

# Statistical Programming Languages – Day 2

SVN-revision: 0

Uwe Ziegenhagen

Institut für Statistik and Ökonometrie  
Humboldt-Universität zu Berlin

<http://www.uweziegenhagen.de>



## Agenda for Today

- Data frames
- Reading and Writing Data
- Exploratory Data Analysis



## Data Frames

- ▣ matrices can only have one datatype
- ▣ data frames: several type allowed
- ▣ equal length of all elements required



## Data Frame Example

```
1 a <- c(10,20,15,43,76,41,25,46) # numeric
2 # Factor 'sex'
3 b <- factor(c("m","f","m","f","m","f","m","f"))
4 # siblings, numeric
5 c <- c(2,5,8,3,6,1,5,6)
6 myframe <- data.frame(a,b,c)
7 myframe
8 colnames(myframe) <- c("Age", "Sex", "Siblings")
```



## Factor variables

- Factor variables: categorical variables (numeric or string )
- advantages:
  - ▶ implemented correctly in statistical modeling
  - ▶ very useful in many different types of graphics
  - ▶ correct number of degrees of freedom



## Addressing Components

```
1 myframe[,1]
2 myframe["Age"]
3 myframe$Age
4 myframe[3,3] <-2 # change value
5 myframe[, -2] # all vars except 2nd
```



## Overview of Objects II

```
1 # add object set to search path  
2 attach(name)  
3 # remove from search path  
4 detach(name)
```



## Subgrouping Data Frames

```
1 > subset(myframe, myframe$Age>30) # 4 entries
2 > mean(subset(myframe$Age, myframe$Sex=="m"))
3 [1] 31.5
4 > mean(subset(myframe$Age, myframe$Sex=="f"))
5 [1] 37.5
6 myframe[(myframe$Sex=="m") & (myframe$Age>30),]
7 # males over 30
8 myframe[(myframe$Sex=="m") | (myframe$Age>30),]
9 # male or over 30
```





## Data Frames - Variables

```
1 > myframe <- cbind (myframe , "Income (USD)"=
2   c(1700,2100,2300,2050,2800,1450,3400,2000))
3 > names(myframe)[names(myframe)=="Income (USD)"] <-
   "IncomeUSD"
```

Task: Add variable IncomeEUR.



## Search and Replace

Use `gsub` to perform replacement of matches determined by regular expressions.

```
1 > names(myframe) <- gsub("In", "Out", names(myframe))
2 > myframe
3   Age Sex Siblings OutcomeUSD
4 1  10  m         2         1700
5 2  20  f         5         2100
```



## Deleting and Sorting

```
1 > myframe$Age <- NULL
2 > myframe
3
4 > myframe[order(myframe$Age),]
5   Age Sex Siblings OutcomeUSD
6 1  10  m         2         1700
7 3  15  m         2         2300
```



## Deleting and Sorting

1. Sortieren nach Sex
2. Sortieren nach Age

```
1 > myframe[order(myframe$Sex,partial=myframe$Age),]  
2   Age Sex Siblings OutcomeUSD  
3  2  20   f         5         2100  
4  6  41   f         1         1450  
5  8  46   f         6         2000  
6  4  43   f         3         2050  
7  1  10   m         2         1700  
8  3  15   m         2         2300  
9  5  76   m         6         2800  
10 7  25   m         5         3400
```



## Short excursion sed & awk

sed:

- ▣ stream editor, rowwise
- ▣ examples
  - ▶ `sed 's/abc/def/' input.txt >output.txt`
  - ▶ `sed 's|/|\|g' input.txt >output.txt`
  - ▶ **example using regular expression**
- ▣ extremely useful to process large amounts of data
- ▣ Tutorial: <http://www.grymoire.com/Unix/Sed.html>



## Short excursion sed & awk

awk:

- **Aho, Weinberger, Kernighan**
- in general used to work on columns
  - ▶ awk 'print 12' concatenates columns 1 and 2
  - ▶ awk 'print 1,3' prints columns 1 and 3
  - ▶ **another example using a sum**
- Tutorial: <http://www.vectorsite.net/tsawk.html>

In general: Avoid data processing inside **R**, try to do it outside.



## Data Management

- Sources of data:
  - ▶ Data in human readable format (CSV, TXT)
  - ▶ Data in binary format (Excel, SPSS, STATA)
  - ▶ Data from relational databases
- **R** has 100 built-in datasets: `objects(package:datasets)`
- many packages bring their own datasets



## Loading data from library

```
1 library("datasets") # loads dataset library
2                       # (automatically loaded)
3 data("pressure")    # loads dataset
4 data(pressure)      # alternative
5 pressure            # output pressure data
```





# Data Management

```
1 objects(package:datasets)
2 help(Titanic)
3 data(Titanic)
4 objects()
```



## Reading & Writing Data

```
1 data <- read.table("filename", header=TRUE)
2 # guess type of variable: int, double, text
3 # header with columnnames is available
4 names(data) # variable names
5 str(data)   # show structure of dataframe
6 head(data)  # show first rows
```



## Reading & Writing Data

```
1 # check if datafile has header
2 # may generate string matrix only
3 # if 1st row less than 2nd assume head
4 data <- read.table("filename")
```



## Reading & Writing Data

```
1 # using wrong separator
2 data <- read.table("file", sep="\t")
3 # assumes tabulator, may read whole
4 # line as one variable
```



## Reading & Writing Data

```
1 # reading NaNs
2 data <- read.table("file", na.strings=".")
3 # assumes NaN to be represented as '.'
```



## Reading & Writing Data

```
1 # reading CSV
2 # decimal sep '.', var. sep ','
3 data <- read.csv("file") #
4 # decimal sep ',', var. sep ';'
5 data2 <- read.csv2("file")
6 # direct import from Excel
7 data <- read.table(file="clipboard")
```



## Reading & Writing Data

```
1 x <- read.csv("beispiel.csv", sep=";")
2 dim(x)
3 names(x)
4 x
5 # write to file
6 write.table(x, file="test.csv", sep=";",
7   row.names= FALSE, quote=FALSE)
```



# Univariate Statistics

*ddistrib* density function

*pdistrib* distribution function

*qdistrib* quantile function

*rdistrib* random numbers





## Univariate Statistics

```
1 dnorm(0) # density value of  $N(0,1)$ 
2 pnorm(0) # cum. density up to 0
3 qnorm(0.5) # quantile for 0.5
4 rnorm(100) # vector with 100 random numbers
```



# Univariate Statistics

*ddistrib* density function

*pdistrib* distribution function

*qdistrib* quantile function

*rdistrib* random numbers



## Distributions in standard *R*

<key>binom	Binomial
<key>chisq	Chi-Squared
<key>exp	Exponential
<key>f	F
<key>hyper	Hypergeometric
<key>multinom	Multinomial
<key>logis	Logistic
<key>norm	Normal
<key>pois	Poisson
<key>t	Student t
<key>unif	Uniform

with <key>= d, p, q or r



## Empirical Distributions in *R*

```
1 density() # KDE using Gaussian kernel  
2 ecdf() # empirical cdf
```



## Sampling in R

```
1 sample(n) # sample 1:n vector
2 sample(x) # shuffle the x vector
3 sample(x, replace=TRUE) bootstrap x vector
4 sample(x, n) # draw sample of size n from x
5 sample(x, n, replace = TRUE) # bootstrap sample from
   x
```

Seed is stored in `.Random.seed`, for simulations use `set.seed()`



## Summary statistics

```
1 mean(x) # mean *
2 median(x) # median
3 var(x) # sample variance
4 sd(x) # sample std. deviation
5 cov(y) # cov of matrix y
6 quantile(x,p) # sample quantile *
7 min(x) # minimum of x *
8 max(x) # maximum of x *
9 range() # range of x *
10 skewness(x) # skewness
11 kurtosis(x) # kurtosis
```

\* can remove NaNs using parameter `na.rm=T`



## Linear Regression

linear regression model

- tries to model relation between dependent variable  $Y$  and  $1 \dots n$  indep. variables  $X_1, \dots, X_n$
- influence of variables is linear, first regressor  $X_1$  usually set to constant
- sample of size  $n$  is fitted to model:

$$y_i = \beta_1 + \beta_2 \cdot x_2 + \dots + \beta_n \cdot x_n + \varepsilon_i$$

$$y_i = x_i^T \beta + \varepsilon_i$$

$$y = X\beta + \varepsilon$$



# Linear Regression

Goals:

- ▣ estimate unknown  $\beta$ s using *least squares*
- ▣ decide if all variables are needed
- ▣ check if resulting model explains data well enough
- ▣ use model to forecast

$$\hat{\beta} = (X^T X)^{-1} X^T y$$





## Linear Regression

```
1 # standard model
2 lm(y~x + z)
3 # no intercept
4 lm(y~x - 1)
5 # using data frame
6 lm(amount ~ price, data = consumption)
7 # using data frame and attach()
8 lm(amount ~ price)
```



## Exercise

Download Hubble data from

<http://lib.stat.cmu.edu/DASL/Datafiles/Hubble.html> and  
estimate the hubble constant  $H$  by the model

$$\text{recession-velocity} = H \cdot \text{distance}$$



```
> summary(lm)
```

```
Call:
```

```
lm(formula = rec.vel ~ distance - 1)
```

```
Residuals:
```

Min	1Q	Median	3Q	Max
-411.544	-191.302	-7.103	127.951	496.063

```
Coefficients:
```

	Estimate	Std. Error	t value	Pr(> t )
distance	423.94	42.15	10.06	6.87e-10 ***

```
---
```

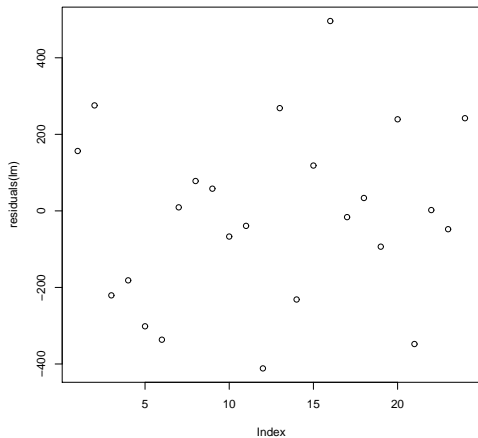
```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 229 on 23 degrees of freedom
```

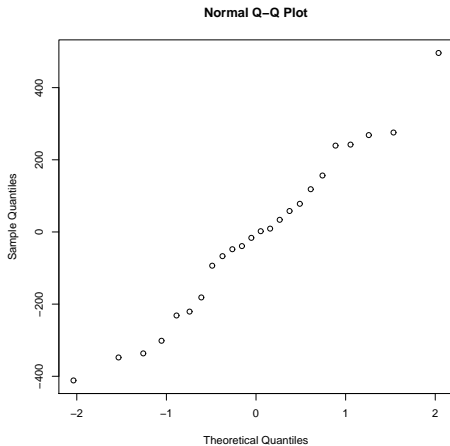
```
Multiple R-squared: 0.8147,    Adjusted R-squared: 0.8067
```

```
F-statistic: 101.1 on 1 and 23 DF,  p-value: 6.869e-10
```

# Residuals



# Residuals



# Multiple Linear Regression

1. Download Cereal data from DASL
2. Read data as dataframe
3. Run linear regression  $\text{rating} = \text{sugars} + \text{fat}$

Call:

```
lm(formula = rating ~ sugars + fat, data = data)
```

Residuals:

	Min	1Q	Median	3Q	Max
	-14.6640	-5.6937	0.2078	4.7660	32.6163

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	61.0886	1.9527	31.284	< 2e-16 ***
sugars	-2.2128	0.2347	-9.428	2.59e-14 ***
fat	-3.0658	1.0365	-2.958	0.00416 **

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 8.755 on 74 degrees of freedom

Multiple R-squared: 0.6218, Adjusted R-squared: 0.6116

F-statistic: 60.84 on 2 and 74 DF, p-value: 2.371e-16

## ***t*-Test**

- checks if certain coefficient  $\beta_j$  is different from 0.
- teststatistics

$$t = \frac{\hat{\beta}_j}{\text{SD}(\hat{\beta}_j)}$$

- under  $H_0 : t \sim t_{n-p}$  with  $p$  as number of independ. variables





## F-Test

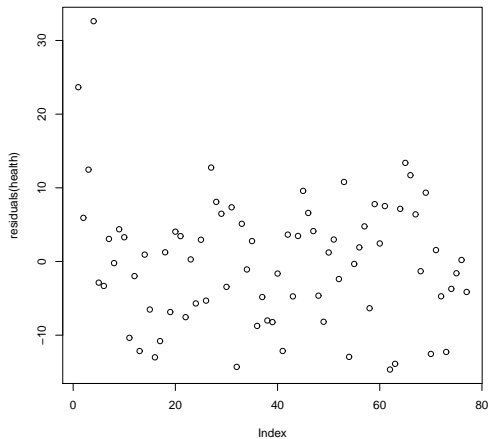
- idea: check if sum of squared residuals is reduced **significantly** if one regressor is added
- add one regressor  $\Rightarrow$  model gets better, but significantly?
- Compute RSS1 for full model with  $k$  parameters, compute RSS2 for simplified model with  $k - q$  parameters
- Compute teststatistics

$$F = \frac{(RSS2 - RSS1)/q}{RSS1/(n - k)}$$

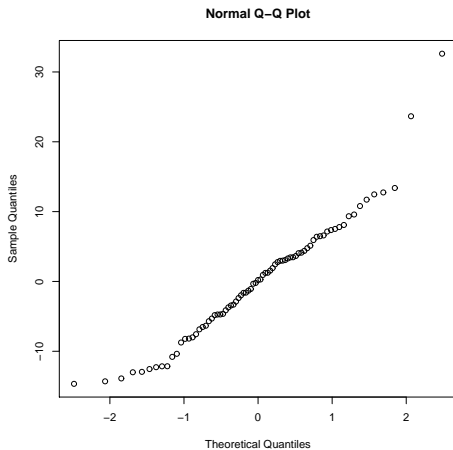
- under  $H_0 : F \sim F_{(n-1, n-q-1)}$
- for  $q = 1$  F-test equivalent to t-Test



# Residuals



# Residuals



# Residuals

