

Course Bio144: Data Analysis in Biology

Stefanie Muff (Lectures) & Owen L. Petchey (Exercises)

Lecture 1: Introduction and Outlook
23./24. February 2017

Organization

All important details, such as testate conditions, exam dates etc. are provided on the OpenEdX course page:

[https://openedx.mnf.uzh.ch/courses/course-v1:
UZH+ BIO144+ FS2017/about](https://openedx.mnf.uzh.ch/courses/course-v1:UZH+ BIO144+ FS2017/about)

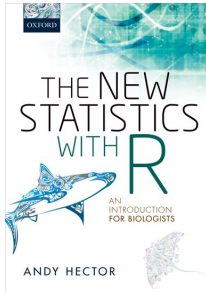
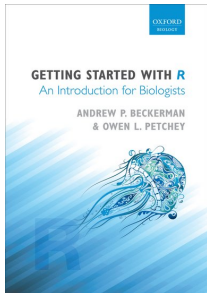
Literature

Compulsory literature:

1. *Lineare Regression* by W. Stahel (pdf on course webpage)
2. *Getting Started with R, An introduction for biologists* (2017, **Second Edition**) Beckerman, Childs & Petchey (DO NOT USE THE FIRST EDITION!).

A UZH library link to the second edition will be added asap, probably mid March.

3. *The New Statistics With R* by A. Hector, Oxford University Press; ISBN 978-0-19-872906-8



Complementary literature:

- *Statistics – An Introduction Using R* by M.J. Crawley (similar to 3.) above)
- *The Analysis of Biological Data* by M.C. Whitlock and D. Schluter
- *Regression - Modelle, Methoden und Anwendungen* by Fahrmeier, Kneib, Lang
- *The Essential Guide to Effect Sizes. Statistical Power, Meta-Analysis, and the Interpretation of Research Results* (2010, First Edition) Ellis.
Ebook via [▶ UZH library](#) .

Overarching goals of the course

- Provide a solid foundation for answering biological questions with quantitative data.
- Help students to understand the language of a statistician.
- Ability to understand and interpret results in research articles.
- Give the students a challenging, engaging, and enjoyable learning experience.

My belief: A solid foundation in statistics makes you independent!

Prerequisite for Bio144

- Mat183 “Stochastik für die Naturwissenschaften” (2nd semester)

Schedule (12 lecture weeks + 2 self-study weeks)

Week 1 Introduction and outlook

Week 2 Simple linear regression

Week 3 Residual analysis, model validation

Week 4 Multiple linear regression

Week 5 ANOVA

Week 6 ANCOVA Matrix Algebra

Week 7 Model selection

Self-study week

Week 8 Interpretation of results, causality

Week 9 Count data (Poisson regression)

Week 10 Binary Data (logistic regression)

Week 11 Measurement error, random effects

Self-study week

Week 12 Selected topics, repetition and outlook

Why is statistical data analysis so relevant for the biological and medical sciences?

What do you think?

Why is statistical data analysis so relevant for the biological and medical sciences?

What do you think?

Awareness that, without a profound knowledge in statistical data analysis, it will be hard to analyze your data from Bachelor, Master or PhD theses....

Examples:

- Medicine: Does a drug have a positive effect? Which factors cause cancer?
- Ecology: What is a suitable habitat for a certain animal? Which resources does it need or prefer?
- Evolutionary biology: Do highly inbred animals have decreased chances to survive or reproduce?

Be careful! "Learning by doing" is not advisable in statistics. Experience is essential, there are many pitfalls.

A good foundation in statistics makes you more independent from consultants or the goodwill of colleagues. Without such a knowledge, you will always need help from others.

Data analysis is itself an exciting part of research!

Data analysis is at the interface between mathematics and biology (and other research fields such as medicine, earth sciences, and so on).

What are the purposes of data analysis?

- To find and quantify associations through graphical representations and modelling.
- To draw conclusion from data.
- To quantify the uncertainty of these conclusions.

Own examples

Otter (*lutra lutra*) (Weinberger et al., 2016)

Research questions: What is the preferred habitat by otters? How do otters adapt to human altered landscapes?

Method: Study in Austria, 9 Otter were radio-tracked and monitored during 2-3 years.

Biological Conservation 199 (2016) 88–95



Contents lists available at [ScienceDirect](https://www.sciencedirect.com)

Biological Conservation

journal homepage: www.elsevier.com/locate/bioc



Flexible habitat selection paves the way for a recovery of otter populations in the European Alps



Irene C. Weinberger ^{a,*}, Stefanie Muff ^{a,b}, Addy de Jongh ^c, Andreas Kranz ^d, Fabio Bontadina ^{e,f}

^a Institute of Ecology and Evolutionary Biology, University of Zurich, Winterthurerstr. 190, 8057 Zurich, Switzerland

^b Epidemiology, Biostatistics and Prevention Institute, University of Zurich, Hirschengraben 84, 8001 Zurich, Switzerland

^c Dutch Otterstation Foundation, Spanjaardslaan 136, 8917 AX Leeuwarden, Netherlands

^d alka-kranz Ingenieurbüro für Wildökologie und Naturschutz, Am Waldgrund 25, 8044 Graz, Austria

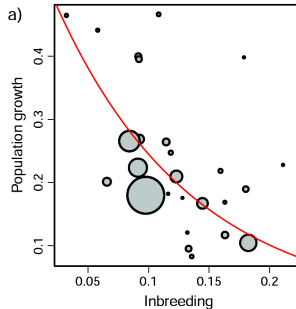
^e SWILD – Urban Ecology & Wildlife Research, Wuhstr. 12, 8003 Zurich, Switzerland

^f Swiss Federal Research Institute WSL, Biodiversity and Conservation Biology, 8903 Birmensdorf, Switzerland

Inbreeding in Alpine ibex

Research question: Does inbreeding in Alpine ibex populations have a negative effect on long-term population growth? Inbreeding depression!

Methods: Genetic information from blood samples allow to quantify the level of inbreeding in each ibex population. In addition, long-term monitoring of population sizes and harvest rates.



Wohnzone im Wallis von Quecksilber vergiftet

Vor über vierzig Jahren hatten 3,1 Tonnen Quecksilber einen Abflusskanal nahe der Walliser Gemeinde Visp verschmutzt. Noch heute müssen die Einwohner mit den Folgen leben.



Artikel zum Thema

Konvention gegen Quecksilber verabschiedet

Ein neues internationales Abkommen schränkt die Verwendung von Quecksilber in der Industrie ein. Massgeblich daran beteiligt war die Schweiz. [Mehr...](#)

19.01.2013

Research question: Is the Hg level in the environment (soil) of people's homes associated to the Hg levels in their bodies (urin, hair)?

Method: Measurements of Hg concentrations on people's properties, as well as measurements and survey of children and their mothers living in these properties.

Highly delicate, emotionally charged, political question!

► Schweiz Aktuell, 20. Juni 2016

Physical activity in children

Research question: Which factors influence physical activity patterns in children aged 2-6 years?

Method: The children had to wear accelerometers for several days. In addition, their parents had to fill in a detailed questionnaire.

Observed variables were, e.g., media consumption, behavior of the parents, age, weight, social structure,...

► [Link to Splashy study](#)

Statistics in the news (April 2016)

MEZ vom Sonntag 3. April 2016

Wissen

Überschätzte Statistiken

Daten-Analysen entscheiden heute darüber, ob ein Medikament als wirksam gilt. Bloss verstehen viele Forscher die Bedeutung dieser Berechnungen gar nicht. **Von Patrick Imhazy**

Karrieren machen können nicht nur Menschen, sondern auch statistische Größen. Das gilt besonders für den sogenannten p-Wert, mit dem jeder Mediziner und jede Statistikerin fortwährend konfrontiert werden muss. Der p-Wert ist ein statistisches Mass für die Wahrscheinlichkeit, dass ein bestimmtes Ergebnis durch Zufall zustande gekommen ist. In der Praxis oftmals zu einem strengen Schlussanlass.

Angibt die statistische Analyse von Daten einen p-Wert < 0,05 (5 Prozent) oder noch kleiner < 0,01 (1 Prozent), gelten diese als hochbedeutend angesehen. Das entscheidet über darüber, ob ein neues Medikament als wirksam eingestuft wird oder ob ein Forscher seine Studie in einem anderen Fachgebiet publizieren kann. Der p-Wert wird aber nie dazu gedacht, wissenschaftlichen Denkanlass zu sein, sondern ist, wie auch bei anderen statistischen Verfahren, ein Werkzeug, das die Ergebnisse der Forschung in Form von Zahlen darstellt.

Wissenschaftler verwenden den p-Wert immer häufiger, ohne zu verstehen, was er ihnen bedeutet. Das fördert schlechte Forschung und ungenutzte die Glaubwürdigkeit der Wissenschaft. Der Mediziner und Epidemiologe John Ioannidis von der Universität Oxford sprach in einem Kommentar von «abgelenkungs» oder «falscher Gewissheit» die p-Werte zu oft missbraucht werden und erfolgt so automatisiert, dass manche nicht

werden danach, vor allem wenn sie mit Forschern und Praktikern besetzt werden. Abgesehen der Missstände, die sich aus der Verallgemeinerung, zum ersten Mal in ihrer 185-jährigen Geschichte, langjährig zu veröffentlichen, wie man mit einer statistischen Größe verfahren sollte.

Wider die Null-Hypothese

Der p-Wert sagt nicht das aus, was man gewöhnlich von ihm erwartet, erklärt der Berner Epidemiologe Peter Hothorn, der seit Jahren an der Angewandten Statistik der Universität Tübingen ist. Das bedeutet: Der p-Wert misst nicht die Wahrscheinlichkeit, ob eine bestimmte Hypothese zutrifft, und auch nicht, ob ein bestimmtes Resultat zufällig zustande gekommen ist, wie die ASA freudlich. Vielmehr misst der p-Wert, das sogenannte Null-Hypothese zu testen und möglichst zu verworfen. Die Hypothese in dieser Permutationstest können zum Beispiel lauten, dass ein Medikament A gegen eine Krankheit besser wirkt als das Medikament B. Die Null-Hypothese ist dann genau das Gegenteil davon, nämlich dass das Medikament A nicht besser wirkt als das Medikament B. Wenn Test bedeutet der Forscher in Prinzip, wie gross die Wahrscheinlichkeit für das Auftreten eines statistisch festgestellten oder nach eigenen Vermessungen das beiden verschiedenen Medikamenten ist, sprich der Annahme, dass das Null-Hypothese stimmt. Diese Wahrscheinlichkeit ist der p-Wert, und je geringer sie ist, desto weniger spricht für die Null-Hypothese. Ein p-Wert von 0,05 bedeutet, dass das festgestellte Resultat ein nach zufälliger Verteilung ist unter den Bedingungen der Null-Hypothese mit einer Wahrscheinlichkeit von lediglich 5 Prozent zustande kommen könnte – und nicht, dass eine bedingte Hypothese mit einer Sicherheit von 95 Prozent wahr ist.

Um die eigentlich interessanten Hypothesen kann der p-Wert nur indirekt etwas aussagen, weil er über zwei Ebenen geschaltet ist. Deronore-Wert liefert das keine statistische Beweis für einen positiven Unterschied oder Zusammenhang. «Der p-Wert ist eine bedingte und nicht eine absolute Wahrscheinlichkeit», erklärt Peter Hothorn. «Auch genau das beschreibt viele Forscher nicht, und es interessiert sie auch nicht.»

Hothorn betont, dass die Signifikanzniveau von 5 Prozent bzw. 1 Prozent fälschlicherweise und insbesondere für defiziente New-Statistiker, Ronald Fisher, der Erfinder des p-Werts, überliefert worden ist und fälschlich die Interpretation des statistischen Prozentsatzes, ab welcher Grösse ein p-Wert in einer Untersuchung ausreicht, haben soll. «Trotzdem haben sich die willkürlich gewählten Signifikanzniveaus in Gefolge von Generationen von Forschern gehalten», sagt Leonhard Held vom Institut für Epidemiologie, Biostatistik und Prävention der Universität Zürich. Die britischen Statistiker Jonathan Sterne und George Davey Smith haben sich vor 13 Jahren im «British Medical Journal» dazu aufgerufen, die Reuse Rate von statistischen Studien nicht mehr als ein Indikator für die Qualität einer Studie zu betrachten, sondern sie nur als ein Indikator für die Wahrscheinlichkeit zu betrachten, dass eine Studie in der Vergangenheit in irgendeiner Weise, die der Qualität nicht entspricht, durchgeführt wurde. Der Test von John Ioannidis ist einerseits ein Indikator für die Qualität einer Studie, andererseits ein Indikator für die Qualität einer Studie, die in der Vergangenheit in irgendeiner Weise, die der Qualität nicht entspricht, durchgeführt wurde. Der Test von John Ioannidis ist einerseits ein Indikator für die Qualität einer Studie, andererseits ein Indikator für die Qualität einer Studie, die in der Vergangenheit in irgendeiner Weise, die der Qualität nicht entspricht, durchgeführt wurde.

5%

Wider die Null-Hypothese

Der p-Wert sagt nicht das aus, was man gewöhnlich von ihm erwartet, erklärt der Berner Epidemiologe Peter Hothorn, der seit Jahren an der Angewandten Statistik der Universität Tübingen ist. Das bedeutet: Der p-Wert misst nicht die Wahrscheinlichkeit, ob eine bestimmte Hypothese zutrifft, und auch nicht, ob ein bestimmtes Resultat zufällig zustande gekommen ist, wie die ASA freudlich. Vielmehr misst der p-Wert, das sogenannte Null-Hypothese zu testen und möglichst zu verworfen. Die Hypothese in dieser Permutationstest können zum Beispiel lauten, dass ein Medikament A gegen eine Krankheit besser wirkt als das Medikament B. Die Null-Hypothese ist dann genau das Gegenteil davon, nämlich dass das Medikament A nicht besser wirkt als das Medikament B. Wenn Test bedeutet der Forscher in Prinzip, wie gross die Wahrscheinlichkeit für das Auftreten eines statistisch festgestellten oder nach eigenen Vermessungen das beiden verschiedenen Medikamenten ist, sprich der Annahme, dass das Null-Hypothese stimmt. Diese Wahrscheinlichkeit ist der p-Wert, und je geringer sie ist, desto weniger spricht für die Null-Hypothese. Ein p-Wert von 0,05 bedeutet, dass das festgestellte Resultat ein nach zufälliger Verteilung ist unter den Bedingungen der Null-Hypothese mit einer Wahrscheinlichkeit von lediglich 5 Prozent zustande kommen könnte – und nicht, dass eine bedingte Hypothese mit einer Sicherheit von 95 Prozent wahr ist.

Um die eigentlich interessanten Hypothesen kann der p-Wert nur indirekt etwas aussagen, weil er über zwei Ebenen geschaltet ist. Deronore-Wert liefert das keine statistische Beweis für einen positiven Unterschied oder Zusammenhang. «Der p-Wert ist eine bedingte und nicht eine absolute Wahrscheinlichkeit», erklärt Peter Hothorn. «Auch genau das beschreibt viele Forscher nicht, und es interessiert sie auch nicht.»

Hothorn betont, dass die Signifikanzniveau von 5 Prozent bzw. 1 Prozent fälschlicherweise und insbesondere für defiziente New-Statistiker, Ronald Fisher, der Erfinder des p-Werts, überliefert worden ist und fälschlich die Interpretation des statistischen Prozentsatzes, ab welcher Grösse ein p-Wert in einer Untersuchung ausreicht, haben soll. «Trotzdem haben sich die willkürlich gewählten Signifikanzniveaus in Gefolge von Generationen von Forschern gehalten», sagt Leonhard Held vom Institut für Epidemiologie, Biostatistik und Prävention der Universität Zürich. Die britischen Statistiker Jonathan Sterne und George Davey Smith haben sich vor 13 Jahren im «British Medical Journal» dazu aufgerufen, die Reuse Rate von statistischen Studien nicht mehr als ein Indikator für die Qualität einer Studie zu betrachten, sondern sie nur als ein Indikator für die Qualität einer Studie, die in der Vergangenheit in irgendeiner Weise, die der Qualität nicht entspricht, durchgeführt wurde. Der Test von John Ioannidis ist einerseits ein Indikator für die Qualität einer Studie, andererseits ein Indikator für die Qualität einer Studie, die in der Vergangenheit in irgendeiner Weise, die der Qualität nicht entspricht, durchgeführt wurde.

<https://www.theguardian.com/science/2016/jul/17/politicians-dodgy-statistics-tricks-guide?&tc=eml>

Data example 1: Prognostic factors for body fat

(From Theo Gasser & Burkhardt Seifert *Grundbegriffe der Biostatistik*)

Body fat is an important indicator for overweight, but difficult to measure.

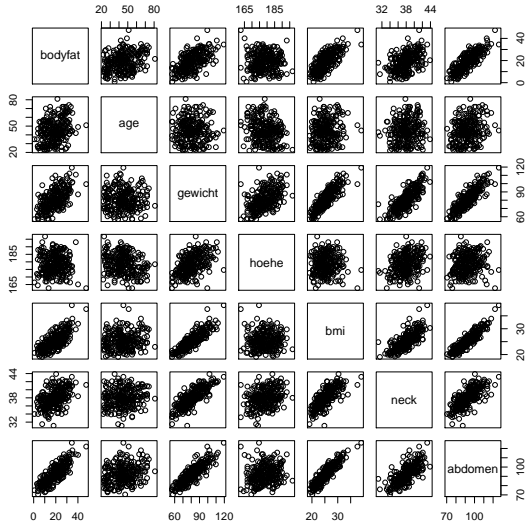
Question: Which factors allow for precise estimation (prediction) of body fat?

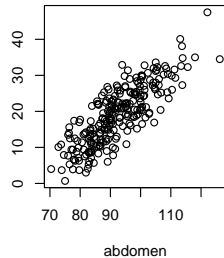
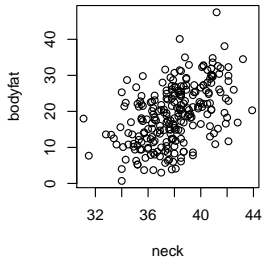
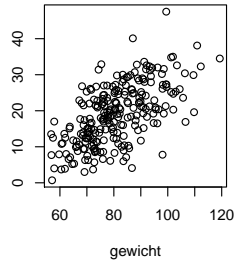
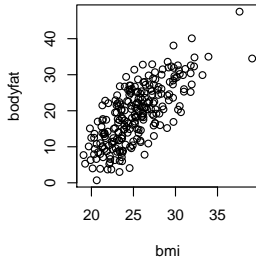
Study with 241 male participants. Measured variable were, among others, body fat (%), age, weight, body size, BMI, neck thickness and abdominal girth.

```
> str(d.bodyfat)
```

```
'data.frame':      243 obs. of  7 variables:
 $ bodyfat: num  12.3 6.1 25.3 10.4 28.7 20.9 19.2 12.4 4.1 11.7 ...
 $ age    : int  23 22 22 26 24 24 26 25 25 23 ...
 $ gewicht: num  70 78.7 69.9 83.9 83.7 ...
 $ hoehe  : num  172 184 168 184 181 ...
 $ bmi    : num  23.6 23.4 24.7 24.9 25.5 ...
 $ neck   : num  36.2 38.5 34 37.4 34.4 39 36.4 37.8 38.1 42.1 ...
 $ abdomen: num  85.2 83 87.9 86.4 100 94.4 90.7 88.5 82.5 88.6 ...
```

```
> pairs(d.bodyfat)
```





We are looking for a *model* that **predicts** body fat as precisely as possible from variables that are easy to measure.

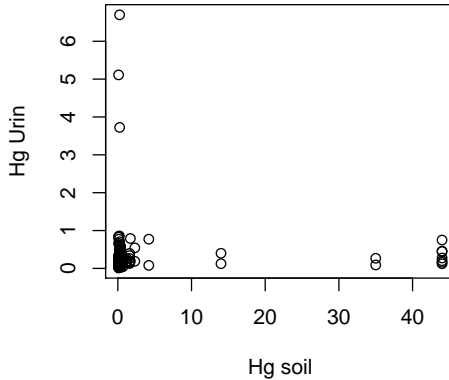
Data example 2: Mercury (Hg) in Valais (Switzerland)

Question: Association between Hg concentrations in the soil and in the urin of the people living in the respective properties. We use a slightly modified data set here.

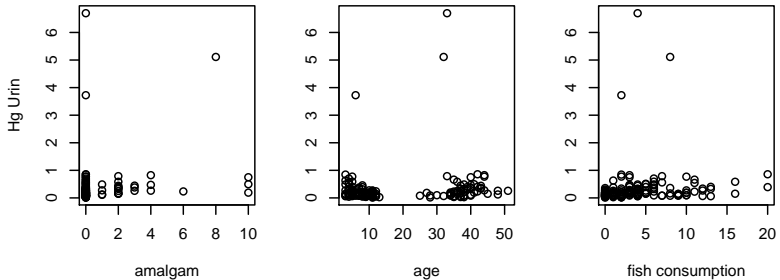
```
> str(d.hg)
```

```
'data.frame':      156 obs. of  10 variables:
 $ Hg_urin      : num  0.258 0.036 0.16 0.314 0.29 ...
 $ Hg_soil      : num  0.49 0.42 0.18 0.49 0.24 0.2 0.1 14 0.1 0.3 ...
 $ veg_garden   : int   1 1 1 1 1 1 1 1 1 1 ...
 $ migration    : int   0 0 0 0 0 0 0 0 0 0 ...
 $ smoking      : int   0 0 0 0 0 0 0 0 0 0 ...
 $ amalgam      : int   3 0 2 0 0 0 0 1 0 0 ...
 $ age          : int  51 11 34 8 6 40 7 48 11 38 ...
 $ fish         : int   3 2 5 4 4 2 2 4 0 7 ...
 $ last_time_fish: int   0 0 0 0 0 0 0 0 0 0 ...
 $ mother       : Factor w/ 2 levels "0","1": 2 1 2 1 1 2 1 2 1 2 ...
```

A first visual inspection is not very informative. There is no association visible by eye:

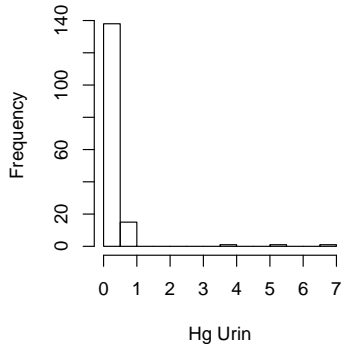
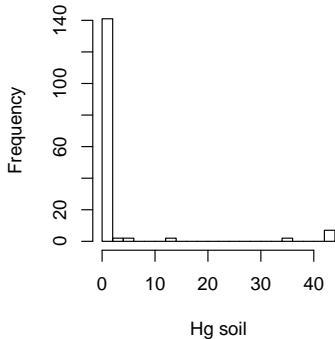


Which other factors might be responsible for high Hg concentrations in urin?



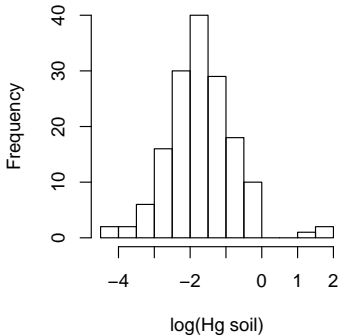
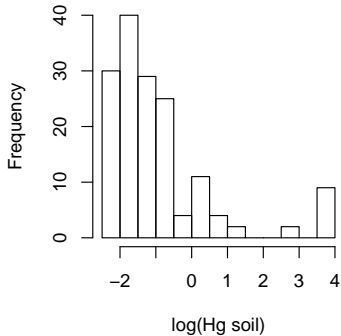
From these plots it is hard to tell which factors exactly influence the Hg pollution in humans.

It is always useful to look at the distribution of the variables in the model.
Let us plot the histogram of Hg concentrations:



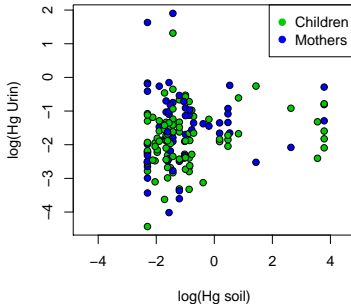
All Hg values seem to “stick” at 0.

In such cases it can help to *log-transform* the respective variables.



The scatterplot does also look much more reasonable with log-transformed values:

```
> plot(log(Hg_urin) ~ log(Hg_soil), data=d.hg, xlab="log(Hg soil)",  
+       ylab = "log(Hg Urin)",pch=21,bg=as.numeric(mother)+2,xlim=c(-4.5,4.5))  
> legend("topright",legend=c("Children","Mothers"),col=c(3,4),pch=21,pt.bg=c(3,4))
```



Remember: The idea to log-transform the variables was mainly obvious thanks to **visual inspection**!

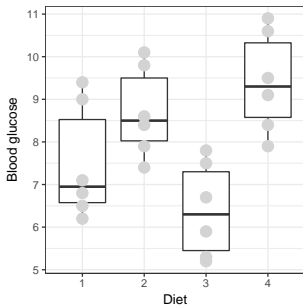
Data example 3: Diet and blood glucose level

(Elpelt and Hartung, 1987, p. 190)

24 persons were split into 4 groups. Each group followed another diet (DIAET). The blood glucose concentrations were measured at the beginning and at the end (after 2 weeks). The difference of these values was stored (BLUTZUCK).

Question: Are there differences among the groups with respect to changes in blood glucose concentrations?

Let's look at the raw data (points and boxplots):



Does this question seem familiar to you?

Hint: what would you do for two groups?

For more than 2 groups we need the *ANOVA* (=ANalysis Of VAriance) approach (see chapter 10.1 in the Mat183 script).

We will see in lecture 5 that there are in fact differences between the diets.

The next question then is: which diets are *pairwise* different.

Data example 4: Blood-screening

(From Hothorn and Everitt, 2014, Chapter 7.1)

Is a high ESR (erythrocyte sedimentation rate) an indicator for certain diseases (rheumatic disease, chronic inflammations)?

Specifically: Is there an association between ESR level $ESR < 20 \text{ mm/hr}$ and the concentrations of the plasma proteins Fibrinogen and Globulin?

Load data from the package that comes with Hothorn and Everitt (2014):

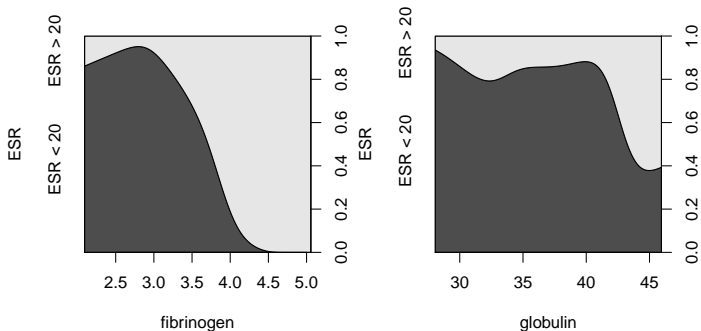
```
> library(HSAUR3)
> data("plasma", package="HSAUR3")
> plasma[c(1,5,9,10,15,29),]
```

	fibrinogen	globulin	ESR
1	2.52	38	ESR < 20
5	3.41	37	ESR < 20
9	3.15	39	ESR < 20
10	2.60	41	ESR < 20
19	2.60	38	ESR < 20
15	2.38	37	ESR > 20

The distinction $ESR < 20\text{mm/hr}$ vs. $ESR \geq 20\text{mm/hr}$ leads to a **binary** response variable.

The relation between the plasmaprotein levels and the binary indicator can be captured by a *conditional density plot*.

```
> par(mfrow=c(1,2))  
> cdplot(ESR ~ fibrinogen, plasma)  
> cdplot(ESR ~ globulin, plasma)
```



What is a model?

A model is an approximation of the reality. Understanding how the real world works is usually only possible thanks to simplifying assumptions. This is exactly the purpose of statistical data analysis.

In 2014, David Hand wrote:

In general, when building statistical models, we must not forget that the aim is to understand something about the real world. Or predict, choose an action, make a decision, summarize evidence, and so on, but always about the real world, not an abstract mathematical world: our models are not the reality – a point well made by George Box in his oft-cited remark that “all models are wrong, but some are useful” (Box, 1979).

Steps in a modelling process

- 1 Formulate a precise question
- 2 Plan your inquiry and the analysis of your data, collect the data (experiments or surveys).
- 3 Tidy and clean the data
- 4 Graphical representation of the data
- 5 Choose an appropriate *model*
- 6 Estimate model parameters and uncertainties
- 7 Check modelling assumptions
- 8 If needed, improve the model; back to step 7
- 9 Interpret your results and compare to step 1
- 10 Communicate results precisely and carefully (publication, articles..)

The scopes of statistical data analysis

- a) **Prediction, interpolation.** Example body fat: use substitute measurements to predict body fat of a person.
- b) **Estimation of parameters.**
- c) **Explanation; determination of important variables.** Example physical activity of children: The study aims to find factors that (positively or negatively) influence the movement behavior of children.
- d) Optimization.
- e) Calibration.

In this course we are concerned with a)-c).

Goals of the course (part 2)

By the end of the course you will be able analyze all data examples introduced here (and of course many more), as well as to draw conclusions from them.

Graphical representation of data

You should remember the following options for graphical data descriptions. Several of them appeared already in previous examples.

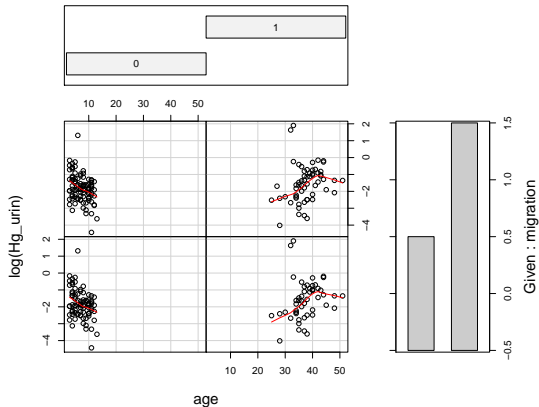
Representation	Useful for
Scatterplots	Pairwise dependency of continuous variables.
Histograms	Distribution of numerical variables.
Boxplots	Distribution of numerical variables, ev. conditionally on categories.
Conditional density plots	Dependency of a binary variable from a continuous variable.
Coplots	Dependencies among multiple variables.

Coplots

Ideal to graphically display dependencies when more than two variables are involved. Very useful for categorical variables. Example: Mercury in Valais.

```
> coplot(log(Hg_urin) ~ age | mother * migration ,d.hg,panel=panel.smooth)
```

Given : mother



There are many “fancy” ways to graphically display data (**nice-to-know**):

- 3D-plots
- Spatial representations (using geodata)
- Interactive graphs and animations

Many R packages are available for various purposes. Interactive apps can, for example, be generated with Shiny (see census app).

Next week: Simple linear regression

It will be partially a repetition of what you heard in Mat183, chapter 10.2.

References:

- Box, G. E. P. (1979). Robustness in the strategy of scientific model building. In R. L. Launer and G. N. Wilkinson (Eds.), *In Robustness in Statistics*, pp. 201–236. New York: Academic Press.
- Elpelt, B. and J. Hartung (1987). *Grundkurs Statistik, Lehr- und Übungsbuch der angewandten Statistik*.
- Hothorn, T. and B. S. Everitt (2014). *A Handbook of Statistical Analyses Using R* (3 ed.). Boca Raton: Chapman & Hall/CRC Press.
- Weinberger, I. C., S. Muff, A. Kranz, and F. Bontadina (2016). Flexible habitat selection paves the way for a recovery of otter populations in the European Alps. *Biological Conservation* 199, 88–95.