

QUANTITATIVE RESEARCH METHODS

DR. MEIKE MORREN



about myself

Background

□ Academic experience

- Bsc Sociology – Msc Social Sciences (University of Amsterdam)
- PhD in methods & statistics (Tilburg University)
 - Latent class modeling
 - Cognitive interviewing

□ Work experience: marketing research agency Centerdata

- Online survey panel - non response
- Project mobile surveys

□ Teaching

- Bachelor Marketing Research/bachelor Business Research Methods

□ Interests

- Modeling approaches & statistics
- Survey responses

Contact data

- Office hours: by appointment

meike.morren@vu.nl

020 5982317

- Secretary:

Sandra van Arendonk (5A64)

sandra.van.arendonk@vu.nl

020-598 7145

COURSE OUTLINE



Course

- Requirements :

 - Attendance lectures (12)

 - Exercises during lecture

 - Study literature => articles & book excerpts

- Grade :

 - Assignments (6, 40%)

 - Exam (60%)

- You are allowed to miss one lecture

- Assignments are made in groups of 2 students, handed in the 11th of May via github

Literature

- Background literature on R (available in ub.vu.nl):
 - ▣ Everitt, B. S. & T. Hothorn (2011). An introduction to applied multivariate analysis with R. New York: Springer

- Articles
 - ▣ Link available via bb, search [google.scholar.com](https://scholar.google.com)

- Book excerpts
 - ▣ See https://github.com/meikemorren/BIS_QRM/issues

Lectures (1)



- Lecture 1: Introduction to R
- Lecture 2: Data cleaning, missing data analysis
- Lecture 3: Writing functions in R
- Lecture 4: Generalized linear framework
- Lecture 5: Simple and multiple regression
- Lecture 6: Categorical data analysis (logit regression)

Lectures (2)



- Lecture 7: PCA
- Lecture 8: Factor analysis
- Lecture 9: Cluster analysis
- Lecture 10: Plots
- Lecture 11: Multilevel analysis
- Lecture 12: webscraping / Q&A

Assignments

1. Analyze missing data
2. Conduct regression analysis
3. Conduct logit regression analysis
4. Analyze multi-item scale(s) using CFA
5. Apply and interpret cluster analysis
6. Compare countries or multilevel regression analysis

Assignment guidelines

- Questions to interpret the results you obtain with programming
- Deadline: 11 May (you can hand in earlier)
- Github

The assignments include:

- Programming (include your code)
- Interpretation of results (focus on argumentation)

Assignment dataset 1

Bring your own data (if suitable for assignment) or use data that is provided for (see bb)

World Values Survey dataset includes:

- Green attitudes
- Work attitudes
- Schwarz values

Across 40+ countries

For assignment 1 (missing data), 2 (regression), 4 (factor analysis) and 6 (multi-level analysis)

Assignment dataset 2

Bring your own data (if suitable for assignment) or use data that is provided for (see bb)

Tablet dataset which includes:

- Sales ranking
- Number of reviews
- Number of characteristics of tablets (screen size, weight, battery life etc)

For assignment 3 (logit regression) and 4 (cluster analysis)

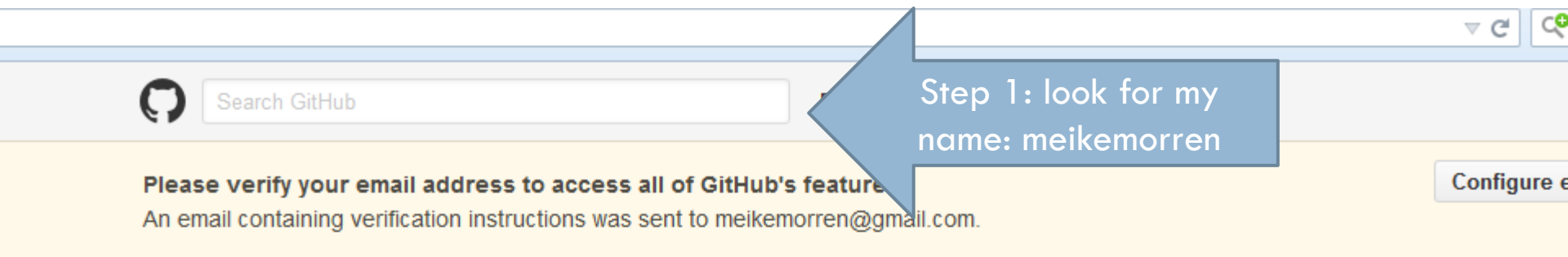
Assignments via github

- Create an account (github.com)
 - ▣ send me your username
- Fork my github repository
 - ▣ (= copy the files to your account)
- Download the files to a location in which you will work on your exercises / assignments
- Upload your answers to github

https://github.com/meikemorren/BIS_QRM.git

GitHub

- ❑ Free
- ❑ Used for coders to work simultaneously on projects



Step 1: look for my
name: meikemorren

Learn Git and GitHub without any code!

Using the Hello World guide, you'll create a repository, start a branch, write comments, and open a pull request.

Let's get started!

GitHub



[Pull requests](#) [Issues](#) [Gist](#)

Please verify your email address to access all of GitHub's features.

An email containing verification instructions was sent to meikemorren@gmail.com.

Search

 **Repositories**

 [Code](#)

2

 [Issues](#)

 [Users](#)

1

[Advanced search](#) [Cheat sheet](#)



We couldn't find any repositories matching

You could try an [advanced search](#)

Step 2: look for users



LECTURE 1



Contents lecture 1

- Introduction to modeling, and the concepts we will use in this lecture:
 - ▣ Introduction to R
 - Data types
 - Explore dataset
 - ▣ Scale types & distributions

R (STUDIO)

Assignment software

- Download R : <http://cran.xl-mirror.nl/>
 - ▣ Install R for the first time
 - ▣ Run exe

- R Studio : <http://www.rstudio.com/products/RStudio/download/>
 - ▣ Choose installers
 - ▣ Select platform

Work environment R studio

The screenshot shows the RStudio interface with the following components:

- Source Editor (Top Left):** Contains R code for creating data structures and matrices. A large blue arrow points to this pane with the text "Text editor Write code".

```
1- #####  
2- ##### INTRODUCTION TO R #####  
3- #####  
4- # assign numbers to words  
5 identity <- 1  
6 identity  
7  
8  
9 # different scale types: http://www.statnet.hawaii.edu/input/dataltypes.html  
10 # vector: a row of numbers, words or i/i statements  
11 x <- c(1,2,3,4,5,6,7,8,9,10) # numeric vector  
12 y <- c("one","two","three") # character vector  
13 z <- c(TRUE,TRUE,TRUE,FALSE,TRUE,FALSE) # logical vector  
14 z  
15 h  
16 c  
17  
18 # refer to elements of the vector  
19 x[2,4] # 2nd and 4th elements of vector  
20  
21 # create matrix  
22 y <- matrix(1:20, nrow=5, ncol=4, byrow=TRUE)  
23 y  
24 y <- matrix(1:20, nrow=5, ncol=4)  
25 y  
26  
27 # identify rows, columns and elements in matrix  
28 y[,4] # 4th column of matrix  
29 y[1,] # 1st row of matrix  
30 y[2:4,1:3] # rows 2,3,4 of columns 1,2,3  
31  
32 # dataframe: more general than matrix  
33 d <- data.frame(x = 1:4)  
34 d  
35 rm(d)
```
- Environment/History (Top Right):** Lists objects in the environment.

| Object | Size |
|-------------|-----------------------------|
| Alpha | 41 obs. of 4 variables |
| ELECT | 6 obs. of 64 variables |
| ENTRY | 141 obs. of 15 variables |
| ENTRYLONG | 87 obs. of 10 variables |
| ENTRYVALUES | 1720 obs. of 87 variables |
| LPI | 177 obs. of 3 variables |
| PIA | 1021 obs. of 22 variables |
| SB | 296779 obs. of 27 variables |
| SG | 24 obs. of 1 variables |
| SG | 29 obs. of 3 variables |
| SGROW | 877 obs. of 1 variables |
| SGROWLONG | 811 obs. of 3 variables |
| SGROWVALUES | 2000 obs. of 20 variables |
- Console (Bottom Left):** Shows the R version and copyright information.

```
R version 3.1.3 (2015-08-09) "Smooth SideWalk"  
Copyright (C) 2015 The R Foundation for Statistical Computing  
Platform: x86_64-w64-mingw32/x86_64 (32-bit)  
  
R is free software and comes with ABSOLUTELY NO WARRANTY.  
You are welcome to redistribute it under certain conditions.  
type 'license()' or 'licence()' for distribution details.  
  
R is a collaborative project with many contributors.  
Type 'contributors()' for more information and  
'citation()' on how to cite R or R packages in publications.  
  
type 'demo()' for some demos, 'help()' for on-line help, or  
'help.start()' for an HTML browser interface to help.  
type 'q()' to quit R.  
  
Error in rgs::readRscript() : attempt to apply non-function  
Workspace loaded from ~/.Rsave/
```
- Help Pane (Bottom Right):** Displays the documentation for the 'fit' function.

Fit Confirmatory Factor Analysis Models

Description

Fit a Confirmatory Factor Analysis (CFA) model.

Usage

```
fit(model = NULL, data = NULL,  
      parameters = "default", fixed = "default",  
      orthogonal = FALSE, random = FALSE, nonrandom = "default",  
      nonrandomization = "default", random = FALSE,  
      missing = "default", ordered = NULL,  
      sample.size = NULL, sample.size.random = "default",  
      sample.mean = NULL, sample.mean.random = NULL,  
      ridge = 1e-06, group = NULL,  
      group.labels = NULL, group.equal = "", group.parameters = "",  
      group.weights = FALSE, cluster = NULL, nonrandomize = "",  
      estimation = "default", likelihood = "default", link = "default",  
      information = "default", se = "default", test = "default",  
      nonrandom = 10000, seed = "default", nonrandomization = "default",  
      do.fis = TRUE, control = list(), MLE.V = NULL, MLEOV = NULL,  
      zero.add = "default", zero.keep.random = "default",  
      zero.call.mean = TRUE,  
      zero = "default", variables = FALSE, warn = TRUE, show = FALSE)
```

Work environment R studio

The screenshot displays the RStudio work environment. The main editor pane on the left contains R code for data manipulation and matrix operations. The console pane at the bottom shows the R version (3.1.3) and copyright information. The environment pane on the right lists loaded packages and data objects. A large blue arrow points from the text 'Output (black) Code (blue)' to the console and editor panes.

```
1- #####
2- ##### INTRODUCTION TO R #####
3- #####
4- # assign numbers to words
5- identity <- R
6- identity
7-
8-
9- # different scale types: http://www.statnet.hawaii.edu/input/dataltypes.html
10- # vector: a row of numbers, words or i/i statements
11- x <- c(1,2,3,4,5,-2,4) # numeric vector
12- b <- c("one","two","three") # character vector
13- i <- c(TRUE,TRUE,TRUE,FALSE,TRUE,FALSE) # logical vector
14- a
15- b
16- c
17-
18- # refer to elements of the vector
19- x[2,4] # 2nd and 4th elements of vector
20-
21- # create matrix
22- y <- matrix(1:20, nrow=5, ncol=4, byrow=TRUE)
23- y
24- y <- matrix(1:20, nrow=5, ncol=4)
25- y
26-
27- # identify rows, columns and elements in matrix
28- y[,4] # 4th column of matrix
29- y[1,] # 1st row of matrix
30- y[2:4,1:3] # rows 2,3,4 of columns 1,2,3
31-
32- # dataframe: more general than matrix
33- d <- data.frame(x = 1:4)
```

R version 3.1.3 (2015-08-09) "Smooth Sailing"
Copyright (C) 2015 The R Foundation for Statistical Computing
Platform: x86_64-w64-mingw32/x86_64 (32-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
type 'license()' or 'licence()' for distribution details.

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

type 'demo()' for some demos, 'help()' for on line help, or
'help.start()' for an HTML browser interface to help.
type 'q()' to quit R.

Error in res.method() : attempt to apply non function
Workspace loaded from ~/.Rsave

Environment: Global Environment

DATA

| Package | Version | Size | Variables |
|------------|---------|-------------|-----------------|
| Alpha | 4.1 | 41 obs. | of 4 variables |
| BIGEST | 6 | 6 obs. | of 64 variables |
| CATRY | 141 | 141 obs. | of 15 variables |
| CATRYLONG | 87 | 87 obs. | of 10 variables |
| CATRYSHORT | 1720 | 1720 obs. | of 87 variables |
| CLM | 177 | 177 obs. | of 3 variables |
| Data | 1021 | 1021 obs. | of 22 variables |
| LSS | 296779 | 296779 obs. | of 27 variables |
| PSM | 29 | 29 obs. | of 1 variables |
| LSSC | 29 | 29 obs. | of 3 variables |
| PSMlong | 107 | 107 obs. | of 1 variables |
| LSSLONG | 811 | 811 obs. | of 3 variables |

Fit Confirmatory Factor Analysis Models

Description

Fit a Confirmatory Factor Analysis (CFA) model.

Usage

```
model = NULL, data = NULL,  
parameters = "default", fixed = "default",  
orthogonal = FALSE, equal = FALSE,  
parameterization = "default", method = "ML",  
fitting = "default", ordered = NULL,  
sample.size = NULL, sample.size.factor = "default",  
sample.mean = NULL, sample.mean = NULL,  
ridge = 1e-05, prior = NULL,  
prior.labels = NULL, prior.equal = "", prior.priors = "",  
prior.weights = FALSE, prior.mean = NULL, prior.weights = "",  
prior.labels = "default", likelihood = "default", link = "default",  
information = "default", se = "default", mean = "default",  
nonlinear = 1000, seed = "default", nonconvergence = "default",  
do.fit = TRUE, control = list(), MEAN = NULL, NONCONV = NULL,  
save.data = "default", save.mean.weights = "default",  
save.all.data = TRUE,  
save = "default", verbose = FALSE, warn = TRUE, show = FALSE
```

Arguments

Work environment R studio

The image shows the RStudio interface. The source editor on the left contains R code for data manipulation. The environment pane on the right lists the objects in the workspace. The console at the bottom shows the R version and license information.

Source Editor:

```
1 # Introduction to R
2 # =====
3 # assign numbers to words
4 identity <- 1
5 identity
6
7 # different scale types: integer
8 # vector: a row of numbers, words
9 # e.g. 1(1,2,3), 1(1,2,3) # matrix
10 # e.g. c("one", "two", "three") # c
11 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
12 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
13 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
14 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
15 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
16 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
17 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
18 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
19 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
20 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
21 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
22 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
23 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
24 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
25 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
26 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
27 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
28 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
29 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
30 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
31 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
32 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
33 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
34 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
35 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
36 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
37 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
38 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
39 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
40 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
41 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
42 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
43 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
44 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
45 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
46 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
47 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
48 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
49 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
50 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
51 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
52 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
53 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
54 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
55 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
56 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
57 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
58 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
59 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
60 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
61 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
62 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
63 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
64 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
65 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
66 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
67 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
68 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
69 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
70 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
71 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
72 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
73 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
74 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
75 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
76 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
77 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
78 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
79 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
80 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
81 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
82 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
83 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
84 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
85 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
86 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
87 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
88 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
89 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
90 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
91 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
92 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
93 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
94 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
95 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
96 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
97 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
98 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
99 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
100 # e.g. 1(TRUE, TRUE, TRUE, FALSE, TRUE)
```

Environment:

| Object | Class | Attributes |
|-------------|--------|-----------------------------|
| Alpha | matrix | 4 obs. of 4 variables |
| BIGOT | matrix | 6 obs. of 64 variables |
| ENTRY | matrix | 141 obs. of 25 variables |
| ENTRYFREQ | matrix | 87 obs. of 10 variables |
| ENTRYWORDS | matrix | 1720 obs. of 16 variables |
| LPI | matrix | 177 obs. of 3 variables |
| DATA | matrix | 1021 obs. of 22 variables |
| LSS | matrix | 206779 obs. of 27 variables |
| ESSE | matrix | 24 obs. of 1 variables |
| LSSC | matrix | 20 obs. of 3 variables |
| ESSEES | matrix | 107 obs. of 1 variables |
| LSSNOVCHONG | matrix | 811 obs. of 3 variables |
| ESSEES | matrix | 107 obs. of 1 variables |

Console:

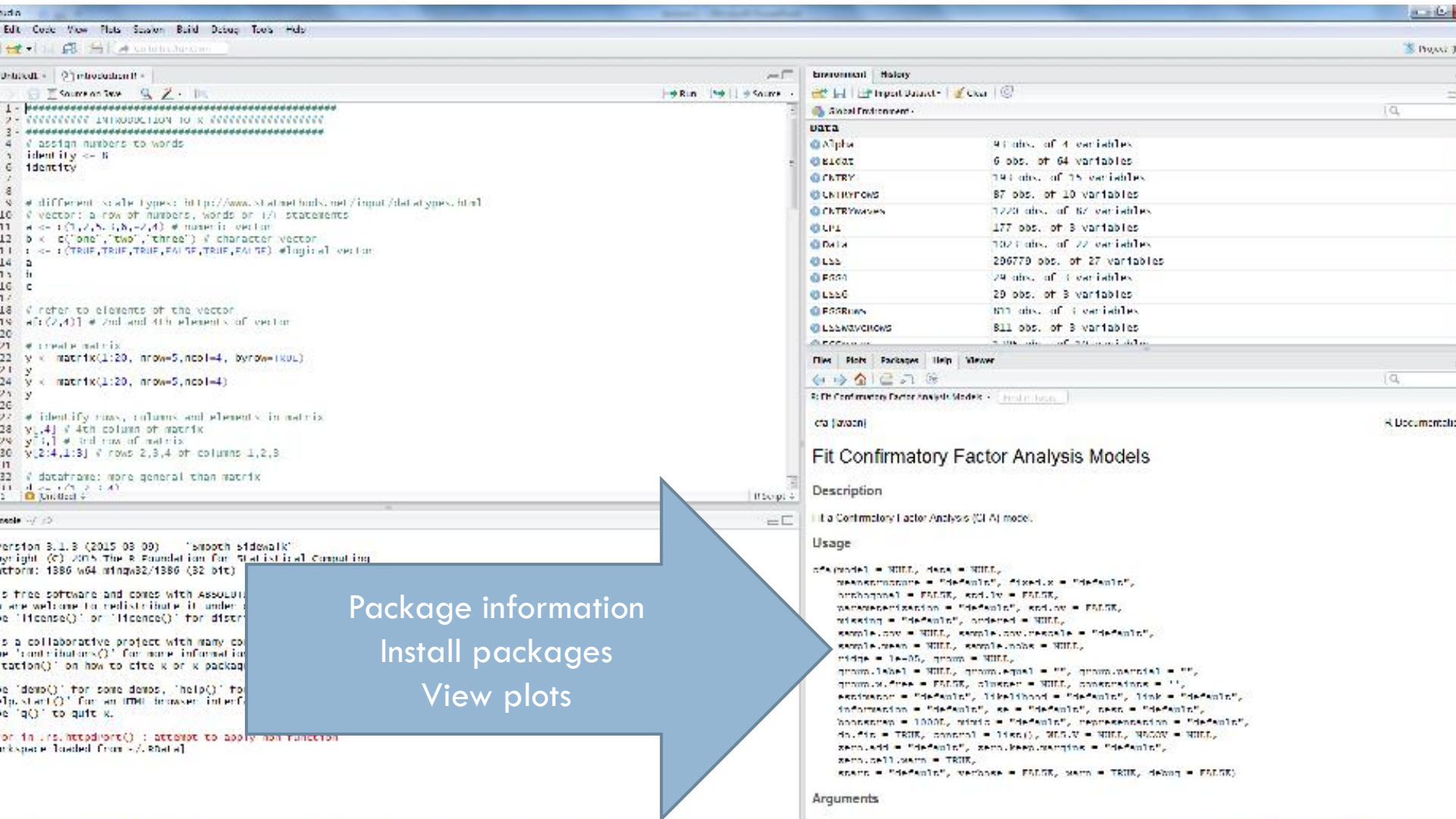
```
R version 3.1.3 (2015-03-09) "smooth sidewalk"
Copyright (C) 2015 The R Foundation for Statistical Computing
Platform: x86_64-w64-mingw32/x86_64 (32-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
For more information on distribution details, see the 'license()' or 'licence()' for distribution details.

R is a collaborative project with many contributors.
For more information on contributors, see 'contributors()'.
For more information on how to cite R or R packages in publications,
see 'citation()' or 'citation()' for distribution details.

Use 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Use 'q()' to quit R.
```


Work environment R studio



READ DATA INTO R

SPSS dataset

- Is considered to be a matrix in R:
 - ▣ Consists of variables (columns)
 - ▣ And observations (rows)
- However, variables can be of different datatypes
 - ▣ Character vectors
 - ▣ Numeric vectors
 - ▣ ...
- Therefore, data is a dataframe

Useful functions

- `read.table`

- `nrow`

- `ncol`

- `colnames`

- `typeof`

DATA TYPES



Data types

- ❑ Integers (no decimals)
- ❑ Numeric (with decimals)
- ❑ Characters (words)
- ❑ Complex
- ❑ Logical (true, false)

DATA OBJECTS



Data objects

- Dataframes
- Factors
- Vectors
- Matrices

The base type in R is vector (not scalar)

These are the abstract data types, also called class

Mode

Data objects are stored in the memory by mode

- Numeric
- Complex
- Character
- Logical
- *List, function*

Changing the mode of an object is often called 'coercion'. The mode of an object can change without necessarily changing the class.

R MANIPULATION

Create vectors, and select from vectors

Create a vector

- Everything is an object
- You can assign a value to an object
 - ▣ `R <- 4`
- Create a vector
 - ▣ `Vec <- rep(1, 4)` # of four 1s
 - ▣ `Vec <- (1:4)` # 1 thru 4
- Assign new value to third element
 - ▣ `Vec[3] <- 5`

Select from vector (matrix)

□ Create matrix

- ▣ `df <- rbind(c(1, 2, 3), c(4, 4, 4))` # create matrix by row
- ▣ `df <- cbind(c(1, 2, 3), c(4, 4, 4))` # create matrix by column
- ▣ `df <- matrix(c(1, 2, 3, 4, 4, 4), nrow=2, ncol=3, byrow=TRUE)`

□ Select element, row, column

- ▣ `df[2, 3]` # select second row, third element
- ▣ `df[, 3]` # select third column
- ▣ `df[2,]` # select second row

Select from data

- `df[df$gender==1,]` # select all data for males
- `df$job[df$gender==1]` # select only job for males
- `df$job[df$edu>3]` # select only job for edu higher than level 3

- For data levels see questionnaire

Accessing attributes

□ Attributes = characteristics of a vector

- ▣ `length(Vec)`
- ▣ `sort(Vec)`
- ▣ `names(Vec)`

□ Assign colnames & select second

- ▣ `colnames(df) <- c("first", "second", "third")`
- ▣ `colnames(df)[2]`

□ Get first impression

- ▣ `summary(Vec)`

R as calculator

□ Basic operations

- ▣ `x <- rnorm(10)` # create vector of 10 random numbers from normal distribution
- ▣ `x + 1`
- ▣ `v <- x * 8 + 1`
- ▣ `v/5`

□ Arithmetic functions

- ▣ `mean(Vec)`
- ▣ `sum(Vec)`
- ▣ `var(Vec) / sd(Vec)`

Some manipulation signs

```
x<-5; y <- 1:10
```

□ Square

- x^2 or x^{**2}

□ Square root

- `sqrt(x)`

□ Equal to

- $y==x$

□ Not equal to

- $y!=x$

□ Not x

- $!x$

□ Smaller than

- $x \leq y$

□ Greater than

- $x \geq y$

□ Equal to

- $y==x$

□ x OR y

- $x \mid y$

□ x AND y

- $x \& y$

Simple calculations

□ Calculate variance on vector

- ▣ `var(x)`

- ▣ `sum((x-mean(x))^2/(length(x)-1))`

□ Calculate z-score

- ▣ `mean(x)`

- ▣ `sum(x)/length(x)`

□ Calculate z-score

- ▣ Standard deviation: `sd(x)`

- ▣ `(x-mean(x))/sd(x)`

- ▣ `(x-mean(x))/sqrt(var(x))`

Algebra (1)

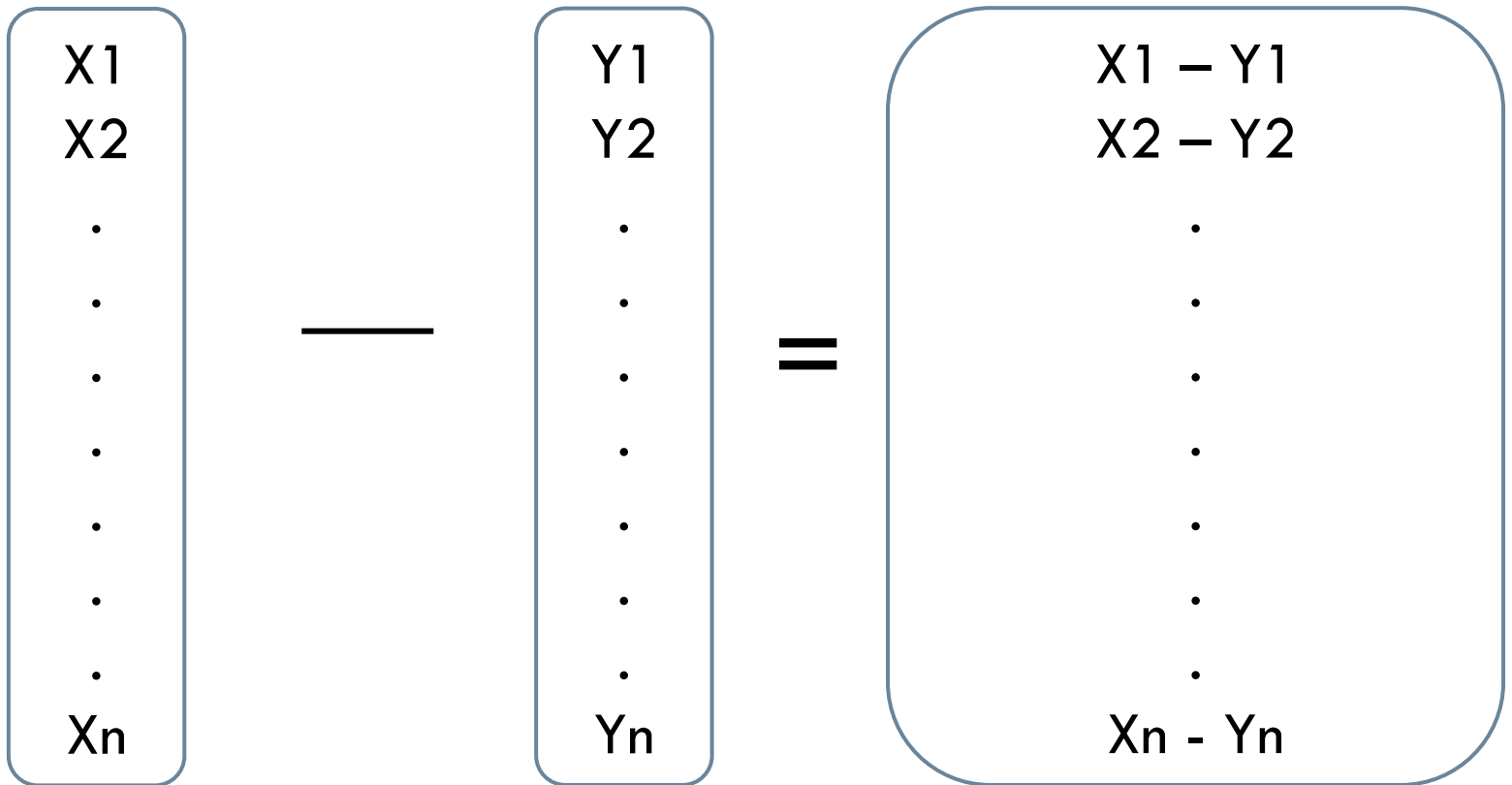
□ $x - \text{mean}(x)$:

$$\begin{array}{c} x_1 \\ x_2 \\ \cdot \\ \cdot \\ \cdot \\ \cdot \\ \cdot \\ \cdot \\ x_n \end{array} - \bar{X} = \begin{array}{c} x_1 - \bar{X} \\ x_2 - \bar{X} \\ \cdot \\ \cdot \\ \cdot \\ \cdot \\ \cdot \\ \cdot \\ x_n - \bar{X} \end{array}$$

Algebra (2)

```
□ y <- rnorm(10)
```

□ $x - y$



R exercise (1)

see blackboard/documents/assignment:

- Open exe1_1.r
- Download data into R
 - ▣ set directory, read.table, give a name to the dataset
- Inspect dataset
 - ▣ Explore elements, rows, columns
 - ▣ Assign names to columns
- Manipulate data
 - ▣ Calculate variance, z-score (first manually, check with R function)

SCALE TYPES



Primary scale types (1)

□ Nominal

- Categories have labels
- Numbers no meaning

□ Ordinal

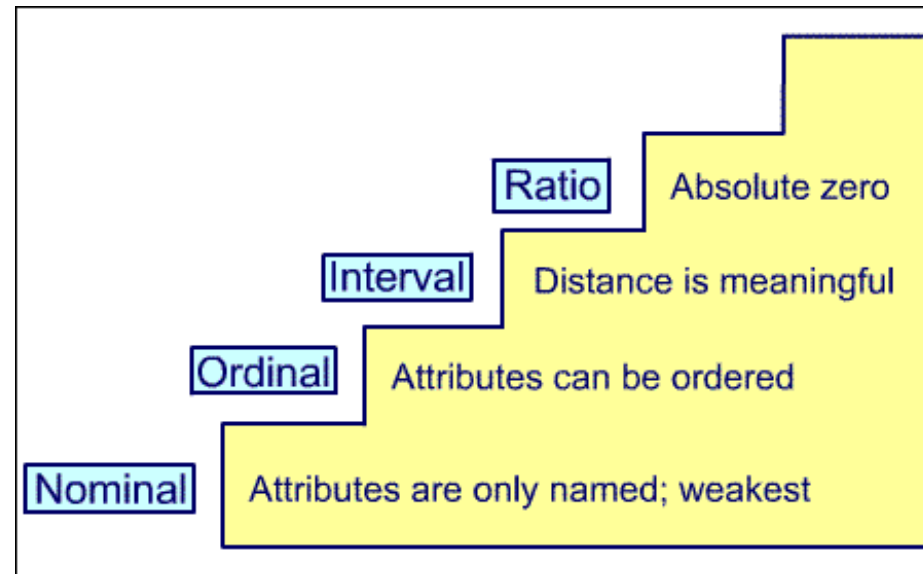
- + ordering
- Numbers arbitrary, only order numbers has meaning

□ Interval

- + equal distances
- Starting point arbitrary, order and distances between numbers have meaning

□ Ratio

- + starting point
- Numbers have meaning in all mathematical senses



Primary scale types (2)

□ Nominal

- Count frequency
- Chi square, binominal test

□ Ordinal

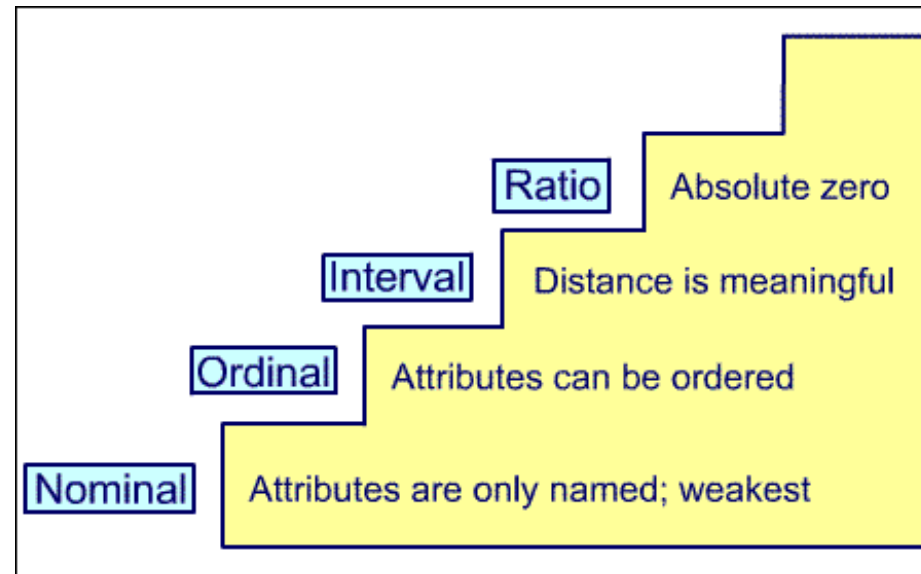
- rank order correlation
- ANOVA (Kruskal-Wallis)
- Kolmogorov-Smirnov test

□ Interval

- T-test
- ANOVA
- Regression
- Factor analysis

□ Ratio

- All mathematical operations are possible



Likert Scale

Itemised scale:

Number or label to describe categories

Ordered categories



| | Strongly disagree | Disagree | Neither agree nor disagree | Agree | Strongly agree |
|---|----------------------------|----------------------------|----------------------------|----------------------------|----------------------------|
| Would you agree to the statement that muffins are more healthy than cupcakes? | <input type="radio"/> 1 | <input type="radio"/> 2 | <input type="radio"/> 3 | <input type="radio"/> 4 | <input type="radio"/> 5 |

R functions to inspect variables

□ Run frequencies

- ▣ `table(df[,1])`
- ▣ `df<-as.data.frame(df)`
- ▣ `table(df$first)`

□ Run descriptives

- ▣ `summary(df$first)`

□ Add variable labels (change into factor)

- ▣ `dataSet$gender <-factor(dataSet$gender, levels = c(-2, 1,2), labels=c("missing", "male", "female"))`
- ▣ `attributes(dataSet$gender)`

R exercise (2): Open exe1_2.r

▣ Run frequencies

- ▣ `table`

▣ Descriptives

- ▣ `summary`, `mean`

▣ Crosstable

- ▣ `table(var1, var2)`

▣ Recode

- ▣ `df[df==4]<-2 # recode all values of 4 to 2`

- ▣ `df[df$first==4,] <-2 # recode values of 4 in first column to 2`

▣ Convert variables to factors

- ▣ `factor(df$first, levels=c(1,2),
labels=c("one", "two"))`

PROBABILITY DENSITY FUNCTIONS



Distributions



Used to

- describe variables
- test hypotheses
- estimate model parameters

Many distributions

- Normal
 - ▣ Continuous variables
- Logistic
 - ▣ Nominal, multi-nominal, or ordinal variables
- Binomial
 - ▣ Nominal variables
- Poisson
 - ▣ Count variables

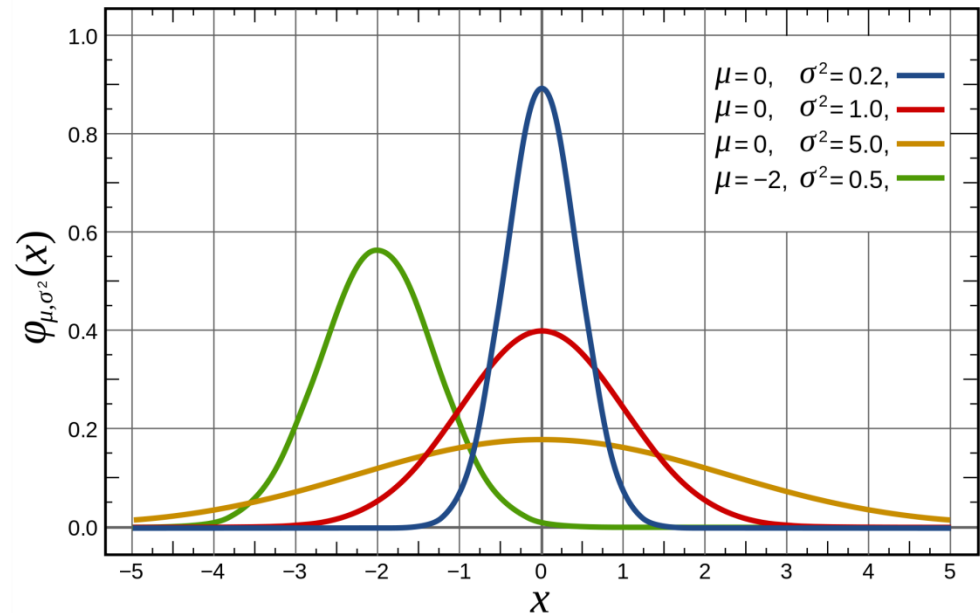
Normal distribution

□ Interval / ratio variable results in normal distribution

μ = mean(location parameter)

σ = standard deviation
(scale parameter)

□ PDF:



$$f(x; \sigma, \mu) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Probability density function (PDF) describes the area under the curve

Central limit theorem

- Random variables
- Independently drawn from
- Independent distributions
- i.i.d.
- Mean of 0
- Standard deviation of 1
- Total area under the curve is 1
- Converge to normal distribution when number of variables is large

Binomial distribution

□ PMF:

$$f(k; n, p) = \Pr(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}$$

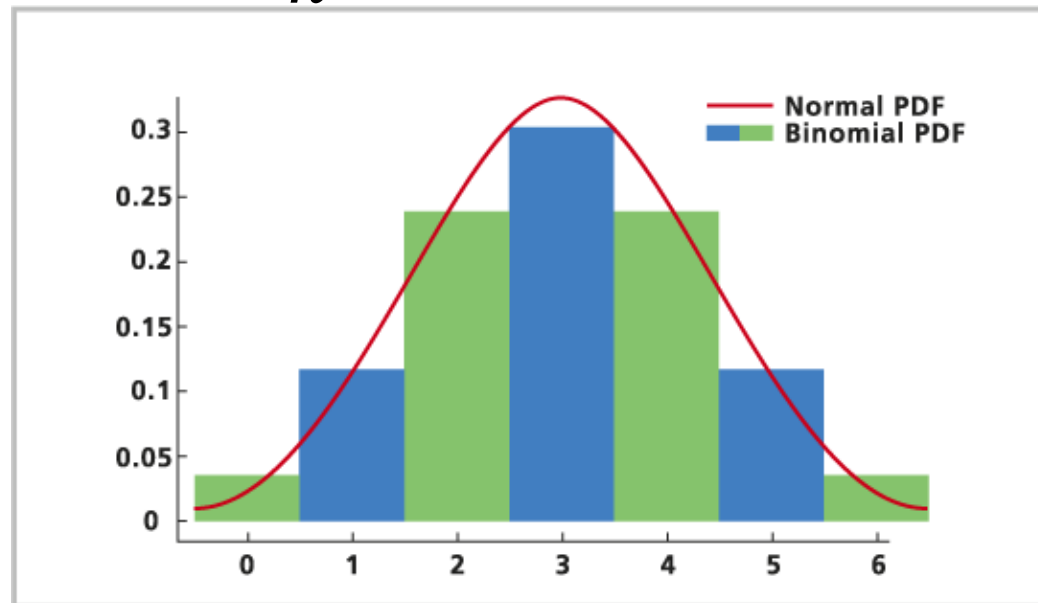
for $x = 1, 2, \dots$

□ $k = \#$ successes

□ $n = \#$ trials

□ $p =$ probability

□ $\binom{n}{k} = \frac{n!}{k!(n-k)!}$



Independent successes/failures (=with replacement)

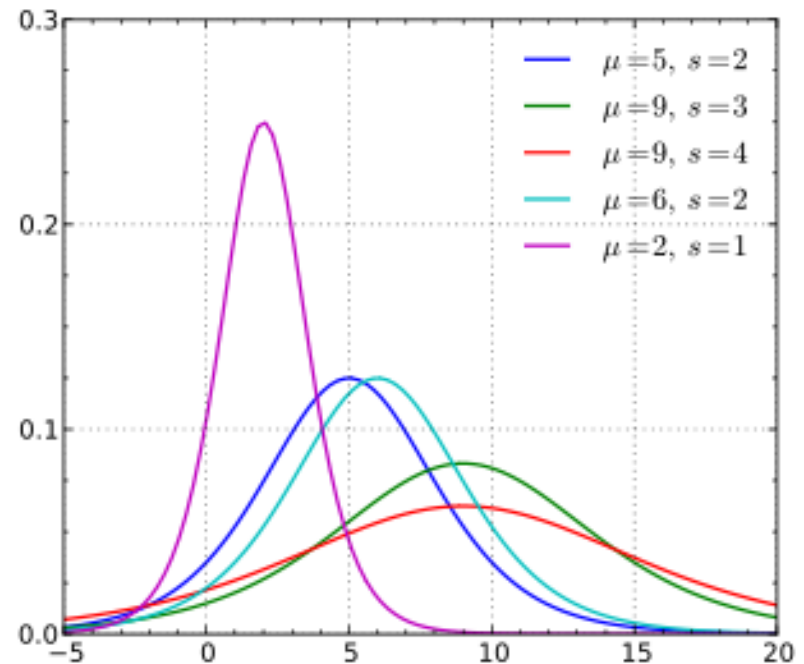
Logistic distribution

- The logistic distribution has slightly longer tails compared to the normal distribution

- μ = mean (location parameter)
- s = variance (scale parameter)

- PDF:

$$f(x; \mu, s) = \frac{e^{-\frac{x-\mu}{s}}}{s(1+e^{-\frac{x-\mu}{s}})^2}$$



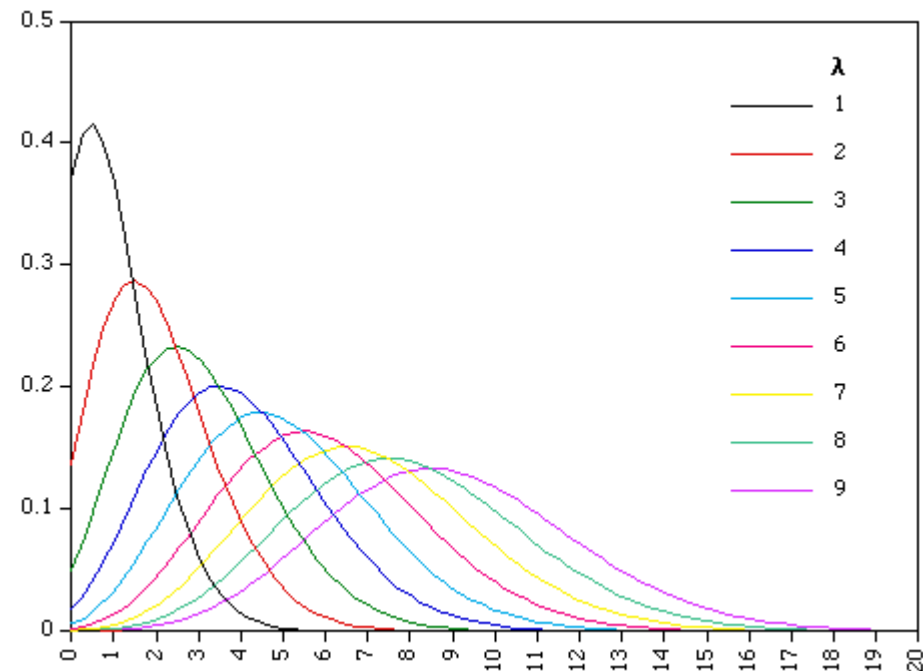
Poisson distribution

□ PDF:

$$f(x; \lambda) = \Pr(X = k) \\ = \frac{\lambda^k e^{-\lambda}}{k!}$$

□ $\lambda = \text{mean}$

□ $\lambda = \text{variance}$



Summary

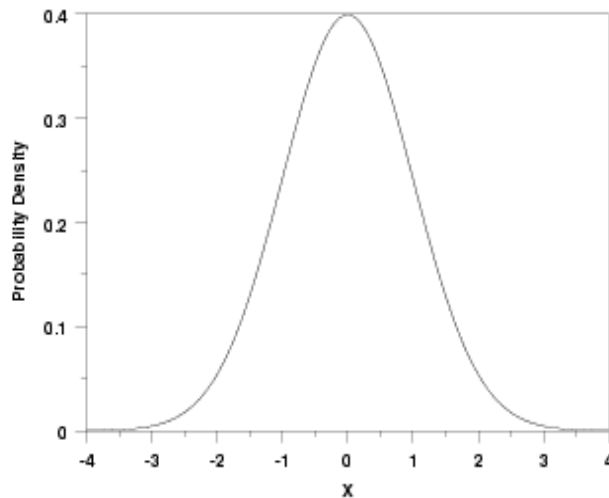
| Distribution | Mean | Variance |
|--------------|-----------|-----------------------|
| Normal | μ | σ^2 |
| Logistic | μ | $s^2 \frac{\pi^2}{3}$ |
| Binomial | np | $np(1-p)$ |
| Bernoulli | p | $p(1-p)$ |
| Poisson | λ | λ |

WHY DISTRIBUTION FUNCTIONS?



Probability density function

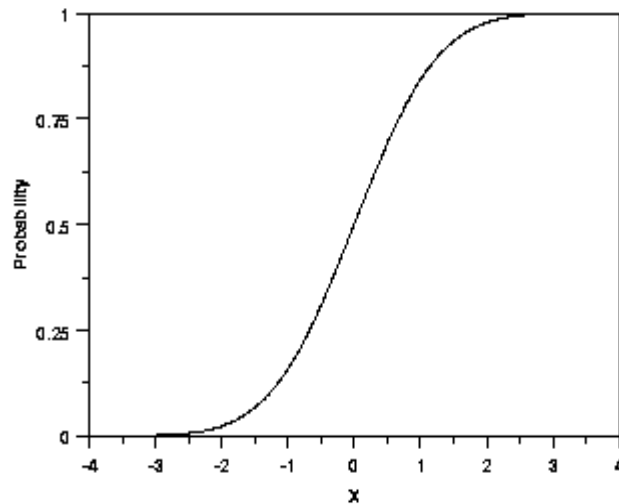
- Assumptions of a distribution allow you to estimate
- the probability of a random variable x



- A large $p(x)$: random variable x is close to μ
- Probability of exact value is zero
- Nonnegative numbers

Cumulative distribution function

- Cumulative distribution allows you to estimate whether random variable x is lower than or equal to value X



Probability distribution in R

| | |
|-----|---|
| “d” | returns the height of the probability density function |
| “p” | returns the cumulative density function |
| “q” | returns the inverse cumulative density function (quantiles) |
| “r” | returns randomly generated numbers |

Various distributions in R

- **Create normally distributed variable :**
 - ▣ `y<-rnorm(n=100)`
- **Create binomially distributed variable:**
 - ▣ `ybi <-rbinom(n=100,size=10,prob=.2)`
- **Create chisquared distributed variable:**
 - ▣ `ychi <- rchisq(n=100, df=2)`
- **Create f distributed variable:**
 - ▣ `yf <- rf(n=100, df1=2, df2=4)`

PLOT YOUR DATA



First data inspection

- Basic plots available in R
- More advanced functions in ggplot2 (we'll discuss later)

Types of plots

□ Barplot

- ▣ a barplot is especially useful for nominal or ordinal variables and shows the frequencies of the various categories

□ Histogram

- ▣ a histogram is useful to show the distribution of a interval or continuous variable. This illustrates whether the variable has a normal or other distribution

□ Boxplot

- ▣ a boxplot is useful to visually detect outliers
- ▣ especially appropriate with interval or continuous variables which is split by groups

R exercise (exe1_3.r)

- Plot (first create variable)
 - ▣ Plot one variable
 - ▣ Plot two variables
- Barplot
 - ▣ Nominal / ordinal/ multinomial variables
- Histogram
 - ▣ Distribution of interval/continuous variables
- Boxplot (first: create variables with certain distribution)
 - ▣ `new <- rnorm(nrow(df))` # create normal variable

Next lecture

- Writing functions in R