ may come from one population

— then