

# LINEAR MODELS IN STATISTICS

LECTURE NOTES

BY

Andrea Kraus, Ph.D.

These lecture notes are intended for students of the 3rd year bachelor course *M5120 Linear Models in Statistics I*. This is a one-semester course offering an overview of linear statistical models as the fundamental tool of statistical analysis. The students encounter theory, software implementation, applications and interpretation. After the course the students are expected to recognize the situations that can be addressed by linear models, formulate and implement the model, and interpret the results. At the same time, the students are made aware of the limitations of the model and should be able to recognize and possibly avoid problems in a given situation.

The lecture notes are based primarily on the following texts:

Julian J. Faraway (2014). *Linear Models with R*. Second edition.

Chapman& Hall/CRC.

Simon N. Wood (2006). *Generalized Additive Models; An introduction with R*.

Chapman& Hall/CRC.

Jiří Anděl (2005). *Základy matematické statistiky*.

Matfyzpress.

I would like to thank Professor Victor Panaretos for kindly sharing his teaching experience and course materials for his course on Regression at the Swiss Federal Institute of Technology in Lausanne. I would also like to thank my students for their feedback on various bits and pieces of these notes.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Motivation . . . . .	1
1.2	Statistics . . . . .	5
1.3	Data analysis in practice . . . . .	10
1.4	Descriptive statistics . . . . .	12
1.4.1	Types of variables . . . . .	12
1.4.2	Relationships between variables . . . . .	15
<b>2</b>	<b>Linear algebra essentials</b>	<b>20</b>
2.1	The problem . . . . .	20
2.1.1	Linear model . . . . .	20
2.1.2	Task for this chapter . . . . .	22
2.2	Linear mapping . . . . .	23
2.2.1	Associated subspaces . . . . .	24
2.2.2	Orthogonality . . . . .	25
2.3	Matrix decompositions . . . . .	26
2.3.1	Eigen-decomposition . . . . .	26
2.3.2	Singular value decomposition . . . . .	28
2.3.3	QR decomposition . . . . .	32
2.4	Pseudoinverse . . . . .	32
2.4.1	Moore–Penrose pseudoinverse . . . . .	32
2.5	Orthogonal projection . . . . .	33
2.6	Application . . . . .	35
<b>3</b>	<b>Normal distribution</b>	<b>37</b>
3.1	The problem . . . . .	37
3.1.1	Linear model . . . . .	37
3.1.2	Task for this chapter . . . . .	38
3.2	Univariate normal distribution . . . . .	39
3.2.1	Definition . . . . .	39
3.2.2	Properties . . . . .	39
3.2.3	Related distributions . . . . .	40
3.3	Multivariate normal distribution . . . . .	42

3.3.1	Definition . . . . .	42
3.3.2	Properties . . . . .	43
3.3.3	Related distributions . . . . .	46
<b>4</b>	<b>Linear model</b>	<b>47</b>
4.1	The problem . . . . .	47
4.1.1	Linear model . . . . .	47
4.1.2	Task for this chapter . . . . .	48
4.2	Estimating $\beta$ . . . . .	49
4.2.1	Orthogonal projection . . . . .	49
4.2.2	Least squares . . . . .	50
4.2.3	Computing $\hat{\beta}$ . . . . .	51
4.3	Quality of estimation . . . . .	53
4.3.1	Gauss–Markov theorem . . . . .	53
4.4	Estimating $\sigma^2$ . . . . .	54
4.4.1	Estimating $\sigma^2$ . . . . .	54
4.5	Quality of model fit . . . . .	54
4.5.1	Coefficient of determination . . . . .	54
<b>5</b>	<b>Normal linear model</b>	<b>56</b>
5.1	The problem . . . . .	56
5.1.1	Normal linear model . . . . .	56
5.1.2	Task for this chapter . . . . .	57
5.2	Estimating $\beta$ and $\sigma^2$ . . . . .	58
5.2.1	Likelihood . . . . .	58
5.2.2	Matrix derivatives . . . . .	59
5.2.3	Maximizing the likelihood . . . . .	60
5.3	Distribution . . . . .	61
5.3.1	Distribution of the MLE . . . . .	61
5.4	Summary . . . . .	63
5.4.1	Estimation in the normal linear model . . . . .	63
<b>6</b>	<b>Inference in normal linear model</b>	<b>64</b>
6.1	The problem . . . . .	64
6.1.1	Normal linear model . . . . .	64
6.1.2	Task for this chapter . . . . .	65
6.2	Estimators & distributions . . . . .	66
6.2.1	Estimators . . . . .	66
6.2.2	Distributions . . . . .	66
6.3	Confidence intervals . . . . .	67
6.4	Prediction . . . . .	68
6.5	Confidence bands . . . . .	70
6.6	Testing hypotheses . . . . .	70

6.6.1	Simple hypothesis . . . . .	70
6.6.2	Composite hypothesis . . . . .	71
6.7	Interpretation . . . . .	72
<b>7</b>	<b>Model selection</b>	<b>75</b>
7.1	The problem . . . . .	75
7.1.1	Normal linear model . . . . .	75
7.1.2	Task for this chapter . . . . .	77
7.2	Why consider various models? . . . . .	77
7.2.1	Should we leave out covariates that appear unnecessary? . . . . .	77
7.2.2	What is the right form of the dependence on covariates? . . . . .	79
7.3	Nested models . . . . .	81
7.4	Selecting the model . . . . .	85
7.4.1	Model selection tools . . . . .	85
7.4.2	Model selection strategies . . . . .	87
<b>8</b>	<b>Model diagnostics</b>	<b>90</b>
8.1	The problem . . . . .	90
8.1.1	Normal linear model . . . . .	90
8.1.2	Task for this chapter . . . . .	91
8.2	Random errors and residuals . . . . .	92
8.3	Checking the assumptions . . . . .	94
8.3.1	General principles . . . . .	94
8.3.2	Assumptions on the expectation . . . . .	95
8.3.3	Assumptions on the variance . . . . .	97
8.3.4	Assumptions on the distribution . . . . .	102
8.4	Influential and unusual observations . . . . .	103
8.4.1	Observations to look at . . . . .	103
<b>9</b>	<b>Reduced-rank design matrix and multicollinearity</b>	<b>106</b>
9.1	The problem . . . . .	106
9.1.1	Normal linear model . . . . .	106
9.1.2	Task for this chapter . . . . .	108
9.2	Rank-deficient design matrix . . . . .	108
9.2.1	Rank-deficient design matrix . . . . .	108
9.2.2	Identifiability . . . . .	111
9.2.3	Choice of the solution . . . . .	112
9.3	Multicollinearity . . . . .	115
<b>10</b>	<b>Miscellanea and recap</b>	<b>119</b>
10.1	The problem . . . . .	119
10.1.1	Normal linear model . . . . .	119
10.1.2	Task for this chapter . . . . .	121

10.2	Linear regression in practice . . . . .	121
10.2.1	Linear regression in practice . . . . .	121
10.3	Notes on interpretation . . . . .	122
10.3.1	Notes on the explanation . . . . .	122
10.3.2	Notes on the prediction . . . . .	124
10.4	Transformations . . . . .	125
10.4.1	Transformations . . . . .	125
10.5	Concluding remarks . . . . .	126
10.5.1	Reflection . . . . .	126

# Chapter 1

## Introduction

### 1.1 Motivation

#### Is eating chocolate good for our health?

#### Effects of chocolate

- it has been suggested that chocolate consumption
  - ▷ is beneficial to cardiovascular health (effects on “bad” cholesterol, blood pressure, stroke, ...)
  - ▷ lowers the risk of diabetes
  - ▷ improves cognitive function & reduces memory decline
  - ▷ ...
- but it has also been suggested that chocolate consumption
  - ▷ leads to obesity (risk for cardiovascular problems, diabetes)
  - ▷ leads to dental problems
  - ▷ decreases bone density
  - ▷ ...
- should be eaten in moderation ...

#### It's an uncertain world ...

- How much of
  - ▷ chocolate and other goodies is good for our health?
  - ▷ levels of bacteria, fertilizers, chemicals, ... is safe?

- What is the right size for
  - ▷ the height of a dam?
  - ▷ insurance premium?
  - ▷ mortgage interest?
- What is
  - ▷ the average salary?
  - ▷ public opinion on ... ?
  - ▷ results in upcoming elections?

### Sources of uncertainty

- we do not fully understand the phenomenon
  - ▷ human body
  - ▷ nature
- we do not know the future
  - ▷ occurrence and size of a flood
  - ▷ occurrence and size of insurance claims
  - ▷ level of inflation
- we do not collect complete data
  - ▷ average salary
  - ▷ public opinion
- measurement error, human factor, ...

### Statistics is all around us

- statistics is used to quantify the uncertainty
- Strategy
  1. build a **mathematical** model, i.e. define
    - ▷ what is known
    - ▷ what is uncertain
  2. build a **probabilistic** model for what is uncertain
  3. use probability calculus to draw conclusions



- 4. “translate” back to the original problem (interpret the results)
- uncertainty at the beginning --> imperfect answers at the end
- statistics is used for quantifying uncertainty,  
not for getting rid of it

## Notation

- random variable  $X, Y$ 
  - ▷ (náhodná veličina)
- random vector/matrix  $\mathbf{X}, \mathbf{Y}$ 
  - ▷ (náhodný vektor/matice)
- density/probability mass function  $f$ 
  - ▷ (hustota/pravděpodobnostní funkce)
- parameters  $\theta, \beta$ , normal distribution  $N(\mu, \sigma^2)$ 
  - ▷ (parametry, normální rozdělení)
- expectation  $E X, E \mathbf{X}$ 
  - ▷ (střední hodnota)
- variance/covariance/variance-covariance matrix  
 $\text{Var } X, \text{Cov}(X, Y), \text{Var } \mathbf{X}$ 
  - ▷ (rozptyl/kovariance/kovarianční matice)

$$\text{Var } \mathbf{X} = \begin{pmatrix} \text{Var } X_1 & \text{Cov}(X_1, X_2) & \dots & \text{Cov}(X_1, X_n) \\ \text{Cov}(X_1, X_2) & \text{Var } X_2 & \dots & \text{Cov}(X_2, X_n) \\ \dots & \dots & \dots & \dots \\ \text{Cov}(X_1, X_n) & \text{Cov}(X_2, X_n) & \dots & \text{Var } X_n \end{pmatrix}$$

## Statistician's TODO list

1. identify right questions
2. collect relevant data  $x_1, \dots, x_n$
3. think of them as realisations of random variables  $X_1, \dots, X_n$   
with distributions (densities/frequency functions)  $f_1, \dots, f_n$

where  $f_i$  is in fact  $f_i(x, \theta)$

4. estimate  $\theta$ /make inference about  $\theta$
5. use the results to answer the questions

### Example

1. Does consuming [amount of] chocolate decrease blood pressure [type, measurement]?
  - ~~Is chocolate good for our health?~~
2. design a trial, collect participants' blood pressures  $x_1, \dots, x_n$
3. suppose e.g. that  $X_i \sim N(\mu_i, \sigma^2)$ 
  - $\mu_i$ : function of eating [amount of] chocolate, age, gender, ...
  - e.g.  $\mu_i = \beta_0 + \beta_1 x_{i,1} + \dots + \beta_k x_{i,k}$
  - $x_{i,1} = \begin{cases} 1 & \text{if the person eats [amount of] chocolate} \\ 0 & \text{otherwise} \end{cases}$
4. test  $H_0 : \beta_1 \geq 0$  versus  $H_1 : \beta_1 < 0$
5. if we reject  $H_0$  in favour of  $H_1$  at  $\alpha\%$  level, we have shown that at  $\alpha\%$  level consuming [amount of] chocolate is associated with a lower blood pressure [type, measurement]
  - if we do not reject  $H_0$  in favour of  $H_1$  at  $\alpha\%$  level, we have not shown that at  $\alpha\%$  level consuming [amount of] chocolate is associated with a lower blood pressure [type, measurement]

### Linear model

- model:  $Y_i = \beta_0 + \beta_1 x_{i,1} + \dots + \beta_k x_{i,k} + \varepsilon_i$ 
  - ▷ matrix notation:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$
  - ▷ assumptions:  $\mathbf{E}\boldsymbol{\varepsilon} = \mathbf{0}$ ,  $\text{Var}\boldsymbol{\varepsilon} = \sigma^2\mathbf{I}$ 
    - \* then  $\mathbf{E}\mathbf{Y} = \mathbf{X}\boldsymbol{\beta}$ ,  $\text{Var}\mathbf{Y} = \sigma^2\mathbf{I}$
  - ▷ we often assume that  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2\mathbf{I})$ 
    - \* then  $\mathbf{Y} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I})$
- parameter:  $\boldsymbol{\theta} = (\beta_0, \dots, \beta_k, \sigma^2)^\top = (\boldsymbol{\beta}^\top, \sigma^2)^\top$ 
  - ▷ estimation (point, interval)
  - ▷ testing
  - ▷ interpretation

## 1.2 Statistics

### Example

1. Does consuming [amount of] chocolate decrease blood pressure [type, measurement]?
  - Is chocolate good for our health?
2. design a trial, collect participants' blood pressures  $x_1, \dots, x_n$
3. suppose e.g. that  $X_i \sim N(\mu_i, \sigma^2)$ 
  - $\mu_i$ : function of eating [amount of] chocolate, age, gender, ...
  - $\mu_i = \beta_0 + \beta_1 x_{i,1} + \dots + \beta_k x_{i,k}$
  - $x_{i,1} = \begin{cases} 1 & \text{if the person eats [amount of] chocolate} \\ 0 & \text{otherwise} \end{cases}$
4. test  $H_0 : \beta_1 \geq 0$  versus  $H_1 : \beta_1 < 0$
5. if we reject  $H_0$  in favour of  $H_1$  at  $\alpha\%$  level, we have shown that at  $\alpha\%$  level consuming [amount of] chocolate is associated with a lower blood pressure [type, measurement]
  - if we do not reject  $H_0$  in favour of  $H_1$  at  $\alpha\%$  level, we have not shown that at  $\alpha\%$  level consuming [amount of] chocolate is associated with a lower blood pressure [type, measurement]

### Linear model

- model:  $Y_i = \beta_0 + \beta_1 x_{i,1} + \dots + \beta_k x_{i,k} + \varepsilon_i$ 
  - ▷ matrix notation:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$
  - ▷ assumptions:  $\mathbf{E}\boldsymbol{\varepsilon} = \mathbf{0}$ ,  $\text{Var}\boldsymbol{\varepsilon} = \sigma^2\mathbf{I}$ 
    - \* then  $\mathbf{E}\mathbf{Y} = \mathbf{X}\boldsymbol{\beta}$ ,  $\text{Var}\mathbf{Y} = \sigma^2\mathbf{I}$
  - ▷ we often assume that  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2\mathbf{I})$ 
    - \* then  $\mathbf{Y} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I})$
- parameter:  $\boldsymbol{\theta} = (\beta_0, \dots, \beta_k, \sigma^2)^\top = (\boldsymbol{\beta}^\top, \sigma^2)^\top$ 
  - ▷ estimation (point, interval)
  - ▷ testing
  - ▷ interpretation

### Parameter estimation

- we observe data with a distribution depending on a parameter
- we would like to use the data to **estimate** the value of the parameter
- **estimator** is a function of data (only!)
- for a one-dimensional parameter  $\theta$ 
  - ▷ point estimator  $\hat{\theta}$
  - ▷ confidence interval  $(\hat{\theta}_L, \hat{\theta}_U)$
- for a vector parameter  $\boldsymbol{\theta}$ 
  - ▷ point estimator  $\hat{\boldsymbol{\theta}}$
  - ▷ confidence region

## Methods of point estimation

### 1. method of moments

- “equate” theoretical and empirical moments
  - ▷  $\widehat{\mathbf{E}Y} = \frac{1}{n} \sum_{i=1}^n Y_i$
  - ▷  $\widehat{\mathbf{E}Y^2} = \frac{1}{n} \sum_{i=1}^n Y_i^2$
  - ▷ ...

### 2. maximum likelihood estimation

- maximize the likelihood with respect to  $\theta$
- likelihood
  - ▷ probability of observing the data at hand under a given model
- very popular thanks to certain asymptotic optimality properties

### 3. other methods exist and we will see some

## Maximum likelihood estimation

- $Y_1, \dots, Y_n \stackrel{\text{ind.}}{\sim} f_i(y, \boldsymbol{\theta})$
- **likelihood**  $L(y_1, \dots, y_n; \boldsymbol{\theta}) = \prod_{i=1}^n f_i(y_i; \boldsymbol{\theta})$
- **log-likelihood**  $\ell(y_1, \dots, y_n; \boldsymbol{\theta}) = \sum_{i=1}^n \log\{f_i(y_i; \boldsymbol{\theta})\}$
- **MLE**  $\hat{\boldsymbol{\theta}}_{\text{MLE}} = \operatorname{argmax}_{\boldsymbol{\theta}} \ell(y_1, \dots, y_n; \boldsymbol{\theta})$

- o usual computation

- ▷ score function  $\mathbf{U}(y_1, \dots, y_n; \boldsymbol{\theta}) = \sum_{i=1}^n \frac{\partial}{\partial \boldsymbol{\theta}} \log\{f_i(y_i; \boldsymbol{\theta})\}$

- ▷ score equation  $\mathbf{U}(y_1, \dots, y_n; \boldsymbol{\theta}) = \mathbf{0}$

- ▷ find the solution  $\hat{\boldsymbol{\theta}}_{\text{MLE}}$  of the score equation

- ▷ observed Fisher information matrix  $\mathbf{J}(y_1, \dots, y_n; \boldsymbol{\theta}) = -\sum_{i=1}^n \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}} \log\{f_i(y_i; \boldsymbol{\theta})\}$

- ▷ show that  $\mathbf{J}(y_1, \dots, y_n; \hat{\boldsymbol{\theta}}_{\text{MLE}})$  is positive definite

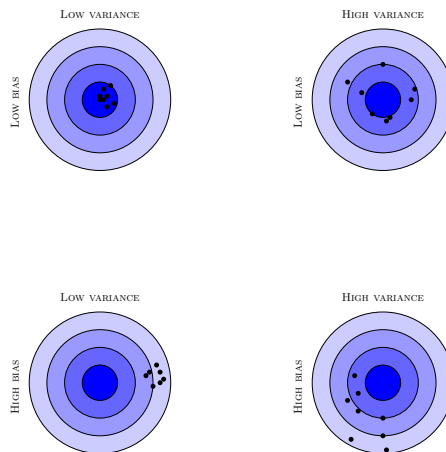
- ▷ Fisher information matrix (under regularity conditions)  $\mathbf{I}(y_1, \dots, y_n; \boldsymbol{\theta}) = \mathbf{E}_{\boldsymbol{\theta}} \mathbf{J}(y_1, \dots, y_n; \boldsymbol{\theta})$

## Properties of estimators

- o parameter  $\theta$  is a number but estimator  $\hat{\theta}$  is a random variable

- ▷  $\hat{\theta}$  has a distribution

- ▷ important distribution summaries:  $\mathbf{E} \hat{\theta}, \text{Var} \hat{\theta}$



- o ideally, estimation improves with sample size

- o let  $\hat{\theta}_n$  be an estimator of  $\theta$  based on  $n$  data points

- o we define desirable properties for the sequence  $\{\hat{\theta}_n\}_{n \in \mathbb{N}}$

- o for a sequence of estimators  $\hat{\theta}_n$  of a parameter  $\theta$

1. unbiasedness

- ▷  $\mathbf{E}_{\boldsymbol{\theta}} \hat{\theta}_n = \boldsymbol{\theta} \forall \boldsymbol{\theta}$

2. consistency

- ▷  $\hat{\theta}_n \rightarrow \boldsymbol{\theta}$  as  $n \rightarrow \infty$  in  $\mathbf{P}_{\boldsymbol{\theta}} \forall \boldsymbol{\theta}$  or a.s.

3. usual asymptotic normality

▷  $\sqrt{n}(\hat{\theta}_n - \theta) \rightarrow N(0, V(\theta))$  as  $n \rightarrow \infty$  in distribution

4. efficiency

▷ “small”  $\text{Var } \hat{\theta}$

**Properties of MLE**○ under regularity conditions

▷ consistency

\*  $\hat{\theta}_{MLE,n} \rightarrow \theta$  a.s. as  $n \rightarrow \infty$

▷ asymptotic normality

\*  $\sqrt{n}(\hat{\theta}_{MLE,n} - \theta) \rightarrow N(0, V(\theta))$  as  $n \rightarrow \infty$  in distribution

▷ asymptotic efficiency

\*  $V(\theta)$  is the smallest possible

▷ bias

\*  $\hat{\theta}_{MLE}$  is often biased, with bias decreasing with  $n$

**Interval estimation**

○ parameter  $\theta$  is a number but estimator  $\hat{\theta}$  is a random variable

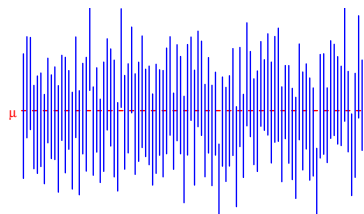
○ confidence interval  $(\hat{\theta}_L, \hat{\theta}_U)$  is a pair of random variables

○  $(1 - \alpha)$  % confidence interval satisfies that

▷  $P_{\theta}\{\theta \in (\hat{\theta}_L, \hat{\theta}_U)\} = 1 - \alpha \forall \theta$

○ note that randomness is in the borders, not in  $\theta$

Confidence interval for  $\mu$  in  $N(\mu, \sigma^2)$  with  $\sigma^2$  unknown



○ properties

▷ coverage  $1 - \alpha$

▷ length  $\hat{\theta}_U - \hat{\theta}_L$

▷ ideally: a short interval with high coverage

## Testing hypotheses

- we observe data with a distribution depending on a parameter
- we would like to use the data to answer questions about the parameter
  - ▷ is  $\theta > 0$ ,  $\theta < 0$ ,  $\theta = 1$ , ...?
- to do so, we can test hypotheses about the parameter
  - ▷  $H_0 : \theta \geq 0$  vs.  $H_1 : \theta < 0$
  - ▷  $H_0 : \theta = 1$  vs.  $H_1 : \theta \neq 1$
  - ▷ ...
- testing has two possible results
  1. we reject  $H_0$  in favour of  $H_1$ 
    - ▷ we can say we have shown  $H_1$  (at the level  $\alpha$ )
  2. we do not reject  $H_0$  in favour of  $H_1$ 
    - ▷ we can say we have not shown  $H_1$  (at the level  $\alpha$ )
    - ▷ !!!we cannot say we have shown  $H_0$ !!!
- the roles of  $H_0$  and  $H_1$  are not symmetric
- testing has two possible results
  1. we reject  $H_0$  in favour of  $H_1$
  2. we do not reject  $H_0$  in favour of  $H_1$
- we can reach a wrong conclusion in two ways
  1. when  $H_0$  is true and we reject  $H_0$  in favour of  $H_1$ 
    - ▷ “ $\mathbf{P}_{H_0}(\text{reject } H_0) = \alpha$ ”
    - ▷  $\alpha$ : type I. error, level of the test
  2. when  $H_1$  is true and we do not reject  $H_0$  in favour of  $H_1$ 
    - ▷ “ $1 - \mathbf{P}_{H_1}(\text{reject } H_0) = \beta$ ”
    - ▷  $\beta$ : type II. error
    - ▷  $1 - \beta$ : power of the test
- often impossible to keep both errors low at the same time

- when choosing a test, we keep the level  $\alpha$  fixed and try to maximize the power  $\beta$
- the roles of  $H_0$  and  $H_1$  are not symmetric
- the roles of  $H_0$  and  $H_1$  are not symmetric
- it is important to choose a good  $H_0, H_1$  pair
- testing
  - ▷  $H_0 : \theta \geq 0$  vs.  $H_1 : \theta < 0$
  - ▷  $H_0 : \theta \leq 0$  vs.  $H_1 : \theta > 0$
  - ▷  $H_0 : \theta = 0$  vs.  $H_1 : \theta \neq 0$

answer different questions

## 1.3 Data analysis in practice

### Example: fev data

- from: <http://www.statsci.org/data/general/fev.html>
- question: association between the FEV[l] and Smoking, corrected for Age[years], Height[cm] and Gender

	FEV	Age	Height	Gender	Smoking
	1.708	9	144.8	Female	Non
	1.724	8	171.5	Female	Non
	1.720	7	138.4	Female	Non
○ data:	1.558	9	134.6	Male	Non
	...	...	...	...	...
	3.727	15	172.7	Male	Current
	2.853	18	152.4	Female	Non
	2.795	16	160.0	Female	Current
	3.211	15	168.9	Female	Non

### Getting the data to R

- data `fev.txt`
- for `*.txt` files:
  - ▷ `read.table(...)`
- for `*.csv` files (from Excel)



```

  ▷ read.csv(...)

○ read the data and look at them

> fev <- read.table("fev.txt", header=TRUE)
>
> class(fev)
[1] "data.frame"
> dim(fev)
[1] 654  6
> names(fev)
[1] "ID"      "Age"      "FEV"      "Height" "Sex"      "Smoker"
>
> fev[1:3, ]
  ID Age  FEV Height  Sex Smoker
1 301  9 1.708  57.0 Female   Non
2 451  8 1.724  67.5 Female   Non
3 501  7 1.720  54.5 Female   Non
>
> fev <- fev[, -1]

```

## Before we fit a model to data

- before we do the analysis, we need to
  - ▷ get to know the variables
  - ▷ get to understand the relationships among the variables
  - ▷ identify possible problems for the analysis
  - ▷ possibly spot obvious mistakes in data
- ⇒ first step in applied data analysis: *descriptive statistics*
  - ▷ informal data descriptions (no model, no inference)
    - \* numerical and graphical
    - \* their choice depends on the type of variable(s) of interest

## First look at the variables

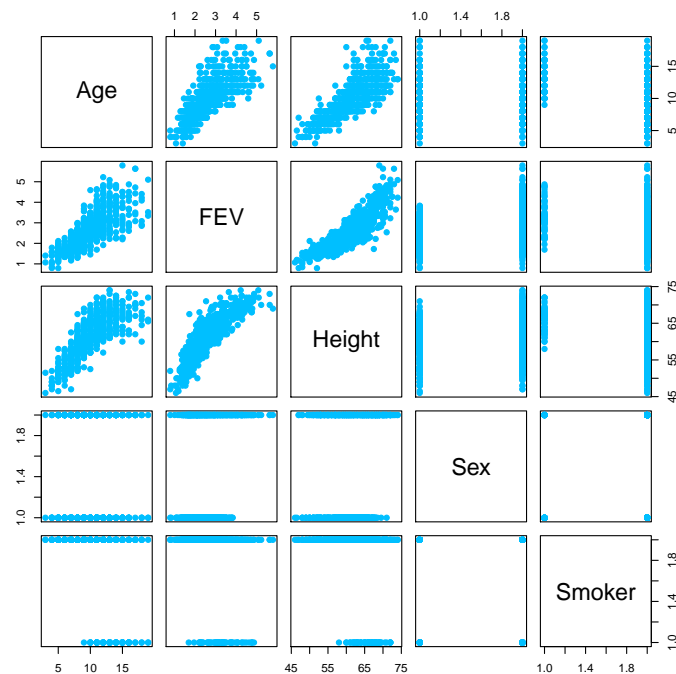
```

> summary(fev)
  Age              FEV              Height              Sex
Min.   : 3.000    Min.   :0.791    Min.   :46.00    Female:318
1st Qu.: 8.000    1st Qu.:1.981    1st Qu.:57.00    Male  :336
Median :10.000    Median :2.547    Median :61.50
Mean   : 9.931    Mean   :2.637    Mean   :61.14
3rd Qu.:12.000    3rd Qu.:3.119    3rd Qu.:65.50
Max.   :19.000    Max.   :5.793    Max.   :74.00
  Smoker
Current: 65
Non     :589

```

## First look at the relationships between the variables


```
> pairs(fev, col="deepskyblue", pch=19)
```



## 1.4 Descriptive statistics

### 1.4.1 Types of variables

#### Types of variables

1. in mathematical statistics:
  - continuous (uncountably many possible values)
  - discrete (at most countably many possible values)
  
2. in applied statistics:
  - quantitative
  - categorical
    - ▷ nominal
    - ▷ ordinal
  
3. in :

- numeric
- ▷ factor
  - ▷ ordered factor

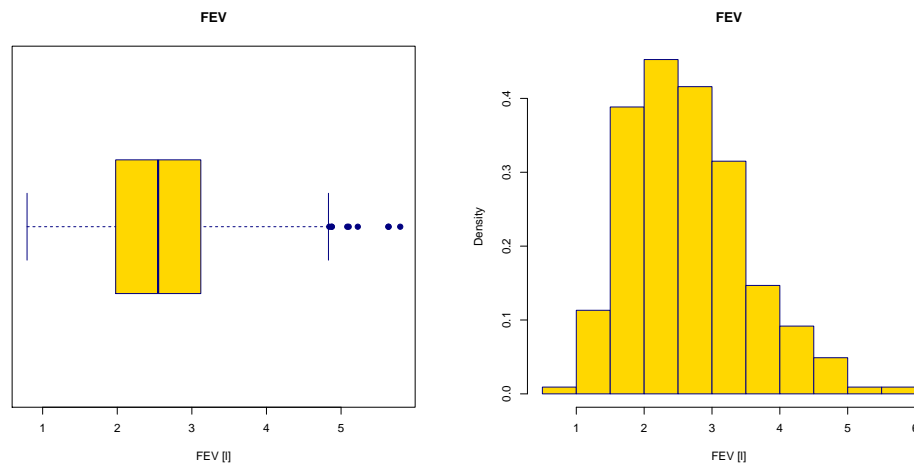
## Quantitative variable

- distribution
  - characteristics of location
    - ▷ mean
    - ▷ maximum, minimum
    - ▷ quantiles, in particular quartiles and median
      - > summary(fev\$FEV)

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
0.791	1.981	2.548	2.637	3.118	5.793
- characteristics of dispersion
  - ▷ standard deviation
  - ▷ interquartile range
    - > sd(fev\$FEV)
    - [1] 0.8670591
    - > IQR(fev\$FEV)
    - [1] 1.1375

## Graphics for quantitative variable

```
> hist(fev$FEV, , freq=FALSE,
+ main="FEV", xlab="FEV [1]",
+ col="gold", border="navyblue")
>
> boxplot(fev$FEV, horizontal=TRUE,
+ main="FEV", xlab="FEV [1]",
+ col="gold", border="navyblue", pch=19)
```



## Categorical variable

- distribution
  - ▷ counts of observations per category
  - ▷ percentage of observations per category
  - ▷ cumulative percentage of observations per category (for ordinal variables)
- characteristics
  - ▷ modus

```
> summary(fev$Sex)
Female  Male
  318   336
> prop.table(table(fev$Sex))
  Female    Male
0.4862385 0.5137615
> cumsum(prop.table(table(fev$Sex)))
  Female    Male
0.4862385 1.0000000
> # not so interesting for a nominal variable
```

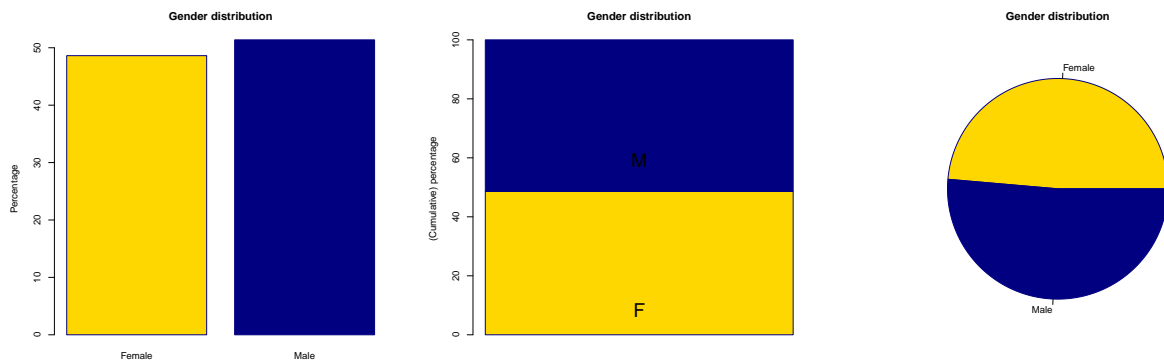
## Graphics for categorical variable

```
> barplot(100*prop.table(table(fev$Sex)),
+         main="Gender distribution", ylab="Percentage",
+         col=c("gold", "navyblue"), border="navyblue")
>
> barplot(100*matrix(prop.table(table(fev$Sex)), nrow=2, ncol=1),
+         main="Gender distribution", ylab="(Cumulative) percentage",
+         col=c("gold", "navyblue"), border="navyblue")
```

```

> text(x=0.7, y=4, labels="F",
+      cex=2, pos=3)
> text(x=0.7, y=55, labels="M",
+      cex=2, pos=3)
>
> pie(summary(fev$Sex),
+     main="Gender distribution",
+     col=c("gold", "navyblue"), border="navyblue")

```



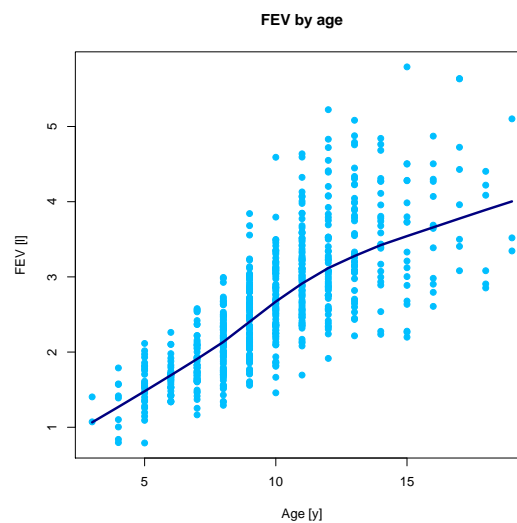
## 1.4.2 Relationships between variables

### Quantitative vs quantitative

```

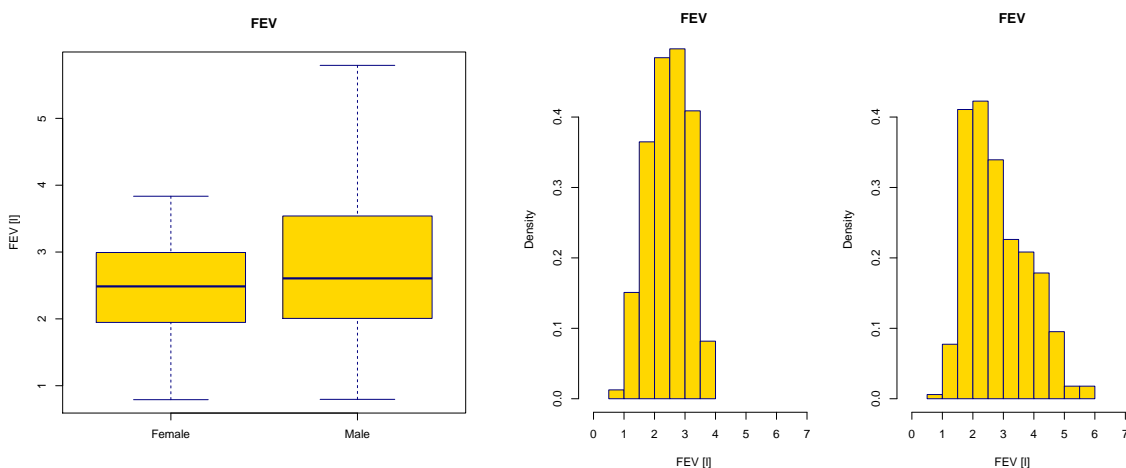
> plot(fev$FEV~fev$Age,
+      main="FEV by age",
+      ylab="FEV [l]", xlab="Age [y]",
+      pch=19, col="deepskyblue")
> lines(lowess(fev$FEV~fev$Age),
+       lwd=3, col="navyblue")

```



## Quantitative vs categorical

```
> boxplot(fev$FEV~fev$Sex,
+         main="FEV", ylab="FEV [l]",
+         col="gold", border="navyblue", pch=19)
>
> par(mfrow=c(1, 2))
> hist(fev$FEV[fev$Sex=="Female"], freq=FALSE,
+      xlim=c(min(fev$FEV)-1, max(fev$FEV)+1), ylim=c(0, 0.48),
+      main="FEV", xlab="FEV [l]",
+      col="gold", border="navyblue")
> hist(fev$FEV[fev$Sex=="Male"], freq=FALSE,
+      xlim=c(min(fev$FEV)-1, max(fev$FEV)+1), ylim=c(0, 0.48),
+      main="FEV", xlab="FEV [l]",
+      col="gold", border="navyblue")
```



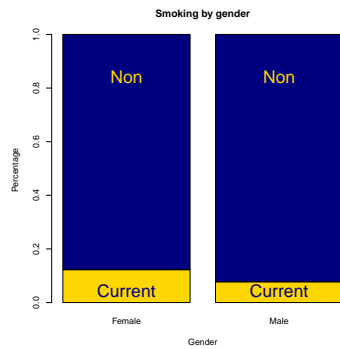
## Categorical vs categorical

```
> table(fev$Smoker, fev$Sex)
      Female Male
Current   39   26
Non      279  310
>
> prop.table(table(fev$Smoker, fev$Sex),
+             margin=1
+             )
      Female      Male
Current 0.600000 0.400000
Non     0.473684 0.5263158
>
> prop.table(table(fev$Smoker, fev$Sex),
+             margin=2
+             )
      Female      Male
Current 0.1226415 0.07738095
Non     0.8773584 0.92261905
```

```

> prop.table(table(fev$Smoker, fev$Sex),
+           margin=2
+           )
           Female      Male
Current 0.12264151 0.07738095
Non     0.87735849 0.92261905
>
> barplot(height=prop.table(table(fev$Smoker, fev$Sex),
+                               margin=2), beside=F,
+         ylab="Percentage", xlab="Gender", main="Smoking by gender",
+         col=c("gold", "navyblue")
+         )
> text(x=1.2*c(0:2)+0.7, y=0, labels="Current",
+      col="navyblue", cex=2, pos=3)
> text(x=1.2*c(0:2)+0.7, y=0.8, labels="Non",
+      col="gold", cex=2, pos=3)

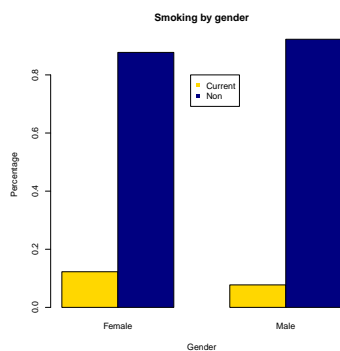
```



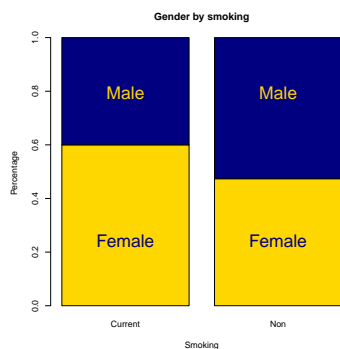
```

> prop.table(table(fev$Smoker, fev$Sex),
+           margin=2
+           )
           Female      Male
Current 0.12264151 0.07738095
Non     0.87735849 0.92261905
>
> barplot(height=prop.table(table(fev$Smoker, fev$Sex),
+                               margin=2), beside=T,
+         ylab="Percentage", xlab="Gender", main="Smoking by gender",
+         col=c("gold", "navyblue")
+         )
>
> legend(x=3.3, y=0.8, legend=c("Current", "Non"),
+       col=c("gold", "navyblue"), pch=15)

```



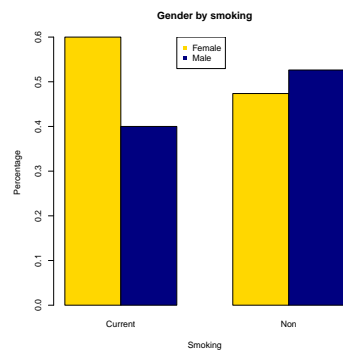
```
> prop.table(table(fev$Sex, fev$Smoker),
+             margin=2
+ )
      Current      Non
Female 0.6000000 0.4736842
Male   0.4000000 0.5263158
>
> barplot(height=prop.table(table(fev$Sex, fev$Smoker),
+                               margin=2), beside=F,
+         ylab="Percentage", xlab="Smoking", main="Gender by smoking",
+         col=c("gold", "navyblue")
+ )
> text(x=1.2*c(0:2)+0.7, y=0.2, labels="Female", col="navyblue",
+      cex=2, pos=3)
> text(x=1.2*c(0:2)+0.7, y=0.75, labels="Male", col="gold",
+      cex=2, pos=3)
```



```
> prop.table(table(fev$Sex, fev$Smoker),
+             margin=2
+ )
      Current      Non
Female 0.6000000 0.4736842
Male   0.4000000 0.5263158
>
> barplot(height=prop.table(table(fev$Sex, fev$Smoker),
```



```
+           margin=2), beside=T,  
+         ylab="Percentage", xlab="Smoking", main="Gender by smoking",  
+         col=c("gold", "navyblue")  
+       )  
>  
> legend(x=3, y=0.6, legend=c("Female", "Male"),  
+       col=c("gold", "navyblue"), pch=15)
```



# Chapter 2

## Linear algebra essentials

### 2.1 The problem

#### 2.1.1 Linear model

##### Chocolatey example

- Does consuming [amount of] chocolate decrease blood pressure [type, measurement]?
- collect blood pressures  $x_1, \dots, x_n$
- suppose that  $X_i \sim N(\mu_i, \sigma^2)$ 
  - ▷  $\mu_i = \beta_0 + \beta_1 x_{i,1} + \dots + \beta_k x_{i,k}$
  - ▷  $x_{i,1} = \begin{cases} 1 & \text{if the person eats [amount of] chocolate} \\ 0 & \text{otherwise} \end{cases}$
  - ▷  $x_{i,j}, j \in \{2, \dots, k\}$ : age, gender, BMI, ...
- test  $H_0 : \beta_1 \geq 0$  versus  $H_1 : \beta_1 < 0$  to answer the question

##### Linear model

- $Y_i = \beta_0 + \beta_1 x_{i,1} + \dots + \beta_k x_{i,k} + \varepsilon_i, i \in \{1, \dots, n\}$ 
  - ▷  $Y_i$ : outcome, response, output, dependent variable
    - \* random variable, we observe a realization  $y_i$
    - \* (odezva, závisle proměnná, regresand)
  - ▷  $x_{i,1}, \dots, x_{i,k}$ : covariates, predictors, explanatory variables,  
input, independent variables
    - \* given, known

- \* (nezávisle proměnné, regresory)
- ▷  $\beta_0, \dots, \beta_k$ : coefficients
  - \* unknown
  - \* (regresní koeficienty)
- ▷  $\varepsilon_i$ : random error
  - \* random variable, unobserved
- $\varepsilon_i \stackrel{\text{iid}}{\sim} (0, \sigma^2)$ ,  $i \in \{1, \dots, n\}$ 
  - ▷  $\mathbb{E} \varepsilon_i = 0$ : no systematic errors
  - ▷  $\text{Var} \varepsilon_i = \sigma^2$ : same precision
- we often assume that  $\varepsilon_i \stackrel{\text{iid}}{\sim} \text{N}(0, \sigma^2)$ ,  $i \in \{1, \dots, n\}$

### Linear model in the matrix form

- $Y_i = \beta_0 + \beta_1 x_{i,1} + \dots + \beta_k x_{i,k} + \varepsilon_i$ ,  $i \in \{1, \dots, n\}$

- let

$$\mathbf{Y} = \begin{pmatrix} Y_1 \\ Y_2 \\ \dots \\ Y_n \end{pmatrix}, \mathbf{X} = \begin{pmatrix} 1 & x_{1,1} & \dots & x_{1,k} \\ 1 & x_{2,1} & \dots & x_{2,k} \\ \dots & \dots & \dots & \dots \\ 1 & x_{n,1} & \dots & x_{n,k} \end{pmatrix}, \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \dots \\ \beta_k \end{pmatrix}, \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \dots \\ \varepsilon_n \end{pmatrix}$$

- then  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon} \sim (\mathbf{0}, \sigma^2 \mathbf{I})$  and often  $\boldsymbol{\varepsilon} \sim \text{N}(\mathbf{0}, \sigma^2 \mathbf{I})$

- ▷  $\mathbf{X}$ : design matrix

- \* (regresní matice, matice plánu)

- let  $p = k + 1$

- then  $\underbrace{\mathbf{Y}}_{n \times 1} = \underbrace{\mathbf{X}}_{n \times p} \underbrace{\boldsymbol{\beta}}_{p \times 1} + \underbrace{\boldsymbol{\varepsilon}}_{n \times 1}$

- we assume that  $n > p$  (and often think about  $n \rightarrow \infty$ ,  $p$  fixed)

### Example: bloodpress data

- from [sites.stat.psu.edu/~lsimon/stat501wc/sp05/data/](https://sites.stat.psu.edu/~lsimon/stat501wc/sp05/data/)

- association between the mean arterial blood pressure[mmHg] and age[years], weight[kg], body surface area[m<sup>2</sup>], duration of hypertension[years], basal pulse[beats/min], stress

◦ data:

BP	Age	Weight	BSA	DoH	Pulse	Stress
105	47	85.4	1.75	5.1	63	33
115	49	94.2	2.10	3.8	70	14
...	...	...	...	...	...	...
110	48	90.5	1.88	9.0	71	99
122	56	95.7	2.09	7.0	75	99

◦ model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$

$$\begin{pmatrix} 105 \\ 115 \\ \dots \\ 110 \\ 122 \end{pmatrix} = \begin{pmatrix} 1 & 47 & 85.4 & 1.75 & 5.1 & 63 & 33 \\ 1 & 49 & 94.2 & 2.10 & 3.8 & 70 & 14 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & 48 & 90.5 & 1.88 & 9.0 & 71 & 99 \\ 1 & 56 & 95.7 & 2.09 & 7.0 & 75 & 99 \end{pmatrix} \times \begin{pmatrix} \beta_0 \\ \dots \\ \beta_6 \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \dots \\ \varepsilon_{19} \\ \varepsilon_{20} \end{pmatrix}$$

## 2.1.2 Task for this chapter

### Design matrix

◦ model:

$$\begin{pmatrix} Y_1 \\ Y_2 \\ \dots \\ Y_n \end{pmatrix} = \begin{pmatrix} 1 & x_{1,1} & \dots & x_{1,k} \\ 1 & x_{2,1} & \dots & x_{2,k} \\ \dots & \dots & \dots & \dots \\ 1 & x_{n,1} & \dots & x_{n,k} \end{pmatrix} \times \begin{pmatrix} \beta_0 \\ \dots \\ \beta_k \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \dots \\ \varepsilon_n \end{pmatrix}$$

◦ design matrix:

$$\mathbf{X} = \begin{pmatrix} 1 & x_{1,1} & \dots & x_{1,k} \\ 1 & x_{2,1} & \dots & x_{2,k} \\ \dots & \dots & \dots & \dots \\ 1 & x_{n,1} & \dots & x_{n,k} \end{pmatrix} = (\mathbf{1} \mid \mathbf{x}_{,1} \mid \mathbf{x}_{,2} \mid \dots \mid \mathbf{x}_{,k}) = \begin{pmatrix} \mathbf{x}_1, \\ \mathbf{x}_2, \\ \dots \\ \mathbf{x}_n, \end{pmatrix}$$

▷  $k$  covariates and  $\mathbf{1}$  are the  $p$  columns of  $\mathbf{X}$

▷  $n$  observations are the  $n$  rows of  $\mathbf{X}$

### Matrix algebra in a linear model

◦ model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$

◦ coefficient vector  $\boldsymbol{\beta}$

▷ fixed but **unknown**

▷  $p \times 1$  matrix

▷  $\boldsymbol{\beta}^\top$  defines a mapping  $\boldsymbol{\beta}^\top : \mathbb{R}^p \mapsto \mathbb{R}$

$$\mathbf{x}_i \in \mathbb{R}^p \rightsquigarrow \mathbf{E} \mathbf{Y}_i \in \mathbb{R}$$

- design matrix  $\mathbf{X}$

- ▷ fixed and known
- ▷  $n \times p$  matrix
- ▷ defines a mapping  $\mathbf{X} : \mathbb{R}^p \mapsto \mathbb{R}^n$

$$\boldsymbol{\beta} \in \mathbb{R}^p \rightsquigarrow \mathbf{E} \mathbf{Y} \in \mathbb{R}^n$$

- ▷ idea: when estimating  $\boldsymbol{\beta}$ , how about choosing  $\hat{\boldsymbol{\beta}}$  so that  $\mathbf{X}$  maps  $\hat{\boldsymbol{\beta}}$  as close to  $\mathbf{Y}$  as possible?

## 2.2 Linear mapping

### Linear mapping from $\mathbb{R}^p$ to $\mathbb{R}^n$

- function  $f : \mathbb{R}^p \mapsto \mathbb{R}^n$  such that
  - ▷  $f(\mathbf{x} + \mathbf{y}) = f(\mathbf{x}) + f(\mathbf{y}) \dots$  additivity
  - ▷  $f(\alpha \mathbf{x}) = \alpha f(\mathbf{x}) \dots$  homogeneity
- described by an  $n \times p$  matrix  $\mathbf{A}$ :  $f(\mathbf{x}) = \mathbf{A}\mathbf{x}$

↔ idea:

- ▷  $\forall \mathbf{x} \in \mathbb{R}^p$  can be written as  $\mathbf{x} = \sum_{i=1}^p c_i \mathbf{v}_i$ ,  
 where  $\mathcal{V} = \{\mathbf{v}_1, \dots, \mathbf{v}_p\}$  is a basis of  $\mathbb{R}^p$

↔  $f(\mathbf{x})$  is determined by  $\{f(\mathbf{v}_1), \dots, f(\mathbf{v}_p)\}$

$$\text{because } f(\mathbf{x}) = f\left(\sum_{i=1}^p c_i \mathbf{v}_i\right) = \sum_{i=1}^p c_i f(\mathbf{v}_i)$$

- ▷  $\forall \mathbf{y} \in \mathbb{R}^n$  can be written as  $\mathbf{y} = \sum_{i=1}^n c_i \mathbf{w}_i$ ,  
 where  $\mathcal{W} = \{\mathbf{w}_1, \dots, \mathbf{w}_n\}$  is a basis of  $\mathbb{R}^n$

↔ just need to write each  $f(\mathbf{v}_i)$  in terms of  $\mathcal{W}$

- ▷ free choice of  $(\mathcal{W}, \mathcal{V}) \rightarrow$  various  $\mathbf{A}$ 's representing the same  $f$

- operations  $f_1 \circ f_2, f_1 + f_2, \alpha f$

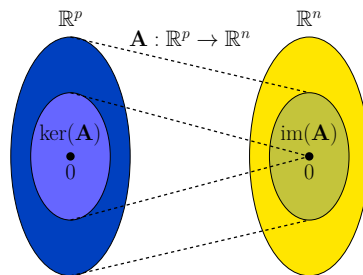
- ▷ result into linear mappings

$$\text{represented by } \mathbf{A}_1 \mathbf{A}_2, \mathbf{A}_1 + \mathbf{A}_2, \alpha \mathbf{A}$$

## 2.2.1 Associated subspaces

### Kernel and image

- kernel (nullspace)
  - ▷  $\ker(\mathbf{A}) = \ker(f) = \{\mathbf{x} \in \mathbb{R}^p; \mathbf{A}\mathbf{x} = \mathbf{0}\}$   
 $= \{\mathbf{x} \in \mathbb{R}^p; f(\mathbf{x}) = \mathbf{0}\}$
  - ▷ subspace of  $\mathbb{R}^p$ ,  $\dim(\ker(\mathbf{A}))$ : nullity of  $\mathbf{A}$
- image (range, column space)
  - ▷  $\text{im}(\mathbf{A}) = \text{im}(f) = \{\mathbf{y} \in \mathbb{R}^n; \exists \mathbf{x} \in \mathbb{R}^p : \mathbf{A}\mathbf{x} = \mathbf{y}\}$   
 $= \{\mathbf{y} \in \mathbb{R}^n; \exists \mathbf{x} \in \mathbb{R}^p : f(\mathbf{x}) = \mathbf{y}\}$
  - ▷ subspace of  $\mathbb{R}^n$ ,  $\dim(\text{im}(\mathbf{A}))$ : rank of  $\mathbf{A}$
- schematically



- rank nullity theorem:  $\dim(\ker(\mathbf{A})) + \dim(\text{im}(\mathbf{A})) = p$

### Four fundamental subspaces associated to $\mathbf{A}$

- kernel and image of  $\mathbf{A}$ 
  - ▷ column space of  $\mathbf{A}$ :  $\text{im}(\mathbf{A}) = \{\mathbf{y} \in \mathbb{R}^n; \exists \mathbf{x} \in \mathbb{R}^p : \mathbf{A}\mathbf{x} = \mathbf{y}\}$   
 \*  $\dim(\text{im}(\mathbf{A})) = \text{rank}(\mathbf{A})$
  - ▷ kernel of  $\mathbf{A}$ :  $\ker(\mathbf{A}) = \{\mathbf{x} \in \mathbb{R}^p; \mathbf{A}\mathbf{x} = \mathbf{0}\}$   
 \*  $\dim(\ker(\mathbf{A})) = p - \text{rank}(\mathbf{A})$
- kernel and image of  $\mathbf{A}^\top$ 
  - ▷ column space of  $\mathbf{A}^\top$ :  $\text{im}(\mathbf{A}^\top)$ : coimage of  $\mathbf{A}$   
 \*  $\dim(\text{im}(\mathbf{A}^\top)) = \text{rank}(\mathbf{A}^\top) = \text{rank}(\mathbf{A})$   
 \* row space of  $\mathbf{A}$
  - ▷ kernel of  $\mathbf{A}^\top$ :  $\ker(\mathbf{A}^\top)$ : cokernel, left nullspace of  $\mathbf{A}$   
 \*  $\dim(\ker(\mathbf{A}^\top)) = n - \text{rank}(\mathbf{A})$ : corank of  $\mathbf{A}$

## 2.2.2 Orthogonality

### Inner product on $\mathbb{R}^n$

◦ dot product:

$$\triangleright \langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{y}^\top \mathbf{x} = \sum_{i=1}^n x_i y_i$$

◦ associated norm (length of  $\mathbf{x}$ ):

$$\triangleright \|\mathbf{x}\| = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle} = \sqrt{\sum_{i=1}^n x_i^2}$$

◦ angle  $\theta$  between  $\mathbf{x}$  and  $\mathbf{y}$ :

$$\triangleright \cos \theta = \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\|\mathbf{x}\| \|\mathbf{y}\|}$$

◦ orthogonality for  $\mathbf{x} \neq \mathbf{0}$  and  $\mathbf{y} \neq \mathbf{0}$ :

$$\triangleright \langle \mathbf{x}, \mathbf{y} \rangle = 0$$

◦ orthogonal complement  $W^\perp$  of a subspace  $W$  of  $\mathbb{R}^n$ :

$$\triangleright W^\perp = \{\mathbf{y} \in \mathbb{R}^n; \langle \mathbf{x}, \mathbf{y} \rangle = 0 \text{ for every } \mathbf{x} \in W\}$$

\*  $W^\perp$  is a subspace of  $\mathbb{R}^n$

\*  $W^\perp \cap W = \{\mathbf{0}\}$

\*  $\dim(W) + \dim(W^\perp) = n$

◦ orthogonality between fundamental subspaces associated to  $\mathbf{A}$

$$\triangleright \ker(\mathbf{A}) = (\text{im}(\mathbf{A}^\top))^\perp \text{ (in } \mathbb{R}^p)$$

$$\triangleright \ker(\mathbf{A}^\top) = (\text{im}(\mathbf{A}))^\perp \text{ (in } \mathbb{R}^n)$$

### Orthogonal columns

◦ matrix with orthogonal columns:

$$\triangleright \mathbf{U} = (\mathbf{u}_1 | \mathbf{u}_2 | \dots | \mathbf{u}_p)$$

$$\langle \mathbf{u}_i, \mathbf{u}_j \rangle = 0 \text{ for } i \neq j$$

◦ matrix with orthonormal columns:

$$\triangleright \mathbf{U} = (\mathbf{u}_1 | \mathbf{u}_2 | \dots | \mathbf{u}_p)$$

$$\langle \mathbf{u}_i, \mathbf{u}_j \rangle = 0 \text{ for } i \neq j$$

$$\|\mathbf{u}_i\| = 1 \text{ for } i \in \{1, \dots, p\}$$

$$\triangleright \mathbf{U}^\top \mathbf{U} = \mathbf{I} \Rightarrow \text{mapping } \mathbf{U} : \mathbf{x} \mapsto \mathbf{U}\mathbf{x} \text{ preserves}$$

- \* inner product
- \* norm
- \* angles
- \* distances

### Orthogonal matrix

- o square matrix  $\mathbf{R}$  with orthonormal columns (and rows)

▷  $\mathbf{R}^\top \mathbf{R} = \mathbf{R} \mathbf{R}^\top = \mathbf{I}$

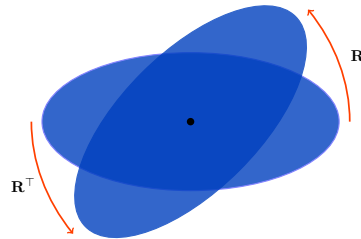
i.e.  $\mathbf{R}^\top = \mathbf{R}^{-1}$

- o  $\mathbf{R}^{-1} = \mathbf{R}^\top$  is also an orthogonal matrix
- o product of orthogonal matrices is also an orthogonal matrix
- o geometrically

- ▷ change of orthonormal basis (coordinate transformation)
- ▷ mapping  $\mathbf{R}: \mathbf{x} \mapsto \mathbf{R}\mathbf{x} \dots$  rotation

- \* preserves the origin
- \* preserves angles
- \* preserves distances
- \* proper rotation

if  $\det \mathbf{R} = 1$

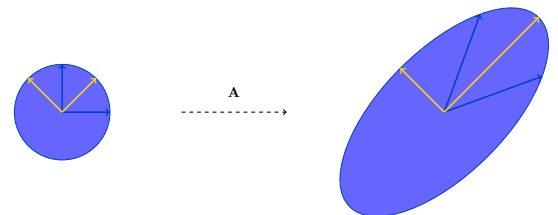


## 2.3 Matrix decompositions

### 2.3.1 Eigen-decomposition

#### Spectral decomposition (eigen-decomposition)

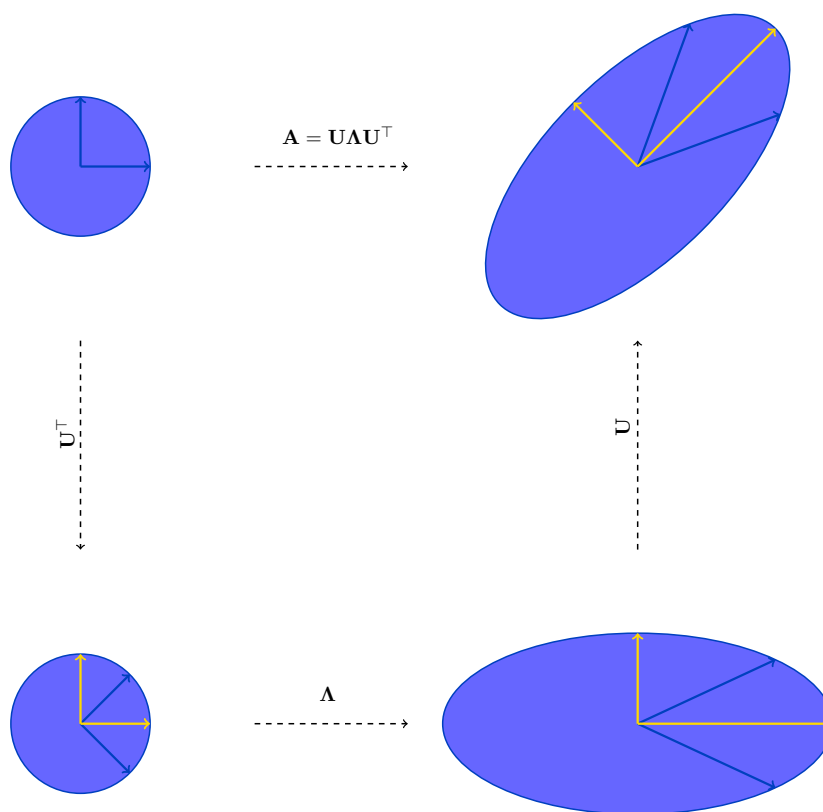
- o let  $\mathbf{A}$  be a symmetric  $p \times p$  matrix
  - ▷ eigenvalues  $\lambda_1, \dots, \lambda_p$
  - ▷ eigenvectors  $\mathbf{u}_{.,1}, \dots, \mathbf{u}_{.,p}$
  - ▷ everything real
  - ▷  $\mathbf{A}\mathbf{u}_{.,i} = \lambda_i \mathbf{u}_{.,i}$
  - ▷  $f$  elongates/shrinks  $\mathbf{u}_{.,i}$  by  $\lambda_i$





- eigen-decomposition:  $\mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top$ , where
  - ▷  $\mathbf{U} = (\mathbf{u}_{\cdot,1} \mid \mathbf{u}_{\cdot,2} \mid \dots \mid \mathbf{u}_{\cdot,p})$ 
    - \*  $\mathbf{U}$  is  $p \times p$  orthogonal matrix
  - ▷  $\mathbf{\Lambda} = \text{diag}\{\lambda_1, \dots, \lambda_p\}$ 
    - \*  $\mathbf{\Lambda}$  is  $p \times p$  diagonal matrix
  - ▷ convention
    - \*  $\lambda_i$  is the  $i^{\text{th}}$  largest eigenvalue of  $\mathbf{A}$
    - \*  $\mathbf{u}_{\cdot,i}$  is the eigenvector corresponding to  $\lambda_i$

### Geometry for $\mathbf{A} \succ 0$

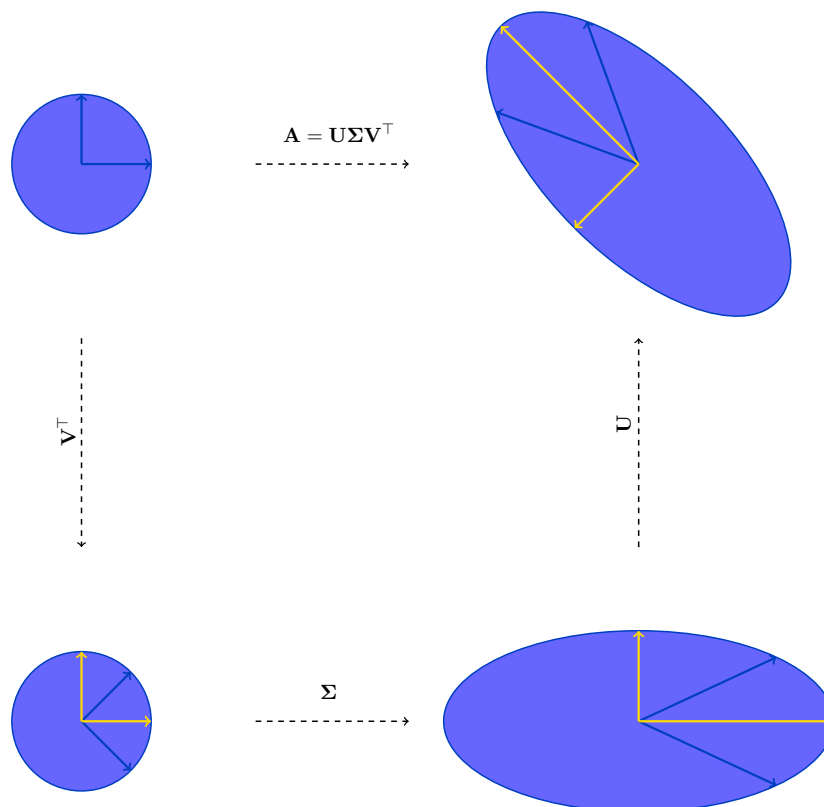


### 2.3.2 Singular value decomposition

#### Singular value decomposition (SVD)

- let  $\mathbf{A}$  be an  $n \times p$  ( $n \geq p$ ) rectangular matrix
  - ▷ singular values  $\sigma_1, \dots, \sigma_p$
  - ▷ left singular vectors  $\mathbf{u}_{\cdot,1}, \dots, \mathbf{u}_{\cdot,p}$
  - ▷ right singular vectors  $\mathbf{v}_{\cdot,1}, \dots, \mathbf{v}_{\cdot,p}$
  - ▷  $\mathbf{A}\mathbf{v}_{\cdot,i} = \sigma_i\mathbf{u}_{\cdot,i}$  &  $\mathbf{A}^\top\mathbf{u}_{\cdot,i} = \sigma_i\mathbf{v}_{\cdot,i}$
  - ▷ everything real, singular values non-negative
- SVD:  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$ , where
  - ▷  $\mathbf{U}$  is  $(n \times n)$  orthogonal
    - \* first  $p$  columns of  $\mathbf{U}$ :  $\mathbf{u}_{\cdot,1}, \dots, \mathbf{u}_{\cdot,p}$
  - ▷  $\mathbf{\Sigma}$   $(n \times p)$  diagonal
    - \*  $\sigma_1, \dots, \sigma_p$  on the diagonal of  $\mathbf{\Sigma}$
  - ▷  $\mathbf{V}$  is  $(p \times p)$  orthogonal
    - \* columns of  $\mathbf{V}$ :  $\mathbf{v}_{\cdot,1}, \dots, \mathbf{v}_{\cdot,p}$
  - ▷ convention : singular values in the descending order

#### Geometry for a square $\mathbf{A}$



## SVD and spectral decomposition

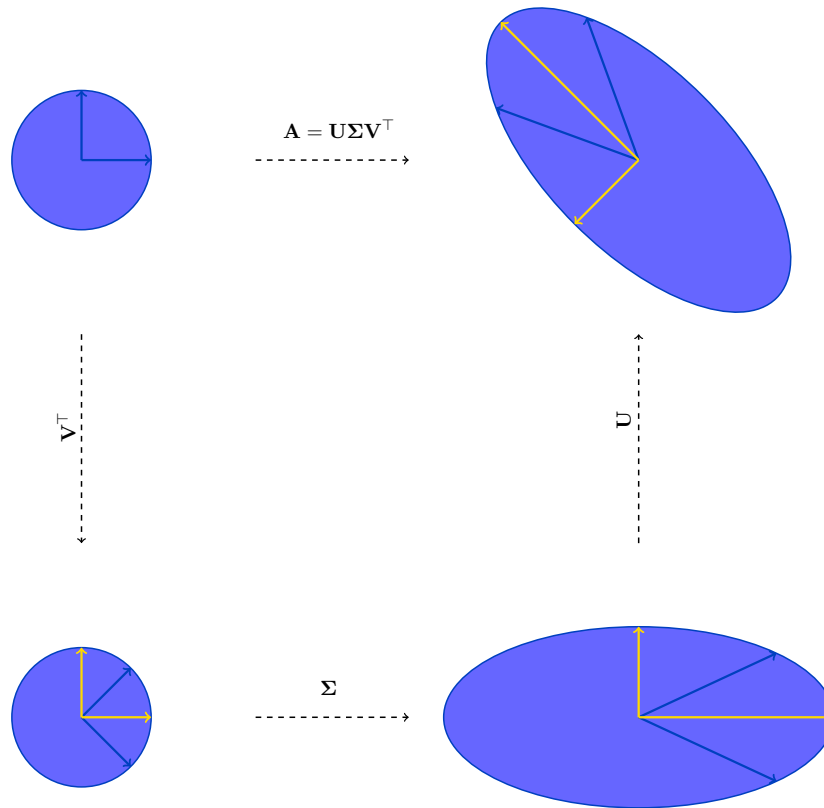
- SVD: rectangular matrix  $\mathbf{A}$  ( $n \times p$ ,  $n \geq p$ )
  - ▷ singular values and vectors  $(\sigma_1, \mathbf{u}_{.,1}, \mathbf{v}_{.,1}), \dots, (\sigma_p, \mathbf{u}_{.,p}, \mathbf{v}_{.,p})$
  - ▷  $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T$ , where
    - \*  $\mathbf{U}$  is  $(n \times n)$  and  $\mathbf{V}$  is  $(p \times p)$ , both orthogonal
    - \*  $\Sigma$  ( $n \times p$ ) diagonal with non-negative diagonal
- Spec. dec.: square symmetric matrix  $\mathbf{A}$  ( $p \times p$ )
  - ▷ eigenvalues and eigenvectors  $(\lambda_1, \mathbf{u}_{.,1}), \dots, (\lambda_p, \mathbf{u}_{.,p})$

- ▷  $\mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top$ , where
  - \*  $\mathbf{U}$  is  $(p \times p)$  orthogonal
  - \*  $\mathbf{\Lambda}$   $(p \times p)$  diagonal
- for a square symmetric  $\mathbf{A}$ ,  $\mathbf{A} \succeq \mathbf{0}$ :  $\mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top$
- for a rectangular matrix  $\mathbf{A}$  ( $n \times p$ ,  $n \geq p$ ),  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$ 
  - ▷  $\mathbf{A}^\top\mathbf{A} = \mathbf{V}\mathbf{\Sigma}^\top\mathbf{\Sigma}\mathbf{V}^\top$  ( $p \times p$ )  $\Rightarrow \mathbf{v}_{\cdot,i}$ 's are eigenvectors of  $\mathbf{A}^\top\mathbf{A}$
  - ▷  $\mathbf{A}\mathbf{A}^\top = \mathbf{U}\mathbf{\Sigma}\mathbf{\Sigma}^\top\mathbf{U}^\top$  ( $n \times n$ )  $\Rightarrow \mathbf{u}_{\cdot,i}$ 's are eigenvectors of  $\mathbf{A}\mathbf{A}^\top$
  - ▷  $\sigma_i$ 's are square roots of non-zero  $\lambda_i$ 's of  $\mathbf{A}^\top\mathbf{A}$  and  $\mathbf{A}\mathbf{A}^\top$

## Reduced SVD's

- SVD: rectangular matrix  $\mathbf{A}$  ( $n \times p$ ,  $n \geq p$ )
  - ▷ singular values and vectors  $(\sigma_1, \mathbf{u}_{\cdot,1}, \mathbf{v}_{\cdot,1}), \dots, (\sigma_p, \mathbf{u}_{\cdot,p}, \mathbf{v}_{\cdot,p})$
  - ▷  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$ , where
    - \*  $\mathbf{U}$  is  $(n \times n)$  and  $\mathbf{V}$  is  $(p \times p)$ , both orthogonal
    - \*  $\mathbf{\Sigma}$   $(n \times p)$  diagonal with non-negative diagonal
  - if  $n > p$ 
    - ▷  $\mathbf{A} = (\mathbf{U}_1 | \mathbf{U}_2) \begin{pmatrix} \mathbf{\Sigma}_1 \\ \mathbf{0} \end{pmatrix} \mathbf{V}^\top$
    - $\Leftrightarrow \mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top = \mathbf{U}_1\mathbf{\Sigma}_1\mathbf{V}^\top$ , where
      - \*  $\mathbf{U}_1$   $(n \times p)$  with orthogonal columns
      - \*  $\mathbf{\Sigma}_1$  is  $(p \times p)$
      - \* **thin SVD**:  $\mathbf{A} = \mathbf{U}_1\mathbf{\Sigma}_1\mathbf{V}^\top$
    - ▷ if  $r = \text{rank}(\mathbf{A}) < p$ ,  $\mathbf{\Sigma}_1$  is  $(r \times r)$  ... **compact SVD**
  - SVD writes  $\mathbf{A}$  as a sum of multiples of rank-one matrices:
 
$$\mathbf{A} = \sum_{i=1}^p \sigma_i \mathbf{u}_{\cdot,i} \mathbf{v}_{\cdot,i}^\top = \sum_{i=1}^r \sigma_i \mathbf{u}_{\cdot,i} \mathbf{v}_{\cdot,i}^\top$$

## Geometry



## SVD and linear mapping

- SVD: rectangular matrix  $\mathbf{A}$  ( $n \times p$ ,  $n \geq p$ )
  - ▷ singular values and vectors  $(\sigma_1, \mathbf{u}_{.,1}, \mathbf{v}_{.,1}), \dots, (\sigma_p, \mathbf{u}_{.,p}, \mathbf{v}_{.,p})$
  - ▷  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ , where
    - \*  $\mathbf{U}$  and  $\mathbf{V}$ : orthonormal bases of  $\mathbb{R}^n$  and  $\mathbb{R}^p$  such that  $\mathbf{A}$  maps the  $i^{\text{th}}$  basis vector of  $\mathbb{R}^p$  to a non-negative multiple of the  $i^{\text{th}}$  basis vector of  $\mathbb{R}^n$ , and sends the left-over basis vectors to zero
- $\ker(\mathbf{A})$ 
  - ▷ spanned by the  $\mathbf{v}_{.,i}$  corresponding to the null  $\sigma_i$

- $\text{im}(\mathbf{A})$ 
  - ▷ spanned by the  $\mathbf{u}_{.,i}$  corresponding to the positive  $\sigma_i$
- $\dim(\ker(\mathbf{A})) + \dim(\text{im}(\mathbf{A})) = p$

### 2.3.3 QR decomposition

#### QR decomposition (factorization)

- let  $\mathbf{A}$  be a  $p \times p$  matrix  $\rightsquigarrow \mathbf{A} = \mathbf{QR}$ , where
  - ▷  $\mathbf{Q}$  is  $p \times p$  orthogonal
  - ▷  $\mathbf{R}$  is  $p \times p$  upper triangular
- let  $\mathbf{A}$  be an  $n \times p$  ( $n \geq p$ ) matrix  $\rightsquigarrow \mathbf{A} = \mathbf{QR}$ , where
  - ▷  $\mathbf{Q}$  is  $n \times n$  orthogonal
  - ▷  $\mathbf{R}$  is  $n \times p$  upper triangular
- if  $n > p$ 
  - ▷  $\mathbf{A} = (\mathbf{Q}_1 | \mathbf{Q}_2) \begin{pmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{pmatrix} = \mathbf{Q}_1 \mathbf{R}_1$ , where
    - \*  $\mathbf{Q}_1$  is  $n \times p$  with orthogonal columns
    - \*  $\mathbf{R}_1$  is  $p \times p$  upper triangular
  - ▷  $\text{rank}(\mathbf{A}) = p \Rightarrow \text{rank}(\mathbf{R}_1) = p$

## 2.4 Pseudoinverse

### 2.4.1 Moore–Penrose pseudoinverse

#### Moore–Penrose pseudoinverse

- let  $\mathbf{A}$  be a  $p \times p$  matrix,  $\text{rank}(\mathbf{A}) = p$
- inverse  $\mathbf{A}^{-1}$  is the  $p \times p$  matrix satisfying
  - ▷  $\mathbf{A}\mathbf{A}^{-1} = \mathbf{I}$
  - ▷  $\mathbf{A}^{-1}\mathbf{A} = \mathbf{I}$
- let  $\mathbf{A}$  be an  $n \times p$  ( $n \geq p$ ) matrix
- Moore–Penrose pseudoinverse  $\mathbf{A}^+$  is the  $p \times n$  matrix satisfying

- ▷  $\mathbf{A}\mathbf{A}^+\mathbf{A} = \mathbf{A}$  (generalized inverse)
- ▷  $\mathbf{A}^+\mathbf{A}\mathbf{A}^+ = \mathbf{A}^+$  (generalized reflexive inverse)
- ▷  $(\mathbf{A}\mathbf{A}^+)^\top = \mathbf{A}\mathbf{A}^+$
- ▷  $(\mathbf{A}^+\mathbf{A})^\top = \mathbf{A}^+\mathbf{A}$
- $\mathbf{A}^+$  exists and is unique

### Construction of $\mathbf{A}^+$

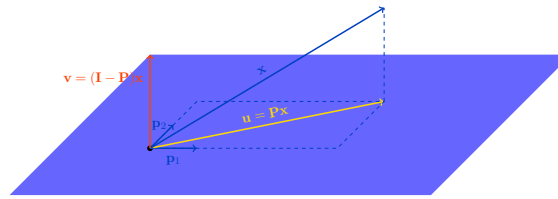
- let  $\mathbf{A}$  be an  $n \times p$  ( $n \geq p$ ) matrix
  - ▷ if  $\text{rank}(\mathbf{A}) = p$ 
    - \*  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$  (thin SVD, i.e.  $\mathbf{U}$  is  $n \times p$  &  $\mathbf{\Sigma}$  is  $p \times p$ )
    - \*  $\mathbf{A}^+ = \mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^\top$
  - ▷ if  $\text{rank}(\mathbf{A}) = r < p$ 
    - \*  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$  (compact SVD, i.e.  $\mathbf{U}$  is  $n \times r$  &  $\mathbf{\Sigma}$  is  $r \times r$ )
    - \*  $\mathbf{A}^+ = \mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^\top$
- let  $\mathbf{A}$  be a  $p \times p$  symmetric matrix
  - ▷  $\mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top$  (spectral decomposition)
  - ▷  $\mathbf{A}^+ = \mathbf{U}\mathbf{\Lambda}^+\mathbf{U}^\top$ , where
    - \*  $\mathbf{\Lambda}^+$ : diagonal with  $1/\lambda_i$  on diagonal if  $\lambda_i \neq 0$ , 0 otherwise

## 2.5 Orthogonal projection

### Orthogonal projection

#### Projection on $\mathbb{R}^n$

- **projection**: linear mapping  $\mathbf{P} : \mathbb{R}^n \mapsto \mathbb{R}^n$  such that  $\mathbf{P}\mathbf{P} = \mathbf{P}$ 
  - ▷  $\mathbf{P}$  idempotent
  - ▷  $\mathbf{P}$  is identity on  $\text{im}(\mathbf{P})$
- $\forall \mathbf{x} \in \mathbb{R}^n : \mathbf{x} = \underbrace{\mathbf{P}\mathbf{x}}_{\mathbf{u}} + \underbrace{(\mathbf{I} - \mathbf{P})\mathbf{x}}_{\mathbf{v}}$ 
  - ▷  $\mathbf{u} \in \text{im}(\mathbf{P})$  &  $\mathbf{v} \in \text{ker}(\mathbf{P})$ : unique decomposition
- $(\mathbf{I} - \mathbf{P})$  is a projection on  $\text{ker}(\mathbf{P})$  and  $\text{ker}(\mathbf{I} - \mathbf{P}) = \text{im}(\mathbf{P})$

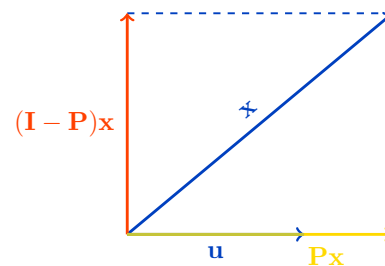


- orthogonal projection: projection with a symmetric  $\mathbf{P}$ 
  - ▷  $\mathbf{P}$  is symmetric iff  $\text{im}(\mathbf{P}) = (\text{im}(\mathbf{I} - \mathbf{P}))^\perp$
  - ▷  $\exists! \mathbf{P}\mathbf{x} \in \text{im}(\mathbf{P})$  and  $\|\mathbf{x} - \mathbf{P}\mathbf{x}\|^2 = \min_{\mathbf{y} \in \text{im}(\mathbf{P})} \|\mathbf{x} - \mathbf{y}\|^2$
  - ▷ if  $\mathbf{P}$  and  $\mathbf{P}_1$  are orthogonal projections and  $\text{im}(\mathbf{P}_1) \leq \text{im}(\mathbf{P})$ 
    - \*  $\mathbf{P}\mathbf{P}_1 = \mathbf{P}_1 = \mathbf{P}_1\mathbf{P}$

### Construction of $\mathbf{P}$

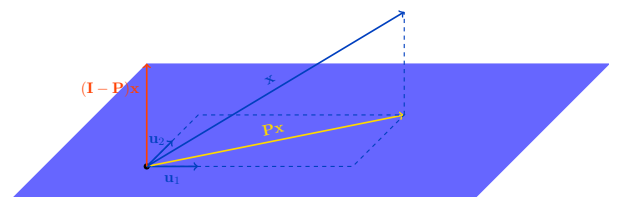
1. orth. projection on  $\text{im}(\mathbf{u})$ :

- $\mathbf{P} = \frac{1}{\|\mathbf{u}\|^2} \mathbf{u}\mathbf{u}^\top$
- $\mathbf{P}\mathbf{x} = \frac{\langle \mathbf{u}, \mathbf{x} \rangle}{\|\mathbf{u}\|^2} \mathbf{u} \in \text{im}(\mathbf{u})$ 
  - ▷ leaves  $c\mathbf{u}$  unchanged
  - ▷ annihilates the complementary basis



2.  $\{\mathbf{u}_1, \dots, \mathbf{u}_p\}$  orthonormal, orth. projection on  $\text{im}(\mathbf{U})$ ,  $\mathbf{U} = (\mathbf{u}_1 | \mathbf{u}_2 | \dots | \mathbf{u}_p)$

- $\mathbf{P} = \mathbf{U}\mathbf{U}^\top$
- $\mathbf{P}\mathbf{x} = \sum_{i=1}^p \langle \mathbf{u}_i, \mathbf{x} \rangle \mathbf{u}_i \in \text{im}(\mathbf{U})$ 
  - ▷ leaves  $\sum_{i=1}^p c_i \mathbf{u}_i$  unchanged
  - ▷ annihilates the complementary basis



### Orthogonal projection onto a column space

- let  $\mathbf{A}$  be an  $n \times p$  ( $n \geq p$ ) matrix

1.  $\text{rank}(\mathbf{A}) = p$

- columns  $\mathbf{a}_{\cdot,1}, \dots, \mathbf{a}_{\cdot,p}$  linearly independent
- $\mathbf{P} = \mathbf{A}(\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top$



- $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$  (thin SVD, i.e.  $\mathbf{U}$  is  $n \times p$  &  $\mathbf{\Sigma}$  is  $p \times p$ )
  - ▷  $(\mathbf{A}^\top \mathbf{A})^{-1} = \mathbf{V}\mathbf{\Sigma}^{-2}\mathbf{V}^\top$
  - ▷  $\mathbf{P} = \mathbf{U}\mathbf{U}^\top$

2.  $\text{rank}(\mathbf{A}) = r < p$

- $\mathbf{P} = \mathbf{A}(\mathbf{A}^\top \mathbf{A})^+ \mathbf{A}^\top$
- $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$  (compact SVD, i.e.  $\mathbf{U}$  is  $n \times r$  &  $\mathbf{\Sigma}$  is  $r \times r$ )
  - ▷  $(\mathbf{A}^\top \mathbf{A})^+ = \mathbf{V}\mathbf{\Sigma}^{-2}\mathbf{V}^\top$
  - ▷  $\mathbf{P} = \mathbf{U}\mathbf{U}^\top$

## 2.6 Application to linear regression

### Application to linear regression

#### Estimation in linear regression (theory)

- linear regression:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ ,  $\mathbb{E}\boldsymbol{\varepsilon} = \mathbf{0}$ ,  $\text{Var}\boldsymbol{\varepsilon} = \sigma^2\mathbf{I}$
- we want to estimate  $\boldsymbol{\beta}$
- start with estimating  $\boldsymbol{\mu} = \mathbb{E}(\mathbf{Y}) = \mathbf{X}\boldsymbol{\beta} \in \text{im}(\mathbf{X})$ 
  - ▷ we look for  $\hat{\boldsymbol{\mu}} \in \text{im}(\mathbf{X})$  closest to  $\mathbf{Y}$ 
    - \* we look to minimize  $\|\hat{\boldsymbol{\mu}} - \mathbf{Y}\|^2$
    - $\Rightarrow \hat{\boldsymbol{\mu}}$  is the orthogonal projection of  $\mathbf{Y}$  onto  $\text{im}(\mathbf{X})$
  - ▷  $\hat{\boldsymbol{\mu}} = \begin{cases} \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} = \mathbf{U}\mathbf{U}^\top \mathbf{Y} & \text{if } \text{rank}(\mathbf{X}) = p, \\ \mathbf{X}(\mathbf{X}^\top \mathbf{X})^+ \mathbf{X}^\top \mathbf{Y} = \mathbf{U}\mathbf{U}^\top \mathbf{Y} & \text{if } \text{rank}(\mathbf{X}) < p, \end{cases}$ 
    - \* where  $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$  (thin/compact SVD)
- $\hat{\boldsymbol{\mu}} \in \text{im}(\mathbf{X}) \Rightarrow \exists \hat{\boldsymbol{\beta}}$  such that  $\hat{\boldsymbol{\mu}} = \mathbf{X}\hat{\boldsymbol{\beta}}$
- if  $\text{rank}(\mathbf{X}) = p \Rightarrow \exists! \hat{\boldsymbol{\beta}}$  such that  $\hat{\boldsymbol{\mu}} = \mathbf{X}\hat{\boldsymbol{\beta}}$

#### Estimation in linear regression (practice in $\mathbb{R}$ )

- aim: minimize  $\|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2$  w.r.t.  $\boldsymbol{\beta}$
- use that  $\mathbf{X} = \mathbf{Q}\mathbf{R}$ 
  - ▷  $\mathbf{Q}$  is  $n \times n$  orthogonal

- ▷  $\mathbf{R}$  is  $n \times p$  upper triangular
- ▷  $\mathbf{Q}$  and  $\mathbf{Q}^\top$  rotations
- $\|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2 = \|\mathbf{Q}^\top(\mathbf{Y} - \mathbf{QR}\boldsymbol{\beta})\|^2$ 

$$= \left\| \begin{pmatrix} \mathbf{Q}_1^\top \\ \mathbf{Q}_2^\top \end{pmatrix} \mathbf{Y} - \begin{pmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{pmatrix} \boldsymbol{\beta} \right\|^2$$

$$= \|\mathbf{Q}_1^\top \mathbf{Y} - \mathbf{R}_1 \boldsymbol{\beta}\|^2 + \|\mathbf{Q}_2^\top \mathbf{Y}\|^2$$
  - ▷ minimize  $\|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2 \Leftrightarrow$  minimize  $\|\mathbf{Q}_1^\top \mathbf{Y} - \mathbf{R}_1 \boldsymbol{\beta}\|^2$
- if  $\text{rank}(\mathbf{X}) = p$ 
  - ▷  $\mathbf{R}_1$  invertible
  - ▷  $\hat{\boldsymbol{\beta}} = \mathbf{R}_1^{-1} \mathbf{Q}_1^\top \mathbf{Y}$

# Chapter 3

## Normal distribution

### 3.1 The problem

#### 3.1.1 Linear model

##### Linear model

- $Y_i = \beta_0 + \beta_1 x_{i,1} + \dots + \beta_k x_{i,k} + \varepsilon_i, i \in \{1, \dots, n\}$ 
  - ▷  $Y_i$ : outcome, response, output, dependent variable
    - \* random variable, we observe a realization  $y_i$
    - \* (odezva, závisle proměnná, regresand)
  - ▷  $x_{i,1}, \dots, x_{i,k}$ : covariates, predictors, explanatory variables, input, independent variables
    - \* given, known
    - \* (nezávisle proměnné, regresory)
  - ▷  $\beta_0, \dots, \beta_k$ : coefficients
    - \* unknown
    - \* (regresní koeficienty)
  - ▷  $\varepsilon_i$ : random error
    - \* random variable, unobserved
- $\varepsilon_i \stackrel{\text{iid}}{\sim} (0, \sigma^2), i \in \{1, \dots, n\}$ 
  - ▷  $E \varepsilon_i = 0$ : no systematic errors
  - ▷  $\text{Var } \varepsilon_i = \sigma^2$ : same precision
- we often assume that  $\varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2), i \in \{1, \dots, n\}$

**Example: bloodpress data**

- from `sites.stat.psu.edu/~lsimon/stat501wc/sp05/data/`
- association between the **mean arterial blood pressure**[mmHg] and `age`[years], `weight`[kg], `body surface area`[ $m^2$ ], `duration of hypertension`[years], `basal pulse`[beats/min], `stress`

○ data:

	BP	Age	Weight	BSA	DoH	Pulse	Stress
	105	47	85.4	1.75	5.1	63	33
	115	49	94.2	2.10	3.8	70	14
	...	...	...	...	...	...	...
	110	48	90.5	1.88	9.0	71	99
	122	56	95.7	2.09	7.0	75	99

- model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$

$$\begin{pmatrix} 105 \\ 115 \\ \dots \\ 110 \\ 122 \end{pmatrix} = \begin{pmatrix} 1 & 47 & 85.4 & 1.75 & 5.1 & 63 & 33 \\ 1 & 49 & 94.2 & 2.10 & 3.8 & 70 & 14 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & 48 & 90.5 & 1.88 & 9.0 & 71 & 99 \\ 1 & 56 & 95.7 & 2.09 & 7.0 & 75 & 99 \end{pmatrix} \times \begin{pmatrix} \beta_0 \\ \dots \\ \beta_6 \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \dots \\ \varepsilon_{19} \\ \varepsilon_{20} \end{pmatrix}$$

**3.1.2 Task for this chapter****Normal distribution in a linear model**

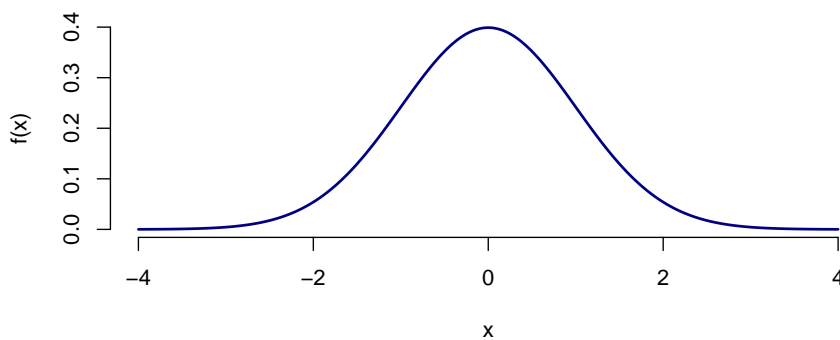
- model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$
- assumptions of the **normal linear model**:
  - ▷  $\mathbf{X}$  fixed and known
  - ▷  $\boldsymbol{\beta}$  fixed unknown
  - ▷  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$   
 $\Rightarrow \mathbf{Y} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$
- estimators of  $\boldsymbol{\beta}$  and  $\sigma^2$ 
  - ▷ functions of  $\mathbf{Y}$
- test statistics concerning  $\boldsymbol{\beta}$  and  $\sigma^2$ 
  - ▷ functions of  $\mathbf{Y}$
- $\Rightarrow$  to make inference in normal linear model, we need to study
  - ▷ multivariate normal distribution  $N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$
  - ▷ distributions of functions of  $N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

## 3.2 Univariate normal distribution

### 3.2.1 Definition

#### Normal distribution $N(\mu, \sigma^2)$

- let  $\mu \in \mathbb{R}$  and  $\sigma^2 > 0$ 
  - ▷ density  $f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(x - \mu)^2\right\}$
  - ▷ for the standard normal distribution ( $\mu = 0, \sigma^2 = 1$ ):

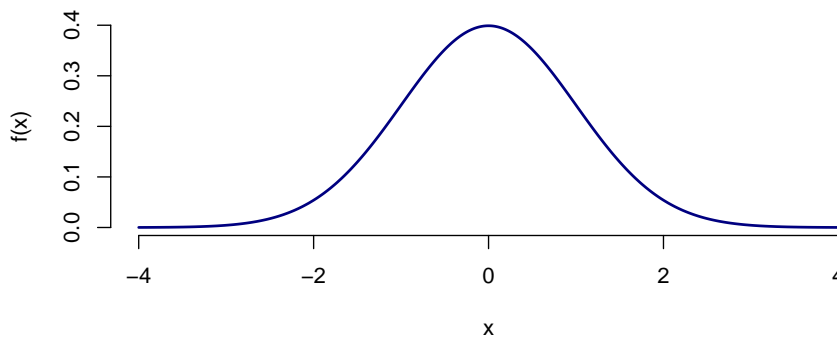


- if  $\sigma^2 = 0$  then  $X = \mu$  a.s.

### 3.2.2 Properties

#### Properties of $N(\mu, \sigma^2)$

- $\mu \in \mathbb{R}$  and  $\sigma^2 > 0$
- Let  $a, b \in \mathbb{R}$ ,  $X \sim N(\mu, \sigma^2)$ . Then  $aX + b \sim N(a\mu + b, a^2\sigma^2)$ .
- Let  $Z \sim N(0, 1)$  and  $X = \mu + \sigma Z$ . Then  $X \sim N(\mu, \sigma^2)$ .



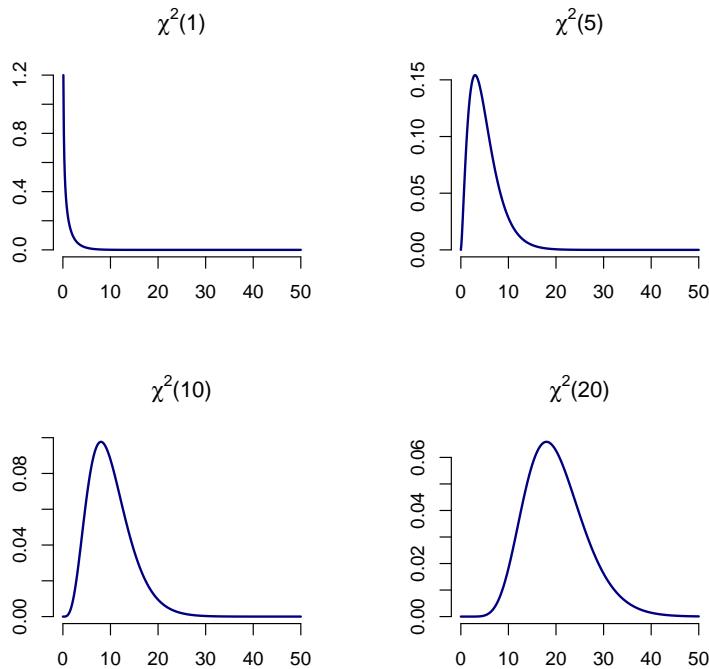
- Let  $a_i, b_i \in \mathbb{R}$ ,  $X_i \stackrel{\text{ind.}}{\sim} N(\mu_i, \sigma_i^2)$  for  $i \in \{1, \dots, n\}$ .

Then  $\sum_{i=1}^n (a_i X_i + b_i) \sim N(\sum_{i=1}^n (a_i \mu_i + b_i), \sum_{i=1}^n a_i^2 \sigma_i^2)$ .

### 3.2.3 Related distributions

#### $\chi^2(n)$ distribution

- let  $Z \sim N(0, 1) \rightsquigarrow Z^2 \sim \chi^2(1)$
- let  $Z_i \stackrel{\text{ind.}}{\sim} N(0, 1)$  for  $i \in \{1, \dots, n\} \rightsquigarrow X = \sum_{i=1}^n Z_i^2 \sim \chi^2(n)$
- density



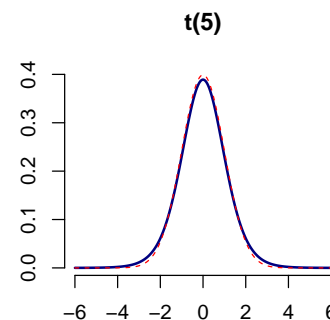
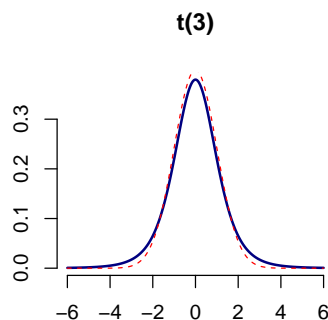
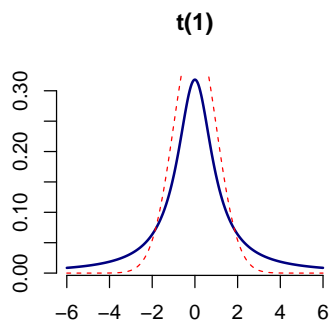
- $E X = n, \text{Var } X = 2n$

#### Student's $t$ -distribution

- let  $Z \sim N(0, 1)$  and  $X \sim \chi^2(n)$ ,  $Z \perp\!\!\!\perp X$

$$\triangleright T = \frac{Z}{\sqrt{X/n}} \sim t(n)$$

- density



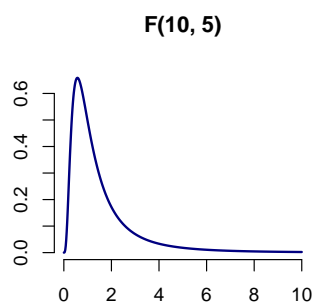
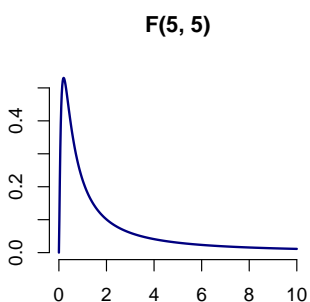
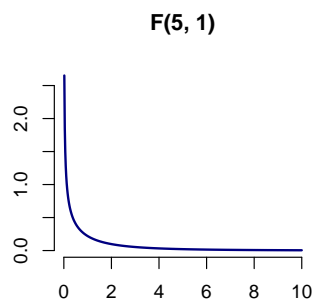
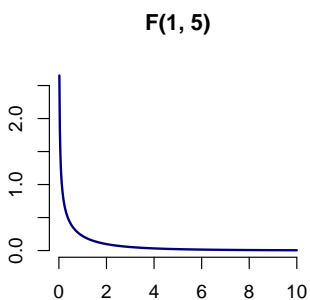
- $ET = 0$  for  $n > 1$ ,  $\text{Var } T = n/(n - 2)$  for  $n > 2$

### Fisher–Snedecor distribution

- let  $X_1 \sim \chi^2(n_1)$  and  $X_2 \sim \chi^2(n_2)$ ,  $X_1 \perp\!\!\!\perp X_2$

$$\triangleright F = \frac{X_1/n_1}{\sqrt{X_2/n_2}} \sim F(n_1, n_2)$$

- density



- $EF = n_2/(n_2 - 2)$  for  $n_2 > 2$

### 3.3 Multivariate normal distribution

#### 3.3.1 Definition

##### Multivariate normal distribution $N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

- $\boldsymbol{\mu} \in \mathbb{R}^n$ ,  $\boldsymbol{\Sigma}$  is an  $n \times n$  positive semidefinite matrix

**Definition.** A random vector  $\mathbf{X} : (\Omega, \mathcal{A}) \mapsto (\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n))$  has *multivariate normal distribution*  $N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  if and only if  $\mathbf{a}^\top \mathbf{X} \sim N(\mathbf{a}^\top \boldsymbol{\mu}, \mathbf{a}^\top \boldsymbol{\Sigma} \mathbf{a})$  for every  $\mathbf{a} \in \mathbb{R}^n$ .

- if  $\text{rank}(\boldsymbol{\Sigma}) = n$  then  $N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  is *non-degenerate*

▷ has density

$$f(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^n \det(\boldsymbol{\Sigma})}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

- if  $\text{rank}(\boldsymbol{\Sigma}) = r < n$  then  $N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  is *degenerate*

▷ a.s. “lives” in a subspace of  $\mathbb{R}^n$  of dimension  $r$

▷ no density w.r.t. Lebesgue measure on  $\mathcal{B}(\mathbb{R}^n)$

##### Non-degenerate multivariate normal distribution

- density

$$f(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^n \det(\boldsymbol{\Sigma})}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

- $\boldsymbol{\Sigma}$ : square symmetric positive definite matrix

▷ spectral decomposition  $\boldsymbol{\Sigma} = \mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^\top$

▷  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n > 0$

▷  $\boldsymbol{\Sigma}^{-1} = \mathbf{U} \boldsymbol{\Lambda}^{-1} \mathbf{U}^\top$

- quadratic form  $(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})$  can be written as

$$\text{▷ } (\mathbf{x} - \boldsymbol{\mu})^\top \mathbf{U} \boldsymbol{\Lambda}^{-1} \mathbf{U}^\top (\mathbf{x} - \boldsymbol{\mu}) = \{\mathbf{U}^\top (\mathbf{x} - \boldsymbol{\mu})\}^\top \boldsymbol{\Lambda}^{-1} \{\mathbf{U}^\top (\mathbf{x} - \boldsymbol{\mu})\}$$

- level sets of  $f(\mathbf{x})$ ,  $I_c = \{\mathbf{x} \in \mathbb{R}^n; f(\mathbf{x}) = c\}$  for  $c > 0$ :

▷ ellipsoids centred at  $\boldsymbol{\mu}$

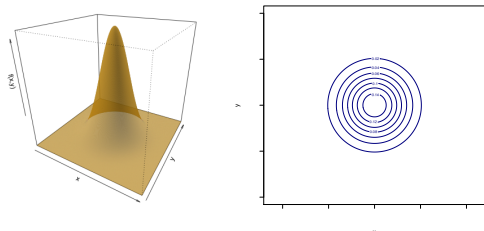
▷ directions of principal axes:  $\mathbf{u}_1, \dots, \mathbf{u}_n$ ,

▷ lengths of principal semi-axes:  $\sqrt{d\lambda_1}, \dots, \sqrt{d\lambda_n}$

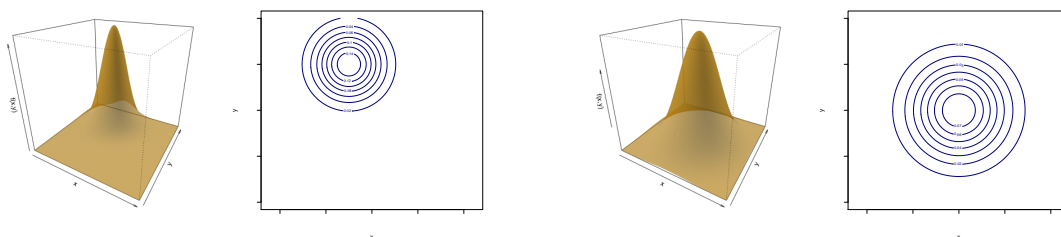


## Non-degenerate bivariate normal distribution

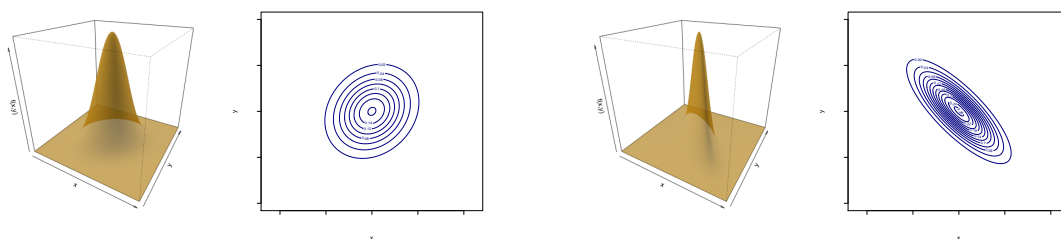
○  $N\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}\right)$



○  $N\left(\begin{pmatrix} -1 \\ 2 \end{pmatrix}, \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}\right)$        $N\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix}\right)$



○  $N\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & 0.2 \\ 0.2 & 1 \end{pmatrix}\right)$        $N\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & -0.8 \\ -0.8 & 1 \end{pmatrix}\right)$



### 3.3.2 Properties

#### Properties of $N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

- $\boldsymbol{\mu} \in \mathbb{R}^n$ ,  $\boldsymbol{\Sigma}$  is an  $n \times n$  symmetric positive semidefinite matrix

**Theorem (MVN 1).** Let  $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ . Then  $E \mathbf{X} = \boldsymbol{\mu}$  and  $\text{Var } \mathbf{X} = \boldsymbol{\Sigma}$ .

**Theorem (MVN 2).** Let  $Z_1, \dots, Z_n \stackrel{\text{iid}}{\sim} N(0, 1)$  and  $\mathbf{Z} = (Z_1, \dots, Z_n)^\top$ . Then  $\mathbf{Z} \sim N(\mathbf{0}, \mathbf{I})$ .

**Theorem (MVN 3).** Let  $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  and let  $\mathbf{A}$  be an  $m \times n$  real matrix and  $\mathbf{b} \in \mathbb{R}^m$ . Then  $\mathbf{AX} + \mathbf{b} \sim N(\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top)$ .

- proofs are given during the lectures and can also be found in Jiří Anděl: Základy matematické statistiky

### $N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ seen through $N(\mathbf{0}, \mathbf{I})$

- $\boldsymbol{\mu} \in \mathbb{R}^n$ ,  $\boldsymbol{\Sigma}$  is an  $n \times n$  symmetric positive semidefinite matrix

#### 1. if $\text{rank}(\boldsymbol{\Sigma}) = n$

- spectral decomposition  $\boldsymbol{\Sigma} = \mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^\top$
- $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n > 0$
- $\boldsymbol{\Sigma} = \mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^\top = \underbrace{\mathbf{U}\boldsymbol{\Lambda}^{1/2}}_{\tilde{\boldsymbol{\Sigma}}} \boldsymbol{\Lambda}^{1/2}\mathbf{U}^\top = \tilde{\boldsymbol{\Sigma}}\tilde{\boldsymbol{\Sigma}}^\top$
- let  $\mathbf{Z} = \tilde{\boldsymbol{\Sigma}}^{-1}(\mathbf{X} - \boldsymbol{\mu}) = \boldsymbol{\Lambda}^{-1/2}\mathbf{U}^\top(\mathbf{X} - \boldsymbol{\mu})$   
 $\Rightarrow \mathbf{Z} \sim N(\mathbf{0}, \mathbf{I})$  ( $n$ -dimensional),  $\mathbf{X} = \boldsymbol{\mu} + \tilde{\boldsymbol{\Sigma}}\mathbf{Z}$  and  $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

#### 2. if $\text{rank}(\boldsymbol{\Sigma}) = r < n$

- spectral decomposition  $\boldsymbol{\Sigma} = \mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^\top$
- $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r > 0, \lambda_{r+1} = \lambda_{r+2} = \dots = \lambda_n = 0$
- $\boldsymbol{\Sigma} = \mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^\top = \underbrace{\mathbf{U}_{n \times r}}_{(\mathbf{u}_{1,1} | \mathbf{u}_{2,1} | \dots | \mathbf{u}_{r,1})} \underbrace{\boldsymbol{\Lambda}_{r \times r}}_{\text{diag}\{\lambda_1, \lambda_2, \dots, \lambda_r\}} \mathbf{U}_{n \times r}^\top = \underbrace{\mathbf{U}_{n \times r} \boldsymbol{\Lambda}_{r \times r}^{1/2}}_{\tilde{\boldsymbol{\Sigma}}} \boldsymbol{\Lambda}_{r \times r}^{1/2} \mathbf{U}_{n \times r}^\top = \tilde{\boldsymbol{\Sigma}}\tilde{\boldsymbol{\Sigma}}^\top$
- let  $\mathbf{Z} = \tilde{\boldsymbol{\Sigma}}^+(\mathbf{X} - \boldsymbol{\mu}) = \boldsymbol{\Lambda}_{r \times r}^{-1/2} \mathbf{U}_{n \times r}^\top(\mathbf{X} - \boldsymbol{\mu})$   
 $\Rightarrow \mathbf{Z} \sim N(\mathbf{0}, \mathbf{I})$  ( $r$ -dimensional),  $\mathbf{X} = \boldsymbol{\mu} + \tilde{\boldsymbol{\Sigma}}\mathbf{Z}$  and  $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

### Density of $N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

- $\boldsymbol{\mu} \in \mathbb{R}^n$ ,  $\boldsymbol{\Sigma}$  is an  $n \times n$  symmetric positive definite matrix

**Theorem (MVN 4).** Let  $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  where  $\text{rank}(\boldsymbol{\Sigma}) = n$ . Then  $\mathbf{X}$  has density  $f(\mathbf{x})$  w.r.t. Lebesgue measure on  $\mathcal{B}(\mathbb{R}^n)$  and

$$f(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^n \det(\boldsymbol{\Sigma})}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right\}.$$

- a proof is given during the lectures and can also be found in Jiří Anděl: Základy matematické statistiky

### Characteristic function (reminder)

**Definition** (Characteristic function of a random variable). *Let  $X$  be a random variable. The function  $\psi_X : \mathbb{R} \mapsto \mathbb{C}$  defined by  $\psi_X(t) = \mathbb{E} \exp\{itX\}$ ,  $t \in \mathbb{R}$ , is the *characteristic function of  $X$* .*

**Definition** (Characteristic function of a random vector). *Let  $\mathbf{X}$  be an  $n$ -dimensional random vector. The function  $\psi_{\mathbf{X}} : \mathbb{R}^n \mapsto \mathbb{C}$  defined by  $\psi_{\mathbf{X}}(\mathbf{t}) = \mathbb{E} \exp\{i\mathbf{t}^\top \mathbf{X}\}$ ,  $\mathbf{t} \in \mathbb{R}^n$ , is the *characteristic function of  $\mathbf{X}$* .*

### Properties of characteristic function (reminder)

**Theorem** (ChF 1). *Let  $X \sim N(\mu, \sigma^2)$ . Then  $\psi_X(t) = \exp\{it\mu - \frac{1}{2}\sigma^2 t^2\}$ .*

**Theorem** (ChF 2). *Let  $\mathbf{X}$  be an  $n$ -dimensional random vector and  $\mathbf{X}_1$  and  $\mathbf{X}_2$  its subvectors such that  $\mathbf{X} = (\mathbf{X}_1^\top, \mathbf{X}_2^\top)^\top$ . Then  $\mathbf{X}_1 \perp\!\!\!\perp \mathbf{X}_2$  iff  $\psi_{\mathbf{X}}(\mathbf{t}) = \psi_{\mathbf{X}_1}(\mathbf{t}_1) \times \psi_{\mathbf{X}_2}(\mathbf{t}_2)$  for every  $\mathbf{t} = (\mathbf{t}_1^\top, \mathbf{t}_2^\top)^\top \in \mathbb{R}^n$ .*

- a proof can be found in *Petr Lachout: Teorie pravděpodobnosti (1998). Nakladatelství Univerzity Karlovy*

### Characteristic function of $N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

- $\boldsymbol{\mu} \in \mathbb{R}^n$ ,  $\boldsymbol{\Sigma}$  is an  $n \times n$  symmetric positive semidefinite matrix

**Theorem** (MVN 5). *Let  $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ . Then*

$$\psi_{\mathbf{X}}(\mathbf{t}) = \exp\left\{i\mathbf{t}^\top \boldsymbol{\mu} - \frac{1}{2}\mathbf{t}^\top \boldsymbol{\Sigma} \mathbf{t}\right\}.$$

- a proof is given during the lectures and can also be found in Jiří Anděl: *Základy matematické statistiky*

### Subvectors of $N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

- $\boldsymbol{\mu} \in \mathbb{R}^n$ ,  $\boldsymbol{\Sigma}$  is an  $n \times n$  symmetric positive semidefinite matrix

**Theorem** (MVN 6). *Let  $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  and let  $k \in \{1, \dots, n\}$ . Then*

$$\begin{pmatrix} X_1 \\ X_2 \\ \dots \\ X_k \end{pmatrix} \sim N\left(\begin{pmatrix} \mu_1 \\ \mu_2 \\ \dots \\ \mu_k \end{pmatrix}, \begin{pmatrix} \sigma_{1,1} & \sigma_{1,2} & \dots & \sigma_{1,k} \\ \sigma_{2,1} & \sigma_{2,2} & \dots & \sigma_{2,k} \\ \dots & \dots & \dots & \dots \\ \sigma_{k,1} & \sigma_{k,2} & \dots & \sigma_{k,k} \end{pmatrix}\right).$$

- a proof is given during the lectures and can also be found in Jiří Anděl: *Základy matematické statistiky*

- analogous statement is true for any sub-vector of  $\mathbf{X}$
- converse is not true

### (In)dependence in $N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

- $\boldsymbol{\mu} \in \mathbb{R}^n$ ,  $\boldsymbol{\Sigma}$  is an  $n \times n$  symmetric positive semidefinite matrix

**Theorem (MVN 7).** Let  $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  and let  $k \in \{1, \dots, n-1\}$ . Denote  $\mathbf{X}_1 = (X_1, \dots, X_k)^\top$ ,  $\mathbf{X}_2 = (X_{k+1}, \dots, X_n)^\top$  and  $\mathbf{X}_1 \sim N(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{1,1})$ ,  $\mathbf{X}_2 \sim N(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_{2,2})$ . If

$$\boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{1,1} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma}_{2,2} \end{pmatrix}$$

then  $\mathbf{X}_1 \perp\!\!\!\perp \mathbf{X}_2$ .

- a proof is given during the lectures and can also be found in Jiří Anděl: *Základy matematické statistiky*
- $\mathbf{A}\mathbf{X} \perp\!\!\!\perp \mathbf{B}\mathbf{X}$  iff  $\mathbf{A}\boldsymbol{\Sigma}\mathbf{B}^\top = \mathbf{0}$

### 3.3.3 Related distributions

#### Quadratic forms

- Let  $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ ,  $\boldsymbol{\mu} \in \mathbb{R}^n$ ,  $\boldsymbol{\Sigma}$  is an  $n \times n$  symmetric positive semidefinite matrix

**Theorem (QF 1).** Let  $\mathbf{Z} \sim N(\mathbf{0}, \mathbf{I})$ . Then  $\mathbf{Z}^\top \mathbf{Z} \sim \chi^2(n)$ .

**Theorem (QF 2).** Let  $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  where  $\text{rank}(\boldsymbol{\Sigma}) = n$ . Then  $(\mathbf{X} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{X} - \boldsymbol{\mu}) \sim \chi^2(n)$ .

**Theorem (QF 3).** Let  $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  where  $\text{rank}(\boldsymbol{\Sigma}) = r < n$ . Then  $(\mathbf{X} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^+(\mathbf{X} - \boldsymbol{\mu}) \sim \chi^2(r)$ .

- proofs are given during the lectures and analogous statements are proved in Jiří Anděl: *Základy matematické statistiky*

#### Quadratic forms

- Let  $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ ,  $\boldsymbol{\mu} \in \mathbb{R}^n$ ,  $\boldsymbol{\Sigma}$  is an  $n \times n$  symmetric positive semidefinite matrix

**Theorem (QF 4).** Let  $\mathbf{Z} \sim N(\mathbf{0}, \mathbf{I})$  and let  $\mathbf{P}$  be an  $n \times n$  projection matrix of rank  $r$ . Then  $\mathbf{Z}^\top \mathbf{P} \mathbf{Z} \sim \chi^2(r)$ .

- a proof is given during the lectures and analogous statements are proved in Jiří Anděl: *Základy matematické statistiky*

# Chapter 4

## Linear model

### 4.1 The problem

#### 4.1.1 Linear model

##### Linear model

- $Y_i = \beta_0 + \beta_1 x_{i,1} + \dots + \beta_k x_{i,k} + \varepsilon_i, i \in \{1, \dots, n\}$ 
  - ▷  $Y_i$ : outcome, response, output, dependent variable
    - \* random variable, we observe a realization  $y_i$
    - \* (odezva, závisle proměnná, regresand)
  - ▷  $x_{i,1}, \dots, x_{i,k}$ : covariates, predictors, explanatory variables, input, independent variables
    - \* given, known
    - \* (nezávisle proměnné, regresory)
  - ▷  $\beta_0, \dots, \beta_k$ : coefficients
    - \* unknown
    - \* (regresní koeficienty)
  - ▷  $\varepsilon_i$ : random error
    - \* random variable, unobserved
- $\varepsilon_i \stackrel{\text{iid}}{\sim} (0, \sigma^2), i \in \{1, \dots, n\}$ 
  - ▷  $E \varepsilon_i = 0$ : no systematic errors
  - ▷  $\text{Var } \varepsilon_i = \sigma^2$ : same precision
- we often assume that  $\varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2), i \in \{1, \dots, n\}$

**Example: bloodpress data**

- from `sites.stat.psu.edu/~lSimon/stat501wc/sp05/data/`
- association between the mean arterial blood pressure[mmHg] and age[years], weight[kg], body surface area[m<sup>2</sup>], duration of hypertension[years], basal pulse[beats/min], stress

	BP	Age	Weight	BSA	DoH	Pulse	Stress
○ data:	105	47	85.4	1.75	5.1	63	33
	115	49	94.2	2.10	3.8	70	14
	...	...	...	...	...	...	...
	110	48	90.5	1.88	9.0	71	99
	122	56	95.7	2.09	7.0	75	99

- model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$

$$\begin{pmatrix} 105 \\ 115 \\ \dots \\ 110 \\ 122 \end{pmatrix} = \begin{pmatrix} 1 & 47 & 85.4 & 1.75 & 5.1 & 63 & 33 \\ 1 & 49 & 94.2 & 2.10 & 3.8 & 70 & 14 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & 48 & 90.5 & 1.88 & 9.0 & 71 & 99 \\ 1 & 56 & 95.7 & 2.09 & 7.0 & 75 & 99 \end{pmatrix} \times \begin{pmatrix} \beta_0 \\ \dots \\ \beta_6 \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \dots \\ \varepsilon_{19} \\ \varepsilon_{20} \end{pmatrix}$$

**4.1.2 Task for this chapter****Estimation in linear model**

- model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ 
  - ▷ outcome  $\mathbf{Y}$ 
    - \* random vector, we observe a realization  $\mathbf{y}$
  - ▷ predictors  $\mathbf{x}_1, \dots, \mathbf{x}_k$ 
    - \* vector of given (known) constants
  - ▷ coefficients  $\boldsymbol{\beta}$ 
    - \* vector of unknown constants
  - ▷ error  $\boldsymbol{\varepsilon}$ 
    - \* unknown random vector, we do not observe its realization
  - ▷ assumptions:  $\boldsymbol{\varepsilon} \sim (\mathbf{0}, \sigma^2 \mathbf{I})$ 
    - \*  $\mathbf{E} \mathbf{Y} = \mathbf{X}\boldsymbol{\beta}$ : the expected value of  $\mathbf{Y}$  is a linear function of  $\boldsymbol{\beta}$
    - \*  $\mathbf{E} \boldsymbol{\varepsilon} = \mathbf{0}$ : no systematic errors
    - \*  $\text{Var} \boldsymbol{\varepsilon} = \sigma^2 \mathbf{I}$ : independence and same precision
- task: given the observed data  $\mathbf{y}$  and known matrix  $\mathbf{X}$ , find estimators  $\hat{\boldsymbol{\beta}}$  (and  $\hat{\sigma}^2$ ) of  $\boldsymbol{\beta}$  (and  $\sigma^2$ ) with desirable properties

## 4.2 Estimating $\beta$

### 4.2.1 Orthogonal projection

#### $\hat{\beta}$ motivated by orthogonal projection

◦ **model:**  $\mathbf{Y} = \mathbf{X}\beta + \varepsilon$ ,  $\varepsilon$  unknown,  $E\varepsilon = \mathbf{0}$

◦ **idea:** set  $\varepsilon \stackrel{!}{=} \mathbf{0}$  and solve  $\mathbf{Y} = \mathbf{X}\beta$  w.r.t.  $\beta$

$$\triangleright \text{then } \underbrace{\mathbf{Y}}_{n \times 1} \stackrel{!}{=} \underbrace{\mathbf{X}}_{n \times p} \underbrace{\beta}_{p \times 1}$$

$\triangleright n$  linear equations with  $p$  unknowns and  $n > p$

$\Rightarrow$  a solution exists only if  $\mathbf{Y} \in \text{im}(\mathbf{X})$

◦ **modified idea:** find  $\hat{\mathbf{Y}} \in \text{im}(\mathbf{X})$  such that  $\|\mathbf{Y} - \hat{\mathbf{Y}}\|^2$  is the smallest possible and solve  $\hat{\mathbf{Y}} = \mathbf{X}\beta$  w.r.t.  $\beta$

$\triangleright$  then  $\hat{\mathbf{Y}}$  is the orthogonal projection of  $\mathbf{Y}$  onto  $\text{im}(\mathbf{X})$

$\triangleright$  projection matrix onto  $\text{im}(\mathbf{X})$  is  $\underbrace{\mathbf{H}}_{\text{hat matrix}} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$

$\triangleright$  solving  $\hat{\mathbf{Y}} = \mathbf{X}\beta$  is solving  $\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} = \mathbf{X}\beta$

$\triangleright$  estimate  $\beta$  by  $\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$

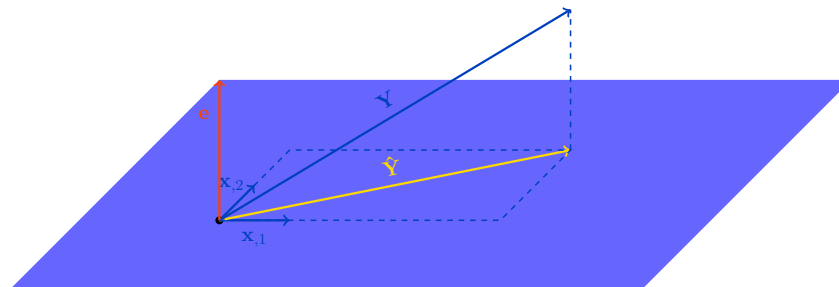
$\triangleright$  but  $\hat{\beta}$  is the unique solution of  $\hat{\mathbf{Y}} = \mathbf{X}\beta$  iff  $\text{rank}(\mathbf{X}) = p$

\* and then  $\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$

#### Geometric intuition

◦ **model:**  $\mathbf{Y} = \mathbf{X}\beta + \varepsilon$ ,  $\varepsilon$  unknown,  $E\varepsilon = \mathbf{0}$

◦ **fitted model:**  $\underbrace{\mathbf{Y}}_{\text{observed value}} = \underbrace{\mathbf{H}\mathbf{Y}}_{\text{fitted value } \hat{\mathbf{Y}}} + \underbrace{(\mathbf{I} - \mathbf{H})\mathbf{Y}}_{\text{residual } \mathbf{e}}$



◦  $\langle \hat{\mathbf{Y}}, \mathbf{e} \rangle = \mathbf{e}^\top \hat{\mathbf{Y}} = \mathbf{Y}^\top (\mathbf{I} - \mathbf{H})^\top \mathbf{H} \mathbf{Y} = 0$ , i.e.  $\hat{\mathbf{Y}} \perp \mathbf{e}$

## 4.2.2 Least squares

### $\hat{\beta}$ as least squares estimator

- **model:**  $\mathbf{Y} = \mathbf{X}\beta + \varepsilon$ ,  $\varepsilon$  unknown,  $E\varepsilon = \mathbf{0}$
- **idea:** make the residuals as small as possible
  - ▷ minimize  $\|\varepsilon\|^2 = \sum_{i=1}^n \varepsilon_i^2$  w.r.t.  $\beta$ 
    - ↪ Least Squares Estimator (LSE)  $\hat{\beta} = \arg \min_{\beta} \sum_{i=1}^n \varepsilon_i^2$
    - ▷ also called the OLS (Ordinary Least Squares) solution
- **computation:**
  - ▷  $\varepsilon = \mathbf{Y} - \mathbf{X}\beta$
  - ▷  $\hat{\beta} = \arg \min_{\beta} \|\mathbf{Y} - \mathbf{X}\beta\|^2 = \arg \min_{\beta} (\mathbf{Y} - \mathbf{X}\beta)^\top (\mathbf{Y} - \mathbf{X}\beta)$
- look for the minimum by differentiating:
 

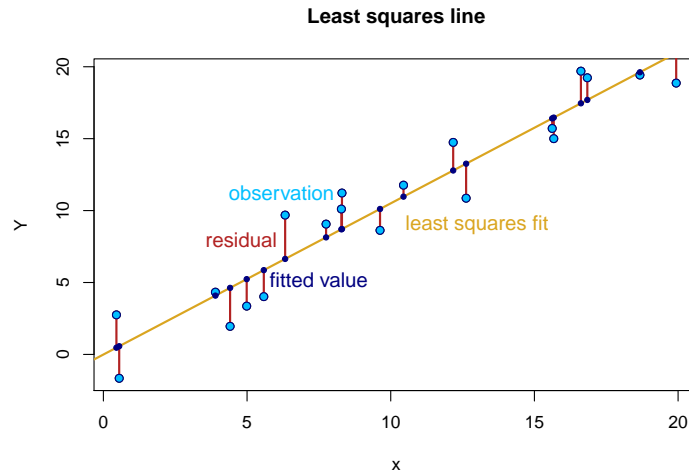
<ul style="list-style-type: none"> <li>▷ <math>\frac{\partial}{\partial \beta} (\mathbf{Y} - \mathbf{X}\beta)^\top (\mathbf{Y} - \mathbf{X}\beta) \stackrel{!}{=} 0</math></li> <li>▷ <math>-2\mathbf{X}^\top \mathbf{Y} + 2\mathbf{X}^\top \mathbf{X}\beta \stackrel{!}{=} 0</math></li> <li>▷ <math>\underbrace{\mathbf{X}^\top \mathbf{X}\beta \stackrel{!}{=} \mathbf{X}^\top \mathbf{Y}}_{\text{normal equations}}</math></li> </ul>	<ul style="list-style-type: none"> <li>▷ <math>\frac{\partial^2}{\partial \beta \partial \beta} (\mathbf{Y} - \mathbf{X}\beta)^\top (\mathbf{Y} - \mathbf{X}\beta) \stackrel{?}{\succ} 0</math></li> <li>at <math>\beta = \hat{\beta}</math></li> <li>▷ <math>2\mathbf{X}^\top \mathbf{X} \succ 0</math> for all <math>\beta</math></li> <li>▷ convex function <math>\Rightarrow</math> minimum</li> </ul>
---	--
- normal equations have unique solution iff  $\text{rank}(\mathbf{X}) = p$ 
  - ▷ and then  $\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$

### Geometric intuition

- **model:**  $\mathbf{Y} = \mathbf{X}\beta + \varepsilon$ ,  $\varepsilon$  unknown,  $E\varepsilon = \mathbf{0}$
- **fitted model:**

$$\underbrace{\mathbf{Y}}_{\text{observed value}} = \underbrace{\mathbf{X}\hat{\beta}}_{\text{fitted value } \hat{\mathbf{Y}}} + \underbrace{(\mathbf{Y} - \mathbf{X}\hat{\beta})}_{\text{residual } \mathbf{e}}$$
- least squares estimator minimizes the sum of squared vertical distances between the fitted and observed values





### 4.2.3 Computing $\hat{\beta}$

$$\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$$

- we have seen two approaches give the same  $\hat{\beta}$
- both approaches give unique  $\hat{\beta}$  iff  $\text{rank}(\mathbf{X}) = p$
- both approaches would give infinitely many  $\hat{\beta}$ s if  $\text{rank}(\mathbf{X}) < p$
- a rank-deficient design matrix means a problem in design/model formulation
- we need to fix that problem to obtain reasonable conclusions from our model
- from now on we assume that  $\text{rank}(\mathbf{X}) = p$
- we will get back to (nearly) rank-deficient  $\mathbf{X}$  in Chapter 9

### $\hat{\beta}$ the way it is computed in $\mathbb{R}$

- model:  $\mathbf{Y} = \mathbf{X}\beta + \varepsilon$ ,  $\varepsilon$  unknown,  $\mathbb{E} \varepsilon = \mathbf{0}$
- $\hat{\beta}$  minimizes  $\|\mathbf{Y} - \mathbf{X}\beta\|^2$  w.r.t.  $\beta$
- $\mathbb{R}$  uses that  $\mathbf{X} = \mathbf{Q}\mathbf{R}$  (QR decomposition from Chapter 2)
  - ▷  $\mathbf{Q}$  ( $n \times n$ ) orthogonal
  - ▷  $\mathbf{R}$  ( $n \times p$ ) upper triangular
  - ▷  $\mathbf{X} = \mathbf{Q}\mathbf{R} = (\mathbf{Q}_1 \mid \mathbf{Q}_2) \begin{pmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{pmatrix} = \mathbf{Q}_1 \mathbf{R}_1$

▷  $\mathbf{R}$  does not allow  $\text{rank}(\mathbf{X}) < p$

$$\Rightarrow \text{rank}(\mathbf{R}_1) = p$$

▷  $\mathbf{Q}$  and  $\mathbf{Q}^\top$  are rotations

$$\begin{aligned} * \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2 &= \|\mathbf{Q}^\top(\mathbf{Y} - \mathbf{Q}\mathbf{R}\boldsymbol{\beta})\|^2 = \left\| \begin{pmatrix} \mathbf{Q}_1^\top \\ \mathbf{Q}_2^\top \end{pmatrix} \mathbf{Y} - \begin{pmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{pmatrix} \boldsymbol{\beta} \right\|^2 \\ &= \|\mathbf{Q}_1^\top \mathbf{Y} - \mathbf{R}_1 \boldsymbol{\beta}\|^2 + \|\mathbf{Q}_2^\top \mathbf{Y}\|^2 \end{aligned}$$

▷ minimize  $\|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2 \Leftrightarrow$  minimize  $\|\mathbf{Q}_1^\top \mathbf{Y} - \mathbf{R}_1 \boldsymbol{\beta}\|^2$

▷  $\hat{\boldsymbol{\beta}} = \mathbf{R}_1^{-1} \mathbf{Q}_1^\top \mathbf{Y}$  (compare with  $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$ )

### Geometric intuition

◦ model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon}$  unknown,  $\mathbf{E}\boldsymbol{\varepsilon} = \mathbf{0}$

$$\underbrace{\mathbf{Y}}_{\text{observed value}} = \underbrace{\mathbf{X}\hat{\boldsymbol{\beta}}}_{\text{fitted value } \hat{\mathbf{Y}}} + \underbrace{(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}})}_{\text{residual } \mathbf{e}}$$

$$\begin{aligned} \triangleright \hat{\mathbf{Y}} &= \mathbf{X}\hat{\boldsymbol{\beta}} \\ &= \mathbf{Q} \left( \mathbf{Q}^\top \mathbf{Q} \mathbf{R} \hat{\boldsymbol{\beta}} \right) \\ &= \mathbf{Q} \left( \begin{pmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{pmatrix} \hat{\boldsymbol{\beta}} \right) \\ &= \mathbf{Q} \begin{pmatrix} \mathbf{R}_1 \hat{\boldsymbol{\beta}} \\ \mathbf{0} \end{pmatrix} \end{aligned}$$

$$\begin{aligned} \triangleright \mathbf{e} &= (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}) \\ &= \mathbf{Q} \left( \mathbf{Q}^\top (\mathbf{Y} - \mathbf{Q} \mathbf{R} \hat{\boldsymbol{\beta}}) \right) \\ &= \mathbf{Q} \left( \begin{pmatrix} \mathbf{Q}_1^\top \\ \mathbf{Q}_2^\top \end{pmatrix} \mathbf{Y} - \begin{pmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{pmatrix} \hat{\boldsymbol{\beta}} \right) \\ &= \mathbf{Q} \begin{pmatrix} \mathbf{0} \\ \mathbf{Q}_2^\top \mathbf{Y} \end{pmatrix} \end{aligned}$$

◦  $\mathbf{Q}^\top$  conveniently rotates  $\mathbf{Y}$  and  $\text{im}(\mathbf{X})$  and  $\mathbf{Q}$  rotates back

◦  $\hat{\mathbf{Y}} \perp \mathbf{e}$

### Geometric intuition

◦ model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon}$  unknown,  $\mathbf{E}\boldsymbol{\varepsilon} = \mathbf{0}$

$$\underbrace{\mathbf{Y}}_{\text{observed value}} = \underbrace{\mathbf{X}\hat{\boldsymbol{\beta}}}_{\text{fitted value } \hat{\mathbf{Y}}} + \underbrace{(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}})}_{\text{residual } \mathbf{e}} = \mathbf{Q} \begin{pmatrix} \mathbf{R}_1 \hat{\boldsymbol{\beta}} \\ \mathbf{0} \end{pmatrix} + \mathbf{Q} \begin{pmatrix} \mathbf{0} \\ \mathbf{Q}_2^\top \mathbf{Y} \end{pmatrix}$$

◦ see Figure 1.5 on page 20 in Simon Wood's *Generalized additive models* for a nice illustration

## 4.3 Quality of estimation

### 4.3.1 Gauss–Markov theorem

#### Linear transformation of a random vector

- we want to study  $\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$

**Theorem.** *Let  $\mathbf{X}$  be an  $n$ -dimensional random vector with a finite variance-covariance matrix and let  $\mathbf{A}$  be an  $m \times n$  matrix. Then*

- $E(\mathbf{A} \mathbf{X}) = \mathbf{A} E \mathbf{X}$ ;
- $\text{Var}(\mathbf{A} \mathbf{X}) = \mathbf{A} (\text{Var} \mathbf{X}) \mathbf{A}^\top$ ;
- $E(\mathbf{X}^\top \mathbf{X}) = (E \mathbf{X})^\top (E \mathbf{X}) + \text{tr}(\text{Var} \mathbf{X})$ .
- proof is a simple exercise

#### Is $\hat{\beta}$ a reasonable estimator?

- model:  $\mathbf{Y} = \mathbf{X}\beta + \varepsilon$ ,  $\varepsilon$  unknown,  $E \varepsilon = \mathbf{0}$
- $\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$
- has a nice motivation but how about properties?

$$\begin{aligned} \triangleright E \hat{\beta} &= E(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top E \mathbf{Y} \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} \beta = \beta \end{aligned}$$

$\Rightarrow$  unbiased

$$\begin{aligned} \triangleright \text{Var} \hat{\beta} &= \text{Var}((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}) \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \text{Var} \mathbf{Y} ((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top)^\top \\ &= \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \\ &= \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} \end{aligned}$$

- how good is  $\text{Var} \hat{\beta}$ ?

$\triangleright \hat{\beta}$  is a linear estimator, i.e.  $\hat{\beta} = \mathbf{A} \mathbf{Y}$  for a matrix  $\mathbf{A}$

$\triangleright \hat{\beta}$  is an unbiased estimator, i.e.  $E_\beta \hat{\beta} = \beta$  for all  $\beta$

$\triangleright$  in fact,  $\hat{\beta}$  is the best linear unbiased estimator of  $\beta$ , i.e.  $\hat{\beta}$  has the smallest variance among all linear unbiased estimators of  $\beta$

#### Gauss–Markov theorem

**Theorem** (Gauss–Markov). Let  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  where  $\mathbf{X}$  is an  $n \times p$  matrix,  $\text{rank}(\mathbf{X}) = p$ ,  $\boldsymbol{\beta} \in \mathbb{R}^p$ , and  $\boldsymbol{\varepsilon}$  is an  $n$ -dimensional random vector with  $\mathbf{E}\boldsymbol{\varepsilon} = \mathbf{0}$  and  $\text{Var}\boldsymbol{\varepsilon} = \sigma^2\mathbf{I}$ . Then  $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\mathbf{Y}$  is the best linear unbiased estimator of  $\boldsymbol{\beta}$ , i.e. if  $\tilde{\boldsymbol{\beta}}$  is a linear unbiased estimator of  $\boldsymbol{\beta}$  then  $\text{Var}\tilde{\boldsymbol{\beta}} - \text{Var}\hat{\boldsymbol{\beta}} \succeq 0$ .

◦ see the blackboard for a proof

▷ main steps

\* show that if  $\tilde{\boldsymbol{\beta}} = \mathbf{A}\mathbf{Y}$  then  $\mathbf{A}\mathbf{X} = \mathbf{I}$

\* show that  $\text{Var}\tilde{\boldsymbol{\beta}} - \text{Var}\hat{\boldsymbol{\beta}} = \sigma^2\mathbf{A}(\mathbf{I} - \mathbf{H})(\mathbf{I} - \mathbf{H})^\top\mathbf{A}^\top$

## 4.4 Estimating $\sigma^2$

### 4.4.1 Estimating $\sigma^2$

#### Estimating $\sigma^2$

◦ model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon}$  unknown,  $\mathbf{E}\boldsymbol{\varepsilon} = \mathbf{0}$ ,  $\text{Var}\boldsymbol{\varepsilon} = \sigma^2\mathbf{I}$

◦ fitted model:  $\mathbf{Y} = \underbrace{\mathbf{H}\mathbf{Y}}_{\hat{\mathbf{Y}}} + \underbrace{(\mathbf{I} - \mathbf{H})\mathbf{Y}}_{\mathbf{e}} = \underbrace{\mathbf{X}\hat{\boldsymbol{\beta}}}_{\hat{\mathbf{Y}}} + \underbrace{(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}})}_{\mathbf{e}}$

◦ idea: estimate  $\boldsymbol{\varepsilon}$  by  $\mathbf{e}$

▷ some care is needed ...

\*  $\mathbf{E}\mathbf{e} = \mathbf{E}(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = \mathbf{E}\mathbf{Y} - \mathbf{X}\mathbf{E}\hat{\boldsymbol{\beta}} = \mathbf{X}\boldsymbol{\beta} - \mathbf{X}\boldsymbol{\beta} = \mathbf{0}$

\*  $\text{Var}\mathbf{e} = \text{Var}((\mathbf{I} - \mathbf{H})\mathbf{Y}) = (\mathbf{I} - \mathbf{H})\text{Var}\mathbf{Y}(\mathbf{I} - \mathbf{H})^\top = \sigma^2(\mathbf{I} - \mathbf{H})$

\*  $\text{rank}((\mathbf{I} - \mathbf{H})) = n - \text{rank}(\mathbf{X}) = n - p < n \Rightarrow$  dependence

▷  $\mathbf{E}(\mathbf{e}^\top\mathbf{e}) = (\mathbf{E}\mathbf{e})^\top(\mathbf{E}\mathbf{e}) + \text{tr}(\text{Var}\mathbf{e})$

$$= \text{tr}(\sigma^2(\mathbf{I} - \mathbf{H}))$$

$$\stackrel{*}{=} \sigma^2(n - \text{rank}(\mathbf{X})) = \sigma^2(n - p)$$

\* \*:  $\text{tr}(\mathbf{P}) = \text{rank}(\mathbf{P})$  for orthogonal projection matrices

▷  $\hat{\sigma}^2 = \frac{1}{n-p}\mathbf{e}^\top\mathbf{e} = \frac{1}{n-p}\sum_{i=1}^n e_i^2 = \frac{1}{n-p}\sum_{i=1}^n (Y_i - \hat{Y}_i)^2$

\* unbiased estimator of  $\sigma^2$

## 4.5 Quality of model fit

### 4.5.1 Coefficient of determination

#### Sums of squares

- for  $\hat{\beta}$  we obtain the minimal  $\|\mathbf{e}\|^2 = \|\mathbf{Y} - \hat{\mathbf{Y}}\|^2 = \|\mathbf{Y} - \mathbf{X}\hat{\beta}\|^2$
- we have seen properties of  $\hat{\beta}$  but how about  $\hat{\mathbf{Y}}$ ?
- a question: how close  $\hat{\mathbf{Y}}$  actually is to  $\mathbf{Y}$ ?
  - ▷ how well do the covariates in  $\mathbf{X}$  explain what we see in  $\mathbf{Y}$ ?
- an answer:
  - ▷ there is some variability in  $Y_i$ s for different  $i$ 
    - \* Total Sum of Squares TSS:  $\sum_{i=1}^n (Y_i - \bar{Y})^2$
    - \* also called SST
  - ▷ the model explains a part of the variability in  $Y_i$ s
    - \* for different  $i$ s there are different  $\mathbf{x}_i$ s and so different  $\hat{Y}_i$ s
    - \* Explained Sum of Squares ESS:  $\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2$
    - \* also called Sum of Squares due to Regression
  - ▷ but some variability remained unexplained by the model
    - \* Residual Sum of Squares RSS:  $\sum_{i=1}^n (Y_i - \hat{Y}_i)^2$
    - \* also called Sum of Squared Residuals or Sum of Squared Errors

## Coefficient of determination $R^2$

- relationship among the sums of squares
  - ▷ TSS = RSS + ESS
    - \*  $\|\mathbf{Y} - \bar{Y}\mathbf{1}\|^2 = \|\mathbf{Y} - \hat{\mathbf{Y}} + \hat{\mathbf{Y}} - \bar{Y}\mathbf{1}\|^2 = \|\mathbf{Y} - \hat{\mathbf{Y}}\|^2 + \|\hat{\mathbf{Y}} - \bar{Y}\mathbf{1}\|^2$
    - \* because  $\langle \mathbf{Y} - \hat{\mathbf{Y}}, \hat{\mathbf{Y}} \rangle = 0 = \langle \mathbf{Y} - \hat{\mathbf{Y}}, \mathbf{1} \rangle$  as  $\hat{\mathbf{Y}}$  is the orthogonal projection of  $\mathbf{Y}$  onto  $\text{im}(\mathbf{X})$  and  $\mathbf{1} \in \text{im}(\mathbf{X})$
  - ▷ variability: total = unexplained + explained
- so how well do the covariates in  $\mathbf{X}$  explain what we see in  $\mathbf{Y}$ ?
  - ▷ coefficient of determination  $R^2 = \frac{ESS}{TSS} = 1 - \frac{RSS}{TSS}$ 
    - \* proportion of variability explained by the model
    - \*  $0 \leq R^2 \leq 1$  and bigger is better
  - ▷ adjusted coefficient of determination  $R_{adj}^2 = 1 - \frac{RSS/(n-p)}{TSS/(n-1)}$ 
    - \* an alternative that takes the number of predictors into account
    - \*  $RSS/(n-p) = \hat{\sigma}^2$  from the linear regression,
    - \*  $TSS/(n-1) = \hat{\sigma}^2$  without the linear regression

# Chapter 5

## Normal linear model

### 5.1 The problem

#### 5.1.1 Normal linear model

##### Normal linear model

- $Y_i = \beta_0 + \beta_1 x_{i,1} + \dots + \beta_k x_{i,k} + \varepsilon_i, i \in \{1, \dots, n\}$ 
  - ▷  $Y_i$ : outcome, response, output, dependent variable
    - \* random variable, we observe a realization  $y_i$
    - \* (odezva, závisle proměnná, regresand)
  - ▷  $x_{i,1}, \dots, x_{i,k}$ : covariates, predictors, explanatory variables, input, independent variables
    - \* given, known
    - \* (nezávisle proměnné, regresory)
  - ▷  $\beta_0, \dots, \beta_k$ : coefficients
    - \* unknown
    - \* (regresní koeficienty)
  - ▷  $\varepsilon_i$ : random error
    - \* random variable, unobserved
- $\varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2), i \in \{1, \dots, n\}$ 
  - ▷  $E \varepsilon_i = 0$ : no systematic errors
  - ▷  $\text{Var } \varepsilon_i = \sigma^2$ : same precision

##### Example: bloodpress data

- from `sites.stat.psu.edu/~lsimon/stat501wc/sp05/data/`
- association between the mean arterial blood pressure[mmHg] and age[years], weight[kg], body surface area[m<sup>2</sup>], duration of hypertension[years], basal pulse[beats/min], stress

	BP	Age	Weight	BSA	DoH	Pulse	Stress
○ data:	105	47	85.4	1.75	5.1	63	33
	115	49	94.2	2.10	3.8	70	14
	...	...	...	...	...	...	...
	110	48	90.5	1.88	9.0	71	99
	122	56	95.7	2.09	7.0	75	99

- model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$

$$\begin{pmatrix} 105 \\ 115 \\ \dots \\ 110 \\ 122 \end{pmatrix} = \begin{pmatrix} 1 & 47 & 85.4 & 1.75 & 5.1 & 63 & 33 \\ 1 & 49 & 94.2 & 2.10 & 3.8 & 70 & 14 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & 48 & 90.5 & 1.88 & 9.0 & 71 & 99 \\ 1 & 56 & 95.7 & 2.09 & 7.0 & 75 & 99 \end{pmatrix} \times \begin{pmatrix} \beta_0 \\ \dots \\ \beta_6 \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \dots \\ \varepsilon_{19} \\ \varepsilon_{20} \end{pmatrix}$$

## 5.1.2 Task for this chapter

### Estimation in normal linear model

- model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$

▷ outcome  $\mathbf{Y}$

\* random vector, we observe a realization  $\mathbf{y}$

▷ predictors  $\mathbf{x}_1, \dots, \mathbf{x}_k$

\* vector of given (known) constants

▷ coefficients  $\boldsymbol{\beta}$

\* vector of unknown constants

▷ error  $\boldsymbol{\varepsilon}$

\* unknown random vector, we do not observe its realization

▷ assumptions:  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$

\*  $E\mathbf{Y} = \mathbf{X}\boldsymbol{\beta}$ : the expected value of  $\mathbf{Y}$  is a linear function of  $\boldsymbol{\beta}$

\*  $E\boldsymbol{\varepsilon} = \mathbf{0}$ : no systematic errors

\*  $\text{Var}\boldsymbol{\varepsilon} = \sigma^2 \mathbf{I}$ : independence and same precision

- task: given the observed data  $\mathbf{y}$  and known matrix  $\mathbf{X}$ , find estimators  $\hat{\boldsymbol{\beta}}$  and  $\hat{\sigma}^2$  of  $\boldsymbol{\beta}$  and  $\sigma^2$  with desirable properties

## 5.2 Estimating $\boldsymbol{\beta}$ and $\sigma^2$

### 5.2.1 Likelihood

#### Likelihood

- **model:**  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2\mathbf{I})$ 
  - ▷  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2\mathbf{I}) \Rightarrow \mathbf{Y} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I})$
  - ▷ density of  $\mathbf{Y}$ :

$$f(\mathbf{y}; \boldsymbol{\beta}, \sigma^2) = \left( \frac{1}{\sqrt{2\pi\sigma^2}} \right)^n \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$$

- ▷ density is a function of  $\mathbf{y}$  (parameters are fixed)
- **likelihood:**

$$L(\boldsymbol{\beta}, \sigma^2; \mathbf{y}) = \left( \frac{1}{\sqrt{2\pi\sigma^2}} \right)^n \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$$

- likelihood is a function of the parameters ( $\mathbf{y}$  is fixed)
- **log-likelihood:**

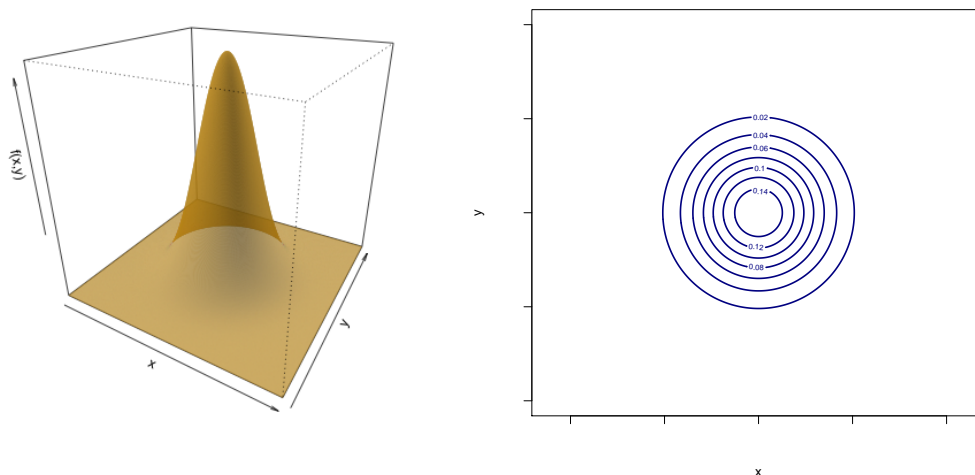
$$\ell(\boldsymbol{\beta}, \sigma^2; \mathbf{y}) = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) - \left\{ \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$$

#### Log-likelihood

- **model:**  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2\mathbf{I})$
- **log-likelihood:**

$$\ell(\boldsymbol{\beta}, \sigma^2; \mathbf{y}) = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) - \left\{ \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$$





## 5.2.2 Matrix derivatives

### Matrix derivatives: definition

- let  $\mathbf{x} \in \mathbb{R}^n$  and  $\mathbf{y} \in \mathbb{R}^m$
- denominator-layout notation:

$$\frac{\partial}{\partial \mathbf{x}} \mathbf{y} = \begin{pmatrix} \frac{\partial}{\partial x_1} y_1 & \cdots & \frac{\partial}{\partial x_1} y_m \\ \cdots & \cdots & \cdots \\ \frac{\partial}{\partial x_n} y_1 & \cdots & \frac{\partial}{\partial x_n} y_m \end{pmatrix}$$

- if  $n = 1$

$$\frac{\partial}{\partial x} \mathbf{y} = \left( \frac{\partial}{\partial x} y_1, \dots, \frac{\partial}{\partial x} y_m \right)$$

- if  $m = 1$

$$\frac{\partial}{\partial \mathbf{x}} y = \begin{pmatrix} \frac{\partial}{\partial x_1} y \\ \cdots \\ \frac{\partial}{\partial x_n} y \end{pmatrix}$$

### Matrix derivatives: useful formulae

- let  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $\mathbf{x} \in \mathbb{R}^n$

$$\triangleright \frac{\partial}{\partial \mathbf{x}} \mathbf{A} \mathbf{x} = \mathbf{A}^\top$$

$$\triangleright \frac{\partial}{\partial \mathbf{x}} \mathbf{x}^\top \mathbf{A} = \mathbf{A}$$

$$\triangleright \frac{\partial}{\partial \mathbf{x}} \mathbf{x}^\top \mathbf{A} \mathbf{x} = (\mathbf{A} + \mathbf{A}^\top) \mathbf{x}$$

### 5.2.3 Maximizing the likelihood

#### Score function

- log-likelihood

$$\ell(\boldsymbol{\beta}, \sigma^2; \mathbf{y}) = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) - \left\{ \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$$

- score function ( $\boldsymbol{\beta}$ -related part):

$$\begin{aligned} \mathbf{U}_{1:p}(\boldsymbol{\beta}, \sigma^2; \mathbf{y}) &= \frac{\partial}{\partial \boldsymbol{\beta}} \left( -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right) \\ &= \frac{\partial}{\partial \boldsymbol{\beta}} \left( -\frac{1}{2\sigma^2} (\mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X}\boldsymbol{\beta} - \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{y} + \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X}\boldsymbol{\beta}) \right) \\ &= -\frac{1}{2\sigma^2} \left( -\mathbf{X}^\top \mathbf{y} - \mathbf{X}^\top \mathbf{y} + (\mathbf{X}^\top \mathbf{X} + \mathbf{X}^\top \mathbf{X}) \boldsymbol{\beta} \right) \\ &= \frac{1}{\sigma^2} \left( \mathbf{X}^\top \mathbf{y} - \mathbf{X}^\top \mathbf{X}\boldsymbol{\beta} \right) \end{aligned}$$

#### Score function ctd.

- log-likelihood

$$\ell(\boldsymbol{\beta}, \sigma^2; \mathbf{y}) = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) - \left\{ \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$$

- score function ( $\sigma^2$ -related part):

$$\begin{aligned} U_{p+1}(\boldsymbol{\beta}, \sigma^2; \mathbf{y}) &= \frac{\partial}{\partial \sigma^2} \left( -\frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right) \\ &= -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \\ &= \frac{1}{2\sigma^2} \left( -n + \frac{1}{\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right) \end{aligned}$$

#### Score equation

- score equation:

$$\mathbf{U}(\boldsymbol{\beta}, \sigma^2; \mathbf{y}) = \left( \begin{array}{c} \frac{1}{\sigma^2} (\mathbf{X}^\top \mathbf{y} - \mathbf{X}^\top \mathbf{X}\boldsymbol{\beta}) \\ \frac{1}{2\sigma^2} \left( -n + \frac{1}{\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right) \end{array} \right) = \mathbf{0}$$

- score equation for  $\beta$

$$\triangleright \frac{1}{\sigma^2} (\mathbf{X}^\top \mathbf{y} - \mathbf{X}^\top \mathbf{X} \beta) \stackrel{!}{=} \mathbf{0}$$

- ▷ actually the normal equations

$$\triangleright \hat{\beta}_{\text{MLE}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$$

- score equation for  $\sigma^2$

$$\triangleright \frac{1}{2\sigma^2} \left( -n + \frac{1}{\sigma^2} (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta) \right) \stackrel{!}{=} 0$$

$$\triangleright \hat{\sigma}_{\text{MLE}}^2 = \frac{1}{n} (\mathbf{Y} - \mathbf{X} \hat{\beta}_{\text{MLE}})^\top (\mathbf{Y} - \mathbf{X} \hat{\beta}_{\text{MLE}})$$

### Fisher information

- observed Fisher information matrix:

$$\mathbf{J}(\beta, \sigma^2; \mathbf{y}) = \frac{1}{\sigma^2} \begin{pmatrix} \mathbf{X}^\top \mathbf{X} & \frac{1}{\sigma^2} (\mathbf{X}^\top \mathbf{y} - \mathbf{X}^\top \mathbf{X} \beta) \\ \frac{1}{\sigma^2} (\mathbf{X}^\top \mathbf{y} - \mathbf{X}^\top \mathbf{X} \beta) & -\frac{1}{\sigma^2} \left( \frac{n}{2} - \frac{1}{\sigma^2} (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta) \right) \end{pmatrix}$$

- $\mathbf{J}(\hat{\beta}_{\text{MLE}}, \hat{\sigma}_{\text{MLE}}^2)$ :

$$\frac{1}{\hat{\sigma}_{\text{MLE}}^2} \begin{pmatrix} \mathbf{X}^\top \mathbf{X} & \mathbf{0} \\ \mathbf{0} & \frac{n}{2\hat{\sigma}_{\text{MLE}}^2} \end{pmatrix} \succ 0$$

- Fisher information matrix:

$$\mathbf{I}(\beta, \sigma^2) = \frac{1}{\sigma^2} \begin{pmatrix} \mathbf{X}^\top \mathbf{X} & \mathbf{0} \\ \mathbf{0} & \frac{n}{2\sigma^2} \end{pmatrix}$$

## 5.3 Distribution

### 5.3.1 Distribution of the MLE

#### Distribution of $\hat{\beta}_{\text{MLE}}$

- model:  $\mathbf{Y} = \mathbf{X}\beta + \varepsilon$ ,  $\varepsilon \sim \text{N}(\mathbf{0}, \sigma^2 \mathbf{I})$
- $\mathbf{Y} \sim \text{N}(\mathbf{X}\beta, \sigma^2 \mathbf{I})$
- distribution of  $\hat{\beta}_{\text{MLE}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$ ?
- MVN 3:

Let  $\mathbf{X} \sim \text{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  and let  $\mathbf{A}$  be an  $m \times n$  real matrix and  $\mathbf{b} \in \mathbb{R}^m$ . Then  $\mathbf{A}\mathbf{X} + \mathbf{b} \sim \text{N}(\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top)$ .

- $\hat{\boldsymbol{\beta}}_{\text{MLE}} \sim N(\boldsymbol{\beta}, \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1})$

### Distribution of $\hat{\sigma}_{\text{MLE}}^2$

- model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$

- $\mathbf{Y} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$

- distribution of  $\hat{\sigma}_{\text{MLE}}^2 = \frac{1}{n}(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}_{\text{MLE}})^\top (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}_{\text{MLE}})$ ?

- recall that

- ▷  $\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{H}\mathbf{Y} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$

- ▷  $(\mathbf{Y} - \hat{\mathbf{Y}}) = \mathbf{e} = (\mathbf{I} - \mathbf{H}) \mathbf{Y}$

- \*  $\mathbf{e} \sim N(\mathbf{0}, \sigma^2(\mathbf{I} - \mathbf{H}))$  (by MVN 3)

- QF 3:

Let  $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  where  $\text{rank}(\boldsymbol{\Sigma}) = r < n$ . Then  $(\mathbf{X} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^+ (\mathbf{X} - \boldsymbol{\mu}) \sim \chi^2(r)$ .

- $(\mathbf{I} - \mathbf{H})^+ = (\mathbf{I} - \mathbf{H})$

- $\Rightarrow \frac{1}{\sigma^2} \mathbf{e}^\top (\mathbf{I} - \mathbf{H}) \mathbf{e} \sim \chi^2(n - p)$

- $\mathbf{e}^\top (\mathbf{I} - \mathbf{H}) \mathbf{e} = \mathbf{Y}^\top (\mathbf{I} - \mathbf{H})^\top (\mathbf{I} - \mathbf{H}) (\mathbf{I} - \mathbf{H}) \mathbf{Y} = \mathbf{e}^\top \mathbf{e}$

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

### Relationship between $\hat{\boldsymbol{\beta}}_{\text{MLE}}$ and $\hat{\sigma}_{\text{MLE}}^2$

- model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$

- $\mathbf{Y} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$

- $\hat{\boldsymbol{\beta}}_{\text{MLE}} \sim N(\boldsymbol{\beta}, \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1})$

- $n \hat{\sigma}_{\text{MLE}}^2 \sim \chi^2(n - p)$

- joint distribution of  $\hat{\boldsymbol{\beta}}_{\text{MLE}}$  and  $\hat{\sigma}_{\text{MLE}}^2$ ?

- recall that

- ▷  $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$

- ▷  $(\mathbf{Y} - \hat{\mathbf{Y}}) = \mathbf{e} = (\mathbf{I} - \mathbf{H}) \mathbf{Y}$

- Corollary of MVN 7:

Let  $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ . Then  $\mathbf{AX} \perp\!\!\!\perp \mathbf{BX}$  iff  $\mathbf{A}\boldsymbol{\Sigma}\mathbf{B}^\top = \mathbf{0}$ .

- 

$$(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{I} - \mathbf{H})^\top = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top = \mathbf{0}$$

- $\hat{\boldsymbol{\beta}} \perp\!\!\!\perp \mathbf{e}$  and  $\hat{\boldsymbol{\beta}} \perp\!\!\!\perp \hat{\sigma}_{\text{MLE}}^2$

## 5.4 Summary

### 5.4.1 Estimation in the normal linear model

#### Estimation in the normal linear model

**Theorem.** Let  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  where  $\mathbf{X}$  is an  $n \times p$  matrix,  $\text{rank}(\mathbf{X}) = p$ ,  $\boldsymbol{\beta} \in \mathbb{R}^p$ , and  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$ .

Then the maximum likelihood estimators of  $\boldsymbol{\beta}$  and  $\sigma^2$  are given by  $\hat{\boldsymbol{\beta}}_{\text{MLE}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$  and  $\hat{\sigma}_{\text{MLE}}^2 = \frac{1}{n} (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}}_{\text{MLE}})^\top (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}}_{\text{MLE}})$ .

Their distributions are  $\hat{\boldsymbol{\beta}}_{\text{MLE}} \sim N(\boldsymbol{\beta}, \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1})$  and  $\frac{n}{\sigma^2} \hat{\sigma}_{\text{MLE}}^2 \sim \chi^2(n-p)$ , and  $\hat{\boldsymbol{\beta}}_{\text{MLE}}$  and  $\hat{\sigma}_{\text{MLE}}^2$  are independent.

- unbiased estimator of  $\sigma^2$ :  $\hat{\sigma}^2 = \frac{1}{n-p} (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}}_{\text{MLE}})^\top (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}}_{\text{MLE}})$
- its distribution:  $\frac{(n-p)}{\sigma^2} \hat{\sigma}^2 \sim \chi^2(n-p)$  and  $\hat{\boldsymbol{\beta}} \perp\!\!\!\perp \hat{\sigma}^2$

# Chapter 6

## Inference in normal linear model

### 6.1 The problem

#### 6.1.1 Normal linear model

##### Normal linear model

- $Y_i = \beta_0 + \beta_1 x_{i,1} + \dots + \beta_k x_{i,k} + \varepsilon_i, i \in \{1, \dots, n\}$ 
  - ▷  $Y_i$ : outcome, response, output, dependent variable
    - \* random variable, we observe a realization  $y_i$
    - \* (odezva, závisle proměnná, regresand)
  - ▷  $x_{i,1}, \dots, x_{i,k}$ : covariates, predictors, explanatory variables, input, independent variables
    - \* given, known
    - \* (nezávisle proměnné, regresory)
  - ▷  $\beta_0, \dots, \beta_k$ : coefficients
    - \* unknown
    - \* (regresní koeficienty)
  - ▷  $\varepsilon_i$ : random error
    - \* random variable, unobserved
- $\varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2), i \in \{1, \dots, n\}$ 
  - ▷  $E \varepsilon_i = 0$ : no systematic errors
  - ▷  $\text{Var } \varepsilon_i = \sigma^2$ : same precision

##### Example: bloodpress data

- from `sites.stat.psu.edu/~lsimon/stat501wc/sp05/data/`
- association between the mean arterial blood pressure[mmHg] and age[years], weight[kg], body surface area[m<sup>2</sup>], duration of hypertension[years], basal pulse[beats/min], stress

	BP	Age	Weight	BSA	DoH	Pulse	Stress
○ data:	105	47	85.4	1.75	5.1	63	33
	115	49	94.2	2.10	3.8	70	14
	...	...	...	...	...	...	...
	110	48	90.5	1.88	9.0	71	99
	122	56	95.7	2.09	7.0	75	99

- model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$

$$\begin{pmatrix} 105 \\ 115 \\ \dots \\ 110 \\ 122 \end{pmatrix} = \begin{pmatrix} 1 & 47 & 85.4 & 1.75 & 5.1 & 63 & 33 \\ 1 & 49 & 94.2 & 2.10 & 3.8 & 70 & 14 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & 48 & 90.5 & 1.88 & 9.0 & 71 & 99 \\ 1 & 56 & 95.7 & 2.09 & 7.0 & 75 & 99 \end{pmatrix} \times \begin{pmatrix} \beta_0 \\ \dots \\ \beta_6 \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \dots \\ \varepsilon_{19} \\ \varepsilon_{20} \end{pmatrix}$$

## 6.1.2 Task for this chapter

### Inference in normal linear model

- model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$

▷ outcome  $\mathbf{Y}$

\* random vector, we observe a realization  $\mathbf{y}$

▷ predictors  $\mathbf{x}_1, \dots, \mathbf{x}_k$

\* vector of given (known) constants

▷ coefficients  $\boldsymbol{\beta}$

\* vector of unknown constants

▷ error  $\boldsymbol{\varepsilon}$

\* unknown random vector, we do not observe its realization

▷ assumptions:  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$

\*  $E\mathbf{Y} = \mathbf{X}\boldsymbol{\beta}$ : the expected value of  $\mathbf{Y}$  is a linear function of  $\boldsymbol{\beta}$

\*  $E\boldsymbol{\varepsilon} = \mathbf{0}$ : no systematic errors

\*  $\text{Var}\boldsymbol{\varepsilon} = \sigma^2 \mathbf{I}$ : independence and same precision

- task: given the observed data  $\mathbf{y}$  and known matrix  $\mathbf{X}$ , draw conclusions about  $\mathbf{Y}$  and the relationship between  $\mathbf{Y}$  and  $\mathbf{X}$

## 6.2 Estimators and distributions

### 6.2.1 Estimators

#### Point estimation in the normal linear model

- model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ 
  - ▷  $\mathbf{X}$  is an  $n \times p$  matrix,  $\text{rank}(\mathbf{X}) = p$
  - ▷  $\boldsymbol{\beta} \in \mathbb{R}^p$
  - ▷  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$
- estimating  $\boldsymbol{\beta}$ 
  - ▷  $\hat{\boldsymbol{\beta}}_{\text{MLE}} = \hat{\boldsymbol{\beta}}_{\text{OLS}} = \hat{\boldsymbol{\beta}}_{\text{MOM}} = \hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$ 
    - \* BLUE
    - \* distribution:  $\hat{\boldsymbol{\beta}} \sim N(\boldsymbol{\beta}, \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1})$
- estimating  $\sigma^2$ 
  - ▷  $\hat{\sigma}_{\text{MLE}}^2 = \frac{1}{n} (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}})^\top (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}})$ 
    - \* distribution:  $\frac{n}{\sigma^2} \hat{\sigma}_{\text{MLE}}^2 \sim \chi^2(n - p)$
  - ▷  $\hat{\sigma}^2 = \frac{1}{n-p} (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}})^\top (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}})$ 
    - \* unbiased
    - \* distribution:  $\frac{(n-p)}{\sigma^2} \hat{\sigma}^2 \sim \chi^2(n - p)$
- $\hat{\boldsymbol{\beta}} \perp\!\!\!\perp \hat{\sigma}_{\text{MLE}}^2$  and  $\hat{\boldsymbol{\beta}} \perp\!\!\!\perp \hat{\sigma}^2$

### 6.2.2 Distributions

#### Distributions in normal linear model

- model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$
- distributions of point estimators
  - ▷  $\hat{\boldsymbol{\beta}} \sim N(\boldsymbol{\beta}, \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1})$
  - ▷  $\frac{(n-p)}{\sigma^2} \hat{\sigma}^2 \sim \chi^2(n - p)$
  - ▷  $\hat{\boldsymbol{\beta}} \perp\!\!\!\perp \hat{\sigma}^2$
- let  $\mathbf{a} \in \mathbb{R}^p$  and  $\mathbf{A} \in \mathbb{R}^{m \times p}$



$$\triangleright \mathbf{a}^\top \widehat{\boldsymbol{\beta}} \sim N(\mathbf{a}^\top \boldsymbol{\beta}, \sigma^2 \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a})$$

$$\triangleright \mathbf{A} \widehat{\boldsymbol{\beta}} \sim N(\mathbf{A} \boldsymbol{\beta}, \sigma^2 \mathbf{A} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{A}^\top)$$

\* proofs: use MVN 3:

Let  $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  and let  $\mathbf{A}$  be an  $m \times n$  real matrix and  $\mathbf{b} \in \mathbb{R}^m$ . Then  $\mathbf{A}\mathbf{X} + \mathbf{b} \sim N(\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top)$ .

### Distributions in normal linear model ctd.

◦ model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$

◦ distributions of statistics

$$\triangleright \mathbf{a}^\top \widehat{\boldsymbol{\beta}} \sim N(\mathbf{a}^\top \boldsymbol{\beta}, \sigma^2 \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a}) \text{ for } \mathbf{a} \in \mathbb{R}^p$$

$$\triangleright \mathbf{A} \widehat{\boldsymbol{\beta}} \sim N(\mathbf{A} \boldsymbol{\beta}, \sigma^2 \mathbf{A} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{A}^\top) \text{ for } \mathbf{A} \in \mathbb{R}^{m \times p}$$

$$\triangleright \frac{(n-p)}{\sigma^2} \widehat{\sigma}^2 \sim \chi^2(n-p)$$

$$\triangleright \widehat{\boldsymbol{\beta}} \perp\!\!\!\perp \widehat{\sigma}^2$$

◦ for  $\mathbf{a} \in \mathbb{R}^p$  and  $\mathbf{A} \in \mathbb{R}^{m \times p}$ ,  $\text{rank}(\mathbf{A}) = m$

$$\triangleright \frac{\mathbf{a}^\top \widehat{\boldsymbol{\beta}} - \mathbf{a}^\top \boldsymbol{\beta}}{\sqrt{\widehat{\sigma}^2 \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a}}} \sim t(n-p)$$

\* proof: verify that the definition of  $t(n-p)$  is satisfied

$$\triangleright \frac{1}{m\widehat{\sigma}^2} (\mathbf{A} \widehat{\boldsymbol{\beta}} - \mathbf{A} \boldsymbol{\beta})^\top (\mathbf{A} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{A}^\top)^{-1} (\mathbf{A} \widehat{\boldsymbol{\beta}} - \mathbf{A} \boldsymbol{\beta}) \sim F(m, n-p)$$

\* proof: use QF2:

$$\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma}), \text{rank}(\boldsymbol{\Sigma}) = n \Rightarrow (\mathbf{X} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{X} - \boldsymbol{\mu}) \sim \chi^2(n)$$

and verify that the definition of  $F(m, n-p)$  is satisfied

## 6.3 Confidence intervals

### Confidence intervals

#### Interval estimation in normal linear model

◦ model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$

◦ let  $\mathbf{a} \in \mathbb{R}^p$

$$\triangleright \frac{\mathbf{a}^\top \hat{\boldsymbol{\beta}} - \mathbf{a}^\top \boldsymbol{\beta}}{\sqrt{\hat{\sigma}^2 \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a}}} \sim t(n-p)$$

$\triangleright (1 - \alpha) \times 100\%$  confidence interval for  $\mathbf{a}^\top \boldsymbol{\beta}$ :

$$\left( \mathbf{a}^\top \hat{\boldsymbol{\beta}} - t_{1-\alpha/2}(n-p) \sqrt{\hat{\sigma}^2 \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a}}, \right. \\ \left. \mathbf{a}^\top \hat{\boldsymbol{\beta}} + t_{1-\alpha/2}(n-p) \sqrt{\hat{\sigma}^2 \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a}} \right)$$

○  $\frac{(n-p)}{\hat{\sigma}^2} \hat{\sigma}^2 \sim \chi^2(n-p)$

$\triangleright (1 - \alpha) \times 100\%$  confidence interval for  $\sigma^2$ :

$$\left( \frac{(n-p) \hat{\sigma}^2}{\chi_{1-\alpha/2}^2(n-p)}, \frac{(n-p) \hat{\sigma}^2}{\chi_{\alpha/2}^2(n-p)} \right)$$

### Confidence intervals for the components of $\boldsymbol{\beta}$

○ model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon} \sim \mathbf{N}(\mathbf{0}, \sigma^2 \mathbf{I})$

○ let  $\mathbf{a} \in \mathbb{R}^p$  such that  $a_i = 1$  and  $a_j = 0$  for  $j \neq i$

$\triangleright (1 - \alpha) \times 100\%$  confidence interval for  $\beta_i$ :

$$\left( \hat{\beta}_i - t_{1-\alpha/2}(n-p) \sqrt{\hat{\sigma}^2 (\mathbf{X}^\top \mathbf{X})_{i,i}^{-1}}, \right. \\ \left. \hat{\beta}_i + t_{1-\alpha/2}(n-p) \sqrt{\hat{\sigma}^2 (\mathbf{X}^\top \mathbf{X})_{i,i}^{-1}} \right)$$

○ let  $\mathbf{a} \in \mathbb{R}^p$  such that  $a_1 = 1$ ,  $a_i = 1$  and  $a_j = 0$  for  $j \neq i$

$\triangleright (1 - \alpha) \times 100\%$  confidence interval for  $\beta_1 + \beta_i$ :

$$\left( \hat{\beta}_1 + \hat{\beta}_i - t_{1-\alpha/2}(n-p) \sqrt{\hat{\sigma}^2 \left( (\mathbf{X}^\top \mathbf{X})_{1,1}^{-1} + 2(\mathbf{X}^\top \mathbf{X})_{1,i}^{-1} + (\mathbf{X}^\top \mathbf{X})_{i,i}^{-1} \right)}, \right. \\ \left. \hat{\beta}_1 + \hat{\beta}_i + t_{1-\alpha/2}(n-p) \sqrt{\hat{\sigma}^2 \left( (\mathbf{X}^\top \mathbf{X})_{1,1}^{-1} + 2(\mathbf{X}^\top \mathbf{X})_{1,i}^{-1} + (\mathbf{X}^\top \mathbf{X})_{i,i}^{-1} \right)} \right)$$

$\triangleright$  and analogously for other sums of components of  $\boldsymbol{\beta}$

## 6.4 Prediction

### Prediction

#### New covariate combinations

- model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon} \sim \mathbf{N}(\mathbf{0}, \sigma^2 \mathbf{I})$
- what can we say about

$$Y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k + \varepsilon ?$$

- let  $\mathbf{x} \in \mathbb{R}^p$  such that  $\mathbf{x} = (1, x_1, \dots, x_k)^\top$
- $Y = \mathbf{x}^\top \boldsymbol{\beta} + \varepsilon$  and  $\mathbf{E}Y = \mathbf{x}^\top \boldsymbol{\beta}$
- we may estimate  $\mathbf{E}Y$  by  $\widehat{\mathbf{E}Y} = \mathbf{x}^\top \widehat{\boldsymbol{\beta}}$
- $(1 - \alpha) \times 100\%$  confidence interval for  $\mathbf{E}Y$ :

$$\left( \mathbf{x}^\top \widehat{\boldsymbol{\beta}} - t_{1-\alpha/2}(n-p) \sqrt{\widehat{\sigma}^2 \mathbf{x}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}}, \right. \\ \left. \mathbf{x}^\top \widehat{\boldsymbol{\beta}} + t_{1-\alpha/2}(n-p) \sqrt{\widehat{\sigma}^2 \mathbf{x}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}} \right)$$

### Prediction in normal linear model

- model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon} \sim \mathbf{N}(\mathbf{0}, \sigma^2 \mathbf{I})$
- how can we estimate

$$Y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k + \varepsilon = \mathbf{x}^\top \boldsymbol{\beta} + \varepsilon ?$$

i.e. how do we predict new  $Y$  for new  $\mathbf{x}$ ?

- prediction  $\hat{Y} = \mathbf{x}^\top \widehat{\boldsymbol{\beta}}$
- $(1 - \alpha) \times 100\%$  confidence interval for  $Y$   
(prediction interval):

$$\left( \mathbf{x}^\top \widehat{\boldsymbol{\beta}} - t_{1-\alpha/2}(n-p) \sqrt{\widehat{\sigma}^2 (1 + \mathbf{x}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x})}, \right. \\ \left. \mathbf{x}^\top \widehat{\boldsymbol{\beta}} + t_{1-\alpha/2}(n-p) \sqrt{\widehat{\sigma}^2 (1 + \mathbf{x}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x})} \right)$$

## 6.5 Confidence bands

### Confidence bands

#### Confidence regions in normal linear model

◦ model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$

◦ let  $\mathbf{A} \in \mathbb{R}^{m \times p}$ ,  $\text{rank}(\mathbf{A}) = m$

$$\triangleright \frac{1}{m \widehat{\sigma}^2} (\mathbf{A} \widehat{\boldsymbol{\beta}} - \mathbf{A} \boldsymbol{\beta})^\top (\mathbf{A} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{A}^\top)^{-1} (\mathbf{A} \widehat{\boldsymbol{\beta}} - \mathbf{A} \boldsymbol{\beta}) \sim F(m, n - p)$$

◦  $(1 - \alpha) \times 100\%$  confidence bands for  $\mathbf{A} \boldsymbol{\beta}$ :

$$\left\{ \mathbf{A} \boldsymbol{\beta}; \frac{1}{m \widehat{\sigma}^2} (\mathbf{A} \widehat{\boldsymbol{\beta}} - \mathbf{A} \boldsymbol{\beta})^\top (\mathbf{A} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{A}^\top)^{-1} (\mathbf{A} \widehat{\boldsymbol{\beta}} - \mathbf{A} \boldsymbol{\beta}) \leq F_{1-\alpha}(m, n - p) \right\}$$

#### Confidence bands in normal linear model

◦ model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$

**Lemma 1.** Let  $\mathbf{B} \in \mathbb{R}^{m \times m}$ ,  $\mathbf{B} \succ 0$ . Then for every  $\mathbf{x} \in \mathbb{R}^m$

$$\mathbf{x}^\top \mathbf{B} \mathbf{x} \leq 1 \Leftrightarrow (\mathbf{b}^\top \mathbf{x})^2 \leq \mathbf{b}^\top \mathbf{B}^{-1} \mathbf{b} \quad \forall \mathbf{b} \in \mathbb{R}^m.$$

◦ a proof can be found in *Jiří Anděl: Základy matematické statistiky (2005). Matfyzpress*; see also multiple comparisons and Scheffé's theorem next semester

◦ for  $\mathbf{A} \in \mathbb{R}^{m \times p}$ ,  $\text{rank}(\mathbf{A}) = m$ :

$$\begin{aligned} 1 - \alpha &= \\ &= P \left( \frac{1}{m \widehat{\sigma}^2} (\mathbf{A} \widehat{\boldsymbol{\beta}} - \mathbf{A} \boldsymbol{\beta})^\top (\mathbf{A} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{A}^\top)^{-1} (\mathbf{A} \widehat{\boldsymbol{\beta}} - \mathbf{A} \boldsymbol{\beta}) \leq F_{1-\alpha}(m, n - p) \right) \\ &= P \left( (\mathbf{b}^\top (\mathbf{A} \widehat{\boldsymbol{\beta}} - \mathbf{A} \boldsymbol{\beta}))^2 \leq m F_{1-\alpha}(m, n - p) \widehat{\sigma}^2 \mathbf{b}^\top (\mathbf{A} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{A}^\top) \mathbf{b}; \forall \mathbf{b} \in \mathbb{R}^m \right) \end{aligned}$$

## 6.6 Testing hypotheses

### 6.6.1 Simple hypothesis

**Testing**  $H_0 : \beta_i = 0$

◦ model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$

◦ for  $\mathbf{a} \in \mathbb{R}^p$

$$\frac{\mathbf{a}^\top \hat{\boldsymbol{\beta}} - \mathbf{a}^\top \boldsymbol{\beta}}{\sqrt{\hat{\sigma}^2 \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a}}} \sim t(n-p)$$

◦ let  $\mathbf{a} \in \mathbb{R}^p$  such that  $a_i = 1$  and  $a_j = 0$  for  $j \neq i$

◦ testing

▷  $H_0 : \beta_i = 0$  vs.

▷  $H_1 : \beta_i \neq 0$

◦ test statistic  $T_i = \frac{\hat{\beta}_i}{\sqrt{\hat{\sigma}^2 (\mathbf{X}^\top \mathbf{X})_{i,i}^{-1}}} \sim t(n-p)$

◦ reject  $H_0$  in favour of  $H_1$  if  $|t_i| > t_{1-\alpha/2}(n-p)$

◦ analogously for linear combinations of elements of  $\boldsymbol{\beta}$

◦ analogously for testing  $H_0 : \beta_i = \beta_{0,i}$

## 6.6.2 Composite hypothesis

**Testing  $H_0 : \boldsymbol{\beta}_{i:p} = \mathbf{0}$**

◦ model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$

◦ for  $\mathbf{A} \in \mathbb{R}^{m \times p}$ ,  $\text{rank}(\mathbf{A}) = m$

$$\frac{1}{m \hat{\sigma}^2} (\mathbf{A} \hat{\boldsymbol{\beta}} - \mathbf{A} \boldsymbol{\beta})^\top (\mathbf{A} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{A}^\top)^{-1} (\mathbf{A} \hat{\boldsymbol{\beta}} - \mathbf{A} \boldsymbol{\beta}) \sim F(m, n-p)$$

◦ testing

▷  $H_0 : \boldsymbol{\beta}_{i:p} = \mathbf{0}$  vs.

▷  $H_1 : \boldsymbol{\beta}_{i:p} \neq \mathbf{0}$

◦ test statistic

$$F_{i:p} = \frac{1}{(p-i+1) \hat{\sigma}^2} \hat{\boldsymbol{\beta}}_{i:p}^\top (\mathbf{X}^\top \mathbf{X})_{i:p,i:p}^{-1} \hat{\boldsymbol{\beta}}_{i:p} \sim F(p-i+1, n-p)$$

- reject  $H_0$  in favour of  $H_1$  if  $f_{i:p} > F_{1-\alpha}(p-i+1, n-p)$

### Testing “the model”

- model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$
- for  $\mathbf{A} \in \mathbb{R}^{m \times p}$ ,  $\text{rank}(\mathbf{A}) = m$

$$\frac{1}{m \hat{\sigma}^2} (\mathbf{A}\hat{\boldsymbol{\beta}} - \mathbf{A}\boldsymbol{\beta})^\top (\mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{A}^\top)^{-1} (\mathbf{A}\hat{\boldsymbol{\beta}} - \mathbf{A}\boldsymbol{\beta}) \sim F(m, n-p)$$

- testing

$$\triangleright H_0 : \boldsymbol{\beta}_{2:p} = \mathbf{0} \text{ vs.}$$

$$\triangleright H_1 : \boldsymbol{\beta}_{2:p} \neq \mathbf{0}$$

- test statistic

$$F = \frac{1}{k \hat{\sigma}^2} \hat{\boldsymbol{\beta}}_{2:p}^\top (\mathbf{X}^\top \mathbf{X})_{2:p, 2:p}^{-1} \hat{\boldsymbol{\beta}}_{2:p} \sim F(k, n-p)$$

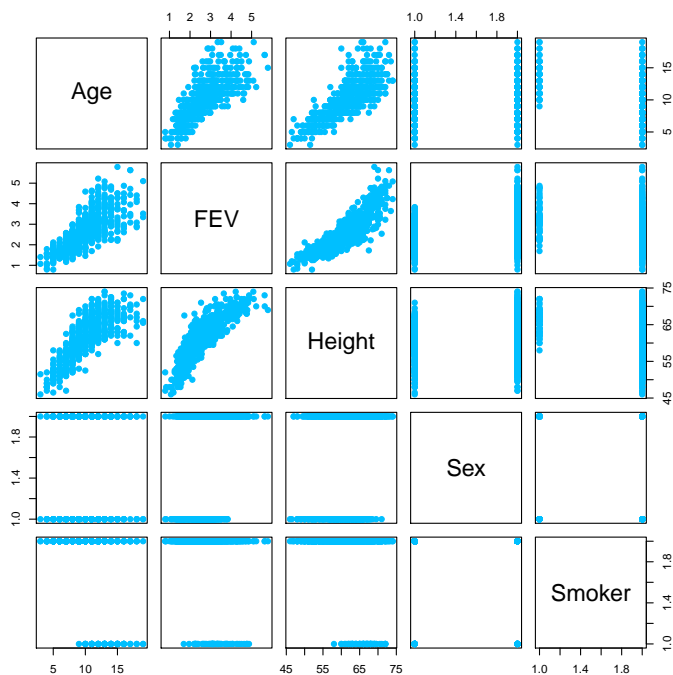
- reject  $H_0$  in favour of  $H_1$  if  $f > F_{1-\alpha}(k, n-p)$

## 6.7 Interpretation

### Interpretation of results for normal linear model

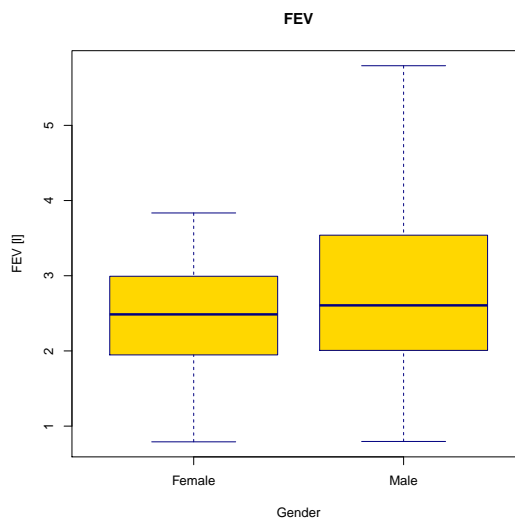
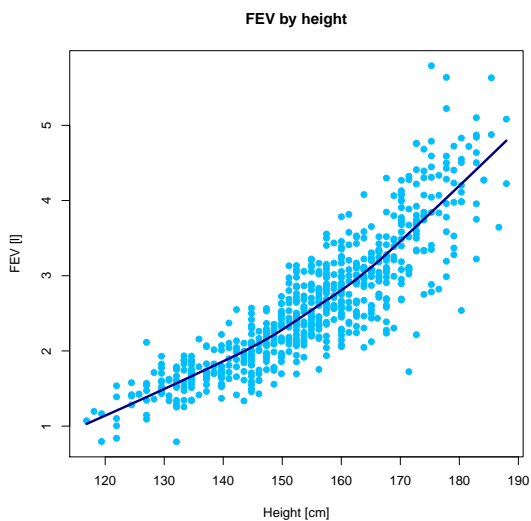
#### A model for fev data

- model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$
- data: fev from <http://www.statsci.org/data/general/fev.html>



**A model for fev data ctd.**

- model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$
- model FEV by Height and Sex



**Fitted model for fev data**

- model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$
- model FEV by Height and Sex

```
> model.simple <- lm(FEV~Height + Sex, data=fev)
> summary(model.simple)

Call:
lm(formula = FEV ~ Height + Sex, data = fev)

Residuals:
    Min       1Q   Median       3Q      Max
-1.6763 -0.2505  0.0001  0.2347  2.0722

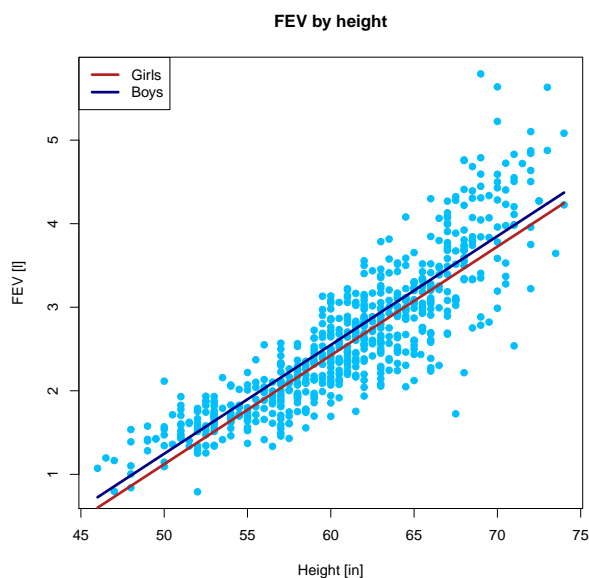
Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) -5.390263   0.180082  -29.932 < 2e-16 ***
Height       0.130231   0.002964   43.933 < 2e-16 ***
SexMale      0.125123   0.033801    3.702 0.000232 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.4265 on 651 degrees of freedom
Multiple R-squared:  0.7587, Adjusted R-squared:  0.758
F-statistic: 1024 on 2 and 651 DF, p-value: < 2.2e-16
```

### Fitted model for fev data ctd.

- model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$
- model FEV by Height and Sex

```
> coefficients(model.simple)
(Intercept)      Height      SexMale
-5.3902632    0.1302305    0.1251234
```





# Chapter 7

## Model selection

### 7.1 The problem

#### 7.1.1 Normal linear model

##### Normal linear model

- $Y_i = \beta_0 + \beta_1 x_{i,1} + \dots + \beta_k x_{i,k} + \varepsilon_i, i \in \{1, \dots, n\}$ 
  - ▷  $Y_i$ : outcome, response, output, dependent variable
    - \* random variable, we observe a realization  $y_i$
    - \* (odezva, závisle proměnná, regresand)
  - ▷  $x_{i,1}, \dots, x_{i,k}$ : covariates, predictors, explanatory variables, input, independent variables
    - \* given, known
    - \* (nezávisle proměnné, regresory)
  - ▷  $\beta_0, \dots, \beta_k$ : coefficients
    - \* unknown
    - \* (regresní koeficienty)
  - ▷  $\varepsilon_i$ : random error
    - \* random variable, unobserved
- $\varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2), i \in \{1, \dots, n\}$ 
  - ▷  $E \varepsilon_i = 0$ : no systematic errors
  - ▷  $\text{Var } \varepsilon_i = \sigma^2$ : same precision

##### Example: bloodpress data

- o from [sites.stat.psu.edu/~lsimon/stat501wc/sp05/data/](http://sites.stat.psu.edu/~lsimon/stat501wc/sp05/data/)
- o association between the mean arterial blood pressure[mmHg] and age[years], weight[kg], body surface area[m<sup>2</sup>], duration of hypertension[years], basal pulse[beats/min], stress

o data:

	BP	Age	Weight	BSA	DoH	Pulse	Stress
	105	47	85.4	1.75	5.1	63	33
	115	49	94.2	2.10	3.8	70	14
	...	...	...	...	...	...	...
	110	48	90.5	1.88	9.0	71	99
	122	56	95.7	2.09	7.0	75	99

- o model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$

$$\begin{pmatrix} 105 \\ 115 \\ \dots \\ 110 \\ 122 \end{pmatrix} = \begin{pmatrix} 1 & 47 & 85.4 & 1.75 & 5.1 & 63 & 33 \\ 1 & 49 & 94.2 & 2.10 & 3.8 & 70 & 14 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & 48 & 90.5 & 1.88 & 9.0 & 71 & 99 \\ 1 & 56 & 95.7 & 2.09 & 7.0 & 75 & 99 \end{pmatrix} \times \begin{pmatrix} \beta_0 \\ \dots \\ \beta_6 \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \dots \\ \varepsilon_{19} \\ \varepsilon_{20} \end{pmatrix}$$

<https://ww2.amstat.org/publications/jse/v13n2/datasets.kahn.html>

### Example: fev data

- o from: <http://www.statsci.org/data/general/fev.html>
- o question: association between the FEV[l] and Smoking,

corrected for Age[years], Height[cm] and Gender

o data:

	FEV	Age	Height	Gender	Smoking
	1.708	9	144.8	Female	Non
	1.724	8	171.5	Female	Non
	1.720	7	138.4	Female	Non
	1.558	9	134.6	Male	Non
	...	...	...	...	...
	3.727	15	172.7	Male	Current
	2.853	18	152.4	Female	Non
	2.795	16	160.0	Female	Current
	3.211	15	168.9	Female	Non

- o model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$

$$\begin{pmatrix} 1.708 \\ 1.724 \\ 1.720 \\ 1.558 \\ \dots \\ 3.727 \\ 2.853 \\ 2.795 \\ 3.211 \end{pmatrix} = \begin{pmatrix} 1 & 9 & 144.8 & 0 & 0 \\ 1 & 8 & 171.5 & 0 & 0 \\ 1 & 7 & 138.4 & 0 & 0 \\ 1 & 9 & 134.6 & 1 & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 1 & 15 & 172.7 & 1 & 1 \\ 1 & 18 & 152.4 & 0 & 0 \\ 1 & 16 & 160.0 & 0 & 1 \\ 1 & 15 & 168.9 & 0 & 0 \end{pmatrix} \times \begin{pmatrix} \beta_0 \\ \dots \\ \beta_5 \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \varepsilon_3 \\ \varepsilon_4 \\ \dots \\ \varepsilon_{651} \\ \varepsilon_{652} \\ \varepsilon_{653} \\ \varepsilon_{654} \end{pmatrix}$$

## 7.1.2 Task for this chapter

### Model building/selection

- o model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$

- ▷ outcome  $\mathbf{Y}$ 
  - \* random vector, we observe a realization  $\mathbf{y}$
- ▷ predictors  $\mathbf{x}_1, \dots, \mathbf{x}_k$ 
  - \* vector of given (known) constants
- ▷ coefficients  $\boldsymbol{\beta}$ 
  - \* vector of unknown constants
- ▷ error  $\boldsymbol{\varepsilon}$ 
  - \* unknown random vector, we do not observe its realization
- ▷ assumptions:  $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ 
  - \*  $\mathbf{E} \mathbf{Y} = \mathbf{X} \boldsymbol{\beta}$ : the expected value of  $\mathbf{Y}$  is a linear function of  $\boldsymbol{\beta}$
  - \*  $\mathbf{E} \boldsymbol{\varepsilon} = \mathbf{0}$ : no systematic errors
  - \*  $\text{Var} \boldsymbol{\varepsilon} = \sigma^2 \mathbf{I}$ : independence and same precision
- **task**: given the observed data  $\mathbf{y}$  and values of potential covariates, construct  $\mathbf{X}$
- Note:  $\mathbf{X}$  should ideally be known a priori based on background knowledge and various optimality considerations but ...

## 7.2 Why consider various models?

### 7.2.1 Should we leave out covariates that appear unnecessary?

#### Testing hypotheses about null coefficients

- model:  $\mathbf{Y} = \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$
- testing
  - ▷  $H_0 : \beta_i = 0$  vs.
  - ▷  $H_1 : \beta_i \neq 0$
- test statistic  $T_i = \frac{\hat{\beta}_i}{\sqrt{\hat{\sigma}^2 (\mathbf{X}^\top \mathbf{X})_{i,i}^{-1}}} \sim t(n-p)$
- reject  $H_0$  in favour of  $H_1$  if  $|t_i| > t_{1-\alpha/2}(n-p)$
- testing
  - ▷  $H_0 : \boldsymbol{\beta}_{i:p} = \mathbf{0}$  vs.

$$\triangleright H_1 : \boldsymbol{\beta}_{i:p} \neq \mathbf{0}$$

- test statistic

$$F_{i:p} = \frac{1}{(p-i+1)\widehat{\sigma}^2} \widehat{\boldsymbol{\beta}}_{i:p}^\top (\mathbf{X}^\top \mathbf{X})_{i:p,i:p}^{-1} \widehat{\boldsymbol{\beta}}_{i:p} \sim F(p-i+1, n-p)$$

- reject  $H_0$  in favour of  $H_1$  if  $f_{i:p} > F_{1-\alpha}(p-i+1, n-p)$

### What if we do not reject?

- model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$

- testing

$$\triangleright H_0 : \beta_i = 0 \text{ vs. } H_1 : \beta_i \neq 0$$

- if  $|t_i| < t_{1-\alpha/2}(n-p)$

$\triangleright$  at  $\alpha\%$  level, we do not reject that  $\beta_i = 0$  in favour of  $\beta_i \neq 0$

- testing

$$\triangleright H_0 : \boldsymbol{\beta}_{i:p} = \mathbf{0} \text{ vs. } H_1 : \boldsymbol{\beta}_{i:p} \neq \mathbf{0}$$

- if  $f_{i:p} < F_{1-\alpha}(p-i+1, n-p)$

$\triangleright$  at  $\alpha\%$  level, we do not reject  $\boldsymbol{\beta}_{i:p} = \mathbf{0}$  in favour of  $\boldsymbol{\beta}_{i:p} \neq \mathbf{0}$

- if we do not reject that some components of  $\boldsymbol{\beta}$  are 0, should we change the model?

$\triangleright$  original model

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$$

$\triangleright$  new model

$$\mathbf{Y} = \mathbf{X}_{,1:(i-1)}\boldsymbol{\beta}_{1:(i-1)} + \boldsymbol{\varepsilon}, \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$$

### Example: bloodpress data

- original model

$$Y_i = \beta_0 + \beta_1 \times \text{Age}_i + \beta_2 \times \text{Weight}_i + \beta_3 \times \text{BSA}_i + \\ + \beta_4 \times \text{Dur}_i + \beta_5 \times \text{Pulse}_i + \beta_6 \times \text{Stress}_i + \varepsilon_i, \quad 1 \leq i \leq 20$$

```

Coefficients:
      Estimate Std. Error t value Pr(>|t|)
(Intercept) -12.870476   2.556650  -5.034 0.000229 ***
Age           0.703259   0.049606  14.177 2.76e-09 ***
Weight       0.969920   0.063108  15.369 1.02e-09 ***
BSA          3.776491   1.580151   2.390 0.032694 *
Dur          0.068383   0.048441   1.412 0.181534
Pulse       -0.084485   0.051609  -1.637 0.125594
Stress       0.005572   0.003412   1.633 0.126491

Residual standard error: 0.4072 on 13 degrees of freedom

> coef.table <- summary(model.full)$coefficients
> V <- vcov(model.full)
> A <- diag(rep(1, 7))[5:7, ]
> F.stat <- t(A*%coef.table[, 1])%*%solve(A*%V*%t(A))%*(A*%coef.table[, 1])/3
> 1-pf(F.stat, df1=3, df2=13)
      [,1]
[1,] 0.1950807

```

- should we rather use the new model?

$$Y_i = \beta_0 + \beta_1 \times \text{Age}_i + \beta_2 \times \text{Weight}_i + \beta_3 \times \text{BSA}_i + \varepsilon_i, \quad 1 \leq i \leq 20$$

### Example: bloodpress data

- original model

$$\triangleright Y_i = \beta_0 + \beta_1 \times \text{Age}_i + \beta_2 \times \text{Weight}_i + \beta_3 \times \text{BSA}_i + \beta_4 \times \text{Dur}_i + \beta_5 \times \text{Pulse}_i + \beta_6 \times \text{Stress}_i + \varepsilon_i, \quad 1 \leq i \leq 20$$

$$\triangleright \mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

$$\begin{pmatrix} 105 \\ 115 \\ \dots \\ 110 \\ 122 \end{pmatrix} = \begin{pmatrix} 1 & 47 & 85.4 & 1.75 & 5.1 & 63 & 33 \\ 1 & 49 & 94.2 & 2.10 & 3.8 & 70 & 14 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & 48 & 90.5 & 1.88 & 9.0 & 71 & 99 \\ 1 & 56 & 95.7 & 2.09 & 7.0 & 75 & 99 \end{pmatrix} \times \begin{pmatrix} \beta_0 \\ \dots \\ \beta_6 \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \dots \\ \varepsilon_{19} \\ \varepsilon_{20} \end{pmatrix}$$

- new model

$$\triangleright Y_i = \beta_0 + \beta_1 \times \text{Age}_i + \beta_2 \times \text{Weight}_i + \beta_3 \times \text{BSA}_i + \varepsilon_i, \quad 1 \leq i \leq 20$$

$$\triangleright \mathbf{Y} = \tilde{\mathbf{X}}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

$$\begin{pmatrix} 105 \\ 115 \\ \dots \\ 110 \\ 122 \end{pmatrix} = \begin{pmatrix} 1 & 47 & 85.4 & 1.75 \\ 1 & 49 & 94.2 & 2.10 \\ \dots & \dots & \dots & \dots \\ 1 & 48 & 90.5 & 1.88 \\ 1 & 56 & 95.7 & 2.09 \end{pmatrix} \times \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \dots \\ \varepsilon_{19} \\ \varepsilon_{20} \end{pmatrix}$$

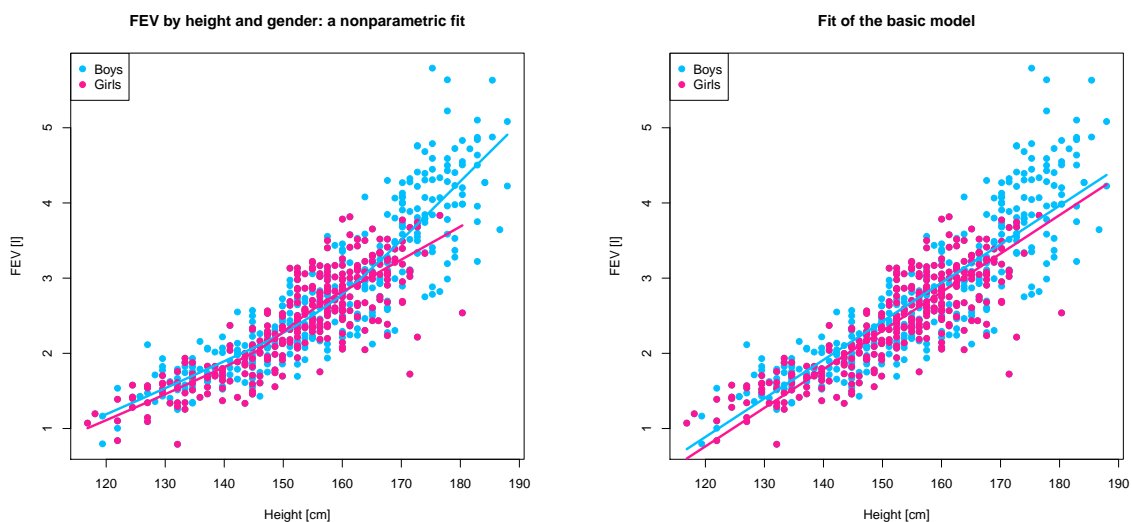
## 7.2.2 What is the right form of the dependence on covariates?

### Specifying the form of dependence in the fev data

- basic model for the dependence of FEV on Height and Sex:

$$Y_i = \beta_0 + \beta_1 \times \text{Height}_i + \beta_2 \times \mathbb{I}\{\text{the } i^{\text{th}} \text{ person is male}\}, \quad 1 \leq i \leq 654$$

- o does the basic model fit the data well enough?



### Example: fev data

- o original model

$$\triangleright Y_i = \beta_0 + \beta_1 \times \text{Height}_i + \beta_2 \times \mathbb{I}\{\text{the } i^{\text{th}} \text{ child is male}\} + \varepsilon_i, 1 \leq i \leq 654$$

$$\begin{pmatrix} 1.708 \\ 1.724 \\ 1.720 \\ 1.558 \\ \dots \\ 3.727 \\ 2.853 \\ 2.795 \\ 3.211 \end{pmatrix} = \begin{pmatrix} 1 & 144.8 & 0 \\ 1 & 171.5 & 0 \\ 1 & 138.4 & 0 \\ 1 & 134.6 & 1 \\ \dots & \dots & \dots \\ 1 & 172.7 & 1 \\ 1 & 152.4 & 0 \\ 1 & 160.0 & 0 \\ 1 & 168.9 & 0 \end{pmatrix} \times \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \varepsilon_3 \\ \varepsilon_4 \\ \dots \\ \varepsilon_{651} \\ \varepsilon_{652} \\ \varepsilon_{653} \\ \varepsilon_{654} \end{pmatrix}$$

- o new model

$$\begin{aligned} \triangleright Y_i = & \beta_0 + \beta_1 \times \text{Height}_i + \beta_2 \times \text{Height}_i^2 + \\ & + \beta_3 \times \mathbb{I}\{\text{the } i^{\text{th}} \text{ child is male}\} + \beta_4 \times \text{Height}_i \mathbb{I}\{\text{the } i^{\text{th}} \text{ child is male}\} + \\ & + \beta_5 \times \text{Height}_i^2 \mathbb{I}\{\text{the } i^{\text{th}} \text{ child is male}\} + \varepsilon_i, 1 \leq i \leq 654 \end{aligned}$$

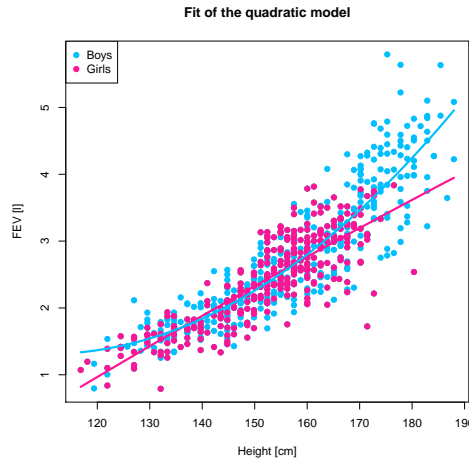
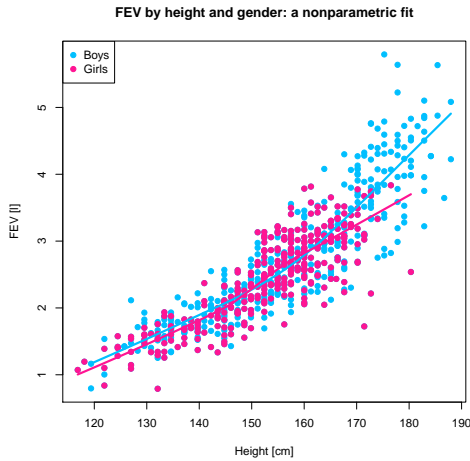
$$\begin{pmatrix} 1.708 \\ 1.724 \\ 1.720 \\ 1.558 \\ \dots \\ 3.727 \\ 2.853 \\ 2.795 \\ 3.211 \end{pmatrix} = \begin{pmatrix} 1 & 144.8 & 20961.3 & 0 & 0 & 0 \\ 1 & 171.5 & 29395.1 & 0 & 0 & 0 \\ 1 & 138.4 & 19162.9 & 0 & 0 & 0 \\ 1 & 134.6 & 18122.5 & 1 & 134.6 & 18122.5 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & 172.7 & 29832.2 & 1 & 172.7 & 29832.2 \\ 1 & 152.4 & 23225.8 & 0 & 0 & 0 \\ 1 & 160.0 & 25606.4 & 0 & 0 & 0 \\ 1 & 168.9 & 28530.6 & 0 & 0 & 0 \end{pmatrix} \times \begin{pmatrix} \beta_0 \\ \dots \\ \beta_5 \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \varepsilon_3 \\ \varepsilon_4 \\ \dots \\ \varepsilon_{651} \\ \varepsilon_{652} \\ \varepsilon_{653} \\ \varepsilon_{654} \end{pmatrix}$$

### Example: fev data

- o fit of the new model

```

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  -5.194e+00  2.740e+00  -1.895  0.0585 .
Height        5.611e-02  3.692e-02   1.520  0.1291
I(Height^2)   -3.977e-05  1.238e-04  -0.321  0.7482
SexMale       1.392e+01  3.423e+00   4.067  5.34e-05 ***
Height:SexMale -1.903e-01  4.545e-02  -4.188  3.20e-05 ***
I(Height^2):SexMale 6.471e-04  1.501e-04   4.310  1.89e-05 ***
    
```



### 7.3 Nested models

#### Submodel

##### Nested models

- Bigger model:  $\mathbf{Y} = \mathbf{X}_b \boldsymbol{\beta}_b + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon} \sim (\mathbf{0}, \sigma^2 \mathbf{I})$ ,

$$\mathbf{X}_b = (\mathbf{1} \mid \mathbf{x}_1 \mid \mathbf{x}_2 \mid \dots \mid \mathbf{x}_{k-1} \mid \mathbf{x}_k)$$

- ▷  $\hat{\boldsymbol{\beta}}_b = (\mathbf{X}_b^\top \mathbf{X}_b)^{-1} \mathbf{X}_b^\top \mathbf{Y}$
- ▷  $\hat{\mathbf{Y}}_b = \mathbf{X}_b \hat{\boldsymbol{\beta}}_b = \mathbf{H}_b \mathbf{Y}$
- ▷  $\mathbf{e}_b = \mathbf{Y} - \hat{\mathbf{Y}}_b = (\mathbf{I} - \mathbf{H}_b) \mathbf{Y}$
- ▷  $\hat{\sigma}_b^2 = \frac{1}{n-p} \|\mathbf{e}_b\|^2$

- Smaller model:  $\mathbf{Y} = \mathbf{X}_s \boldsymbol{\beta}_s + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon} \sim (\mathbf{0}, \sigma^2 \mathbf{I})$ ,

$$\mathbf{X}_s = (\mathbf{1} \mid \mathbf{x}_1 \mid \mathbf{x}_2 \mid \dots \mid \mathbf{x}_{k-r})$$

- ▷  $\hat{\boldsymbol{\beta}}_s = (\mathbf{X}_s^\top \mathbf{X}_s)^{-1} \mathbf{X}_s^\top \mathbf{Y}$
- ▷  $\hat{\mathbf{Y}}_s = \mathbf{X}_s \hat{\boldsymbol{\beta}}_s = \mathbf{H}_s \mathbf{Y}$
- ▷  $\mathbf{e}_s = \mathbf{Y} - \hat{\mathbf{Y}}_s = (\mathbf{I} - \mathbf{H}_s) \mathbf{Y}$

$$\triangleright \hat{\sigma}_s^2 = \frac{1}{n-p+r} \|\mathbf{e}_s\|^2$$

- more generally, any  $\mathbf{X}_s$  such that  $\text{im}(\mathbf{X}_s) \leq \text{im}(\mathbf{X}_b)$ 
  - ▷  $\exists \mathbf{A} \in \mathbb{R}^{p \times (p-r)}$  such that  $\mathbf{X}_s = \mathbf{X}_b \mathbf{A}$
  - ▷  $\mathbf{X}_s = \left( \sum_{i=1}^p a_{i,1} \mathbf{x}_{i,1} \mid \dots \mid \sum_{i=1}^p a_{i,p-r} \mathbf{x}_{i,p-r} \right)$

### Relationship between the two models

- if the smaller model holds, so does the bigger one
- $\exists \mathbf{A} \in \mathbb{R}^{p \times (p-r)}$  such that  $\mathbf{X}_s = \mathbf{X}_b \mathbf{A}$ 
  - ▷  $\mathbf{X}_s = \left( \sum_{i=1}^p a_{i,1} \mathbf{x}_{i,1} \mid \dots \mid \sum_{i=1}^p a_{i,p-r} \mathbf{x}_{i,p-r} \right)$
- bigger model:  $\mathbf{Y} = \mathbf{X}_b \boldsymbol{\beta}_b + \boldsymbol{\varepsilon}$
- smaller model:  $\mathbf{Y} = \mathbf{X}_s \boldsymbol{\beta}_s + \boldsymbol{\varepsilon} = \mathbf{X}_b \underbrace{\mathbf{A} \boldsymbol{\beta}_s}_{\boldsymbol{\beta}_b} + \boldsymbol{\varepsilon}$
- smaller model is the bigger model with a condition on  $\boldsymbol{\beta}_b$

$$\triangleright \underbrace{\boldsymbol{\beta}_b}_{p \times 1} = \underbrace{\mathbf{A}}_{p \times (p-r)} \underbrace{\boldsymbol{\beta}_s}_{(p-r) \times 1} = \begin{pmatrix} \sum_{j=1}^{n-p} a_{1,j} \beta_{s,j} \\ \dots \\ \sum_{j=1}^{n-p} a_{p,j} \beta_{s,j} \end{pmatrix}$$

$$\triangleright \exists \mathbf{B} \in \mathbb{R}^{r \times p} \text{ such that } \mathbf{B} \boldsymbol{\beta}_b = \mathbf{0}$$

- in the bigger normal linear model we may test for the validity of the smaller model by testing whether  $\mathbf{B} \boldsymbol{\beta}_b = \mathbf{0}$  (see Week 7)

### Relationship between the fits of the two models

- difference between the fits
  - ▷  $\hat{\mathbf{Y}}_b - \hat{\mathbf{Y}}_s = (\mathbf{H}_b - \mathbf{H}_s) \mathbf{Y}$
- difference between the residuals
  - ▷  $\mathbf{e}_s - \mathbf{e}_b = (\mathbf{I} - \mathbf{H}_s) \mathbf{Y} - (\mathbf{I} - \mathbf{H}_b) \mathbf{Y} = (\mathbf{H}_b - \mathbf{H}_s) \mathbf{Y}$
- comparison of the nested models' fits
  - ▷  $\|\mathbf{e}_s\|^2 = \|\mathbf{e}_b\|^2 + \|\mathbf{e}_s - \mathbf{e}_b\|^2$ 
    - \* proof: realize that  $\langle \mathbf{e}_b, \mathbf{e}_s - \mathbf{e}_b \rangle = 0$
    - \* note:  $\|\mathbf{e}_s\|^2 \geq \|\mathbf{e}_b\|^2 \Rightarrow$  the fit of the bigger model is closer to the observed data



\* note:  $\|\mathbf{e}_s - \mathbf{e}_b\|^2 = \|\mathbf{e}_s\|^2 - \|\mathbf{e}_b\|^2$

- in the normal linear model ( $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$ )

▷  $\mathbf{e}_b \perp\!\!\!\perp (\mathbf{e}_s - \mathbf{e}_b)$

\* proof: Corollary of MVN 7:

Let  $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ . Then  $\mathbf{A}\mathbf{X} \perp\!\!\!\perp \mathbf{B}\mathbf{X}$  iff  $\mathbf{A}\boldsymbol{\Sigma}\mathbf{B}^\top = \mathbf{0}$ .

### Does the bigger model fit significantly better?

- assume that both models hold (i.e. the smaller model holds) and that  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$  (normal linear model)

- $\frac{1}{\sigma^2} \|\mathbf{e}_b\|^2 \sim \chi_{n-p}^2$

▷ proof: see Week 6

- $\frac{1}{\sigma^2} \|\mathbf{e}_s - \mathbf{e}_b\|^2 \sim \chi_r^2$

▷ proof: MVN 3:

Let  $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  and let  $\mathbf{A}$  be an  $m \times n$  real matrix and  $\mathbf{b} \in \mathbb{R}^m$ . Then  $\mathbf{A}\mathbf{X} + \mathbf{b} \sim N(\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top)$ .

▷ and QF 4:

Let  $\mathbf{Z} \sim N(\mathbf{0}, \mathbf{I})$  and let  $\mathbf{P}$  be an  $n \times n$  projection matrix of rank  $r$ . Then  $\mathbf{Z}^\top \mathbf{P} \mathbf{Z} \sim \chi^2(r)$ .

- $\|\mathbf{e}_b\|^2 \perp\!\!\!\perp \|\mathbf{e}_s - \mathbf{e}_b\|^2$

▷ proof: see the previous slide

- $\frac{\|\mathbf{e}_s - \mathbf{e}_b\|^2/r}{\|\mathbf{e}_b\|^2/(n-p)} = \frac{(\|\mathbf{e}_s\|^2 - \|\mathbf{e}_b\|^2)/r}{\|\mathbf{e}_b\|^2/(n-p)} \sim F_{r, n-p}$

▷ proof: verify that the definition of  $F_{r, n-p}$  is satisfied

### More submodels

#### Several models nested within one another

- Big model:  $\mathbf{Y} = \mathbf{X}_b \boldsymbol{\beta}_b + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon} \sim (\mathbf{0}, \sigma^2 \mathbf{I})$ ,

▷  $\hat{\mathbf{Y}}_b = \mathbf{X}_b \hat{\boldsymbol{\beta}}_b = \mathbf{H}_b \mathbf{Y}$

▷  $\mathbf{e}_b = \mathbf{Y} - \hat{\mathbf{Y}}_b = (\mathbf{I} - \mathbf{H}_b) \mathbf{Y}$

- Small model:  $\mathbf{Y} = \mathbf{X}_s \boldsymbol{\beta}_s + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon} \sim (\mathbf{0}, \sigma^2 \mathbf{I})$ ,

- ▷  $\hat{\mathbf{Y}}_s = \mathbf{X}_s \hat{\boldsymbol{\beta}}_s = \mathbf{H}_s \mathbf{Y}$
- ▷  $\mathbf{e}_s = \mathbf{Y} - \hat{\mathbf{Y}}_s = (\mathbf{I} - \mathbf{H}_s) \mathbf{Y}$
- ▷  $\hat{\sigma}_s^2 = \frac{1}{n-p+r} \|\mathbf{e}_s\|^2$
- Super-small model:  $\mathbf{Y} = \mathbf{X}_{ss} \boldsymbol{\beta}_{ss} + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon} \sim (\mathbf{0}, \sigma^2 \mathbf{I})$ ,
  - ▷  $\hat{\mathbf{Y}}_{ss} = \mathbf{X}_{ss} \hat{\boldsymbol{\beta}}_{ss} = \mathbf{H}_{ss} \mathbf{Y}$
  - ▷  $\mathbf{e}_{ss} = \mathbf{Y} - \hat{\mathbf{Y}}_{ss} = (\mathbf{I} - \mathbf{H}_{ss}) \mathbf{Y}$
  - ▷  $\hat{\sigma}_{ss}^2 = \frac{1}{n-p+s} \|\mathbf{e}_{ss}\|^2$
- $\text{im}(\mathbf{X}_{ss}) \leq \text{im}(\mathbf{X}_s) \leq \text{im}(\mathbf{X}_b)$

### Relationship between the fits of the models

- difference between the fits
  - ▷  $\hat{\mathbf{Y}}_b - \hat{\mathbf{Y}}_s = (\mathbf{H}_b - \mathbf{H}_s) \mathbf{Y}$
  - ▷  $\hat{\mathbf{Y}}_b - \hat{\mathbf{Y}}_{ss} = (\mathbf{H}_b - \mathbf{H}_{ss}) \mathbf{Y}$
  - ▷  $\hat{\mathbf{Y}}_s - \hat{\mathbf{Y}}_{ss} = (\mathbf{H}_s - \mathbf{H}_{ss}) \mathbf{Y}$
- difference between the residuals
  - ▷  $\mathbf{e}_s - \mathbf{e}_b = (\mathbf{I} - \mathbf{H}_s) \mathbf{Y} - (\mathbf{I} - \mathbf{H}_b) \mathbf{Y} = (\mathbf{H}_b - \mathbf{H}_s) \mathbf{Y}$
  - ▷  $\mathbf{e}_{ss} - \mathbf{e}_b = (\mathbf{I} - \mathbf{H}_{ss}) \mathbf{Y} - (\mathbf{I} - \mathbf{H}_b) \mathbf{Y} = (\mathbf{H}_b - \mathbf{H}_{ss}) \mathbf{Y}$
  - ▷  $\mathbf{e}_{ss} - \mathbf{e}_s = (\mathbf{I} - \mathbf{H}_{ss}) \mathbf{Y} - (\mathbf{I} - \mathbf{H}_s) \mathbf{Y} = (\mathbf{H}_s - \mathbf{H}_{ss}) \mathbf{Y}$
- in the normal linear model ( $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$ )
  - ▷  $\mathbf{e}_b \perp\!\!\!\perp (\mathbf{e}_s - \mathbf{e}_b)$
  - ▷  $\mathbf{e}_b \perp\!\!\!\perp (\mathbf{e}_{ss} - \mathbf{e}_b)$
  - ▷  $\mathbf{e}_b \perp\!\!\!\perp (\mathbf{e}_{ss} - \mathbf{e}_s)$
  - \* proof: Corollary of MVN 7:  
Let  $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ . Then  $\mathbf{A}\mathbf{X} \perp\!\!\!\perp \mathbf{B}\mathbf{X}$  iff  $\mathbf{A}\boldsymbol{\Sigma}\mathbf{B}^\top = \mathbf{0}$ .

### How about the super-small model's fit?

- assume that all three models hold (i.e. the super-small model holds) and that  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$  (normal linear model)
- $\frac{1}{\sigma^2} \|\mathbf{e}_b\|^2 \sim \chi_{n-p}^2$
- $\frac{1}{\sigma^2} \|\mathbf{e}_{ss} - \mathbf{e}_s\|^2 \sim \chi_{s-r}^2$

▷ proof: MVN 3:

Let  $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  and let  $\mathbf{A}$  be an  $m \times n$  real matrix and  $\mathbf{b} \in \mathbb{R}^m$ . Then  $\mathbf{AX} + \mathbf{b} \sim N(\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top)$ .

▷ and QF 4:

Let  $\mathbf{Z} \sim N(\mathbf{0}, \mathbf{I})$  and let  $\mathbf{P}$  be an  $n \times n$  projection matrix of rank  $r$ . Then  $\mathbf{Z}^\top \mathbf{P} \mathbf{Z} \sim \chi^2(r)$ .

○  $\|\mathbf{e}_b\|^2 \perp\!\!\!\perp \|\mathbf{e}_{ss} - \mathbf{e}_s\|^2$

○  $\frac{\|\mathbf{e}_{ss} - \mathbf{e}_s\|^2 / (s - r)}{\|\mathbf{e}_b\|^2 / (n - p)} = \frac{(\|\mathbf{e}_{ss}\|^2 - \|\mathbf{e}_s\|^2) / (s - r)}{\|\mathbf{e}_b\|^2 / (n - p)} \sim F_{s-r, n-p}$

▷ proof: verify that the definition of  $F_{s-r, n-p}$  is satisfied

## 7.4 Selecting the model

### 7.4.1 Model selection tools

#### Model selection based on sequential testing

○ statistical tests

▷  $t$  test for testing  $\beta_i = 0$  vs.  $\beta_i \neq 0$

▷  $F$  test for testing  $\mathbf{A}\boldsymbol{\beta} = \mathbf{0}$  vs.  $\mathbf{A}\boldsymbol{\beta} \neq \mathbf{0}$

▷ likelihood ratio test

- \*  $2(\max_{\boldsymbol{\theta}_b} \ell(\text{Big model}) - \max_{\boldsymbol{\theta}_s} \ell(\text{Small model})) \stackrel{as.}{\sim} \chi^2_{|\boldsymbol{\theta}_b| - |\boldsymbol{\theta}_s|}$
- \* details next semester

○ we may start with a big model and sequentially leave out terms that do not appear significant

▷ multiple testing  $\Rightarrow$  we do not keep the overall  $\alpha$

- \* often  $\alpha > 0.05$  is used at this stage (even  $\alpha \approx 0.2$ )
- \* the procedure is an ad-hoc one (rather than valid testing)
- \* “clean” ways exist (e.g. error-spending function)

▷ an approach of this kind is often applied when the interest is in  $\boldsymbol{\beta}$  and the model is there to explain the phenomenon

○ words of caution

▷  $p > 0.05$  does not guarantee the absence of the relationship

- ▷ significance of the terms in the final model may be amplified

### Model selection based on “criteria”

- model selection criterion

- ▷ a number that describes the overall fit of the model

- often applied when the interest is in prediction

- focus is on  $\|\mathbf{e}\|^2 = \|\mathbf{Y} - \hat{\mathbf{Y}}\|^2$

- already seen

- ▷ coefficient of determination

$$R^2 = 1 - \frac{\|\mathbf{e}\|^2}{\|\mathbf{Y} - \bar{Y}\mathbf{1}\|^2}$$

- \* always bigger for a bigger model

- \* bigger model is not necessarily better, so is the difference big enough to justify the use of the bigger model?

- ▷ adjusted coefficient of determination

$$R_{adj}^2 = 1 - \frac{\|\mathbf{e}\|^2/(n-p)}{\|\mathbf{Y} - \bar{Y}\mathbf{1}\|^2/(n-1)}$$

- \* penalizes for the model complexity

### Likelihood-based information criteria

- model fit versus model complexity trade-off

- Akaike information criterion

- ▷  $AIC = -2 \max_{\boldsymbol{\theta}} \ell(\text{model}) + 2 \times |\boldsymbol{\theta}|$

- ▷ motivation

- \* information theory

- \* prediction

- ▷ favours bigger models

- Bayesian information criterion

- ▷  $BIC = -2 \max_{\boldsymbol{\theta}} \ell(\text{model}) + \log(n) \times |\boldsymbol{\theta}|$

- ▷ motivation

- \* Bayesian model comparison

- \* selection of covariates
  - ▷ favours smaller models
- smaller is better
- can be used to compare non-nested models
- can be used for more general models (cf. next semester)

### Mallows's $C_P$

- criterion specific for linear regression:
  - ▷ suppose that the full model has  $\boldsymbol{\beta}$  of length  $p$
  - ▷ describe the fit (focus on prediction) of its submodel with  $\tilde{\boldsymbol{\beta}}$  of length  $P$
  - ▷ estimate the average mean square error of prediction  $\frac{1}{\sigma^2} \sum_{i=1}^n \mathbb{E}(\hat{Y}_i - \mathbb{E}Y_i)^2$  by
 
$$\frac{1}{\hat{\sigma}_b^2} \sum_{i=1}^n (\hat{Y}_i - Y_i)^2 = \frac{\|\mathbf{e}_s\|^2}{\|\mathbf{e}_b\|^2/(n-p)}$$
- $C_P = \frac{\|\mathbf{e}_s\|^2}{\|\mathbf{e}_b\|^2/(n-p)} - n + 2P$ 
  - ▷ for the full model:  $C_p = p$
  - ▷ models with  $C_P \approx P$  are considered good
  - ▷ we may plot  $C_P$  against  $P$  and choose a small model that has  $C_P \approx P$  (if small is preferred)
  - ▷ related to the AIC

## 7.4.2 Model selection strategies

### To leave out or not to leave out?

- setting  $\beta_i = 0$  if the true  $\beta_i \neq 0$ 
  - i.e. leaving out a covariate that should have been kept
    - ▷ possible bias in the estimators of  $\beta_j$  for  $i \neq j$
    - ▷ possible invalidity of the resulting model (cf. Week 10)
- allowing  $\beta_i \neq 0$  if the true  $\beta_i \approx 0$ 
  - i.e. keeping unnecessary covariates in the model

- ▷ possibly worse estimation of  $\beta_j$  for  $i \neq j$  and larger confidence intervals (cf. Week 11)
- ▷ possibility of overfitting
- ▷ sometimes/often simple explanations are preferable
- conclusion
  - ▷ avoid blind automatic model selection procedures if possible

### Model selection strategies

- step-wise procedures based on p-values of the  $t/F$  test
  - ▷ backward selection
    - \* start with a biggest model, leave out the covariate with the largest p-value, end when p-values for all included covariates are smaller than  $\alpha_{\text{crit}}$
  - ▷ forward selection
    - \* start with a smallest model, add the covariate with the smallest p-value, end when p-values of all non-included covariates are larger than  $\alpha_{\text{crit}}$
  - ▷ step-wise selection
    - \* a combination of forward and backward selection
  - ▷ issues
    - \* non-exhaustive search
    - \* multiple testing; tests invalid unless the smaller model is true
    - \* not recommended for prediction
- step-wise procedures with a model selection criterion
- exhaustive search with a model selection criterion
  - ▷ e.g. plot  $C_P$  or  $R^2$  against the number of predictors

### Notes on model selection

- hierarchical modelling
  - ▷ powers of lower order should be kept in the model if powers of higher order are present
  - ▷ main terms and interactions of lower order should be kept in the model if interactions of higher order are present
  - ▷ there may be a good reason for a non-hierarchical model but such a model is not invariant to affine transformations and rotations of covariates

- several models may fit equally well
  - ▷ if they give qualitatively different answers, reconsider the use of the data to answer the question
- avoid blind automatic model selection procedures if possible
  - ▷ if impossible, choose a selection procedure to fit the purpose of the modelling and carefully examine the final model
- make sure that the models you considered were fitted to the same data

### Concluding notes

- there is no best/foolproof way to do the model selection except for common sense and sound understanding of the phenomenon
- *A model should be as simple as possible but no simpler.*

Albert Einstein

- *All models are wrong but some are useful.*

George Box

# Chapter 8

## Model diagnostics

### 8.1 The problem

#### 8.1.1 Normal linear model

##### Normal linear model

- $Y_i = \beta_0 + \beta_1 x_{i,1} + \dots + \beta_k x_{i,k} + \varepsilon_i, i \in \{1, \dots, n\}$ 
  - ▷  $Y_i$ : outcome, response, output, dependent variable
    - \* random variable, we observe a realization  $y_i$
    - \* (odezva, závisle proměnná, regresand)
  - ▷  $x_{i,1}, \dots, x_{i,k}$ : covariates, predictors, explanatory variables, input, independent variables
    - \* given, known
    - \* (nezávisle proměnné, regresory)
  - ▷  $\beta_0, \dots, \beta_k$ : coefficients
    - \* unknown
    - \* (regresní koeficienty)
  - ▷  $\varepsilon_i$ : random error
    - \* random variable, unobserved
- $\varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2), i \in \{1, \dots, n\}$ 
  - ▷  $E \varepsilon_i = 0$ : no systematic errors
  - ▷  $\text{Var } \varepsilon_i = \sigma^2$ : same precision

##### Example: bloodpress data



- from [sites.stat.psu.edu/~lsimon/stat501wc/sp05/data/](http://sites.stat.psu.edu/~lsimon/stat501wc/sp05/data/)
- association between the mean arterial blood pressure[mmHg] and age[years], weight[kg], body surface area[m<sup>2</sup>], duration of hypertension[years], basal pulse[beats/min], stress

○ data:

	BP	Age	Weight	BSA	DoH	Pulse	Stress
	105	47	85.4	1.75	5.1	63	33
	115	49	94.2	2.10	3.8	70	14
	...	...	...	...	...	...	...
	110	48	90.5	1.88	9.0	71	99
	122	56	95.7	2.09	7.0	75	99

- model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$

$$\begin{pmatrix} 105 \\ 115 \\ \dots \\ 110 \\ 122 \end{pmatrix} = \begin{pmatrix} 1 & 47 & 85.4 & 1.75 & 5.1 & 63 & 33 \\ 1 & 49 & 94.2 & 2.10 & 3.8 & 70 & 14 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & 48 & 90.5 & 1.88 & 9.0 & 71 & 99 \\ 1 & 56 & 95.7 & 2.09 & 7.0 & 75 & 99 \end{pmatrix} \times \begin{pmatrix} \beta_0 \\ \dots \\ \beta_6 \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \dots \\ \varepsilon_{19} \\ \varepsilon_{20} \end{pmatrix}$$

### Example: fev data

- from: <http://www.statsci.org/data/general/fev.html>
- question: association between the FEV[l] and Smoking, corrected for Age[years], Height[cm] and Gender

○ data:

	FEV	Age	Height	Gender	Smoking
	1.708	9	144.8	Female	Non
	1.724	8	171.5	Female	Non
	1.720	7	138.4	Female	Non
	1.558	9	134.6	Male	Non
	...	...	...	...	...
	3.727	15	172.7	Male	Current
	2.853	18	152.4	Female	Non
	2.795	16	160.0	Female	Current
	3.211	15	168.9	Female	Non

- model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$

$$\begin{pmatrix} 1.708 \\ 1.724 \\ 1.720 \\ 1.558 \\ \dots \\ 3.727 \\ 2.853 \\ 2.795 \\ 3.211 \end{pmatrix} = \begin{pmatrix} 1 & 9 & 144.8 & 0 & 0 \\ 1 & 8 & 171.5 & 0 & 0 \\ 1 & 7 & 138.4 & 0 & 0 \\ 1 & 9 & 134.6 & 1 & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 1 & 15 & 172.7 & 1 & 1 \\ 1 & 18 & 152.4 & 0 & 0 \\ 1 & 16 & 160.0 & 0 & 1 \\ 1 & 15 & 168.9 & 0 & 0 \end{pmatrix} \times \begin{pmatrix} \beta_0 \\ \dots \\ \beta_5 \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \varepsilon_3 \\ \varepsilon_4 \\ \dots \\ \varepsilon_{651} \\ \varepsilon_{652} \\ \varepsilon_{653} \\ \varepsilon_{654} \end{pmatrix}$$

## 8.1.2 Task for this chapter

### Checking the model assumptions

- model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ 
  - ▷ outcome  $\mathbf{Y}$

- \* random vector, we observe a realization  $\mathbf{y}$
- ▷ predictors  $\mathbf{x}_1, \dots, \mathbf{x}_k$ 
  - \* vector of given (known) constants
- ▷ coefficients  $\boldsymbol{\beta}$ 
  - \* vector of unknown constants
- ▷ error  $\boldsymbol{\varepsilon}$ 
  - \* unknown random vector, we do not observe its realization
- ▷ assumptions:  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$ 
  - \*  $E \mathbf{Y} = \mathbf{X}\boldsymbol{\beta}$ : the expected value of  $\mathbf{Y}$  is a linear function of  $\boldsymbol{\beta}$
  - \*  $E \boldsymbol{\varepsilon} = \mathbf{0}$ : no systematic errors
  - \*  $\text{Var } \boldsymbol{\varepsilon} = \sigma^2 \mathbf{I}$ : independence and same precision
- **task**: do the assumptions appear to be satisfied?
- Note: if they are not, inference is not valid ...

## 8.2 Random errors and residuals

### Random errors

#### Random errors in the normal linear model

- model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$
- assumptions
  - ▷  $E \mathbf{Y} = \mathbf{X}\boldsymbol{\beta}$ : the expected value of  $\mathbf{Y}$  is a linear function of  $\boldsymbol{\beta}$
  - ▷  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$ 
    - \*  $E \boldsymbol{\varepsilon} = \mathbf{0}$ : no systematic errors
    - \*  $\text{Var } \boldsymbol{\varepsilon} = \sigma^2 \mathbf{I}$ : independence and the same precision
- we need to verify the assumptions on
  - ▷ expectation:  $E \mathbf{Y} = \mathbf{X}\boldsymbol{\beta}$ , i.e.  $E \boldsymbol{\varepsilon} = \mathbf{0}$
  - ▷ variance:  $\text{Var } \boldsymbol{\varepsilon} = \sigma^2 \mathbf{I}$
  - ▷ distribution:  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$
- all assumptions are made on unobserved random errors  $\boldsymbol{\varepsilon}$
- fitted model:  $\mathbf{Y} = \hat{\mathbf{Y}} + (\mathbf{Y} - \hat{\mathbf{Y}}) = \mathbf{X}\hat{\boldsymbol{\beta}} + \mathbf{e}$

- residuals  $\mathbf{e}$  sometimes seen as “estimates” of  $\boldsymbol{\varepsilon}$ 
  - ▷  $\boldsymbol{\varepsilon}$  is an unobserved random vector, not a parameter (constant)
  - ▷  $\mathbf{e}$  are not estimates in the usual sense

## Residuals

### Residuals in the normal linear model

- model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$
- fitted model:  $\mathbf{Y} = \hat{\mathbf{Y}} + \mathbf{e} = \mathbf{H}\mathbf{Y} + (\mathbf{I} - \mathbf{H})\mathbf{Y}$
- $\mathbf{e} \sim N(\mathbf{0}, \sigma^2 (\mathbf{I} - \mathbf{H}))$ 
  - ▷ proof: cf. Week 6 or use MVN 3
  - ▷  $\text{rank}(\mathbf{I} - \mathbf{H}) = n - p$  if  $\text{rank}(\mathbf{X}) = p$ 
    - \*  $\mathbf{e} \stackrel{d}{=} \mathbf{A}\mathbf{Z}$  for an  $(n - p)$ -dimensional  $\mathbf{Z} \sim N(\mathbf{0}, \mathbf{I})$
    - $\mathbf{I} - \mathbf{H} = \mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^\top$  (spec. dec.)  $\Rightarrow \mathbf{A} = \mathbf{U}_{n \times (n-p)} \boldsymbol{\Lambda}_{(n-p) \times (n-p)}^{1/2}$
    - (cf. Week 4 or use MVN 3)
- if the assumptions are satisfied, residuals are
  - ▷ zero-mean
  - ▷ with unequal variances:  $\text{Var } e_i = \sigma^2 (1 - h_{i,i})$
  - ▷ with a degenerate normal distribution
  - ▷ correlated:  $\text{Cor}(e_i, e_j) = -\frac{h_{i,j}}{\sqrt{(1-h_{i,i})(1-h_{j,j})}}$
- compare to  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}) \dots$

### Standardized residuals in the normal linear model

- model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ 
  - ▷  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$
- fitted model:  $\mathbf{Y} = \hat{\mathbf{Y}} + \mathbf{e} = \mathbf{H}\mathbf{Y} + (\mathbf{I} - \mathbf{H})\mathbf{Y}$ 
  - ▷  $\mathbf{e} \sim N(\mathbf{0}, \sigma^2 (\mathbf{I} - \mathbf{H}))$
- to check the assumptions, we often use

standardized residuals  $r_i = \frac{e_i}{\sqrt{\hat{\sigma}^2 (1 - h_{i,i})}}$ ,  $1 \leq i \leq n$

- if the assumptions are satisfied
  - ▷ we expect that  $r_i \approx N(0, 1)$ 
    - \* it can be shown that  $E r_i = 0$  and  $\text{Var } r_i = 1$   
(some technical work needed to prove this)
    - \* we did not derive the distribution of  $r_i$ 's
    - \* we did not try to get rid of the correlation
- compare to  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}) \dots$

## 8.3 Model diagnostics I: checking the assumptions

### 8.3.1 General principles

#### Checking the assumptions

- Specifying the possible departures
  - ▷ need to specify in what sense the assumption might be violated
  - ▷ if the assumption is  $H_0$ , need to specify  $H_1$
- 1. Graphical checking
  - plots that allows us to “see” departures from the assumptions
  - based on residuals ( $\mathbf{e}$  or  $\mathbf{r}$ )
- 2. Testing the validity of assumptions
  - usually by fitting a more general model that allows them not to be satisfied and testing whether the generalization is needed
  - useful as numerical indications BUT
  - we cannot “prove the null hypothesis”
  - problems with the validity of inference:
    - ▷ chains of tests and multiple testing
    - ▷ assumptions on assumptions
    - ▷ we should \*know\* in advance they are satisfied

#### Overall check: residuals versus fitted values

- $\mathbf{e} \perp \hat{\mathbf{Y}}$  by definition

- no systematic patterns should appear between  $\mathbf{e}$  and  $\hat{\mathbf{Y}}$

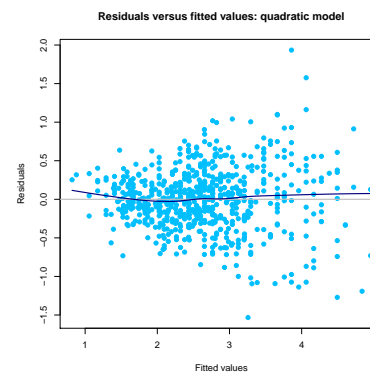
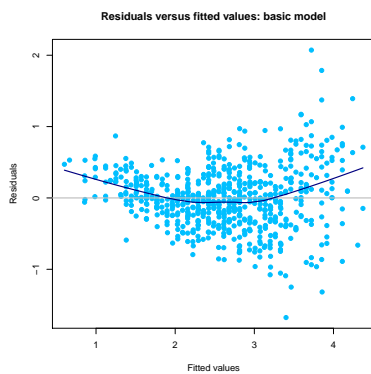
- example: fev data

- ▷ basic model:

$$Y_i = \beta_0 + \beta_1 \times \text{Height}_i + \beta_2 \times \mathbb{I}\{\text{the } i^{\text{th}} \text{ child is male}\} + \varepsilon_i, \quad 1 \leq i \leq 654$$

- ▷ quadratic model:

$$Y_i = \beta_0 + \beta_1 \times \text{Height}_i + \beta_2 \times \text{Height}_i^2 + \\ + \beta_3 \times \mathbb{I}\{\text{the } i^{\text{th}} \text{ child is male}\} + \beta_4 \times \text{Height}_i \mathbb{I}\{\text{the } i^{\text{th}} \text{ child is male}\} + \\ + \beta_5 \times \text{Height}_i^2 \mathbb{I}\{\text{the } i^{\text{th}} \text{ child is male}\} + \varepsilon_i, \quad 1 \leq i \leq 654$$



### 8.3.2 Assumptions on the expectation

#### Checking $E\varepsilon = 0$ , i.e. $EY = X\beta$

- suspected departures from the assumption

- ▷ incorrectly specified form of dependence

- \* plot  $\mathbf{e}$  against the included covariates

- \*  $\mathbf{e} \perp \mathbf{x}_i, 1 \leq i \leq p$ , by definition

- \* no systematic patterns should appear between  $\mathbf{e}$  and  $\mathbf{x}_i$

- \* a trend indicates a dependence not captured by the model

- \* a formal test: fit a more complicated dependence and test against the original model

- ▷ missing covariates

- \* plot  $\mathbf{e}$  against covariates that are not included in the model

- \* no systematic patterns should appear

- \* a trend indicates a dependence not captured by the model

- \* a formal test: fit a larger model and test the effect of the additional covariate

## Incorrectly specified form of dependence

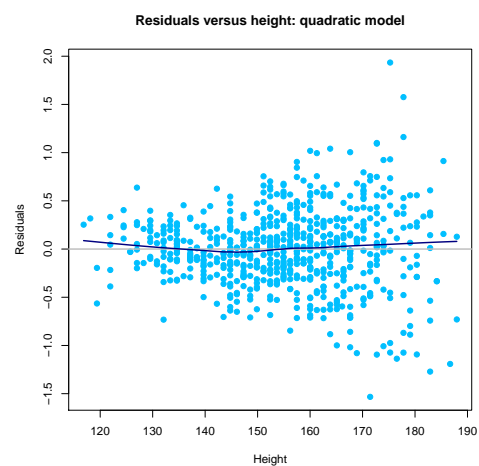
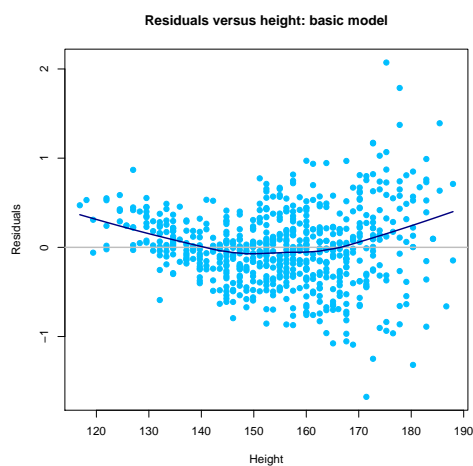
○ example: fev data

▷ basic model:

$$Y_i = \beta_0 + \beta_1 \times \text{Height}_i + \beta_2 \times \mathbb{I}\{\text{the } i^{\text{th}} \text{ child is male}\} + \varepsilon_i, \quad 1 \leq i \leq 654$$

▷ quadratic model:

$$\begin{aligned} Y_i = & \beta_0 + \beta_1 \times \text{Height}_i + \beta_2 \times \text{Height}_i^2 + \\ & + \beta_3 \times \mathbb{I}\{\text{the } i^{\text{th}} \text{ child is male}\} + \beta_4 \times \text{Height}_i \mathbb{I}\{\text{the } i^{\text{th}} \text{ child is male}\} + \\ & + \beta_5 \times \text{Height}_i^2 \mathbb{I}\{\text{the } i^{\text{th}} \text{ child is male}\} + \varepsilon_i, \quad 1 \leq i \leq 654 \end{aligned}$$

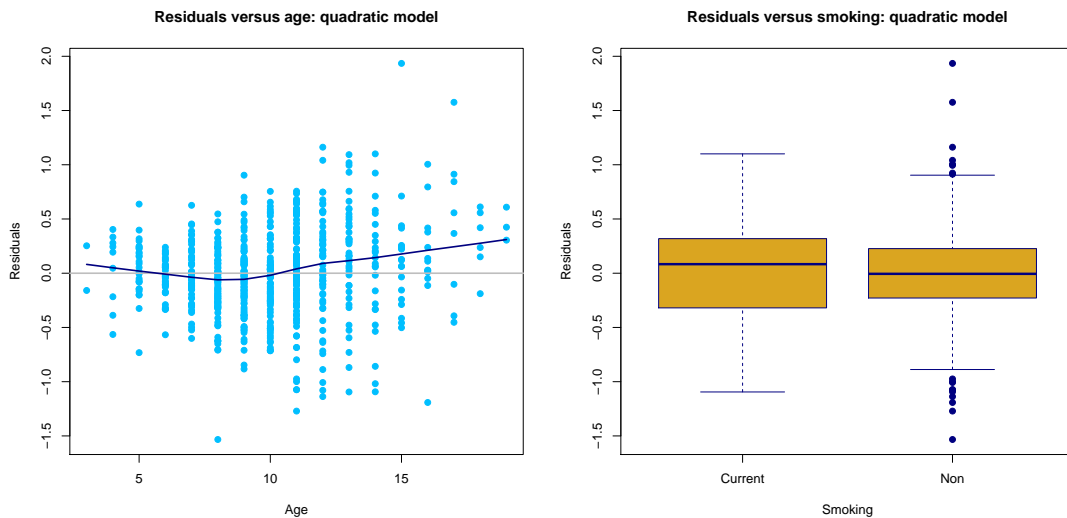


## Missing covariates

○ example: fev data

○ quadratic model:

$$\begin{aligned} Y_i = & \beta_0 + \beta_1 \times \text{Height}_i + \beta_2 \times \text{Height}_i^2 + \\ & + \beta_3 \times \mathbb{I}\{\text{the } i^{\text{th}} \text{ child is male}\} + \beta_4 \times \text{Height}_i \mathbb{I}\{\text{the } i^{\text{th}} \text{ child is male}\} + \\ & + \beta_5 \times \text{Height}_i^2 \mathbb{I}\{\text{the } i^{\text{th}} \text{ child is male}\} + \varepsilon_i, \quad 1 \leq i \leq 654 \end{aligned}$$



### 8.3.3 Assumptions on the variance

#### Checking $\text{Var } \varepsilon = \sigma^2 \mathbf{I}$ : homoskedasticity

- suspected departures from the assumption
  - ▷ variance changing with fitted values (usually increasing)
    - \* plot standardized residuals (usually square root of the absolute value) against fitted values
    - \* no pattern should appear
  - ▷ variance changing with covariates
    - \* plot standardized residuals (usually square root of the absolute value) against covariates
    - \* no pattern should appear
    - \* a formal test: studentized Breusch–Pagan test
  - ▷ subgroups with the same within-group variance
    - \* plot boxplots of standardized residuals by groups
    - \* boxes should be of approximately equal sizes
    - \* a formal test: fit a more general model and test against the original model

#### Breusch–Pagan test

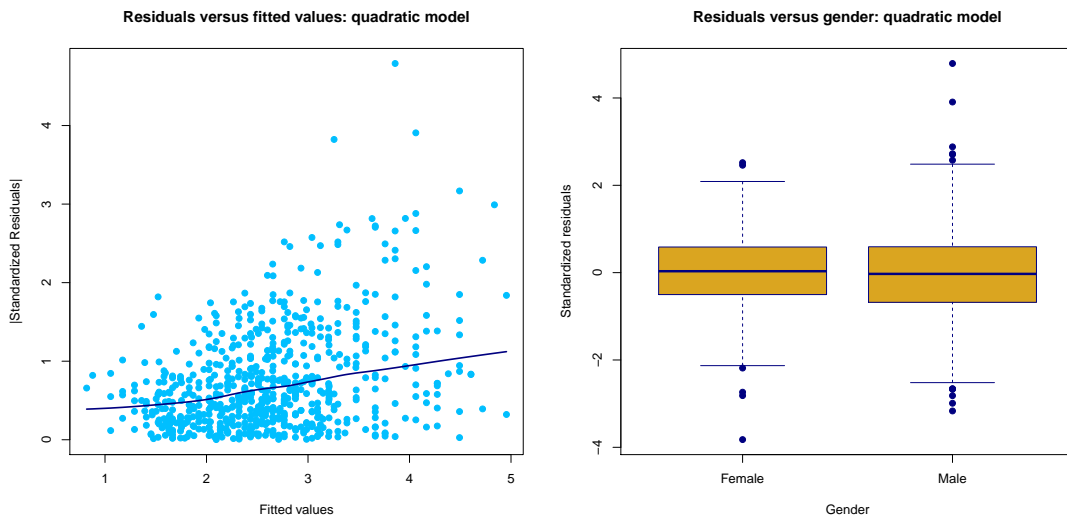
- original model
  - ▷  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$

- ▷  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$
- more general model
  - ▷  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$
  - ▷  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \text{diag}(\sigma_1^2, \dots, \sigma_n^2))$
  - ▷  $\sigma^2 = \mathbf{X}\boldsymbol{\alpha}$
- Breusch–Pagan test: test  $\boldsymbol{\alpha}_{2:p} = \mathbf{0}$  in the more general model
- studentized Breusch–Pagan test less sensitive to the assumption of normality
- more general versions of the Breusch–Pagan test and more general tests exist

### Checking $\text{Var } \boldsymbol{\varepsilon} = \sigma^2 \mathbf{I}$ : homoskedasticity

- example: fev data
- quadratic model:

$$\begin{aligned}
 Y_i = & \beta_0 + \beta_1 \times \text{Height}_i + \beta_2 \times \text{Height}_i^2 + \\
 & + \beta_3 \times \mathbb{I}\{\text{the } i^{\text{th}} \text{ child is male}\} + \beta_4 \times \text{Height}_i \mathbb{I}\{\text{the } i^{\text{th}} \text{ child is male}\} + \\
 & + \beta_5 \times \text{Height}_i^2 \mathbb{I}\{\text{the } i^{\text{th}} \text{ child is male}\} + \varepsilon_i, \quad 1 \leq i \leq 654
 \end{aligned}$$



### Checking $\text{Var } \boldsymbol{\varepsilon} = \sigma^2 \mathbf{I}$ : independence

- suspected departures from the assumption
  - ▷ clustering



- \* suspected e.g. when several data points collected from one individual (e.g. same individuals followed over time)
  - \* plot boxplots of residuals by the suspected groups
  - \* no pattern should appear
  - \* a formal test: fit a more general model allowing for the within-group dependence and test against the original model
- ▷ serial correlation
- \* suspected when data collected over time or space
  - \* plot  $e_i$  against  $e_{i-1}$
  - \* no pattern should appear
  - \* plot the (partial) autocorrelation function
  - \* a formal test: fit a more general model and test against the original model
  - \* a formal test: Durbin–Watson test

### Durbin–Watson test

- original model
  - ▷  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$
  - ▷  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$
- more general model
  - ▷  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$
  - ▷  $\varepsilon_i = \rho \varepsilon_{i-1} + w_i, w_i \stackrel{\text{iid}}{\sim} (0, \sigma^2), |\rho| < 1$   
(autoregression of the first order on the error terms)
- Durbin–Watson test: test  $\rho = 0$  against  $\rho > 0$  in the more general model
- also possible to test  $\rho = 0$  against  $\rho < 0$  and  $\rho = 0$  against  $\rho \neq 0$
- more general tests available

### Time series models

- **time series** is a random sequence  $\{X_t, t \in \mathbb{Z}\}$ 
  - ▷ **stationary** if  $E X_t = \mu, \text{Var } X_t = \sigma^2, \text{Cov}(X_t, X_{t+s}) = \gamma(s)$
- The **autocovariance function** of a stationary random sequence  $\{X_t, t \in \mathbb{Z}\}$  is defined as  $\gamma(h) = \text{Cov}(X_t, X_{t+h}), h \in \mathbb{Z}$ .

- The **autocorrelation function** (ACF) is defined as  $\rho(h) = \text{Cor}(X_t, X_{t+h}) = \gamma(h)/\gamma(0)$ ,  $h \in \mathbb{Z}$ .
- The **partial autocorrelation function** (PACF) is defined as  $\alpha(1) = \text{Cor}(X_t, X_{t+1}) = \rho(1)$  and  $\alpha(h) = \text{Cor}(X_t - \hat{X}_t, X_{t+h} - \hat{X}_{t+h})$ ,  $h = 2, 3, \dots$ , where  $\hat{X}_t$  and  $\hat{X}_{t+h}$  are the fitted values from the linear regressions  $X_t \sim X_{t+1}, \dots, X_{t+h-1}$  and  $X_{t+h} \sim X_{t+1}, \dots, X_{t+h-1}$ .

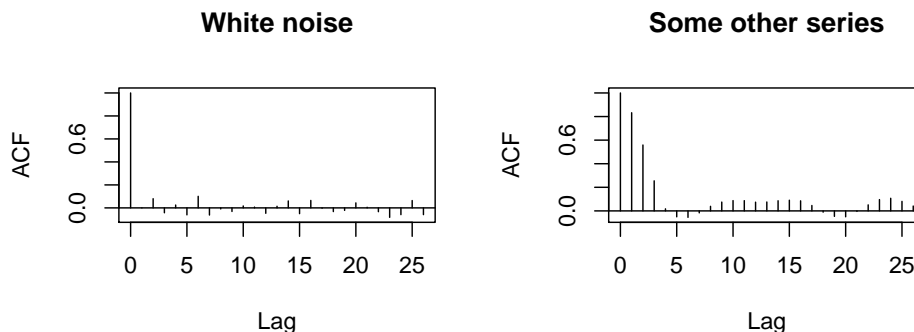
### ACF and PACF for ARMA models

- special time series models
- Let  $\{\epsilon_t\} \stackrel{iid}{\sim} (0, \sigma^2)$ . Then  $\{X_t, t \in \mathbb{Z}\}$  is
  - ▷ **AR(p)** if
    - \*  $X_t = \phi_1 X_{t-1} + \dots + \phi_p X_{t-p} + \epsilon_t$ ;
  - ▷ **MA(q)** if
    - \*  $X_t = \epsilon_t + \theta_1 \epsilon_{t-1} + \dots + \theta_q \epsilon_{t-q}$ ;
  - ▷ **ARMA(p, q)** if
    - \*  $X_t = \phi_1 X_{t-1} + \dots + \phi_p X_{t-p} + \epsilon_t + \theta_1 \epsilon_{t-1} + \dots + \theta_q \epsilon_{t-q}$ .
- ACF and PACF for AR/MA

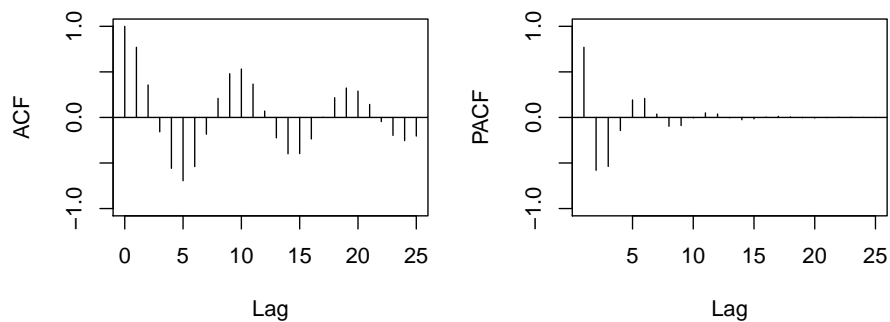
	ACF	PACF
AR(p)	Exponential decay	Cuts off after lag p
MA(q)	Cuts off after lag q	Exponential decay
ARMA(p, q)	Exponential decay	Exponential decay

### ACFs and PACFs

- simulated ACFs

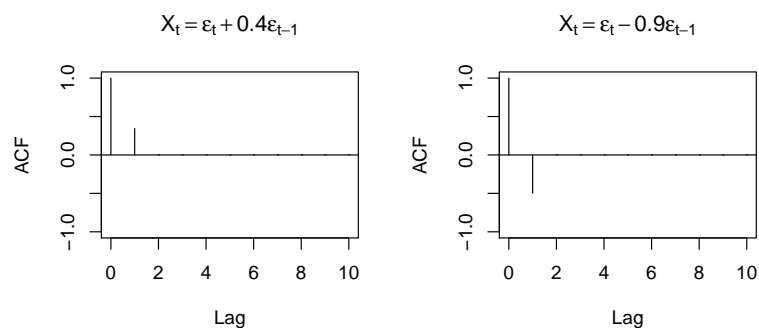


- theoretical ACF and PACF for an ARMA

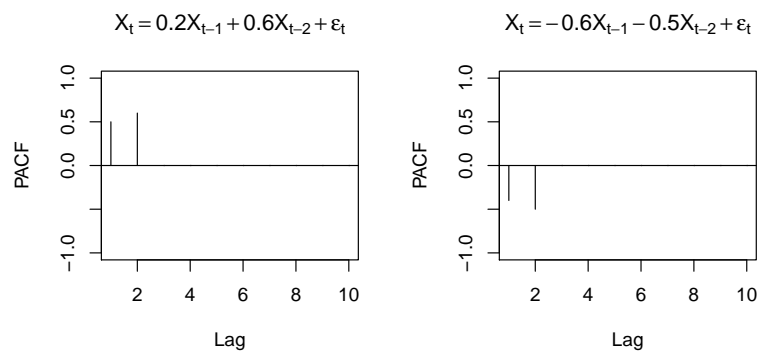


### ACFs and PACFs

- theoretical ACF for MA(1)



- theoretical PACF for AR(2)



### Other types of dependence

- spatial correlation diagnosed via **semivariogram**

▷ for a stationary isotropic random field  $\{Z(\mathbf{x}); \mathbf{x} \in \mathbb{R}^2\}$ , **semivariogram** is

$$* \gamma(\mathbf{x}, \mathbf{y}) = \frac{1}{2} \text{Var}(Z(\mathbf{x}) - Z(\mathbf{y})) = \frac{1}{2} \text{E}(Z(\mathbf{x}) - Z(\mathbf{y}))^2 = \gamma(h), \text{ where } h = \|\mathbf{x} - \mathbf{y}\|^2$$

- clustering (boxplots of residuals by group)

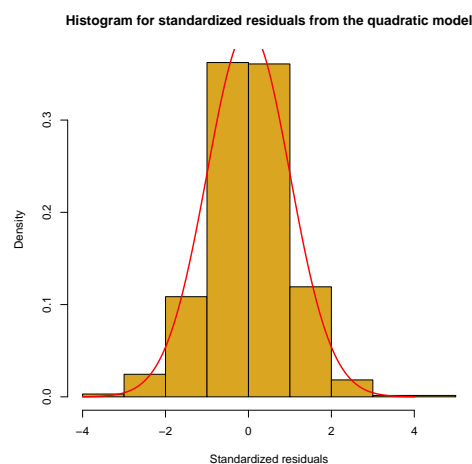
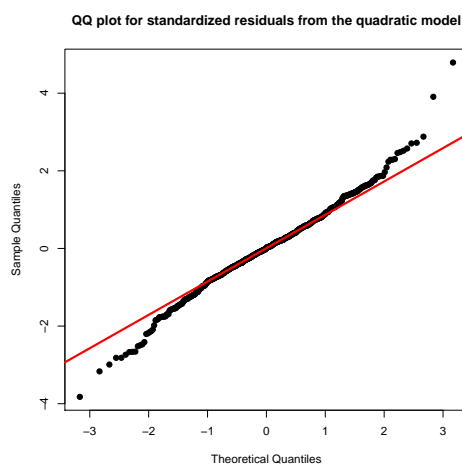
### 8.3.4 Assumptions on the distribution

#### Checking $\varepsilon \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$

- suspected departures from the assumption
  - ▷ non-normal distribution
    - \* skewed distribution
    - \* heavy-tailed distribution
- plot a QQ plot for (standardized) residuals
- plot a histogram for (standardized) residuals
- formal tests: Shapiro–Wilk test, Kolmogorov–Smirnov test
  - ▷ warning: valid for iid’s (and residuals are not iid’s)

#### QQ plot and histogram

- QQ plot (preferred)
  - ▷ quantiles of  $N(0, 1)$  against empirical quantiles
  - ▷ should be near a straight line
  - ▷ problems to look for
    - \* S shape (heavy tails)
    - \* an arc (skewness)



#### Shapiro–Wilk test and Kolmogorov–Smirnov test

- valid for iid’s (and residuals are not iid’s)

## ▷ Shapiro–Wilk test

- \* can be seen as a numerical summary of the QQ plot
- \* rather a strong one

```
> shapiro.test(rstandard(model.basic.quad))

Shapiro-Wilk normality test

data:  rstandard(model.basic.quad)
W = 0.9865, p-value = 9.713e-06

> shapiro.test(rstandard(model.basic.quad)[sample(1:654, 50)])

Shapiro-Wilk normality test

data:  rstandard(model.basic.quad)[sample(1:654, 50)]
W = 0.97011, p-value = 0.2338
```

## ▷ Kolmogorov–Smirnov test

- \* rather a weak one

**Importance of the assumption**○ large-sample distribution of  $\hat{\beta}$ 

- ▷ Let  $\mathbf{X}_n$ ;  $n \in \mathbb{N}$  be a sequence of  $n \times p$  design matrices of full rank defining a sequence of linear models  $\mathbf{Y}_n = \mathbf{X}_n \boldsymbol{\beta} + \boldsymbol{\varepsilon}_n$  with  $\boldsymbol{\varepsilon}_n \sim (\mathbf{0}, \sigma^2 \mathbf{I}_n)$ . If  $\max_{1 \leq i \leq n} \mathbf{x}_i^\top (\mathbf{X}_n^\top \mathbf{X}_n)^{-1} \mathbf{x}_i \xrightarrow[n \rightarrow \infty]{} 0$  then

$$(\mathbf{X}_n^\top \mathbf{X}_n)^{1/2} (\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}) \xrightarrow[n \rightarrow \infty]{d} \mathbf{N}(\mathbf{0}, \sigma^2 \mathbf{I}),$$

where  $\hat{\boldsymbol{\beta}}_n = (\mathbf{X}_n^\top \mathbf{X}_n)^{-1} \mathbf{X}_n^\top \mathbf{Y}_n$ .

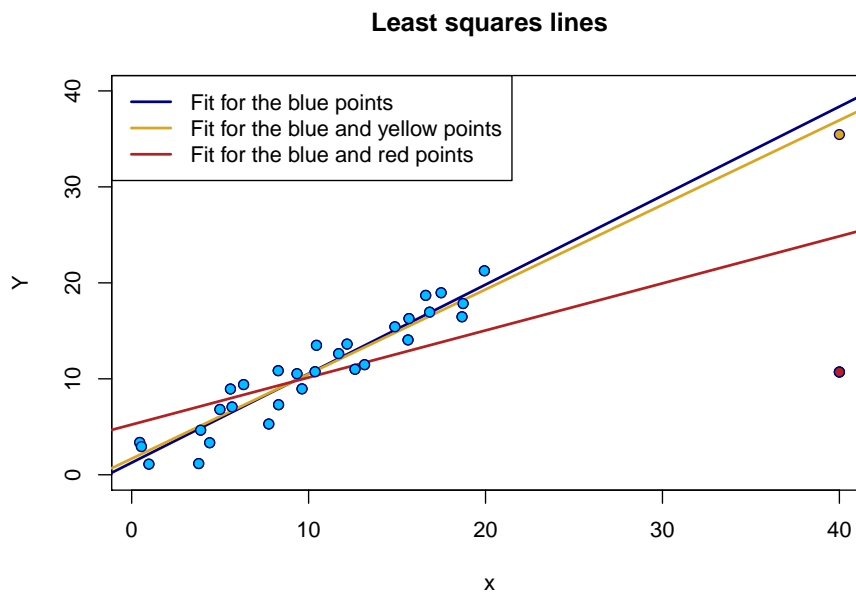
- normality not crucial in large samples unless there are special observations

**8.4 Model diagnostics II: influential and unusual observations****8.4.1 Observations to look at****Leverage**

- $\mathbf{Y} = \hat{\mathbf{Y}} + \mathbf{e} = \mathbf{H}\mathbf{Y} + (\mathbf{I} - \mathbf{H})\mathbf{Y}$
- $\text{Var } \hat{Y}_i = h_{i,i} \dots$  leverage
- $\mathbf{H} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$  and  $\text{rank}(\mathbf{H}) = \text{tr}(\mathbf{H}) = p$
- $h_{i,i} = \mathbf{x}_i^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_i$ , and  $\sum_{i=1}^n h_{i,i} = p$

- the variance of  $\hat{Y}_i$  determined by the corresponding covariates
- we want all observations to contribute  $\approx$  equally to the fit
  - ▷ we want that  $h_{i,i} \approx \frac{p}{n}$
- if  $h_{i,i}$  much larger for some  $i$ , the fit may be influenced by  $(Y_i, \mathbf{x}_i)$  much more than by the other observations
- observations with  $h_{i,i} > \frac{2p}{n}$  should be checked

### Potentially influential and influential observations



- both points have a high leverage, but only one is influential

### Model with an excluded observation

- consider a model  $\mathbf{Y}_{[-i]} = \mathbf{X}_{[-i]} \boldsymbol{\beta} + \boldsymbol{\varepsilon}_{[-i]}$  without the  $i^{\text{th}}$  observation
- fit the model
  - ▷ compute  $\hat{\boldsymbol{\beta}}_{[-i]}$  and  $\hat{\sigma}^2_{[-i]}$
- compute  $\hat{y}_{[-i]} = \mathbf{x}_i^\top \hat{\boldsymbol{\beta}}_{[-i]}$ 
  - ▷ prediction of  $y_i$  based on the model without the  $i^{\text{th}}$  observation
- if  $y_i - \hat{y}_{[-i]}$  is large, the  $i^{\text{th}}$  observation is an outlier
  - ▷ how large is “too large?”

- ▷  $\text{Var}(y_i - \hat{y}_{[-i]}) = \sigma^2(1 + \mathbf{x}_i^\top (\mathbf{X}_{[-i]}^\top \mathbf{X}_{[-i]})^{-1} \mathbf{x}_i)$
- ▷ define jackknife residuals  $t_i = \frac{y_i - \hat{y}_{[-i]}}{\sqrt{\hat{\sigma}_{[-i]}^2(1 + \mathbf{x}_i^\top (\mathbf{X}_{[-i]}^\top \mathbf{X}_{[-i]})^{-1} \mathbf{x}_i)}$
- ▷ there is a simpler equivalent formula that does not require fitting  $n$  models with excluded observations

## Influential and unusual observations

- in the normal linear model:
  - ▷  $t_i \sim t_{n-p-1}$
- we can test whether an observation is an **outlier**
  - ▷ heavy multiple testing  $\rightsquigarrow$  Bonferroni correction
    - \* use  $t_{n-p-1}(1 - \alpha/(2n))$  instead of  $t_{n-p-1}(1 - \alpha/2)$
- to evaluate whether the observation is **influential**
  - ▷ **Cook's distance**:  $d_i = \frac{1}{p \hat{\sigma}^2} \|\hat{\mathbf{Y}} - \hat{\mathbf{Y}}_{[-i]}\|^2 = \frac{1}{p} r_i^2 \frac{h_{i,i}}{1 - h_{i,i}}$
  - ▷ how large is “too large”?
    - \* rule of thumb:  $d_i \geq 0.5$  deserve some attention
    - $d_i \geq 1 \rightsquigarrow$  highly influential observation

# Chapter 9

## Reduced-rank design matrix and multicollinearity

### 9.1 The problem

#### 9.1.1 Normal linear model

##### Normal linear model

- $Y_i = \beta_0 + \beta_1 x_{i,1} + \dots + \beta_k x_{i,k} + \varepsilon_i, i \in \{1, \dots, n\}$ 
  - ▷  $Y_i$ : outcome, response, output, dependent variable
    - \* random variable, we observe a realization  $y_i$
    - \* (odezva, závisle proměnná, regresand)
  - ▷  $x_{i,1}, \dots, x_{i,k}$ : covariates, predictors, explanatory variables, input, independent variables
    - \* given, known
    - \* (nezávisle proměnné, regresory)
  - ▷  $\beta_0, \dots, \beta_k$ : coefficients
    - \* unknown
    - \* (regresní koeficienty)
  - ▷  $\varepsilon_i$ : random error
    - \* random variable, unobserved
- $\varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2), i \in \{1, \dots, n\}$ 
  - ▷  $E \varepsilon_i = 0$ : no systematic errors
  - ▷  $\text{Var } \varepsilon_i = \sigma^2$ : same precision



### Example: bloodpress data

- o from [sites.stat.psu.edu/~lmsimon/stat501wc/sp05/data/](http://sites.stat.psu.edu/~lmsimon/stat501wc/sp05/data/)
- o association between the mean arterial blood pressure[mmHg] and age[years], weight[kg], body surface area[m<sup>2</sup>], duration of hypertension[years], basal pulse[beats/min], stress

o data:

	BP	Age	Weight	BSA	DoH	Pulse	Stress
	105	47	85.4	1.75	5.1	63	33
	115	49	94.2	2.10	3.8	70	14
	...	...	...	...	...	...	...
	110	48	90.5	1.88	9.0	71	99
	122	56	95.7	2.09	7.0	75	99

- o model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$

$$\begin{pmatrix} 105 \\ 115 \\ \dots \\ 110 \\ 122 \end{pmatrix} = \begin{pmatrix} 1 & 47 & 85.4 & 1.75 & 5.1 & 63 & 33 \\ 1 & 49 & 94.2 & 2.10 & 3.8 & 70 & 14 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & 48 & 90.5 & 1.88 & 9.0 & 71 & 99 \\ 1 & 56 & 95.7 & 2.09 & 7.0 & 75 & 99 \end{pmatrix} \times \begin{pmatrix} \beta_0 \\ \dots \\ \beta_6 \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \dots \\ \varepsilon_{19} \\ \varepsilon_{20} \end{pmatrix}$$

<https://ww2.amstat.org/publications/jse/v13n2/datasets.kahn.html>

### Example: fev data

- o from: <http://www.statsci.org/data/general/fev.html>
- o question: association between the FEV[l] and Smoking, corrected for Age[years], Height[cm] and Gender

o data:

	FEV	Age	Height	Gender	Smoking
	1.708	9	144.8	Female	Non
	1.724	8	171.5	Female	Non
	1.720	7	138.4	Female	Non
	1.558	9	134.6	Male	Non
	...	...	...	...	...
	3.727	15	172.7	Male	Current
	2.853	18	152.4	Female	Non
	2.795	16	160.0	Female	Current
	3.211	15	168.9	Female	Non

- o model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$

$$\begin{pmatrix} 1.708 \\ 1.724 \\ 1.720 \\ 1.558 \\ \dots \\ 3.727 \\ 2.853 \\ 2.795 \\ 3.211 \end{pmatrix} = \begin{pmatrix} 1 & 9 & 144.8 & 0 & 0 \\ 1 & 8 & 171.5 & 0 & 0 \\ 1 & 7 & 138.4 & 0 & 0 \\ 1 & 9 & 134.6 & 1 & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 1 & 15 & 172.7 & 1 & 1 \\ 1 & 18 & 152.4 & 0 & 0 \\ 1 & 16 & 160.0 & 0 & 1 \\ 1 & 15 & 168.9 & 0 & 0 \end{pmatrix} \times \begin{pmatrix} \beta_0 \\ \dots \\ \beta_5 \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \varepsilon_3 \\ \varepsilon_4 \\ \dots \\ \varepsilon_{651} \\ \varepsilon_{652} \\ \varepsilon_{653} \\ \varepsilon_{654} \end{pmatrix}$$

## 9.1.2 Task for this chapter

### Rank-deficiency/near-rank deficiency of $\mathbf{X}$

- model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ 
  - ▷ outcome  $\mathbf{Y}$ 
    - \* random vector, we observe a realization  $\mathbf{y}$
  - ▷ predictors  $\mathbf{x}_1, \dots, \mathbf{x}_k$ 
    - \* vector of given (known) constants
  - ▷ coefficients  $\boldsymbol{\beta}$ 
    - \* vector of unknown constants
  - ▷ error  $\boldsymbol{\varepsilon}$ 
    - \* unknown random vector, we do not observe its realization
  - ▷ assumptions:  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$ 
    - \*  $E \mathbf{Y} = \mathbf{X}\boldsymbol{\beta}$ : the expected value of  $\mathbf{Y}$  is a linear function of  $\boldsymbol{\beta}$
    - \*  $E \boldsymbol{\varepsilon} = \mathbf{0}$ : no systematic errors
    - \*  $\text{Var } \boldsymbol{\varepsilon} = \sigma^2 \mathbf{I}$ : independence and same precision
- task: so far we have assumed that  $\text{rank}(\mathbf{X}) = p$ 
  - What happens if  $\text{rank}(\mathbf{X}) < p$  or “nearly so”?

## 9.2 Rank-deficient design matrix

### 9.2.1 Rank-deficient design matrix

#### Full-rank design matrix $\mathbf{X}$

- design matrix  $\mathbf{X}$  is  $n \times p$ ,  $n > p$
- SVD:  $\mathbf{X} = \underbrace{\mathbf{U}}_{n \times n} \underbrace{\boldsymbol{\Sigma}}_{n \times p} \underbrace{\mathbf{V}^\top}_{p \times p}$
- if all covariates are linearly independent
  - ▷  $\text{rank}(\mathbf{X}) = p$
  - ▷  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p > 0$
  - ▷ thin SVD:  $\mathbf{X} = \underbrace{\mathbf{U}_1}_{n \times p} \underbrace{\boldsymbol{\Sigma}_1}_{p \times p} \underbrace{\mathbf{V}^\top}_{p \times p}$
  - ▷ the columns generate a  $p$ -dimensional space  $\text{im}(\mathbf{X})$

- \*  $\{\mathbf{x}_1, \dots, \mathbf{x}_p\}$  is a basis of  $\text{im}(\mathbf{X})$
- \*  $\{\mathbf{u}_1, \dots, \mathbf{u}_p\}$  is an orthonormal basis of  $\text{im}(\mathbf{X})$
- \*  $\mathbf{H} = \mathbf{U}_1 \mathbf{U}_1^\top = \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$  is a projection matrix on  $\text{im}(\mathbf{X})$

### Rank-deficient design matrix $\mathbf{X}$

- o design matrix  $\mathbf{X}$  is  $n \times p$ ,  $n > p$
- o SVD:  $\mathbf{X} = \underbrace{\mathbf{U}}_{n \times n} \underbrace{\mathbf{\Sigma}}_{n \times p} \underbrace{\mathbf{V}^\top}_{p \times p}$
- o if covariates are not linearly independent
  - ▷  $\text{rank}(\mathbf{X}) = r < p$
  - ▷  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0 = \sigma_{r+1} = \dots = \sigma_p$
  - ▷ compact SVD:  $\mathbf{X} = \underbrace{\mathbf{U}_1}_{n \times r} \underbrace{\mathbf{\Sigma}_1}_{r \times r} \underbrace{\mathbf{V}^\top}_{r \times r}$
  - ▷ the columns generate an  $r$ -dimensional space  $\text{im}(\mathbf{X})$ 
    - \*  $\{\mathbf{u}_1, \dots, \mathbf{u}_r\}$  is an orthonormal basis of  $\text{im}(\mathbf{X})$
    - \*  $\mathbf{H} = \mathbf{U}_1 \mathbf{U}_1^\top = \mathbf{X} (\mathbf{X}^\top \mathbf{X})^+ \mathbf{X}^\top$  is a projection matrix on  $\text{im}(\mathbf{X})$

### $\hat{\boldsymbol{\beta}}$ motivated by orthogonal projection (reminder)

- o model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon}$  unknown,  $\mathbf{E}\boldsymbol{\varepsilon} = \mathbf{0}$
- o idea: set  $\boldsymbol{\varepsilon} \stackrel{!}{=} \mathbf{0}$  and solve  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta}$  w.r.t.  $\boldsymbol{\beta}$ 
  - ▷ then  $\underbrace{\mathbf{Y}}_{n \times 1} \stackrel{!}{=} \underbrace{\mathbf{X}}_{n \times p} \underbrace{\boldsymbol{\beta}}_{p \times 1}$
  - ▷  $n$  linear equations with  $p$  unknowns and  $n > p$   
 $\Rightarrow$  a solution exists only if  $\mathbf{Y} \in \text{im}(\mathbf{X})$
- o modified idea: find  $\hat{\mathbf{Y}} \in \text{im}(\mathbf{X})$  such that  $\|\mathbf{Y} - \hat{\mathbf{Y}}\|^2$  is the smallest possible and solve  $\hat{\mathbf{Y}} = \mathbf{X}\boldsymbol{\beta}$  w.r.t.  $\boldsymbol{\beta}$ 
  - ▷ then  $\hat{\mathbf{Y}}$  is the orthogonal projection of  $\mathbf{Y}$  onto  $\text{im}(\mathbf{X})$
  - ▷ projection matrix onto  $\text{im}(\mathbf{X})$  is  $\underbrace{\hat{\mathbf{H}}}_{\text{hat matrix}} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^+ \mathbf{X}^\top$
  - ▷ solving  $\hat{\mathbf{Y}} = \mathbf{X}\boldsymbol{\beta}$  is solving  $\mathbf{X}(\mathbf{X}^\top \mathbf{X})^+ \mathbf{X}^\top \mathbf{Y} = \mathbf{X}\boldsymbol{\beta}$
  - ▷ estimate  $\boldsymbol{\beta}$  by  $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^+ \mathbf{X}^\top \mathbf{Y}$
  - ▷ but  $\hat{\boldsymbol{\beta}}$  is the unique solution of  $\hat{\mathbf{Y}} = \mathbf{X}\boldsymbol{\beta}$  iff  $\text{rank}(\mathbf{X}) = p$

\* and then  $\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$

### $\hat{\beta}$ as least squares estimator (reminder)

- **model:**  $\mathbf{Y} = \mathbf{X}\beta + \varepsilon$ ,  $\varepsilon$  unknown,  $E\varepsilon = \mathbf{0}$
- **idea:** make the residuals as small as possible
  - ▷ minimize  $\|\varepsilon\|^2 = \sum_{i=1}^n \varepsilon_i^2$  w.r.t.  $\beta$ 
    - ↪ Least Squares Estimator (LSE)  $\hat{\beta} = \arg \min_{\beta} \sum_{i=1}^n \varepsilon_i^2$
    - ▷ also called the OLS (Ordinary Least Squares) solution
- **computation:**
  - ▷  $\varepsilon = \mathbf{Y} - \mathbf{X}\beta$
  - ▷  $\hat{\beta} = \arg \min_{\beta} \|\mathbf{Y} - \mathbf{X}\beta\|^2 = \arg \min_{\beta} (\mathbf{Y} - \mathbf{X}\beta)^\top (\mathbf{Y} - \mathbf{X}\beta)$
- look for the minimum by differentiating:
  - ▷  $\frac{\partial}{\partial \beta} (\mathbf{Y} - \mathbf{X}\beta)^\top (\mathbf{Y} - \mathbf{X}\beta) \stackrel{!}{=} 0$
  - ▷  $-2 \mathbf{X}^\top \mathbf{Y} + 2 \mathbf{X}^\top \mathbf{X} \beta \stackrel{!}{=} 0$
  - ▷  $\mathbf{X}^\top \mathbf{X} \beta \stackrel{!}{=} \mathbf{X}^\top \mathbf{Y}$ : normal equations
- normal equations have unique solution iff  $\text{rank}(\mathbf{X}) = p$ : then
  - ▷ the solution is  $(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$
  - ▷  $\frac{\partial^2}{\partial \beta \partial \beta} (\mathbf{Y} - \mathbf{X}\beta)^\top (\mathbf{Y} - \mathbf{X}\beta) = 2 \mathbf{X}^\top \mathbf{X} \succ 0$  for all  $\beta$ 
    - ⇒ the solution is the minimum ⇒  $\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$

### If $\text{rank}(\mathbf{X}) = r < p$

- orthogonal projection approach
  - ▷  $\hat{\mathbf{Y}}$  exists and is unique
  - ▷  $\hat{\beta}$  such that  $\hat{\mathbf{Y}} = \mathbf{X}\hat{\beta}$  is a vector of coordinates of  $\hat{\mathbf{Y}} \in \text{im}(\mathbf{X})$  w.r.t.  $\{\mathbf{x}_1, \dots, \mathbf{x}_p\}$ 
    - \* if  $\{\mathbf{x}_1, \dots, \mathbf{x}_p\}$  is not a basis of  $\text{im}(\mathbf{X})$ ,  $\hat{\beta}$  is not unique
  - ▷  $\{\hat{\beta}; \hat{\mathbf{Y}} = \mathbf{X}\hat{\beta}\}$  is a linear subspace of  $\mathbb{R}^p$  of dimension  $p - r$
  - ▷ neither  $\hat{\mathbf{Y}}$  nor  $\|\mathbf{Y} - \hat{\mathbf{Y}}\|^2$  depend on the choice of  $\hat{\beta}$
- ordinary least squares approach
  - ▷ normal equations  $\mathbf{X}^\top \mathbf{X} \beta = \mathbf{X}^\top \mathbf{Y}$  are consistent

- \*  $\text{rank}(\mathbf{X}^\top \mathbf{X}) = \text{rank}((\mathbf{X}^\top \mathbf{Y}, \mathbf{X}^\top \mathbf{X}))$
- ▷ normal equations have infinitely many solutions
  - \* the linear subspace of  $\mathbb{R}^p$  of dimension  $p - r$
- ▷ the minimum  $\min_{\boldsymbol{\beta}} \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2$  is attained for each of the solutions and its value is the same for all the solutions
  - \* the  $\|\mathbf{Y} - \hat{\mathbf{Y}}\|^2$
- ▷ proofs can be found in *Anděl: Základy matematické statistiky*

## 9.2.2 Identifiability

### Identifiable parameters

- $\hat{\mathbf{Y}}$  and  $\|\mathbf{Y} - \hat{\mathbf{Y}}\|^2 = \min_{\boldsymbol{\beta}} \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2$  does not depend on  $\hat{\boldsymbol{\beta}}$
- any other quantities with such properties?

**Theorem.** Let  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  where  $\mathbf{X}$  is an  $n \times p$  matrix,  $\text{rank}(\mathbf{X}) = r < p$ ,  $\boldsymbol{\beta} \in \mathbb{R}^p$ , and  $\boldsymbol{\varepsilon}$  is an  $n$ -dimensional random vector with  $\mathbf{E}\boldsymbol{\varepsilon} = \mathbf{0}$  and  $\text{Var}\boldsymbol{\varepsilon} = \sigma^2 \mathbf{I}$ . Let  $\mathbf{c} \in \mathbb{R}^p$  and  $\theta = \mathbf{c}^\top \boldsymbol{\beta}$ .

If  $\theta \in \text{im}((\mathbf{X}\boldsymbol{\beta})^\top)$ , equivalently if  $\mathbf{c} \in \text{im}(\mathbf{X}^\top)$ , then

- (i) the value of  $\hat{\theta} = \mathbf{c}^\top \hat{\boldsymbol{\beta}}$  where  $\hat{\boldsymbol{\beta}}$  is a solution to the normal equations does not depend on the choice of the solution;
- (ii)  $\exists$  a linear unbiased estimator of  $\theta$ ;
- (iii)  $\hat{\theta} = \mathbf{c}^\top \hat{\boldsymbol{\beta}}$  is BLUE for  $\theta$ .

- parameter  $\theta$  that is a linear combination of  $\mathbf{E}\mathbf{Y}$  is **identifiable**
- a proof can be found in *Jiří Anděl: Základy matematické statistiky*

### Inference for identifiable parameters

- model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon} = \mathbf{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ ,  $\text{rank}(\mathbf{X}) = r < p$
- $\mathbf{E}\mathbf{Y}$  is identifiable
- $\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$  is BLUE for  $\mathbf{E}\mathbf{Y}$  for any  $\hat{\boldsymbol{\beta}}$  that solves the normal equations
- it can be shown that
  - ▷  $\frac{n-r}{\sigma^2} \hat{\sigma}^2 \sim \chi_{n-r}^2$
  - ▷  $\hat{\sigma}^2 = \frac{1}{n-r} \|\mathbf{Y} - \hat{\mathbf{Y}}\|^2$  is an unbiased estimator of  $\sigma^2$

- ▷  $\hat{\sigma}^2 \perp\!\!\!\perp \hat{\beta}$  for any  $\hat{\beta}$  that solves the normal equations
- ▷ proofs are similar to the full-rank case
  - \* can be found in *Jiří Anděl: Základy matematické statistiky (2005). Mat-fyzpress.*
- inference for identifiable parameters and vectors is as in the full-rank model but we need to adjust the degrees of freedom
  - ▷  $n - r$  instead of  $n - p$

### 9.2.3 Choice of the solution

#### Choice of $\hat{\beta}$

- $\hat{Y}$  and  $\hat{\sigma}^2$  do not depend on the choice of  $\hat{\beta}$
- $\{\hat{\beta}; \hat{Y} = X\hat{\beta}\}$  is a linear subspace of  $\mathbb{R}^p$  of dimension  $p - r$ 
  - ▷ we can choose  $\hat{\beta}$  by specifying  $p - r$  linear constraints
    - \* choose an  $(p - r) \times p$  matrix  $D$ ,  $\text{rank}(D) = p - r$
    - \* require that  $D\beta = 0$
- for a given  $D$ 
  - ▷ QR decompose  $D^T = (Q_1 | Q_2) \begin{pmatrix} R_1 \\ 0 \end{pmatrix} = Q_1 R_1$
  - ▷  $C_D = Q_2$  is a  $p \times r$  matrix,  $\text{rank}(C_D) = r$
  - ▷  $X_D = X C_D$  is an  $n \times r$  matrix,  $\text{rank}(X_D) = r$
  - ▷ fit the (full-rank) model  $Y = X_D \beta_D + \varepsilon$
  - ▷  $\hat{\beta} = C_D \hat{\beta}_D$  is the solution to the original normal equations satisfying the constraints given by  $D$

#### Common example: factor variables (fev data)

- basic model: FEV  $\sim$  Height + Gender
  - ▷ naïve parametrization

$$Y_i = \beta_0 + \beta_H \times \text{Height}_i +$$

$$+ \beta_M \times \mathbb{I}\{\text{the } i^{\text{th}} \text{ child is male}\} + \beta_F \times \mathbb{I}\{\text{the } i^{\text{th}} \text{ child is female}\} +$$

$$+ \varepsilon_i, 1 \leq i \leq 654$$

$$\begin{pmatrix} 1.708 \\ 1.724 \\ 1.720 \\ 1.558 \\ \dots \\ 3.211 \end{pmatrix} = \begin{pmatrix} 1 & 144.8 & 0 & 1 \\ 1 & 171.5 & 0 & 1 \\ 1 & 138.4 & 0 & 1 \\ 1 & 134.6 & 1 & 0 \\ \dots & \dots & \dots & \dots \\ 1 & 168.9 & 0 & 1 \end{pmatrix} \times \begin{pmatrix} \beta_0 \\ \beta_H \\ \beta_M \\ \beta_F \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \varepsilon_3 \\ \varepsilon_4 \\ \dots \\ \varepsilon_{654} \end{pmatrix}$$

▷ standard parametrization

$$Y_i = \beta_0 + \beta_H \times \text{Height}_i + \\ + \beta_M \times \mathbb{I}\{\text{the } i^{\text{th}} \text{ child is male}\} + \\ + \varepsilon_i, \quad 1 \leq i \leq 654$$

○ basic model with interaction:  $FEV \sim \text{Height} * \text{Gender}$

▷ standard parametrization

$$Y_i = \beta_0 + \beta_H \times \text{Height}_i + \\ + \beta_M \times \mathbb{I}\{\text{the } i^{\text{th}} \text{ child is male}\} + \\ + \beta_{H:M} \times \mathbb{I}\{\text{the } i^{\text{th}} \text{ child is male}\} \times \text{Height}_i + \varepsilon_i, \quad 1 \leq i \leq 654$$

## One-way ANOVA

○  $Y_{i,j} = \mu + \alpha_i + \varepsilon_{i,j}, \varepsilon_{i,j} \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$

$$i \in \{1, \dots, I\}, j \in \{1, \dots, n_i\}$$

○ matrix form  $\mathbf{Y} = \mathbf{X}(\mu, \boldsymbol{\alpha})^\top + \boldsymbol{\varepsilon}, \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$

$$\begin{pmatrix} Y_{1,1} \\ \dots \\ Y_{1,n_1} \\ Y_{2,1} \\ \dots \\ Y_{2,n_2} \\ \dots \\ \dots \\ Y_{I,1} \\ \dots \\ Y_{I,n_I} \end{pmatrix} = \begin{pmatrix} 1 & 1 & 0 & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & 1 & 0 & 0 & \dots & 0 \\ 1 & 0 & 1 & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & 0 & 1 & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & 0 & 0 & 0 & \dots & 1 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & 0 & 0 & 0 & \dots & 1 \end{pmatrix} \begin{pmatrix} \mu \\ \alpha_1 \\ \alpha_2 \\ \dots \\ \alpha_I \end{pmatrix} + \begin{pmatrix} \varepsilon_{1,1} \\ \dots \\ \varepsilon_{1,n_1} \\ \varepsilon_{2,1} \\ \dots \\ \varepsilon_{2,n_2} \\ \dots \\ \dots \\ \varepsilon_{I,1} \\ \dots \\ \varepsilon_{I,n_I} \end{pmatrix}$$

▷  $\mathbf{X}$  is an  $n \times (I + 1)$  matrix with  $\text{rank}(\mathbf{X}) = I$

## ANOVA

○ one-way ANOVA

▷  $Y_{i,j} = \mu + \alpha_i + \varepsilon_{i,j}, \varepsilon_{i,j} \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$

$$i \in \{1, \dots, I\}, j \in \{1, \dots, n_i\}$$

▷ matrix form  $\mathbf{Y} = (\mathbf{1} | \mathbf{X}_\alpha)(\mu, \boldsymbol{\alpha})^\top + \boldsymbol{\varepsilon}, \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$

\*  $\mathbf{X}$  is an  $n \times (I + 1)$  matrix with  $\text{rank}(\mathbf{X}) = I$

○ two-way ANOVA

- ▷  $Y_{i,j,k} = \mu + \alpha_i + \beta_j + \varepsilon_{i,j,k}$ ,  $\varepsilon_{i,j,k} \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$   
 $i \in \{1, \dots, I\}$ ,  $j \in \{1, \dots, J\}$ ,  $k \in \{1, \dots, n_{i,j}\}$
- ▷ matrix form  $\mathbf{Y} = (\mathbf{1} | \mathbf{X}_\alpha | \mathbf{X}_\beta) (\mu, \boldsymbol{\alpha}, \boldsymbol{\beta})^\top + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$   
 \*  $\mathbf{X}$  is an  $n \times (I + J + 1)$  matrix with  $\text{rank}(\mathbf{X}) = I + J - 1$

○ two-way ANOVA with interactions

- ▷  $Y_{i,j,k} = \mu + \alpha_i + \beta_j + \gamma_{i,j} + \varepsilon_{i,j,k}$ ,  $\varepsilon_{i,j,k} \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$   
 $i \in \{1, \dots, I\}$ ,  $j \in \{1, \dots, J\}$ ,  $k \in \{1, \dots, n_{i,j}\}$
- ▷  $\mathbf{Y} = (\mathbf{1} | \mathbf{X}_\alpha | \mathbf{X}_\beta | \mathbf{X}_\alpha \cdot \mathbf{X}_\beta) (\mu, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\gamma})^\top + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$   
 (. denotes component-wise multiplication in the  $n \times (I \times J)$  matrix)  
 \*  $\mathbf{X}$  is an  $n \times (I + J + (I \times J) + 1)$  matrix,  $\text{rank}(\mathbf{X}) = I \times J$

## ANOVA parametrizations

○ one-way ANOVA

- ▷  $Y_{i,j} = \mu + \alpha_i + \varepsilon_{i,j}$ ,  $\varepsilon_{i,j} \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$   
 $i \in \{1, \dots, I\}$ ,  $j \in \{1, \dots, n_i\}$   
 \*  $\mathbb{R}$  parametrization:  $\alpha_1 = 0$   
 \* other parametrizations: e.g.  $\sum_{i=1}^I n_i \alpha_i = 0$

○ two-way ANOVA

- ▷  $Y_{i,j,k} = \mu + \alpha_i + \beta_j + \varepsilon_{i,j,k}$ ,  $\varepsilon_{i,j,k} \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$   
 $i \in \{1, \dots, I\}$ ,  $j \in \{1, \dots, J\}$ ,  $k \in \{1, \dots, n_{i,j}\}$   
 \*  $\mathbb{R}$  parametrization:  $\alpha_1 = 0$ ,  $\beta_1 = 0$   
 \* other: e.g.  $\sum_{i=1}^I \alpha_i \sum_{j=1}^J n_{i,j} = 0$ ,  $\sum_{j=1}^J \beta_j \sum_{i=1}^I n_{i,j} = 0$

○ two-way ANOVA with interactions

- ▷  $Y_{i,j,k} = \mu + \alpha_i + \beta_j + \gamma_{i,j} + \varepsilon_{i,j,k}$ ,  $\varepsilon_{i,j,k} \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$   
 $i \in \{1, \dots, I\}$ ,  $j \in \{1, \dots, J\}$ ,  $k \in \{1, \dots, n_{i,j}\}$   
 \*  $\mathbb{R}$  parametrization:  $\alpha_1 = 0$ ,  $\beta_1 = 0$ ,  $\gamma_{1,j} = 0 \forall j$ ,  $\gamma_{i,1} = 0 \forall i$   
 \* other: e.g.  $\sum_{i=1}^I \alpha_i \sum_{j=1}^J n_{i,j} = 0$ ,  $\sum_{j=1}^J \beta_j \sum_{i=1}^I n_{i,j} = 0$ ,  $\sum_{i=1}^I n_{i,j} \gamma_{i,j} = 0 \forall j$ ,  $\sum_{j=1}^J n_{i,j} \gamma_{i,j} = 0 \forall i$

## ANOVA parametrizations via matrices of contrasts



- one-way ANOVA
  - ▷  $\mathbf{Y} = (\mathbf{1} | \mathbf{X}_\alpha) (\mu, \boldsymbol{\alpha})^\top + \boldsymbol{\varepsilon}, \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$
  - ▷ replace  $(\mathbf{1} | \mathbf{X}_\alpha)$  by  $(\mathbf{1} | \mathbf{X}_\alpha \mathbf{C}_\alpha)$ 
    - \*  $\mathbf{C}_\alpha \in \mathbb{R}^{I \times (I-1)}, \text{rank}((\mathbf{1} | \mathbf{X}_\alpha \mathbf{C}_\alpha)) = I$
  - ▷ estimate  $\boldsymbol{\alpha}$  by  $\mathbf{C}_\alpha \hat{\boldsymbol{\alpha}}$  from the fitted model
- two-way ANOVA
  - ▷  $\mathbf{Y} = (\mathbf{1} | \mathbf{X}_\alpha | \mathbf{X}_\beta) (\mu, \boldsymbol{\alpha}, \boldsymbol{\beta})^\top + \boldsymbol{\varepsilon}, \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$
  - ▷ replace  $(\mathbf{1} | \mathbf{X}_\alpha | \mathbf{X}_\beta)$  by  $(\mathbf{1} | \mathbf{X}_\alpha \mathbf{C}_\alpha | \mathbf{X}_\beta \mathbf{C}_\beta)$ 
    - \*  $\mathbf{C}_\alpha \in \mathbb{R}^{I \times (I-1)}, \mathbf{C}_\beta \in \mathbb{R}^{J \times (J-1)}$
    - \*  $\text{rank}((\mathbf{1} | \mathbf{X}_\alpha \mathbf{C}_\alpha | \mathbf{X}_\beta \mathbf{C}_\beta)) = I + J - 1$
  - ▷ estimate  $\boldsymbol{\alpha}$  and  $\boldsymbol{\beta}$  by  $\mathbf{C}_\alpha \hat{\boldsymbol{\alpha}}$  and  $\mathbf{C}_\beta \hat{\boldsymbol{\beta}}$  from the fitted model
- two-way ANOVA with interactions
  - ▷  $\mathbf{Y} = (\mathbf{1} | \mathbf{X}_\alpha | \mathbf{X}_\beta | \mathbf{X}_\alpha \cdot \mathbf{X}_\beta) (\mu, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\gamma})^\top + \boldsymbol{\varepsilon}, \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$
  - ▷ replace  $(\mathbf{1} | \mathbf{X}_\alpha | \mathbf{X}_\beta | \mathbf{X}_\alpha \cdot \mathbf{X}_\beta)$  by  $(\mathbf{1} | \mathbf{X}_\alpha \mathbf{C}_\alpha | \mathbf{X}_\beta \mathbf{C}_\beta | \mathbf{X}_\alpha \mathbf{C}_\alpha \cdot \mathbf{X}_\beta \mathbf{C}_\beta)$ 
    - \*  $\mathbf{C}_\alpha \in \mathbb{R}^{I \times (I-1)}, \mathbf{C}_\beta \in \mathbb{R}^{J \times (J-1)}$
    - \*  $\text{rank}((\mathbf{1} | \mathbf{X}_\alpha \mathbf{C}_\alpha | \mathbf{X}_\beta \mathbf{C}_\beta | \mathbf{X}_\alpha \mathbf{C}_\alpha \cdot \mathbf{X}_\beta \mathbf{C}_\beta)) = I J$
  - ▷ estimate  $\boldsymbol{\alpha}, \boldsymbol{\beta}$  and  $\boldsymbol{\gamma}$  by  $\mathbf{C}_\alpha \hat{\boldsymbol{\alpha}}, \mathbf{C}_\beta \hat{\boldsymbol{\beta}}$  and  $(\mathbf{C}_\alpha \otimes \mathbf{C}_\beta) \hat{\boldsymbol{\gamma}}$

## 9.3 Multicollinearity

### Multicollinearity

#### Multicollinearity

- we have seen that if  $\text{rank}(\mathbf{X}) = r < p$ , we do not lose anything by leaving out  $p - r$  columns
- but what if  $\text{rank}(\mathbf{X}) = p$  but “only nearly so”?
  - ▷ the columns of  $\mathbf{X}$  linearly independent BUT
  - ▷  $\frac{\langle \mathbf{x}_i, \mathbf{x}_j \rangle}{\|\mathbf{x}_i\| \|\mathbf{x}_j\|} \approx \pm 1$  for some  $(i, j)$   
and/or for some linear combinations of the columns
- we would lose information by leaving out columns but keeping them all is a problem as well

- ▷  $\mathbf{X}^\top \mathbf{X}$  is ill-conditioned
  - \*  $\hat{\boldsymbol{\beta}}$  solves  $(\mathbf{X}^\top \mathbf{X}) \boldsymbol{\beta} = \mathbf{X}^\top \mathbf{Y}$
  - \* small change in  $\mathbf{Y} \Rightarrow$  large change in  $\hat{\boldsymbol{\beta}}$
  - \* fit extremely sensitive to errors  $\boldsymbol{\varepsilon}$
- ▷ large  $\text{Var} \hat{\boldsymbol{\beta}}$ 
  - \* imprecise estimation of  $\boldsymbol{\beta}$
  - \* wide confidence intervals for  $\beta$ 's
  - \* large p-values of the t-tests  
(not necessarily of the overall F-test)

### Detecting multicollinearity

- pairwise relationships
  - ▷ graphically: plot pairs of covariates one against another
  - ▷ numerically: compute pairwise correlations
- pairwise and/or higher-order relationships
  - ▷ regressing each covariate in turn on all the others
    - \* large values of the corresponding  $R^2$  problematic
  - ▷ compute eigenvalues of  $\mathbf{X}^\top \mathbf{X}$ 
    - \* large values of  $\sqrt{\lambda_1/\lambda_j}$  problematic
- other indications
  - ▷ large  $p$ -values of the individual  $t$ -tests but a small  $p$ -value of the overall  $F$ -test
  - ▷ estimates of  $\boldsymbol{\beta}$  and  $\text{Var}(\hat{\boldsymbol{\beta}})$  very sensitive to adding/leaving out covariates and/or perturbing  $\mathbf{Y}$

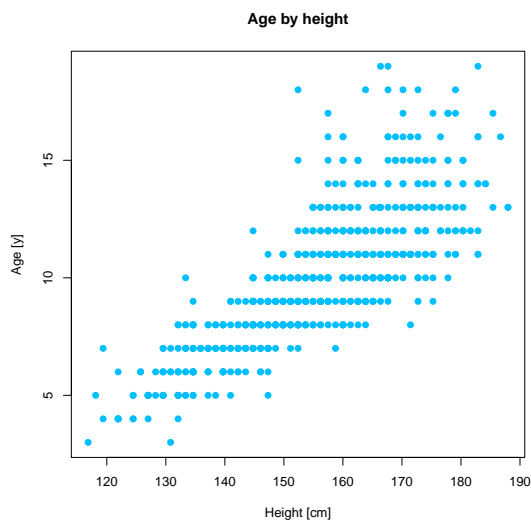
### Variance inflation factors

- fit  $\text{lm}(\mathbf{X}_{,j} \sim \mathbf{X}_{,1} + \dots + \mathbf{X}_{,j-1} + \mathbf{X}_{,j+1} + \dots + \mathbf{X}_{,p})$ 
  - ▷  $R_j^2$  ... the corresponding coefficient of determination
- it can be shown that  $\text{Var}(\hat{\beta}_j) = \frac{s^2}{(n-1)s_{X_{,j}}^2} \times \frac{1}{1-R_j^2}$   
in  $\text{lm}(\mathbf{Y} \sim \mathbf{X}_{,1} + \dots + \mathbf{X}_{,p})$
- variance inflation factor  $\text{VIF}_j = \frac{1}{1-R_j^2}$

- ▷ measures linear dependence of the  $j^{\text{th}}$  covariate on the other covariates
- ▷ interpretation
  - \* standard error of  $\hat{\beta}_j$  is  $\approx \sqrt{\text{VIF}_j} \times$  larger than it would be were the  $j^{\text{th}}$  covariate independent of the other covariates
- ▷ = 1 for orthogonal covariates, large values indicate problems
- ▷ how big is “too big”?
  - \* some consider  $\text{VIF} > 5$  problematic
  - \*  $\text{VIF} > 10$  is definitely considered problematic
- a generalization **gVIF** exists for categorical variables

**Example: fev data**

- $\text{Cor}(\text{Age}, \text{Height}) = 0.79$



- $R^2_{\text{Age}} = 0.69$
- $\text{VIF}_{\text{Age}} = 3.24$

**Ill-conditioned  $\mathbf{X}^T \mathbf{X}$**

- linear model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$
- model fitting: 
$$\underbrace{(\mathbf{X}^T \mathbf{X})}_{(p \times p)} \underbrace{\hat{\boldsymbol{\beta}}}_{(p \times 1)} = \underbrace{\mathbf{X}^T \mathbf{Y}}_{(p \times 1)} \dots \underbrace{\mathbf{A}}_{(p \times p)} \underbrace{\mathbf{x}}_{(p \times 1)} = \underbrace{\mathbf{b}}_{(p \times 1)}$$

- solving for  $\hat{\beta}$  with machine precision
  - ▷ if the error in  $\mathbf{b}$  is  $\epsilon$ , the error in the solution  $\mathbf{A}^{-1}\mathbf{b}$  is  $\mathbf{A}^{-1}\epsilon$
  - ▷ relative error in the solution divided by the relative error in  $\mathbf{b}$ :
    - \*  $\frac{\|\mathbf{A}^{-1}\epsilon\|/\|\mathbf{A}^{-1}\mathbf{b}\|}{\|\epsilon\|/\|\mathbf{b}\|}$  for some norm  $\|\bullet\|^2$
    - \* maximal value:  $\frac{\|\mathbf{A}^{-1}\|}{\|\mathbf{A}\|}$
  - ▷ for Euclidean/spectral norm:  $\frac{\|\mathbf{A}^{-1}\|}{\|\mathbf{A}\|} = \sqrt{\frac{\lambda_1}{\lambda_p}}$ :  $\sqrt{\phantom{x}}$  of the ratio of the smallest and largest eigenvalue: **condition number**
    - \* some consider  $\geq 30$  problematic
    - \* the condition number depends also on the scales of covariates (not only on their relationships)
    - \* can improve a lot if all covariates are on similar scales

### Tackling multicollinearity

- having independent covariates helps a lot but inherent relationships cannot be circumvented
- with collinear covariates, information does not increase as we would expect with the number of covariates
- “solutions”
  - ▷ excluding covariates
    - \* we avoid “repeating the same thing” but lose information
    - \* keep covariates that are of interest and/or are easy to measure
    - \* do not misinterpret leaving out a covariate as implying that it has no significant influence on the outcome
  - ▷ orthogonalizing and/or standardizing the predictors
    - \* more complicated interpretation
    - \* not a problem for prediction (but then multicollinearity might not have been a big issue unless extrapolation was planned)
  - ▷ a different method for estimation (e.g. ridge regression)
    - \* we lose some nice properties of the estimators

# Chapter 10

## Miscellanea and recap

### 10.1 The problem

#### 10.1.1 Normal linear model

##### Normal linear model

- $Y_i = \beta_0 + \beta_1 x_{i,1} + \dots + \beta_k x_{i,k} + \varepsilon_i, i \in \{1, \dots, n\}$ 
  - ▷  $Y_i$ : outcome, response, output, dependent variable
    - \* random variable, we observe a realization  $y_i$
    - \* (odezva, závisle proměnná, regresand)
  - ▷  $x_{i,1}, \dots, x_{i,k}$ : covariates, predictors, explanatory variables, input, independent variables
    - \* given, known
    - \* (nezávisle proměnné, regresory)
  - ▷  $\beta_0, \dots, \beta_k$ : coefficients
    - \* unknown
    - \* (regresní koeficienty)
  - ▷  $\varepsilon_i$ : random error
    - \* random variable, unobserved
- $\varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2), i \in \{1, \dots, n\}$ 
  - ▷  $E \varepsilon_i = 0$ : no systematic errors
  - ▷  $\text{Var } \varepsilon_i = \sigma^2$ : same precision

##### Example: bloodpress data

- o from [sites.stat.psu.edu/~lsimon/stat501wc/sp05/data/](http://sites.stat.psu.edu/~lsimon/stat501wc/sp05/data/)
- o association between the mean arterial blood pressure[mmHg] and age[years], weight[kg], body surface area[m<sup>2</sup>], duration of hypertension[years], basal pulse[beats/min], stress

o data:

	BP	Age	Weight	BSA	DoH	Pulse	Stress
	105	47	85.4	1.75	5.1	63	33
	115	49	94.2	2.10	3.8	70	14
	...	...	...	...	...	...	...
	110	48	90.5	1.88	9.0	71	99
	122	56	95.7	2.09	7.0	75	99

- o model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$

$$\begin{pmatrix} 105 \\ 115 \\ \dots \\ 110 \\ 122 \end{pmatrix} = \begin{pmatrix} 1 & 47 & 85.4 & 1.75 & 5.1 & 63 & 33 \\ 1 & 49 & 94.2 & 2.10 & 3.8 & 70 & 14 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & 48 & 90.5 & 1.88 & 9.0 & 71 & 99 \\ 1 & 56 & 95.7 & 2.09 & 7.0 & 75 & 99 \end{pmatrix} \times \begin{pmatrix} \beta_0 \\ \dots \\ \beta_6 \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \dots \\ \varepsilon_{19} \\ \varepsilon_{20} \end{pmatrix}$$

<https://ww2.amstat.org/publications/jse/v13n2/datasets.kahn.html>

### Example: fev data

- o from: <http://www.statsci.org/data/general/fev.html>
- o question: association between the FEV[l] and Smoking,

corrected for Age[years], Height[cm] and Gender

o data:

	FEV	Age	Height	Gender	Smoking
	1.708	9	144.8	Female	Non
	1.724	8	171.5	Female	Non
	1.720	7	138.4	Female	Non
	1.558	9	134.6	Male	Non
	...	...	...	...	...
	3.727	15	172.7	Male	Current
	2.853	18	152.4	Female	Non
	2.795	16	160.0	Female	Current
	3.211	15	168.9	Female	Non

- o model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$

$$\begin{pmatrix} 1.708 \\ 1.724 \\ 1.720 \\ 1.558 \\ \dots \\ 3.727 \\ 2.853 \\ 2.795 \\ 3.211 \end{pmatrix} = \begin{pmatrix} 1 & 9 & 144.8 & 0 & 0 \\ 1 & 8 & 171.5 & 0 & 0 \\ 1 & 7 & 138.4 & 0 & 0 \\ 1 & 9 & 134.6 & 1 & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 1 & 15 & 172.7 & 1 & 1 \\ 1 & 18 & 152.4 & 0 & 0 \\ 1 & 16 & 160.0 & 0 & 1 \\ 1 & 15 & 168.9 & 0 & 0 \end{pmatrix} \times \begin{pmatrix} \beta_0 \\ \dots \\ \beta_5 \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \varepsilon_3 \\ \varepsilon_4 \\ \dots \\ \varepsilon_{651} \\ \varepsilon_{652} \\ \varepsilon_{653} \\ \varepsilon_{654} \end{pmatrix}$$

## 10.1.2 Task for this chapter

### Miscellanea & recap

- **model:**  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ 
  - ▷ outcome  $\mathbf{Y}$ 
    - \* random vector, we observe a realization  $\mathbf{y}$
  - ▷ predictors  $\mathbf{x}_1, \dots, \mathbf{x}_k$ 
    - \* vector of given (known) constants
  - ▷ coefficients  $\boldsymbol{\beta}$ 
    - \* vector of unknown constants
  - ▷ error  $\boldsymbol{\varepsilon}$ 
    - \* unknown random vector, we do not observe its realization
  - ▷ assumptions:  $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ 
    - \*  $\mathbf{E} \mathbf{Y} = \mathbf{X}\boldsymbol{\beta}$ : the expected value of  $\mathbf{Y}$  is a linear function of  $\boldsymbol{\beta}$
    - \*  $\mathbf{E} \boldsymbol{\varepsilon} = \mathbf{0}$ : no systematic errors
    - \*  $\text{Var} \boldsymbol{\varepsilon} = \sigma^2 \mathbf{I}$ : independence and same precision
- **task:** miscellanea & recap

## 10.2 Linear regression in practice

### 10.2.1 Linear regression in practice

#### Statistical analysis with linear regression

1. build a mathematical model, i.e. define
  - what is known
  - what is uncertain

linear regression example:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$
2. build a probabilistic model for what is uncertain
 

linear regression example:  $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$
3. use probability calculus to draw conclusions
 

linear regression example:

- $\widehat{\boldsymbol{\beta}} \sim N(\boldsymbol{\beta}, \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1})$
  - $\frac{n-p}{\sigma^2} \widehat{\sigma}^2 \sim \chi_{n-p}^2$
  - $\widehat{\boldsymbol{\beta}} \perp \widehat{\sigma}^2$
- $\rightsquigarrow$  confidence intervals  
 & hypotheses testing

4. “translate back” to the original problem (interpret the results)

linear regression example:

- $\widehat{\boldsymbol{\beta}}, \mathbf{a}^\top \widehat{\boldsymbol{\beta}}, \mathbf{A} \widehat{\boldsymbol{\beta}},$
  - $\widehat{\mathbf{E}} \mathbf{Y}, \mathbf{a}^\top (\widehat{\mathbf{E}} \mathbf{Y}), \mathbf{A} (\widehat{\mathbf{E}} \mathbf{Y})$
- confidence intervals  
 ◦ hypotheses testing

## Usual additions to the basic analysis

1. find a suitable mathematical model

- propose a suitable functional dependence of  $\mathbf{Y}$  on  $\mathbf{X}$
- propose a suitable model for the error

linear regression example: model selection

2. build a probabilistic model for what is uncertain

linear regression example: check the normality, potentially propose a different error distribution

3. use probability calculus to draw conclusions

- might need to adjust for multiple testing, post-hoc testing, poor design, ...

4. “translate back” to the original problem (interpret the results)

linear regression example:

- explanation
- prediction

## 10.3 Notes on interpretation

### 10.3.1 Notes on the explanation

#### Explanation using linear regression

- model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$



- estimate  $\beta$  by  $\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$
- estimate  $\mathbf{a}^\top \beta$  by  $\mathbf{a}^\top \hat{\beta}$
- $(1 - \alpha) \times 100\%$  confidence interval for  $\mathbf{a}^\top \beta$ :

$$\left( \mathbf{a}^\top \hat{\beta} - t_{1-\alpha/2}(n-p) \sqrt{\widehat{\sigma}^2 \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a}}, \right. \\ \left. \mathbf{a}^\top \hat{\beta} + t_{1-\alpha/2}(n-p) \sqrt{\widehat{\sigma}^2 \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a}} \right)$$

- estimate  $\mathbf{A} \beta$  by  $\mathbf{A} \hat{\beta}$
- $(1 - \alpha) \times 100\%$  confidence bands for  $\mathbf{A} \beta$ :

$$\left\{ \mathbf{A} \beta; \frac{1}{m \widehat{\sigma}^2} (\mathbf{A} \hat{\beta} - \mathbf{A} \beta)^\top (\mathbf{A} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{A}^\top)^{-1} (\mathbf{A} \hat{\beta} - \mathbf{A} \beta) \leq F_{1-\alpha}(m, n-p) \right\}$$

## Interpretation

- “keeping the values of all the other covariates fixed, a unit increase in  $x_i$  is associated with a  $\hat{\beta}_i$  increase in  $\mathbf{E} Y$ ”
  - ▷ suitably adapted for categorical predictors and potentially interactions, and depends on the choice of the identifiability conditions
  - ▷ polynomials need a more complex interpretation
- is it meaningful to imagine that a covariate changes while all the other remain fixed?

## Be careful with

- confounding: suppose that
  - ▷ the truth is  $Y_i = \beta_0 + \beta_E E_i + \beta_C C_i + \varepsilon_i$
  - ▷ we do not know about  $C$  and use  $Y_i = \beta_0 + \beta_E E_i + \varepsilon_i$  instead
  - ▷  $C$  and  $E$  are connected, e.g.  $E_i = \gamma_0 + \gamma_C C_i + \tilde{\varepsilon}_i$
  - ▷ then if  $C$  has an effect on  $Y$ , we will (erroneously) attribute an effect on  $Y$  to  $E$
  - ▷ may be solved by multiple regression model, provided the confounders and the form of their association to the outcome are known
- causality
  - ▷ very hard to be confident about a causal relationship rather than the “association”
- both can be helped by a sound design

### 10.3.2 Notes on the prediction

#### Prediction from linear regression

- model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon} \sim \mathbf{N}(\mathbf{0}, \sigma^2 \mathbf{I})$
- what can we say about  $Y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k + \varepsilon$   
for a new  $\mathbf{x} = (1, x_1, \dots, x_k)^\top$ ?

- $Y = \mathbf{x}^\top \boldsymbol{\beta} + \varepsilon$  and  $\mathbf{E}Y = \mathbf{x}^\top \boldsymbol{\beta}$
- estimate  $\mathbf{E}Y$  and  $Y$  by  $\mathbf{x}^\top \widehat{\boldsymbol{\beta}}$
- $(1 - \alpha) \times 100\%$  confidence interval for  $\mathbf{E}Y$ :

$$\left( \mathbf{x}^\top \widehat{\boldsymbol{\beta}} - t_{1-\alpha/2}(n-p) \sqrt{\widehat{\sigma}^2 \mathbf{x}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}}, \right. \\ \left. \mathbf{x}^\top \widehat{\boldsymbol{\beta}} + t_{1-\alpha/2}(n-p) \sqrt{\widehat{\sigma}^2 \mathbf{x}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}} \right)$$

- $(1 - \alpha) \times 100\%$  confidence interval for  $Y$

$$\left( \mathbf{x}^\top \widehat{\boldsymbol{\beta}} - t_{1-\alpha/2}(n-p) \sqrt{\widehat{\sigma}^2 (1 + \mathbf{x}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x})}, \right. \\ \left. \mathbf{x}^\top \widehat{\boldsymbol{\beta}} + t_{1-\alpha/2}(n-p) \sqrt{\widehat{\sigma}^2 (1 + \mathbf{x}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x})} \right)$$

#### Be careful with

- extrapolation
  - ▷ predicting  $Y$  for  $\mathbf{x}$  that is far from the  $\mathbf{x}_i$ 's in  $\mathbf{X}$
  - ▷ predicting for different situations/populations than the one satisfying  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$
- overfitting
  - ▷ fitting a model  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  that is “too close to the data”
  - ▷ estimated  $\sigma^2$  is small
- having seen enough data

## 10.4 Transformations

### 10.4.1 Transformations

#### Transformations of variables

- model:  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$
- we have seen transformations of predictors to find a suitable functional dependence of  $Y$  on  $x$
- how about transforming  $Y$ ?
  - ▷ done in practice to improve the functional dependence or fix heteroskedasticity
  - ▷ most common are  $\log(Y)$ ,  $\sqrt{Y}$ , some use other powers of  $Y$
  - ▷ this is a fundamental change to the model
    - \* leaving the simple linear regression framework ...

#### Log-transformation of the response

- original model:  $Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$
- model after the log transform:  $\log(Y_i) = \beta_0 + \beta_1 x_i + \varepsilon_i$ 
  - ▷ on the original scale:  $Y_i = \exp\{\beta_0\} \times \exp\{\beta_1 x_i\} \times \exp\{\varepsilon_i\}$
  - ▷ the effects of covariates are on the multiplicative scale
  - ▷ the error enters multiplicatively and the multiplicative error has log-normal distribution
    - \*  $\exp\{x\} \approx 1 + x$  for small  $x$
    - \*  $\Rightarrow Y_i = \exp\{\beta_0\} \times \exp\{\beta_1 x_i\} \times (1 + \varepsilon_i)$  for small  $\varepsilon_i$
    - \* non-linear regression model with non-constant variance
    - \*  $Y_i = \exp\{\beta_0\} \times \exp\{\beta_1\} \exp\{x_i\} + \sigma_i^2 \varepsilon_i$  for small  $\varepsilon_i$  ...
  - ▷ prediction on the original scale
    - \* predict by  $\exp\{\hat{Y}\}$  with CI ( $\exp\{L\}, \exp\{U\}$ )
  - ▷ interpretation of  $\boldsymbol{\beta}$  on the log-scale
  - ▷ problems with interpretation on the original scale
    - \*  $\log(\mathbf{E}Y) \neq \mathbf{E} \log(Y)$  but the median is preserved
    - \*  $\log(1 + x) \approx x$  for small  $x$  ...
    - \* e.g.  $\hat{\beta}_1 = 0.09$  can be interpreted as a 9% increase in  $\text{med}Y$  associated with a unit increase in  $x$

## Box–Cox transformation of the response

- original model:  $Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$
- looking for a more general transform. . .
- suppose that  $Y > 0$
- **Box–Cox transformation:**
  - ▷  $g_\lambda(y) = \begin{cases} \frac{y^\lambda - 1}{\lambda} & \lambda \neq 0 \\ \log(y) & \lambda = 0 \end{cases}$  (a continuous function of  $\lambda$ )
  - ▷  $\lambda$  can be viewed as a parameter and  $\hat{\lambda}$  found by MLE
    - \* also gives a CI
  - ▷ for prediction, you may use  $y^{\hat{\lambda}}$
  - ▷ for interpretation, you had better round  $\hat{\lambda}$  to the nearest interpretable value (check the CI)
  - ▷ use CI to see if you need a transform at all

## 10.5 Concluding remarks

### 10.5.1 Reflection

#### It's an uncertain world . . . use statistics to decide

- How much of
  - ▷ chocolate and other goodies is good for our health?
  - ▷ levels of bacteria, fertilizers, chemicals, . . . is safe?
- What is the right size for
  - ▷ the height of a dam?
  - ▷ insurance premium?
  - ▷ mortgage interest?
- What is
  - ▷ the average salary?
  - ▷ public opinion on . . . ?
  - ▷ results in upcoming elections?

- uncertainty at the beginning --> imperfect answers at the end
- statistics is used for quantifying uncertainty,  
not for getting rid of it

### **Statistics is collaboration**

- *The best thing about being a statistician is that you get to play at everyone's backyard.*

John Tukey

### **Statistics does not guarantee the right answers**

- if there is no uncertainty, there is no need for statistics
- ↔ statistics might give a wrong answer
- !!!but we should not abuse this!!!
- only incompetent statisticians do not know how to lie with statistics
  - good statisticians know the pitfalls and know they must be cautious

### **Ingredients of a statistical analysis**

- mathematics, programming, communication . . .
- but above all: **COMMON SENSE**