

# Course Bio144:

## Data analysis in biology

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

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

# Organization

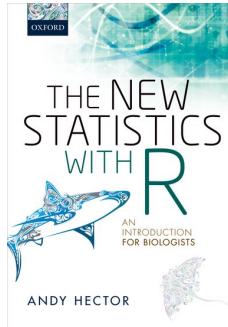
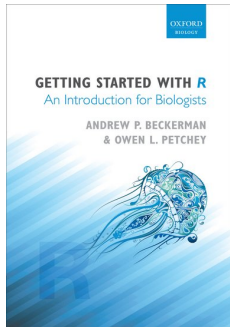
To do:

- Testate conditions
- Examination date
- Link to course webpage

# Literature

Compulsory literature:

1. *Lineare Regression* by W. Stahel (pdf on course webpage)
2. *Getting Started With R* by A.P. Beckerman und O.L. Petchey, Oxford University Press; ISBN 978-0-19-960162-2
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

# 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 (14 weeks, 12 lectures)

1. Introduction and outlook
2. Simple linear regression
3. Residual analysis, model validation
4. Multiple linear regression
5. ANOVA and ANCOVA
6. Matrix Algebra
7. Model choice
8. Self-study week
9. Interpretation of the results
10. Binary Data (logistic regression)
11. Count data (Poisson regression)
12. Measurement error, random effects
13. Self-study week
14. Selected topics, repetition and outlook

**Short-dated changes possible!**

# 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](http://www.elsevier.com/locate/bioc)



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



Irene C. Weinberger <sup>a,\*</sup>, Stefanie Muff <sup>a,b</sup>, Addy de Jongh <sup>c</sup>, Andreas Kranz <sup>d</sup>, Fabio Bontadina <sup>e,f</sup>

<sup>a</sup> Institute of Ecology and Evolutionary Biology, University of Zurich, Winterthurerstr. 190, 8057 Zurich, Switzerland

<sup>b</sup> Epidemiology, Biostatistics and Prevention Institute, University of Zurich, Hirschengraben 84, 8001 Zurich, Switzerland

<sup>c</sup> Dutch Otterstation Foundation, Spanjaardslaan 136, 8917 AX Leeuwarden, Netherlands

<sup>d</sup> alka-kranz Ingenieurbüro für Wildökologie und Naturschutz, Am Waldgrund 25, 8044 Graz, Austria

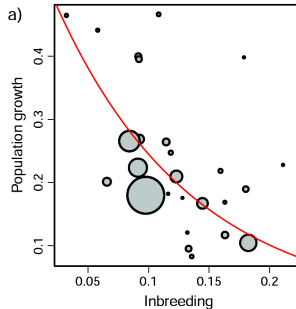
<sup>e</sup> SWILD – Urban Ecology & Wildlife Research, Wuhstr. 12, 8003 Zurich, Switzerland

<sup>f</sup> 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-ten Evening 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**



5%

Wahrscheinlich ist das ein Zufall, dass die Daten aus einer Studie ein wenig anders aussehen, als sie in der Praxis tatsächlich sein werden.

**K**arrieren machen können nicht nur Menschen, sondern auch statistische Größen. Das gilt besonders für den sogenannten p-Wert, mit dem jeder Mensch (oder auch jede Studie) in Formidabelkeit vor sich hat. Denn wenn die p-Wert-Werte zu klein sind, ist das ein Zeichen dafür, dass die Daten aus einer Studie ein wenig anders aussehen, als sie in der Praxis tatsächlich sein werden.

Die statistische Analyse von Daten ist ein p-Wert (0,05) (5 Prozent) oder noch besser (0,01) (1 Prozent), geben diese als Maß für die Wahrscheinlichkeit an, dass die Daten aus einer Studie ein wenig anders aussehen, als sie in der Praxis tatsächlich sein werden.

Die statistische Analyse von Daten ist ein p-Wert (0,05) (5 Prozent) oder noch besser (0,01) (1 Prozent), geben diese als Maß für die Wahrscheinlichkeit an, dass die Daten aus einer Studie ein wenig anders aussehen, als sie in der Praxis tatsächlich sein werden.

Die statistische Analyse von Daten ist ein p-Wert (0,05) (5 Prozent) oder noch besser (0,01) (1 Prozent), geben diese als Maß für die Wahrscheinlichkeit an, dass die Daten aus einer Studie ein wenig anders aussehen, als sie in der Praxis tatsächlich sein werden.

Das Problem ist dabei nicht nur, dass die p-Werte ein wenig anders aussehen, sondern auch, dass die Daten aus einer Studie ein wenig anders aussehen, als sie in der Praxis tatsächlich sein werden.

Die statistische Analyse von Daten ist ein p-Wert (0,05) (5 Prozent) oder noch besser (0,01) (1 Prozent), geben diese als Maß für die Wahrscheinlichkeit an, dass die Daten aus einer Studie ein wenig anders aussehen, als sie in der Praxis tatsächlich sein werden.

Die statistische Analyse von Daten ist ein p-Wert (0,05) (5 Prozent) oder noch besser (0,01) (1 Prozent), geben diese als Maß für die Wahrscheinlichkeit an, dass die Daten aus einer Studie ein wenig anders aussehen, als sie in der Praxis tatsächlich sein werden.

Die statistische Analyse von Daten ist ein p-Wert (0,05) (5 Prozent) oder noch besser (0,01) (1 Prozent), geben diese als Maß für die Wahrscheinlichkeit an, dass die Daten aus einer Studie ein wenig anders aussehen, als sie in der Praxis tatsächlich sein werden.

<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.

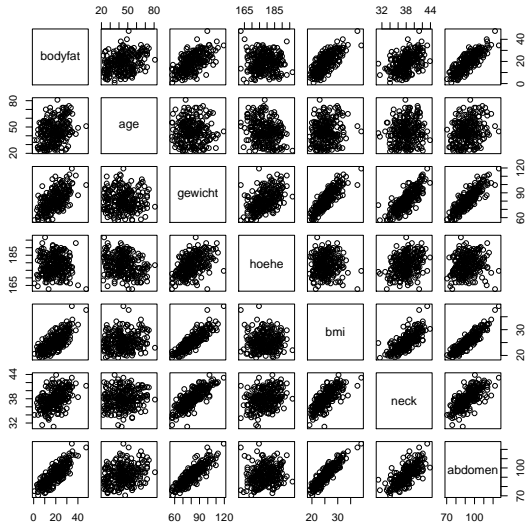
**Question:** Which factors allow for precise estimation (prediction) of body fat?

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

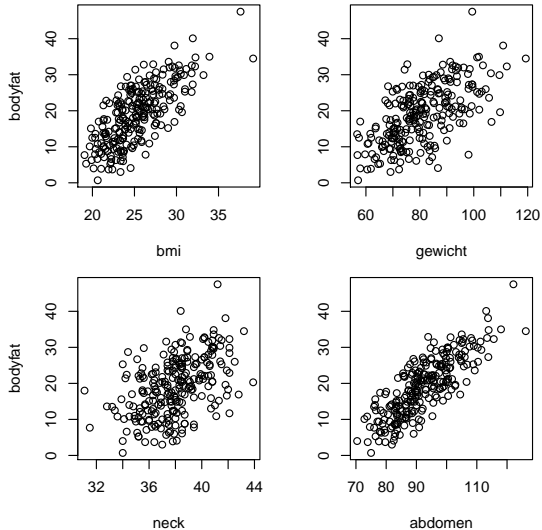
```
> str(d.bodyfat)
```

```
'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)
```



`pairs()` gives us the scatterplots of all against all variables.



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 in Valais (Switzerland)

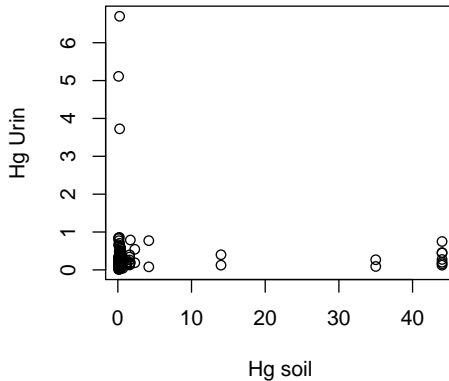
**Question:** Association between Hg concentrations in the soil and in the urin of the people living in the respective properties. We use a slightly modified data set here.

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

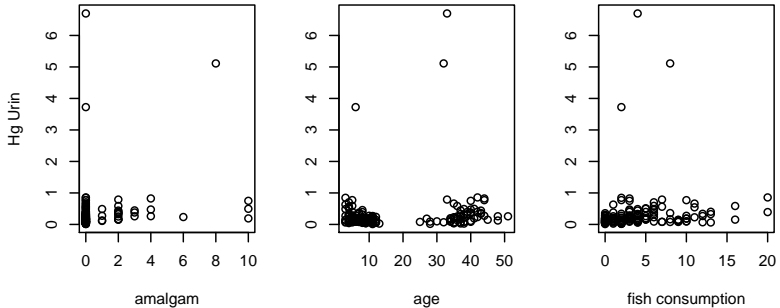
```
'data.frame':      156 obs. of  10 variables:
 $ Hg_urin      : num  0.258 0.036 0.16 0.314 0.29 ...
 $ Hg_soil      : num  0.49 0.42 0.18 0.49 0.24 0.2 0.1 14 0.1 0.3 ...
 $ veg_garden   : int  1 1 1 1 1 1 1 1 1 1 ...
 $ migration    : int  0 0 0 0 0 0 0 0 0 0 ...
 $ smoking      : int  0 0 0 0 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:

```
> plot(Hg_urin ~ Hg_soil, data=d.hg, xlab="Hg soil", ylab = "Hg Urin")
```



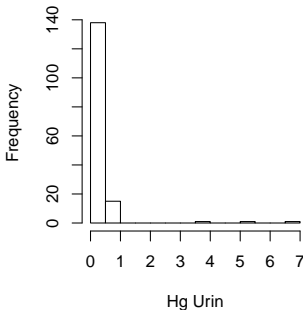
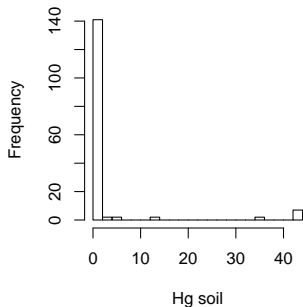
Which other factors might be responsible for high Hg concentrations in urin?



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:

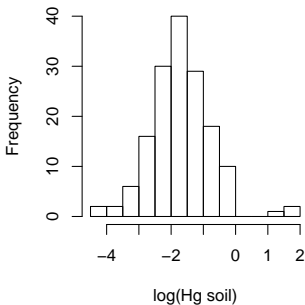
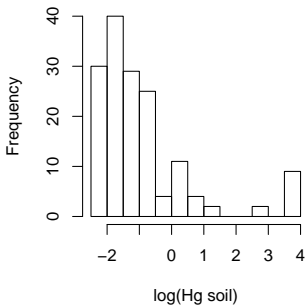
```
> par(mfrow=c(1,2))  
> hist(d.hg$Hg_soil,xlab="Hg soil",nclass=20,main="")  
> hist(d.hg$Hg_urin,xlab="Hg Urin",nclass=20,main="")
```



All Hg values seem to “stick” at 0.

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

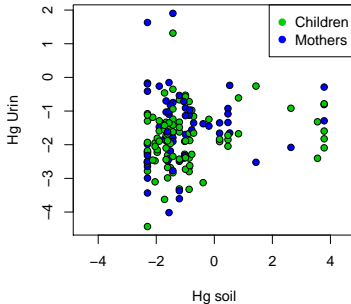
```
> par(mfrow=c(1,2))  
> hist(log(d.hg$Hg_soil),xlab="log(Hg soil)",nclass=20,main="")  
> hist(log(d.hg$Hg_urin),xlab="log(Hg soil)",nclass=20,main="")
```





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

```
> plot(log(Hg_urin) ~ log(Hg_soil), data=d.hg, xlab="Hg soil",  
+       ylab = "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**!

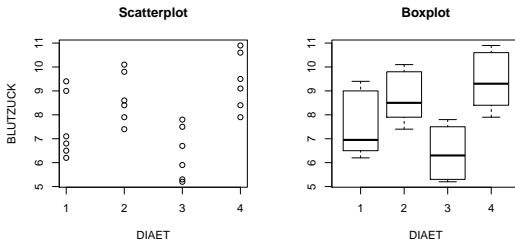
## Data example 3: Diet and blood glucose level

(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?

```
> par(mfrow=c(1,2))
> plot(BLUTZUCK ~ DIAET,d.blz,xaxt="n",main="Scatterplot")
> axis(1,1:4)
> boxplot(BLUTZUCK ~ DIAET,d.blz,xaxt="n",xlab="DIAET",main="Boxplot")
> axis(1,1:4)
```



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.

We will see in week 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):

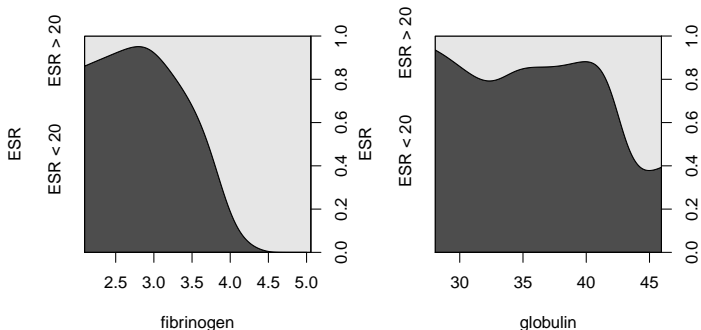
```
> library(HSAUR3)
> data("plasma", package="HSAUR3")
> plasma[c(1,5,9,10,15,29),]
```

	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** 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) **Determine 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.

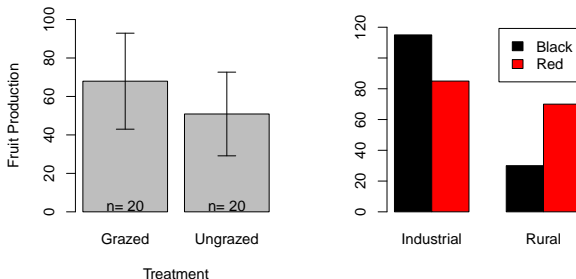
Representation	Useful for
Scatterplots	Pairwise dependency of continuous variables.
Histograms	Distribution of numerical variables.
Boxplots	Distribution of numerical variables, ev. conditionally on categories.
Conditional density plots	Dependency of a binary variable from a continuous variable.
Barplots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

# Barplots

Examples from Beckerman & Petchey:

**Left:** Fruit production in grazed and ungrazed areas. Grazing reduces above-ground biomass. How is fruit production affected by this?

**Right:** Number of birds of certain colors in two environments.

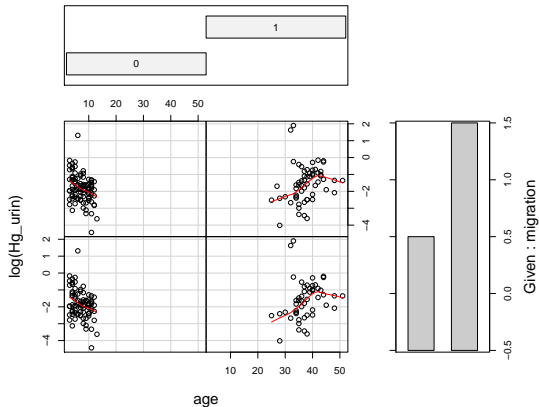


# Coplots

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).

# Experimental vs observational data

To do, ev mention only later in weeks 7/8

## Next week: Simple linear regression

## 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.