

## Arthur Charpentier

e-mail : [charpentier.arthur@uqam.ca](mailto:charpentier.arthur@uqam.ca)

url : <http://freakonometrics.hypotheses.org/>

Data Science pour l'Actuariat, Mars - Juin 2015

DataMining & R

avec Stéphane Tufféry

**A Brief Introduction To More Advanced R**

*“An expert is a man who has made all the mistakes which can be made, in a narrow field”* N. Bohr

## References

Wickam, H. [Advanced R](#). CRC Press, 2014.

(cf <http://adv-r.had.co.nz/>)

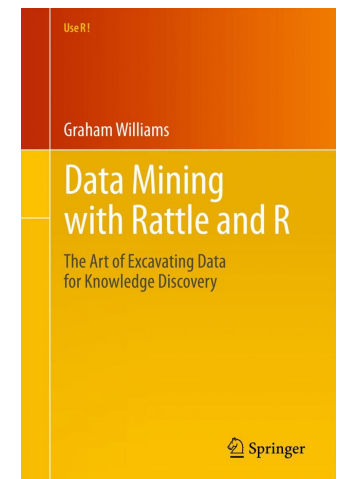
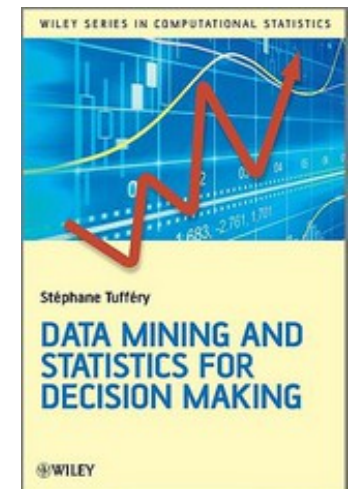
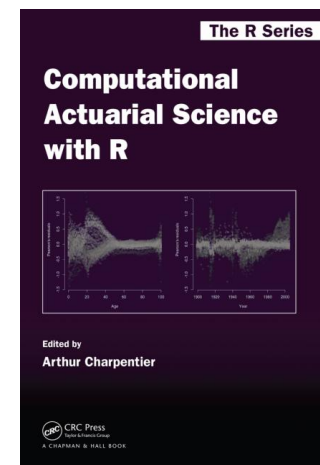
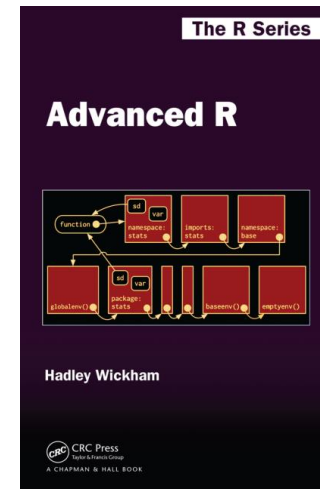
Tufféry, S. [Data Mining and Statistics for Decision Making](#). Wiley, 2013.

Charpentier, A. [Computational Actuarial Science with R](#). CRC Press. 2014.

Williams, G. [Data Mining with Rattle and R](#). Springer. 2011.

Zhao, Y. [R and Data Mining : Examples and Case Studies](#). 2013

(cf <http://cran.r-project.org/contrib/>)



# Additional Information, R in Insurance, 2015

## Amsterdam School of Economics

[Home](#) [Education](#) [Research](#) [News & Events](#) [Sections](#) [People](#) [About the ASE](#) [Alumni](#) [Contact](#)

### ▼ News & Events

- › Events
- › Archive

### ▼ R in Insurance

- › Talk proposal submission
- › Registration and fees
- › Programme
- › Travel and accommodation
- › Contact
- › Rob in Insurance

## R in Insurance 2015

Following two successful events at Cass in London, R in Insurance will travel across the Channel to Amsterdam in 2015. The third conference on R in Insurance will be held at the Amsterdam School of Economics in Amsterdam, The Netherlands, on Monday, 29 June 2015.

The intended audience of the conference includes both academics and practitioners who are interested in practical aspects of insurance and more specifically the applications of R in Insurance.

### Topics

This one-day conference will focus on all subdisciplines of actuarial science where professionals can use R, the lingua franca for statistical computation. Topics covered include risk management, reserving, pricing, loss modelling and the use of R in a production environment.

The central theme is how one can use R as a primary tool for insurance risk management, analysis and modelling.

# Additional Information, R in Insurance, 2015

## Amsterdam School of Economics

[Home](#) [Education](#) [Research](#) [News & Events](#) [Sections](#) [People](#) [About the ASE](#) [Alumni](#) [Contact](#)

### ▼ News & Events

- › Events
- › Archive
- › R in Insurance
- ▼ **Rob in Insurance**
  - › Programme
  - › Location
  - › Contact

## Rob in Insurance

After a distinguished career, Rob Kaas, Professor of Actuarial Statistics at the University of Amsterdam, will retire in the summer of 2015. He will be granted emeritus status at the University of Amsterdam, where he started working as early as 1969.

Among many other aspects, Prof. Kaas is author of *Modern Actuarial Risk Theory - Using R*, one of the main textbooks in Actuarial Science that has been translated into many languages, and of over 130 scholarly research papers in Actuarial Science; Program Director of the renowned Actuarial Science Education Programs at the University of Amsterdam; Managing Editor of *Insurance: Mathematics and Economics*, the leading journal in Actuarial Science; and organizer of the first IME conference in Amsterdam in 1997, starting the series of what are considered the major academic conferences in Actuarial Science today.



## R

R is an *interpreted language* where expressions are entered into the R console, and a program within the R system (called the interpreter) executes the code, unlike C but like Javascript.

```
1 > 2+3
2 [1] 5
```

R is an Object-Oriented Programming language. The idea is to create various **objects** that contain useful information, and that could be called by other functions (e.g. graphs).

```
1 > a <- 2+3
2 > print(a)
3 [1] 5
```

## R

A **package** is a related set of functions, including help files, and data files, that have been bundled together and is shared among the R community

```
1 > install.packages("quantreg", dependencies=TRUE)
2 > library(quantreg)
```

## S3 and S4 classes

*“Everything in S is an object. Every object in S has a class.”*

S3 is a primitive concept of classes, in R To define a class, use

```

1 > person3 <- function(name, age,
2   weight, height) {
3   crct <- list(name=name, age=age,
4     weight=weight, height=height)
5   class(crct) <- "person3"
6   return(crct)
7 }
8
9 > JohnDoe3 <- person3(name="John",
10   age=28, weight=76, height=182)
11
12 > JohnDoe3
13 $name
14 [1] "John"
15 $age
16 [1] 28
17 $weight
18 [1] 76
19 $height
20 [1] 182
21 attr(,"class")
22 [1] "person3"

```

To create a person, use

## S3 and S4 classes

If we want a function that returns the BMI (Body Mass Index), use

```
1 > BMI3 <- function(object, ...) {  
  return(object$weight*1e4/  
  object$height^2)}
```

e.g.

```
1 > BMI3(JohnDoe3)  
2 [1] 22.94409
```



## S3 and S4 classes

The analogous S4 version is

```

1 > setClass("person4",
  representation(name="
    character", age="numeric",
    weight="numeric", height="
    numeric"))
2 > JohnDoe4 <- new("person4", name
  ="John", age=28, weight=76,
  height=182)
1 > JohnDoe4
2 An object of class "person"
3 Slot "name":
4 [1] "John"
5 Slot "age":
6 [1] 28
7 Slot "weight":
8 [1] 76
9 Slot "height":
10 [1] 182

```

Here we have

## S3 and S4 classes

Observe that

```
1 > JohnDoe3$age
2 [1] 28
3 > JohnDoe4@age
4 [1] 28
```

```
1 > setGeneric("BMI4", function(
   object, separator) return(
   standardGeneric("BMI")))
2
3 > setMethod("BMI4", "person4",
4 function(object){return(object
   weight*1e4/objectheight^2)})
5
6 > BMI4(JohnDoe)
7 [1] 22.94409
```

To create our BMI function, use

## Numbers, in R

```
1 > x <- exp(1)
```

```
2 > x
```

```
3 [1] 2.718282
```

About large numbers,

```
4 > 1/0
```

```
5 [1] Inf
```

```
6 > .Machine$double.xmax
```

```
7 [1] 1.797693e+308
```

```
8 > 2e+307 < Inf
```

```
9 [1] TRUE
```

```
10 > 2e+308 < Inf
```

```
11 [1] FALSE
```

R has a [recycling rule](#)

```
11 > x <- c
    (100, 200, 300, 400, 500, 600, 700, 800)
```

```
12 > y <- c(1, 2, 3, 4)
```

```
13 > x+y
```

```
14 [1] 101 202 303 404 501 602 703
    804
```

```
4 > for(i in 1:2){
5 +   nom_var <- paste0("x", i)
6 +   assign(nom_var, rpois(5, 7))
7 + }
```

```
8 > x1
```

```
9 [1] 4 6 5 8 9
```

```
10 > x2
```

```
11 [1] 6 9 5 8 5
```

## Numbers, in R

About naming components of a vector

```

4 > x <- 1:6
5 > names(x) <- letters[1:6]
6 > x
7 a b c d e f
8 1 2 3 4 5 6
9 > x[2:4]
10 b c d
11 2 3 4
12 > x[c("b", "c", "d")]
13 b c d
14 2 3 4

```

It is possible to consider some **active building** variables

```

4 > x <- runif(4)
5 > x
6 [1] 0.8018263 0.1685260
   0.5900765 0.8230110
7 > x
8 [1] 0.8018263 0.1685260
   0.5900765 0.8230110
9 > library(pryr)
10 > x %<a-% runif(1)
11 > x
12 [1] 0.08434417
13 > x
14 [1] 0.9253761
15 > x <- 1:2
16 > x
17 [1] 0.1255551

```

## Numbers ( $\in \mathbb{R}$ ), in R

```

1 > (3/10-1/10)
2 [1] 0.2
3 > (3/10-1/10)==(7/10-5/10)
4 [1] FALSE
5 > (3/10-1/10)-(7/10-5/10)
6 [1] 2.775558e-17
7 > all.equal((3/10-1/10),(7/10-5/
8 [1] TRUE
9 > (eps<-Machine$double.eps)
10 [1] 2.220446e-16

```

```

11 > set.seed(1)
12 > U <- runif(20)
13 > U[1:4]
14 [1] 0.2655087 0.3721239
15 [1] 0.5728534 0.9082078
16 > options(digits = 3)
17 > U[1:4]
18 [1] 0.266 0.372 0.573 0.908
19 > options(digits = 22)
20 > U[1:4]
21 [1] 0.2655086631420999765396
22 [1] 0.3721238996367901563644
23 [1] 0.5728533633518964052200
24 [1] 0.9082077899947762489319

```

## Matrices, in R

```

1 > (N<-rpois(20,5))
2 [1] 9 3 6 3 4 4 1 4 8 4 5 5 5 3
   7 6 7 2 6 4
3 > (M=matrix(N,4,5))
4 [ ,1] [ ,2] [ ,3] [ ,4] [ ,5]
5 [1 ,] 9 4 8 5 7
6 [2 ,] 3 4 4 3 2
7 [3 ,] 6 1 5 7 6
8 [4 ,] 3 4 5 6 4
9 > dim(N)=c(4,5)
10 > N
11 [ ,1] [ ,2] [ ,3] [ ,4] [ ,5]
12 [1 ,] 9 4 8 5 7
13 [2 ,] 3 4 4 3 2
14 [3 ,] 6 1 5 7 6
15 [4 ,] 3 4 5 6 4
16 > M>.6
17 [ ,1] [ ,2] [ ,3] [ ,4] [ ,5]
18 [1 ,] FALSE FALSE TRUE TRUE TRUE
19 [2 ,] FALSE TRUE FALSE FALSE TRUE
20 [3 ,] FALSE TRUE FALSE TRUE FALSE
21 [4 ,] TRUE TRUE FALSE FALSE TRUE
22 > M[1]
23 [1] 9
24 > M[1,1]
25 [1] 9
26 > M[1,]
27 [1] 9 4 8 5 7
28 > M[,1]
29 [1] 9 3 6 3

```

## Matrices, in R

```

1 > u<-1:24
2 > dim(u)=c(6,4)
3 > u
4      [,1] [,2] [,3] [,4]
5 [1,]    1    7   13   19
6 [2,]    2    8   14   20
7 [3,]    3    9   15   21
8 [4,]    4   10   16   22
9 [5,]    5   11   17   23
10 [6,]    6   12   18   24
11 > str(u)
12 int [1:6, 1:4] 1 2 3 4 5 6 7 8
13      9 10 ...
13 > colnames(u)=letters[1:4]
14 > rownames(u)=LETTERS[10:15]

```

```

1 > u[c("K", "N"), ]
2      a    b    c    d
3 K 2    8   14   20
4 N 5   11   17   23
5 > u[c(1,2), ]
6      a    b    c    d
7 J 1    7   13   19
8 K 2    8   14   20
9 > u[c("K",5), ]
10 Error in u[c("K", 5), ] :
      subscript out of bounds

```

## Memory Issues, in R

R holds all the objects in 'memory', and only limited amount of memory can be used.

classical R message : **cannot allocate vector of size \_\_\_\_ MB**

big datasets are often larger than the size of the RAM that is available

on Windows, the limits are 2Gb and 4Gb for 32-bit and 64-bit respectively

```
1 > rm(list=ls())
2 > memory.size()
3 [1] 15.05
4 > memory.limit()
5 [1] 3583
6 > memory.limit(size=8000)
7 Error in memory.limit(size = 8000) :
8   don't be silly!: your machine has a 4Gb address limit
9 > memory.limit(size=4000)
10 [1] 4000
```



## Memory Issues, in R

```
1 > x3 <- 1:1e3
2 > x4 <- 1:1e4
3 > x5 <- 1:1e5
4 > x6 <- 1:1e6
```

```
1 > object.size(x3)
2 4024 bytes
3 > object.size(x4)
4 40024 bytes
5 > object.size(x5)
6 400024 bytes
7 > object.size(x6)
8 4000024 bytes
```

```
1 > z1 <- matrix(0,1e+06,3)
2 > z2 <- matrix(as.integer(0),1e
+06,3)
```

```
1 > object.size(z1)
2 24000112 bytes
3 > object.size(z2)
4 12000112 bytes
```

```
1 > z3 <- matrix(0,1e+07,30)
2 Error: cannot allocate vector of
size 2.2 Gb
```

## Memory Issues, in R

Alternative : matrices are stored on disk, rather than in RAM, and only elements needed are read from the disk, and cached in RAM

```

1 > z1 <- matrix(as.integer(5), 3,
2               4)
3 > object.size(z1)
3 248 bytes
4 > z1 <- matrix(as.integer(5),
5               30, 40)
6 > object.size(z1)
6 5000 bytes
7 > z1 <- matrix(as.integer(5),
8               300, 400)
9 > object.size(z1)
9 480200 bytes

```

```

1 > z2 <- big.matrix(3, 4, type="
2               integer", init=5)
3 > object.size(z2)
3 664 bytes
4 > z2 <- big.matrix(300, 400,
5               type="integer", init=5)
6 > object.size(z2)
6 664 bytes
7 > z2
8 An object of class "big.matrix"
9 Slot "address":
10 <pointer: 0x1aa0e220>

```

## Integers, in R

```
1 > (x_num=c(1,6,10))
2 [1] 1 6 10
3 > (x_int=c(1L,6L,10L))
4 [1] 1 6 10
5 > object.size(x_num)
6 72 bytes
7 > object.size(x_int)
8 56 bytes
9 > typeof(x_num)
10 [1] "double"
11 > typeof(x_int)
12 [1] "integer"
```

```
13 > is.integer(x_num)
14 [1] FALSE
15 > is.integer(x_int)
16 [1] TRUE
17 > str(x_num)
18 num [1:3] 1 6 10
19 > str(x_int)
20 int [1:3] 1 6 10

19 > c(1,c(2,c(3,c(4,5))))
20 [1] 1 2 3 4 5
```

## Factors, in R

```

1 > (x <- c("A", "A", "B", "B", "C"))
2 [1] "A" "A" "B" "B" "C"
3 > (x <- factor(x))
4 [1] A A B B C
5 Levels: A B C
6 > unclass(x)
7 [1] 1 1 2 2 3
8 attr(,"levels")
9 [1] "A" "B" "C"
10 > model.matrix(~0+x)
11   xA  xB  xC
12 1  1  0  0
13 2  1  0  0
14 3  0  1  0
15 4  0  1  0
16 5  0  0  1
17 > x[1]
18 [1] A
19 Levels: A B C
20 > x[1, drop=TRUE]
21 [1] A
22 Levels: A
23 > x <- factor(x, labels=c("Young",
24   "Adult", "Senior"))
25 > x
26 [1] Young Young Adult Adult
27   Senior
28 Levels: Young Adult Senior
29 > relevel(x, "Senior")
30 [1] Young Young Adult Adult
31   Senior
32 Levels: Senior Young Adult

```

## Factors, in R

```

1 > cut(U, breaks=2, labels=c("small
    ", "large"))
2 [1] small small large large
    small large large large
    large small small small
    large small large small
    large large small large
3 Levels: small large
4 > table(cut(U, breaks=c
    (0, .3, .8, 1), labels=c("small"
    , "medium", "large")))
5 small medium large
6      5      11      4

```

```

15 > x <- factor(c("b", "a", "b"))
16 > levels(x)
17 [1] "a" "b"
18 > x[3] = "c"
19 Warning Message :
20 In '[<- .factor' ('*tmp*', 3,
    value = "c") :
21   invalid factor level, NA
    generated
22 > x
23 [1] b      a      <NA>
24 Levels: a b

```

## Factors, in R

```
15 > X <- runif(20)
16 > group <- rep(c("A", "B"), c(8,
17   12))
18 > mean(X[group=="A"])
19 [1] 0.526534
20 > mean(X[group=="B"])
21 [1] 0.5185935
22 > tapply(X, group, mean)
23      A      B
24 0.5265340 0.5185935
25 > sapply(split(X, group), mean)
26      A      B
27 0.5265340 0.5185935
```

## Lists, in R

```

1 > x <- list(1:5, c(1,2,3,4,5),
2 + a="test", b=c(TRUE,FALSE),
3 + rpois(5,8))
4 > x
5 [[1]]
6 [1] 1 2 3 4 5
7 [[2]]
8 [1] 1 2 3 4 5
9 $a
10 [1] "test"
11 $b
12 [1] TRUE FALSE
13 [[5]]
14 [1] 11 12 5 6 3

```

```

15 > str(x)
16 List of 5
17 $ : int [1:5] 1 2 3 4 5
18 $ : num [1:5] 1 2 3 4 5
19 $ a: chr "test"
20 $ b: logi [1:2] TRUE FALSE
21 $ : int [1:5] 11 12 5 6 3

22 > names(x)
23 [1] "" "" "a" "b" ""

```

## Lists, in R

```

15 > x <- list(list(1:5, c
    (1,2,3,4,5)), a="test", b=c(
    TRUE,FALSE), rpois(5,8))
16 > x
17 [[1]]
18 [[1]][[1]]
19 [1] 1 2 3 4 5
20 [[1]][[2]]
21 [1] 1 2 3 4 5
22 $a
23 [1] "test"
24 $b
25 [1] TRUE FALSE
26 [[4]]
27 [1] 10 3 8 10 9

```

```

1 > is.list(x)
2 > is.recursive(x)
3 > x[[3]]
4 [1] TRUE FALSE
5 > x[["b"]]
6 [1] TRUE FALSE
7 > x[[1]]
8 [[1]]
9 [1] 1 2 3 4 5
10 [[2]]
11 [1] 1 2 3 4 5
12 > x[[1]][[1]]
13 [1] 1 2 3 4 5

```



## Lists, in R

```

1 > list(abc=1)$a
2 [1] 1
3 > list(abc=1,a=2)$a
4 [1] 2
5 > list(bac=1)$a
6 NULL
7 > list(abc=1,b=2)$a
8 [1] 1
9 > list(bac=1,a=2)$a
10 [1] 2
11 > list(bac=1)$a
12 NULL
13 > list(bac=1)$b
14 [1] 1

```

Lists are standard with S3 functions

```

1 > f <- function(x) { return(x*
   (1-x)) }
2 > optim.f <- optimize(f,
   interval=c(0, 1), maximum=
   TRUE)
3 > names(optim.f)
4 [1] "maximum" "objective"
5 > optim.f$maximum
6 [1] 0.5

```

## Characters and Strings, in R

```

1 > cities <- c("New York, NY", "
  Los Angeles, CA", "Boston,
  MA")
2 > substr(cities, nchar(cities)
  -1, nchar(cities))
3 [1] "NY" "CA" "MA"
4 > unlist(strsplit(cities, ", "))
  [seq(2,6,by=2)]
5 [1] "NY" "CA" "MA"

```

```

1 "Be carefull of 'quotes'"
2 'Be carefull of "quotes"'

```

```

15 > library(stringr)
16 > tweet <- "Emerging #climate
  change, environment and
  sustainability concerns for
  #actuaries Apr 9 #Toronto.
  Register TODAY http://bit.ly/CIAClimateForum"
17 > hash <- "#[a-zA-Z]{1,}"
18 > str_extract(tweet, hash)
19 [1] "#climate"
20 > str_extract_all(tweet, hash)
21 [[1]]
22 [1] "#climate" "#actuaries" "#
  Toronto"

```

## Characters and Strings, in R

```

1 > str_locate(tweet, hash)
2   start end
3 [1,]   10  17
4 > str_locate_all(tweet, hash)
5 [[1]]
6   start end
7 [1,]   10  17
8 [2,]   71  80
9 [3,]   88  95
10 > urls <- "http://([\\da-z
    \\. -]+)\\.([a-z\\.]{2,6})/[a
    -zA-Z0-9]{1,}"
11 > str_extract(tweet, urls)
12 [1] "http://bit.ly/
    CIAClimateForum"

```

```

1 > email="^([a-z0-9_\\.-]+)@([\\
    da-z\\. -]+)\\.([a-z
    \\.]{2,6})$"
2 > grep(pattern = email, x = "
    charpentier.arthur@uqam.ca")
3 [1] 1
4 > grepl(pattern = email, x = "
    charpentier.arthur@uqam.ca")
5 [1] TRUE
6 > str_detect(pattern = email,
    string=c("charpentier.
    arthur@uqam.ca", "
    @freakonometrics"))
7 [1] TRUE FALSE

```

## Characters and Strings, in R

Consider a simple sentence

```
1 > ex_sentence = "This is 1 simple sentence, just to play with, then  
we'll play with 4, and that will be more difficult "  
2 > ex_sentence  
3 [1] "This is 1 simple sentence, just to play with, then we'll play  
with 4, and that will be more difficult "
```

The first step is to create a **corpus**

```
1 > library(tm)  
2 > ex_corpus <- Corpus(VectorSource(ex_sentence))  
3 > ex_corpus  
4 <<VCorpus (documents: 1, metadata (corpus/indexed): 0/0)>>  
5 > inspect(ex_corpus)  
6 [[1]]  
7 <<PlainTextDocument (metadata: 7)>>  
8 This is 1 simple sentence, just to play with, then we'll play with 4,  
and that will be more difficult
```

## Characters and Strings, in R

Here we have one *document* in that corpus. We see if some documents do contain some specific words

```
1 > grep("hard", ex_sentence)
2 integer(0)
3 > grep("difficult", ex_sentence)
4 [1] 1
```

Since here we do not need the corpus structure (we have only one sentence) we can use more basic functions

```
1 > library(stringr)
2 > word(ex_sentence, 4)
3 [1] "simple"
```

## Characters and Strings, in R

To get the list of all the words

```

1 > word(ex_sentence, 1:20)
2 [1] "This" "is" "1" "simple" "sentence," "
   just" "to" "play" "with," "then" "we'll"
3 [12] "play" "with" "4," "and" "that" "
   will" "be" "more" "difficult"
4 > ex_words <- strsplit(ex_sentence, split=" ") [[1]]
5 > ex_words
6 [1] "This" "is" "1" "simple" "sentence," "
   just" "to" "play" "with," "then" "we'll"
7 [12] "play" "with" "4," "and" "that" "
   will" "be" "more" "difficult"
8 > grep(pattern="w", ex_words, value=TRUE)
9 [1] "with," "we'll" "with" "will"

```

## Characters and Strings, in R

We can count the occurrence of w's or i's in each word

```

1 > str_count(ex_words, "w")
2 [1] 0 0 0 0 0 0 0 0 1 0 1 0 1 0 0 0 1 0 0 0
3 > str_count(ex_words, "i")
4 [1] 1 1 0 1 0 0 0 0 1 0 0 0 1 0 0 0 1 0 0 2

```

or get all the words with a l

```

1 > grep(pattern="l", ex_words, value=TRUE)
2 [1] "simple"      "play"      "we'll"     "play"      "will"      "
      difficult "
3 > grep(pattern="l{2}", ex_words, value=TRUE)
4 [1] "we'll" "will"

```

or get all the words with an a or an i

```

1 > grep(pattern="[ai]", ex_words, value=TRUE)
2 [1] "This"      "is"      "simple"   "play"   "with,"   "
      play"

```

## Characters and Strings, in R

or a punctuation symbol

```
1 > grep(pattern=" [[:punct:]] ", ex_words, value=TRUE)
2 [1] "sentence , " "with , " "we' ll " "4 , "
```

It is possible, here, to create some WordCloud, e.g.

```
1 > require(wordcloud)
2 > wordcloud(ex_corpus)
3 > cols<-colorRampPalette(c("lightgrey",
  "blue"))(40)
4 > wordcloud(words = ex_corpus, max.
  words = 40, random.order=FALSE,
  scale = c(5, 0.5), colors=cols)
```





## Characters and Strings, in R

The corpus can be used to generate a list of words along with counts of their occurrence.

```

1 > tdm <- TermDocumentMatrix(ex_
  corpus)
2 > inspect(tdm)
3 <<TermDocumentMatrix (terms: 14,
  documents: 1)>>
4 Non-/sparse entries: 14/0
5 Sparsity           : 0%
6 Maximal term length: 9
7 Weighting          : term
  frequency (tf)

```

	Docs
1	
2	Terms 1
3	and 1
4	difficult 1
5	just 1
6	more 1
7	play 2
8	sentence , 1
9	simple 1
10	that 1
11	then 1
12	this 1
13	we 'll 1
14	will 1
15	with 1
16	with , 1

## Characters and Strings, in R

Note that the Corpus should be cleaned. This involves the following steps :

- convert all text to lowercase
- expand all contractions
- remove all punctuation
- remove all *noise words*

We start with

```
1 > inspect(ex_corpus)
2 <<VCorpus (documents: 1, metadata (corpus/indexed): 0/0)>>
3
4 [[1]]
5 <<PlainTextDocument (metadata: 7)>>
6 This is 1 simple sentence, just to play with, then we'll play with 4,
   and that will be more difficult
```

## Characters and Strings, in R

The first step might be to fix contractions

```
1 > fix_contractions <- function(doc) {
2 +   doc <- gsub("won't", "will not", doc)
3 +   doc <- gsub("n't", " not", doc)
4 +   doc <- gsub("'ll", " will", doc)
5 +   doc <- gsub("'re", " are", doc)
6 +   doc <- gsub("'ve", " have", doc)
7 +   doc <- gsub("'m", " am", doc)
8 +   doc <- gsub("'s", "", doc)
9 +   return(doc)
10 + }
11 > ex_corpus <- tm_map(ex_corpus, fix_contractions)
12 > inspect(ex_corpus)
13 [[1]]
14 [1] This is 1 simple sentence, just to play with, then we will play
    with 4, and that will be more difficult
```

## Characters and Strings, in R

Then we can remove numbers

```
1 > ex_corpus <- tm_map(ex_corpus, removeNumbers)
2 > inspect(ex_corpus)
3 [[1]]
4 [1] This is simple sentence, just to play with, then we will play
    with , and that will be more difficult
```

as well as punctuation

```
1 > gsub("[:punct:]", "", ex_sentence)
2 [1] "This is 1 simple sentence just to play with then well play with
    4 and that will be more difficult"
3 > ex_corpus <- tm_map(ex_corpus, gsub, pattern="[:punct:]",
    replacement = " ")
4 > inspect(ex_corpus)
5 [[1]]
6 [1] This is simple sentence just to play with then we will play
    with and that will be more difficult
```

## Characters and Strings, in R

Then, we usually remove stop words,

```

1 > stopwords("en")[sample(1:length(stopwords("en")),size=10)]
2 [1] "can't" "could" "because" "i've" "there's" "who's"
   "for" "couldn't" "we've" "him"
3 > ex_corpus <- tm_map(ex_corpus, removeWords, words=stopwords("en"))
4 > inspect(ex_corpus)
5 [[1]]
6 [1] This simple sentence just play will play will
   difficult

```

We should also convert all words to lower case

```

1 > ex_corpus <- tm_map(ex_corpus, tolower)
2 > inspect(ex_corpus)
3 [[1]]
4 [1] this simple sentence just play will play will
   difficult

```

## Characters and Strings, in R

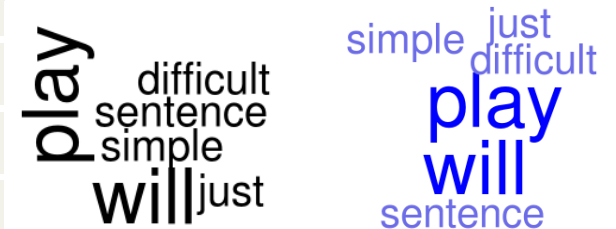
And finally, we stem the text

```
1 > library(SnowballC)
2 > ex_corpus <- tm_map(ex_corpus, stemDocument)
3 > inspect(ex_corpus)
4 [[1]]
5 [1] this simple sentence just play will play will
    difficult
```

## Characters and Strings, in R

We now have a clean list of words, it is possible to create some WordCloud

```
1 > wordcloud(ex_corpus[[1]])  
2 > wordcloud(words = ex_corpus[[1]], max.words =  
40, random.order=FALSE, scale = c(5, 0.5),  
colors=cols)
```



## Dates, in R

```
1 > (some.dates <- as.Date(c("16/10/12", "19/11/12"),
2   format="%d/%m/%y"))
3 [1] "2012-10-16" "2012-11-19"
4 > (sequence.date <- seq(from=some.dates[1], to=some.
5   dates[2], by=7))
6 [1] "2012-10-16" "2012-10-23" "2012-10-30" "2012-11-06"
7   "2012-11-13"
8 > format(sequence.date, "%b")
9 [1] "oct" "oct" "oct" "nov" "nov"
10 > weekdays(some.dates)
11 [1] "Tuesday" "Monday"
12 > Sys.setlocale("LC_TIME", "fr_FR")
13 [1] "fr_FR"
14 > weekdays(some.dates)
15 [1] "Mardi" "Lundi"
```



## Symbolic Expressions, in R

Consider a regression model,  $Y_i = \beta_0 + \beta_1 X_{1,i} + \beta_2 X_{2,i} + \beta_3 X_{3,i} + \varepsilon_i$ . The code to fit such a model is based on

```
1 > fit <- lm(formula = y ~ x1 + x2 + x3, data=df)
```

In a formula, + stands for inclusion (not for summation), and - for exclusion.

To illustrate the use of categorical variables, consider

```
1 > set.seed(123)
```

```
2 > df <- data.frame(Y=rnorm(50), X1=as.factor(sample(LETTERS[1:4], size
  =50, replace=TRUE)), X2=as.factor(sample(1:3, size=50, replace=TRUE)
  ))
```

```
3 > tail(df, 3)
```

```
4      Y X1 X2
5 48 -0.557 B  2
6 49  0.950 C  2
7 50 -0.498 A  3
```

## Symbolic Expressions, in R

```
1 > reg <- lm(Y~X1+X2, data=df)
```

```
2 > model.matrix(reg)[47:50,]
```

```
3 (Intercept) X1B X1C X1D X22 X23
4 47          1  0  0  0  1  0
5 48          1  0  0  0  1  0
6 49          1  0  0  0  0  0
7 50          1  0  1  0  1  0
```

```
1 > reg <- lm(Y~X1*X2, data=df)
```

```
2 > model.matrix(reg)[47:50,]
```

```
3 (Intercept) X1B X1C X1D X22 X23 X1B:X22 X1C:X22 X1D:X22 X1B:X23
4 47          1  1  0  0  0  1          0          0          0          1
5 48          1  1  0  0  1  0          1          0          0          0
6 49          1  0  1  0  1  0          0          1          0          0
7 50          1  0  0  0  0  1          0          0          0          0
```

## Functions, in R

```

1 > factorial
2 function (x)
3 gamma(x + 1)
4 <bytecode: 0x1708aa7c>
5 <environment: namespace:base>

```

```

1 > gamma
2 function (x) .Primitive("gamma")

```

```

1 > x<-rexp(6)
2 > sum(x)
3 [1] 5.553364
4 > .Primitive("sum")(x)
5 [1] 5.553364
6 > cppFunction('double sum_C(
      NumericVector x) {
7 +   int n = x.size();
8 +   double total = 0;
9 +   for(int i = 0; i < n; ++i) {
10 +     total += x[i];
11 +   }
12 +   return total;
13 + }')
14 > sum_C(x)
15 [1] 5.553364

```

# Functions, in R

```

1 > f <- function(x) x^2
2 > f
3 function(x) x^2
4 > formals(f)
5 $x
6 > body(f)
7 x^2

```

```

1 > x <- 10
2 > f <- function() x<-5
3 > f()
4 > x
5 [1] 10

```

```

1 > f <- function() {
2 +   x<-5
3 +   return(x)
4 + }
5 > f()
6 [1] 5
7 > x
8 [1] 10
9 > f <- function() {
10 +   x<<-5
11 +   return(x)
12 + }
13 > f()
14 [1] 5
15 > x
16 [1] 5

```

## Functions, in R

```

1 > x=function(y) y/2
2 > x
3 function(y) y/2
4 > x <- 5
5 > x(x)
6 Error: could not find function "
  x"
7
8 > x=function(y) y/2
9 > y=function() {
10 +   x <- 10
11 +   x(x)
12 + }
13 > y()
14 [1] 5

```

```

1 > names_list <- function(...) {
2 +   names(list(...))
3 + }
4 > names_list(a=5,b=7)
5 [1] "a" "b"

```

Replacement functions act like they modify their arguments in place

```

1 > 'second<-' <- function(x,
2 +   value) {
3 +   return(x)
4 + }
5 > x <- 1:8
6 > second(x) <- 5
7 > x
8 [1] 1 5 3 4 5 6 7 8

```

## Functions, in R

```
1 > f <- function(x,m=0,s=1){
2 + H<-function(t) 1-pnorm(t,m,s)
3 + integral<-integrate(H,lower=x,upper=Inf)$value
4 + res<-H(x)/integral
5 + return(res)
6 + }
7 > f(0:1)
8 [1] 1.2533141 0.3976897
9 Warning :
10 In if (is.finite(lower)) { :
11 the condition has length > 1 and only the first
    element will be used
12 > Vectorize(f)(0:1)
13 [1] 1.253314 1.904271
14 > sapply(0:1,"f")
15 [1] 1.253314 1.904271
```

## Functions, in R

```
1 > fibonacci <- function(n) {  
2 +   if(n<2) return(1)  
3 +   fibonacci(n-2)+fibonacci(n-1)  
4 + }  
5 > fibonacci(20)  
6 [1] 10946  
7 > system.time(fibonacci(30))  
8   user  system elapsed  
9 3.687    0.000    3.719
```

## Functions, in R

It is possible to use Memoisation : all previous inputs are stored... tradeoff speed and memory

```
1 > library(memoise)
2 > fibonacci <- memoise(function(n) {
3 +   if(n<2) return(1)
4 +   fibonacci(n-2)+fibonacci(n-1)
5 + })
6 > system.time(fibonacci(30))
7   user  system elapsed
8 0.004   0.000   0.004
```



## Functions, in R

```

1 > binorm <- function(x1, x2, r=0){
2 + exp(-(x1^2+x2^2-2*r*x1*x2)/(2*(1-r^2)))/(2*pi*sqrt(1-r^2))
3 + }

```

```

1 > u <- seq(-2,2)
2 > binorm(u, u)
3 [1] 0.00291 0.05854 0.15915 0.05854 0.00291
4 > outer(u, u, binorm)
5           [,1]  [,2]  [,3]  [,4]  [,5]
6 [1,] 0.00291 0.0130 0.0215 0.0130 0.00291
7 [2,] 0.01306 0.0585 0.0965 0.0585 0.01306
8 [3,] 0.02153 0.0965 0.1591 0.0965 0.02153
9 [4,] 0.01306 0.0585 0.0965 0.0585 0.01306
10 [5,] 0.00291 0.0130 0.0215 0.0130 0.00291

```

```
11 > (uv<-expand.grid(u,u))
12 > matrix(binorm(uv$Var1,uv$Var2),5,5)

1 > "%pm%" <- function(x,s) x + c(qnorm(.05),qnorm(.95))*s
2 > 100 %pm% 10
3 [1] 83.55146 116.44854

1 > f(0:1)
2 [1] 1.2533141 0.3976897
3 Warning :
4 In if (is.finite(lower)) { :
5 the condition has length > 1 and only the first element will be used
6 > Vectorize(f)(0:1)
7 [1] 1.253314 1.904271
8 > sapply(0:1,"f")
9 [1] 1.253314 1.904271
```

## Functions, in R

```
1 > f<-function(x) dlnorm(x)
2 > integrate(f,0,Inf)
3 1 with absolute error < 2.5e-07
4 > integrate(f,0,1e5)
5 1.819813e-05 with absolute error < 3.6e-05
6 > integrate(f,0,1e3)$value+integrate(f,1e3,1e5)$value
7 [1] 1
```

## Functions, in R

```
1 > set.seed(1)
2 > u <- runif(1)
3 > if(u>.5) {"greater than 50%"} else {"smaller than 50%"}
4 [1] "smaller than 50%"
5 > ifelse(u>.5,("greater than 50%"),("smaller than 50%"))
6 [1] "smaller than 50%"
7 > u
8 [1] 0.2655087
9
10 > v_x <- runif(1e5)
11 > sqrt_x <- NULL
12 > system.time(for(x in v_x) sqrt_x <- c(sqrt_x, sqrt(x)))
13      user      system    elapsed
14 23.774      0.224    24.071
```

## Functions, in R

```

1 > sqrt_x <- numeric(length(v_x))
2 > system.time(for(x in seq_along(v_x)) sqrt_x[i] <- sqrt(x[i]))
3           user      system      elapsed
4           1.320      0.000      1.328
5 >
6 > system.time(Vectorize(sqrt)(v_x))
7           user      system      elapsed
8           0.008      0.000      0.009
9 >
10 > sqrt_x <- sapply(v_x, sqrt)
11 > system.time(unlist(lapply(v_x, sqrt)))
12           user      system      elapsed
13           0.300      0.000      0.299

```

## Functions, in R

```
1 > library(parallel)
2 > (allcores <- detectCores(all.tests=TRUE))
3 [1] 4
4 > system.time(unlist(mclapply(v_x, sqrt, mc.cores=4)))
5      user      system    elapsed
6 0.396    0.224    0.362
```

## Functions, in R

Write a function to generate random numbers drawn from a **compound Poisson**,  $X = Y_1 + \dots + Y_N$  with  $N \sim \mathcal{P}(\lambda)$  and  $Y_i$  i.i.d.  $\mathcal{E}(\alpha)$ .

```
1 > rN.Poisson <- function(n) rpois(n,5)
```

```
2 > rX.Exponential <- function(n) rexp(n,2)
```

```
1 > rcpd1 <- function(n,rN=rN.Poisson ,rX=rX.Exponential){
```

```
2   V <- rep(0,n)
```

```
3   for(i in 1:n){
```

```
4     N <- rN(1)
```

```
5     if(N>0){V[i] <- sum(rX(N))}
```

```
6   }
```

```
7   return(V)}
```

## Functions, in R

```
1 > rcpd1 <- function(n, rN=rN.Poisson, rX=rX.Exponential) {  
2   V <- rep(0, n)  
3   for (i in 1:n) V[i] <- sum(rX(rN(1)))  
4   return(V)}
```

```
1 > rcpd2 <- function(n, rN=rN.Poisson, rX=rX.Exponential) {  
2   N <- rN(n)  
3   X <- rX(sum(N))  
4   I <- factor(rep(1:n, N), levels=1:n)  
5   return(as.numeric(xtabs(X ~ I)))}
```



## Functions, in R

```
1 > rcpd3 <- function(n, rN=rN.Poisson, rX=rX.Exponential) {  
2   N <- rN(n)  
3   X <- rX(sum(N))  
4   I <- factor(rep(1:n, N), levels=1:n)  
5   V <- tapply(X, I, sum)  
6   V[is.na(V)] <- 0  
7   return(as.numeric(V)) }
```

## Functions, in R

```
1 > rcpd4 <- function(n, rN=rN.Poisson, rX=rX.Exponential) {  
2   return(sapply(rN(n), function(x) sum(rX(x))))}
```

```
1 > rcpd5 <- function(n, rN=rN.Poisson, rX=rX.Exponential) {  
2   return(sapply(Vectorize(rX)(rN(n)), sum))}
```

```
1 > rcpd6 <- function(n, rN=rN.Poisson, rX=rX.Exponential) {  
2   return(unlist(lapply(lapply(t(rN(n)), rX), sum)))}
```

## Functions, in R

```

1 > n <- 100
2 > library(microbenchmark)
3 > options(digits=1)
4 > microbenchmark(rcpd1(n), rcpd2(n), rcpd3(n), rcpd4(n), rcpd5(n), rcpd6(n)
   ))
5 Unit: microseconds
6      expr   min    lq  mean  median    uq   max  neval  cld
7  rcpd1(n)  756   794   875   829   884  1624   100   c
8  rcpd2(n) 1292  1394  1656  1474  1578  6799   100   d
9  rcpd3(n)  629   677   741   707   756  2079   100  bc
10 rcpd4(n)  482   515   595   540   561  5095   100  ab
11 rcpd5(n)  613   663   864   694   744  9943   100   c
12 rcpd6(n)  426   453   496   469   492  1020   100   a

```

## Functions, in R

```

1 > n <- 50
2 > X <- matrix(rnorm(m * n, mean = 10, sd = 3), nrow = m)
3 > group <- rep(c("A", "B"), each = n / 2)
4 >
5 > system.time(for(i in 1:m) t.test(X[i, ] ~ group)$stat)
6 utilisateur      syst me      coul
7      1.028      0.000      1.030
8 > system.time(for(i in 1:m) t.test(X[i, group == "A"], X[i, group ==
9      "B"])$stat)
9 utilisateur      syst me      coul
10      0.224      0.000      0.222

```

## Functions, in R

```
1 > f <- function(x) log(x)
2 > f("x")
3 Error in log(x) : non-numeric argument to mathematical function
4 > try(f("x"))
5 Error in log(x) : non-numeric argument to mathematical function
6 > inherits(try(f("x"), silent=TRUE), "try-error")
7 [1] TRUE
8 > x=2:4
9 > a=0
10 > try(a<-f("x"), silent=TRUE)
11 > a
12 [1] 0
13 > try(a<-f(x), silent=TRUE)
14 > a
15 [1] 0.6931472 1.0986123 1.3862944
```

## Functions, in R

```

1 > power <- function(exponent) {
2   +   function(x) {
3     +     x ^ exponent
4     +   }
5   + }
6 > square <- power(2)
7 > square(4)
8 [1] 16
9 > cube <- power(3)
10 > cube(4)
11 [1] 64
12 > cube
13 function(x) {
14   +   x ^ exponent
15   + }
16 <environment: 0x9c810a0>

```

```

1 > x=1:10
2 > g=function(f) f(x)
3 > g(mean)
4 [1] 5.5

```

## Progress Bar, in R

```

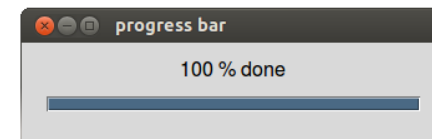
1 > v_x <- runif(2e4)
2 > sqrt_x <- NULL
3 > pb <- txtProgressBar(min = 0,
  max = total, style = 3)
4 | 0%
5 > for(i in seq_along(v_x)){
6 + sqrt_x <- c(sqrt_x, sqrt(x[i
  ]))
7 + if(i %% 1e3==0)
  setTxtProgressBar(pb, i%%1
  e3)
8 + }
9 |=====| 100%

```

```

1 > library(tcltk)
2 > total <- 20
3 > pb <- tkProgressBar(title = "
  progress bar", min = 0, max =
  total, width = 300)
4 > for(i in seq_along(v_x)){
5 + sqrt_x <- c(sqrt_x, sqrt(x[i]
  ))
6 + if(i %% 1e3==0)
  setTkProgressBar(pb, i%%1e3
  , label=paste(round(i%%1e3
  /total*100, 0), "% done"))
7 + }

```



# Data Frames, in R

```

1 > df <- data.frame(x=1:3,y=letters[1:3])
2 > str(df)
3 'data.frame': 3 obs. of 2 variables:
4 $ x: int 1 2 3
5 $ y: Factor w/ 3 levels "a","b","c": 1 2 3
6 > class(df)
7 [1] "data.frame"

```

```

1 > cbind(df, z=9:7)
2   x y z
3 1 1 a 9
4 2 2 b 8
5 3 3 c 7

```

```

1 > df$z <- 5:3
2 > df
3   x y z
4 1 1 a 5
5 2 2 b 4
6 3 3 c 3

```



## Data Frames, in R

```

1 > cbind(df, z=9:7)
2   x y z z
3   1 1 a 5 9
4   2 2 b 4 8
5   3 3 c 3 7
6 > df$z<-5:3
7 > df
8   x y z
9   1 1 a 5
10  2 2 b 4
11  3 3 c 3

```

```

1 > df <- data.frame(x=1:3,y=
2   letters[1:3],xy=19:17)
3 > df[1]
4   x
5   1 1
6   2 2
7   3 3
8 > df[,1,drop=FALSE]
9   x
10  1 1
11  2 2

```

## Data Frames, in R

```

1 > df[,1,drop=TRUE]
2 [1] 1 2 3
3 > df[[1]]
4 [1] 1 2 3
5 > df[[1]]
6 [1] 1 2 3
7 > df$x
8 [1] 1 2 3
9 > df[, "x"]
10 [1] 1 2 3
11 > df[["x"]]
12 [1] 1 2 3
13 > df[["x",exact=FALSE]]
14 [1] 1 2 3

1 > set.seed(1)
2 > df[sample(nrow(df)),]
3   x y xy
4 1 1 a 19
5 3 3 c 17
6 2 2 b 18
7 > set.seed(1)
8 > df[sample(nrow(df),nrow(df)*2,
9           replace=TRUE),]
9   x y xy
10 1 1 a 19
11 2 2 b 18
12 2.1 2 b 18
13 3 3 c 17
14 1.1 1 a 19
15 3.1 3 c 17

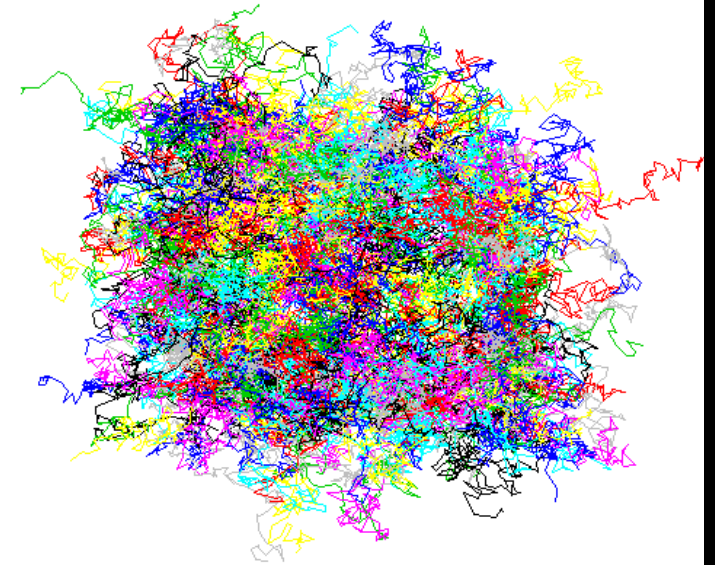
```

## Data Frames, in R

```
1 > rm(list=ls())
2 > library(RCurl)
3 > dropbox_ldf <- "https://dl.dropboxusercontent.com/s/l7rolinwojn37e2
  /ldf_json_2.RData"
4 > dropbox_df <- "https://dl.dropboxusercontent.com/s/wzr19v9pyl7ah0j
  /df_json_2.RData"
5 > dropbox_dt <- "https://dl.dropboxusercontent.com/s/nlwi1mmsxr5f23k
  /dt_json_2.RData"
6 > source_https <- function(loc, ...) {
7 +   curl=getCurlHandle()
8 +   curlSetOpt(cookiejar="cookies.txt", useragent="Mozilla/5.0",
9 +   followlocation=TRUE)
9 +   tmp <- getURLContent(loc, .opts=list(ssl.verifypeer=FALSE), curl=
10 +   curl)
10 +   fch <- source(tmp)
11 +   fch
12 + }
```

## Data Frames, in R

```
1 > source_https(dropbox_ldf)
2 > load("df_json_2.RData")
3 > tail(df)
4           Pers_Id Traj_Id      lat      lon
5 159996158      10000 2000091 3.860666 -2.6781690
6 159996159      10000 2000091 3.983418 -2.2454256
7 159996160      10000 2000091 3.929773 -2.0908522
8 159996161      10000 2000091 3.967067 -1.8922986
9 159996162      10000 2000091 3.881188 -2.1948032
10 159996163      10000 2000091 2.989197  0.0869032
```



## Data Frames, in R

```

1 > n <- nrow(df)
2 > system.time(df$first<-c(1,df$
3   Traj_Id[2:n] !=
4   df$Traj_Id[1:(n-1)]))
5   user      system    elapsed
6   3.12      1.23      4.37
7 > system.time(df$last<-c(df$Traj
8   _Id[2:n] !=
9   df$Traj_Id[1:(n-1)],1))
10  user      system    elapsed
11  2.90      6.23      31.58
12 > object.size(df)
13 6399847720 bytes

1 > lat_0=0
2 > lon_0=0
3 > system.time(df$test <- (((lat_
4   0-df$lat)^2+(lon_0-df$lon)^2)
5   <=1)&(df$last==1))
6 Error: cannot allocate vector of
7   size 1.2 Gb

1 > df$first <- NULL
2 > df$last <- NULL
3 > object.size(df)
4 3839908904 bytes

```

## Data Frames, in R

```

1 > system.time(df$test <- ((
2 (lat_0-df$lat)^2+(lon_0-df$lon)^2)
3 <=1)
4 &(c(df$Traj_Id[2:n] !=df$Traj_Id
5 [1:(n-1)],1)==1))
6 user      system      elapsed
7 9.39      11.10      41.17
8 > system.time(list_Traj <-
9 unique(df$Traj_Id[df$test==
10 TRUE]))
11 user      system      elapsed
12 0.72      0.25      0.98

```

```

1 > system.time(base <- df[(df$
2 Traj_Id %in% list_Traj)&(c
3 (1,df$Traj_Id[2:n] !=df$Traj_
4 Id[1:(n-1)])==1),c("lat",
5 "lon")])
6 user      system      elapsed
7 11.7      2.7      14.4
8 > head(base)
9           lat      lon
10 556 -0.9597891 2.469243
11 2866 -0.9597891 2.469243
12 4321 -0.9597891 2.469243
13 5677 -0.9597891 2.469243
14 9403 -0.9597891 2.469243
15 10432 -0.9597891 2.469243
16 > nrow(base)
17 [1] 63453

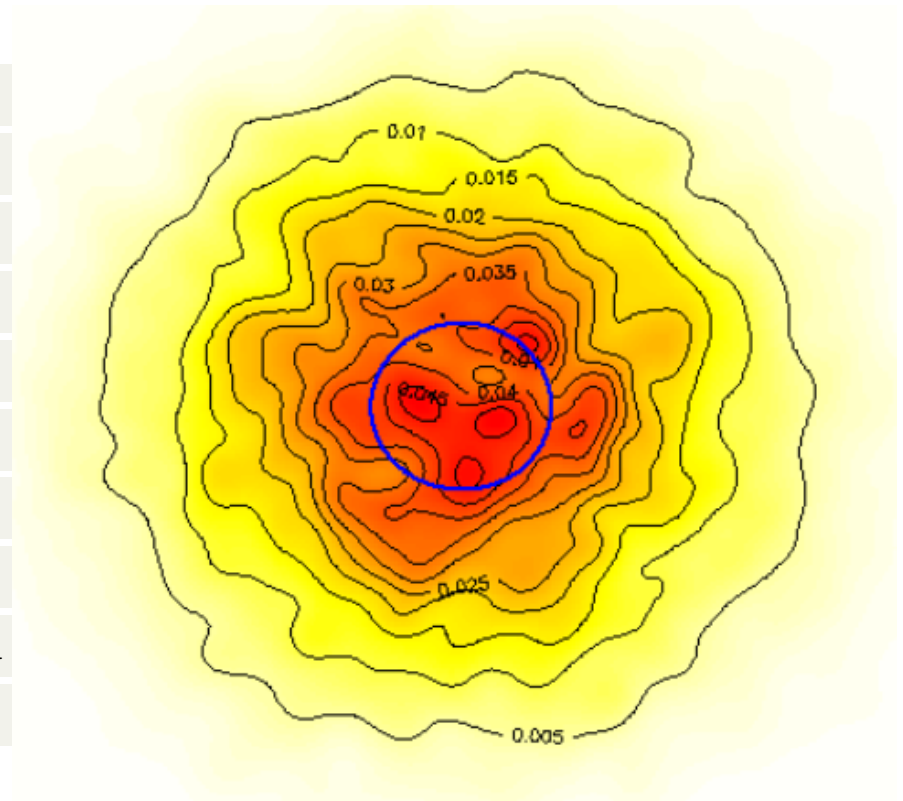
```

## Data Frames, in R

```

1 > X <- base[,c("lon", "lat")]
2 > library(KernSmooth)
3 > kde2d <- bkde2D(X, bandwidth=c(bw.ucv(
  X[,1]), bw.ucv(X[,2])), gridsize = c
  (251L, 251L))
4 > image(x=kde2d$x1, y=kde2d$x2, z=kde2d$
  fhat, col=
5 rev(heat.colors(100)))
6 > contour(x=kde2d$x1, y=kde2d$x2, z=kde2d
  $fhat, add=TRUE)

```



## Databases, in R

Consider the `gapminderDataFiveYear.txt` dataset, inspired from [stat545-ubc](#)

```

1 > gdf <- read.delim("gapminderDataFiveYear.txt")
2 > head(gdf,4)
3      country year      pop continent lifeExp gdpPercap
4 1 Afghanistan 1952  8425333      Asia   28.801   779.4453
5 2 Afghanistan 1957  9240934      Asia   30.332   820.8530
6 3 Afghanistan 1962 10267083      Asia   31.997   853.1007
7 4 Afghanistan 1967 11537966      Asia   34.020   836.1971
8 > str(gdf)
9 'data.frame': 1704 obs. of 6 variables:
10 $ country : Factor w/ 142 levels "Afghanistan",...: 1 1 1 1 1 1 ...
11 $ year    : int  1952 1957 1962 1967 1972 1977 1982 1987 1992 ...
12 $ pop     : num  8425333 9240934 10267083 11537966 13079460 ...
13 $ continent: Factor w/ 5 levels "Africa","Americas",...: 3 3 3 3 ...
14 $ lifeExp : num  28.8 30.3 32 34 36.1 ...
15 $ gdpPercap: num  779 821 853 836 740 ...

```



## Databases, in R

One can consider `tbl_df()` to get an improved data frame (called `local dataframe`)

```

1 > gtbl <- tbl_df(gdf)
2 > gtbl
3 Source: local data frame [1,704 x 6]
4
5   country year      pop continent lifeExp gdpPercap
6 1  Afghanistan 1952  8425333      Asia    28.801   779.4453
7 2  Afghanistan 1957  9240934      Asia    30.332   820.8530
8 3  Afghanistan 1962 10267083      Asia    31.997   853.1007
9 4  Afghanistan 1967 11537966      Asia    34.020   836.1971
10 5  Afghanistan 1972 13079460      Asia    36.088   739.9811
11 6  Afghanistan 1977 14880372      Asia    38.438   786.1134
12 ..          ...      ...          ...          ...          ...

```

## Databases, in R

For instance, to reproduce

```
1 > subset(gdf, lifeExp < 30)
2           country year      pop continent  lifeExp  gdpPercap
3 1  Afghanistan 1952 8425333      Asia   28.801   779.4453
4 1293      Rwanda 1992 7290203      Africa 23.599   737.0686
```

use

```
1 > filter(gtbl, lifeExp < 30)
2 Source: local data frame [2 x 6]
3
4           country year      pop continent  lifeExp  gdpPercap
5 1 Afghanistan 1952 8425333      Asia   28.801   779.4453
6 2      Rwanda 1992 7290203      Africa 23.599   737.0686
```

## Databases, in R

The `%>%` operator can be used to generate (conveniently) datasets

```
1 > gtbl %>%  
2 +   filter(country == "Italy") %>%  
3 +   select(year, lifeExp)  
4 Source: local data frame [12 x 2]  
5  
6   year lifeExp  
7 1  1952  65.940  
8 2  1957  67.810  
9 3  1962  69.240  
  
17 11 2002  80.240  
18 12 2007  80.546
```

which is (almost) the same as

```
19 > gdf[gdf$country == "Italy", c("year", "lifeExp")]
```

## Local Data Frames, in R

```

1 > load("ldf_json_2.RData")
2 > system.time( ldepart <- ldf
  %>% group_by(Traj_Id) %>%
3 + summarise( first_lat=head(lat
  ,1) ,
4 first_lon=head(lon,1) ) )
5   user      system    elapsed
6  60.82      1.46      80.54

```

```

1 > system.time( larrive <- ldf
  %>% group_by(Traj_Id) %>%
2 + summarise( last_lat=tail(lat,1)
  , last_lon=tail(lon,1) ) )
3   user      system    elapsed
4  60.81      0.31      62.15
5 > lat_0=0
6 > lon_0=0
7 > system.time( system.time(
  larrive <- mutate(larrive ,
  dist=(last_lat-lat_0)2+(last
  _lon-lon_0)2 ) )
8   user      system    elapsed
9  0.08      0.00      0.09

```

## Local Data Frames, in R

```

1 > system.time( lfin <- filter(larrive , dist <=1) )
2   user      system    elapsed
3   0.05      0.00      0.04
4 > system.time( lbase <- left_join(lfin , ldepart) )
5 Joining by: "Traj_Id"
6   user      system    elapsed
7   0.53      0.05      0.66
8 > head(lbase)
9 Source: local data frame [63,453 x 6]
10
11   Traj_Id    last_lat    last_lon    dist    first_lat    first_lon
12 1         8    0.41374639  0.491248980  0.412511638 -0.9597891  2.469243
13 2        36    0.58774352  0.003360806  0.345453735 -0.9597891  2.469243
14 3        54    0.34479069 -0.358867800  0.247666719 -0.9597891  2.469243
15 4        71   -0.04341135  0.014686416  0.002100236 -0.9597891  2.469243
16 5       117   -0.05103682 -0.070353141  0.007554322 -0.9597891  2.469243

```

17 6 130 -0.56196768 -0.715720445 0.828063425 -0.9597891 2.469243

## Databases, in R

As in [stat545-ubc.github.io](https://github.com/stat545-ubc) consider the two following datasets (from `superheroes.RData`)

```
1 > load("superheroes.RData")
2 > superheroes
3      name alignment gender publisher
4 1  Magneto      bad  male      Marvel
5 2   Storm     good female      Marvel
6 3  Mystique     bad female      Marvel
7 4   Batman     good  male         DC
8 5    Joker     bad  male         DC
9 6 Catwoman     bad female         DC
10 7 Hellboy     good  male Dark Horse Comics
```

for the superheroes,

## Databases, in R

and for the publishers, consider

```
1 > publishers
2   publisher yr_founded
3 1         DC      1934
4 2      Marvel      1939
5 3      Image      1992
```

There are many ways to merge those databases.



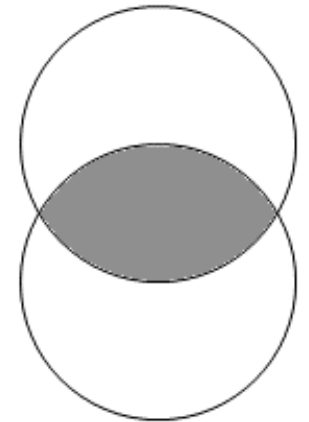
## Databases, in R

Function `inner_join(x, y)` return all rows from `x` where there are matching values in `y`

```

1 > inner_join(superheroes, publishers)
2 Joining by: "publisher"
3   publisher      name alignment  gender yr_founded
4 1    Marvel  Magneto      bad    male    1939
5 2    Marvel   Storm      good  female    1939
6 3    Marvel  Mystique      bad  female    1939
7 4         DC   Batman      good    male    1934
8 5         DC    Joker      bad    male    1934
9 6         DC Catwoman      bad  female    1934

```



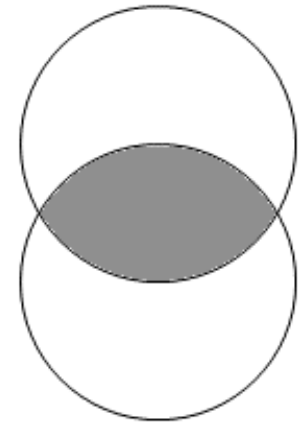
## Databases, in R

Function `semi_join(x, y)` return all rows from `x` where there are matching values in `y`, but only columns from `x` are kept,

```

1 > semi_join(superheroes, publishers)
2 Joining by: "publisher"
3   name alignment gender publisher
4 1  Batman      good  male      DC
5 2   Joker      bad   male      DC
6 3 Catwoman    bad  female    DC
7 4  Magneto    bad   male    Marvel
8 5   Storm    good  female    Marvel
9 6  Mystique    bad  female    Marvel

```



## Databases, in R

```

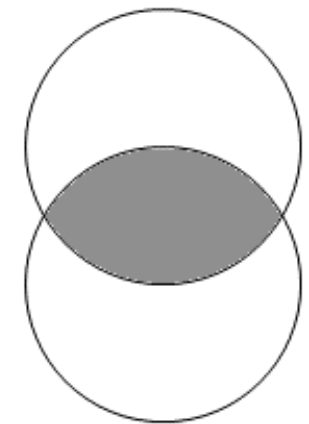
1 > inner_join(publishers, superheroes)
2 Joining by: "publisher"
3   publisher yr_founded      name alignment gender
4 1    Marvel    1939   Magneto      bad   male
5 2    Marvel    1939    Storm      good  female
6 3    Marvel    1939  Mystique      bad  female
7 4      DC     1934   Batman      good   male
8 5      DC     1934    Joker      bad   male
9 6      DC     1934  Catwoman      bad  female

```

```

1 > semi_join(publishers, superheroes)
2 Joining by: "publisher"
3   publisher yr_founded
4 1    Marvel    1939
5 2      DC     1934

```



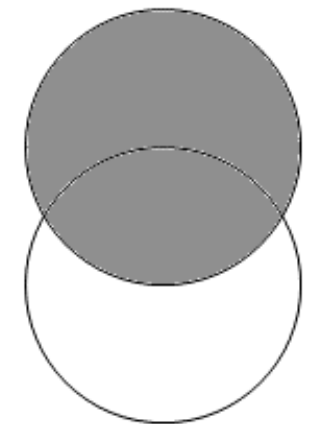
## Databases, in R

Function `left_join(x, y)` return all rows from `x` and all columns from `x` and `y`

```

1 > left_join(superheroes, publishers)
2 Joining by: "publisher"
3   publisher      name alignment gender yr_founded
4 1      Marvel  Magneto      bad   male    1939
5 2      Marvel   Storm     good  female  1939
6 3      Marvel  Mystique     bad  female  1939
7 4         DC   Batman     good   male   1934
8 5         DC    Joker     bad   male   1934
9 6         DC  Catwoman     bad  female  1934
10 7 Dark Horse Comics  Hellboy     good   male    NA

```



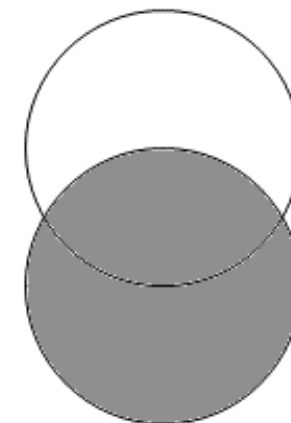
## Databases, in R

There is no `right_join(x, y)` so we have to permutate x and y

```

1 > left_join(publishers, superheroes)
2 Joining by: "publisher"
3   publisher yr_founded      name alignment gender
4 1         DC      1934   Batman      good   male
5 2         DC      1934    Joker      bad    male
6 3         DC      1934 Catwoman      bad   female
7 4    Marvel      1939   Magneto      bad   male
8 5    Marvel      1939    Storm      good   female
9 6    Marvel      1939 Mystique      bad   female
10 7    Image      1992    <NA>    <NA>   <NA>

```



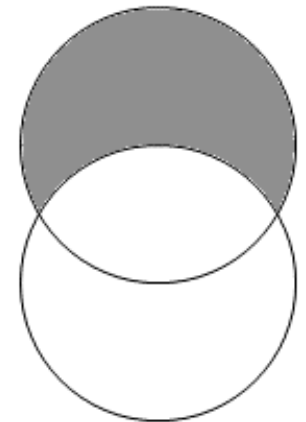
## Databases, in R

One can use `anti_join(x, y)` for rows of `x` that have no match in `y`

```
1 > anti_join(superheroes, publishers)
2 Joining by: "publisher"
3   name alignment gender publisher
4 1 Hellboy      good  male Dark Horse Comics
```

and conversely

```
1 > anti_join(publishers, superheroes)
2 Joining by: "publisher"
3 publisher yr_founded
4 1      Image      1992
```



## Databases, in R

Note that it is possible to use a standard `merge()` function

```

1 > merge(superheroes , publishers , all = TRUE)
2           publisher      name alignment gender yr_founded
3 1 Dark Horse Comics Hellboy      good  male      NA
4 2           DC      Batman      good  male    1934
5 3           DC      Joker      bad   male    1934
6 4           DC Catwoman      bad  female    1934
7 5           Marvel Magneto      bad   male    1939
8 6           Marvel  Storm      good  female    1939
9 7           Marvel Mystique      bad  female    1939
10 8           Image  <NA>      <NA>  <NA>    1992

```

but it is much slower (in `dplyr` integrates R with C++)

There is also a `sql_join` for more advanced SQL requests).

## Data Tables, in R

```

1 > system.time( load("dt_json_2.RData") )
2           user           system           elapsed
3      21.53             1.33             27.71
4 > system.time( setkey(dt, Traj_Id) )
5           user           system           elapsed
6       0.38             0.09             0.47
7 > system.time( depart <- dt[J(unique(Traj_Id)), mult = "first"] )
8           user           system           elapsed
9       3.82             8.02            101.77
10 > system.time( arrivee <- dt[J(unique(Traj_Id)), mult = "last"] )
11           user           system           elapsed
12      2.62             0.52             3.39

```



## Data Tables, in R

```

1 > lat_0=0
2 > lon_0=0
3 > system.time( arrivee[, dist:=(lat-lat_0)2+(lon-lon_0)2] )
4           user           system           elapsed
5           0.03            0.08            1.60
6 > system.time( fin <- subset( arrivee, dist <= 1) )
7           user           system           elapsed
8           0.04            0.00            0.78
9 > system.time( fin[, Pers_Id:=NULL] )
10          user           system           elapsed
11          0.00            0.00            0.07
12 > system.time( fin[, lat:=NULL] )
13          user           system           elapsed
14          0.0             0.0             0.2
15 > system.time( fin[, lon:=NULL] )
16          user           system           elapsed
17          0              0              0

```

## Data Tables, in R

```

1 > system.time( base <- merge(fin , depart , all.x=TRUE) )
2           user           system           elapsed
3           0.02            0.00            0.17
4 > system.time( base <- merge(fin , depart , all.x=TRUE) )
5           user           system           elapsed
6           0.02            0.00            0.17
7 >
8 > head(base)
9   Traj_Id      dist Pers_Id      lat      lon
10 1:         8 0.41251163         1 -0.9597891 2.469243
11 2:        36 0.34545373         1 -0.9597891 2.469243
12 3:        54 0.24766671         1 -0.9597891 2.469243
13 4:        71 0.00210023         1 -0.9597891 2.469243
14 5:       117 0.00755432         1 -0.9597891 2.469243
15 6:       130 0.82806342         1 -0.9597891 2.469243

```

## Memory and Datasets, in R

Instead of loading the complete dataset in the RAM, it is also possible to load it by chunks. Consider e.g. the ‘[Death Master File](#)’ .info,

```

1 > cols <- c(1,9,20,4,15,15,1,2,2,4,2,2,4,2,5,5,7)
2 > noms_col <- c("code", "ssn", "last_name", "name_suffix", "first_name",
  "middle_name", "VorPCode", "date_death_m", "date_death_d", "date_death
  _y", "date_birth_m", "date_birth_d", "date_birth_y", "state", "zip_
  resid", "zip_payment", "blanks")
3 > library(LaF)
4 > temp <- "ssdm3"
5 > ssn <- laf_open_fwf( temp, column_widths = cols, column_types=rep("
  character", length(cols)), column_names = noms_col, trim = TRUE)
6 > object.size(ssn)
7 3544 bytes
8 > go_through <- seq(1, nrow(ssn), by = 1e05)
9 > if(go_through[ length(go_through)] != nrow(ssn)) go_through <- c(
  go_through, nrow(ssn))

```

## Memory and Datasets, in R

```

1 > go_through <- cbind(go_through[-length(go_through)], c(go_through[-c
  (1, length(go_through)) ]-1, go_through [ length(go_through)]))
2 > go_through
3           [,1]      [,2]
4 [1,]          1  100000
5 [2,]    100001  200000
6 [3,]    200001  300000
7
8 [286,] 28500001 28600000
9 [287,] 28600001 28607398
10 >
11 > pb <- txtProgressBar(min = 0, max = nrow( go_through), style = 3)
12 > count_birthday <- function(s){
13 +   #print(s)
14 +   setTxtProgressBar(pb, s)

```

## Memory and Datasets, in R

```

1 + data <- ssn[ go_through[s,1]:go_through[s,2], c("date_death_m", "
   date_death_d",
2 +           "date_birth_m", "date_birth_d")]
3 + sum((data$date_death_m==data$date_birth_m) &
4 +     (data$date_death_d==data$date_birth_d) )
5 + }
6 >
7 > system.time( data <- lapply( seq_len(nrow( go_through) ), count_
   birthday) )
8 |=====| 100
9 %utilisateur      syst me      coul
10      19.48          1.37       37.31
11 > sum( unlist(data) ) /nrow(ssn)
12 [1] 0.001753847

```

## Environments, in R

An environment is a collection of names, and each name points to an object stored somewhere

```
1 > a <- 1
2 > ls(globalenv())
3 [1] "a"
4 > environment(sd)
5 <environment: namespace:stats>
6 > find("pi")
7 [1] "package:base"
1 > e <- new.env()
2 > e$d <- 1
3 > e$f <- 1:5
4 > e$g <- 1:5
5 > ls(e)
6 [1] "d" "f" "g"
7 > str(e)
8 <environment: 0x8b14918>
```

## Environments, in R

```
1 > identical(globalenv(), e)
2 [1] FALSE
3 > search()
4 [1] ".GlobalEnv" "package: memoise"
5 [3] "package: microbenchmark" "package: Rcpp"
6 [5] "package: lubridate" "package: pryr"
7 [7] "package: parallel" "package: sp"
8 [9] "tools: rstudio" "package: stats"
9 [11] "package: graphics" "package: grDevices"
10 [13] "package: utils" "package: datasets"
11 [15] "package: methods" "Autoloads"
12 [17] "package: base"
```

# Filling Forms & Web Scrapping



Mein Telefonbuch Anmelden

As in Munzert *et al.* (2014, <http://eu.wiley.com>) Consider here all people in Germany with the name Feuerstein,

```
1 > tb <- getForm("http://www.dastelefonbuch.de/", .params = c(kw = "
  Feuerstein", cmd = "search", ao1 = "1", reccount = "2000"))
```

Let us store that webpage on our computer

```
1 > write(tb, file = "phonebook_feuerstein.html")
2 > tb_parse <- htmlParse("phonebook_feuerstein.html", encoding = "UTF
  -8")

1 > xpath <- "//ul/li/a[contains(text(), 'Privat')]"
```





```

11 > names_vec <- xpathSApply(tb_parse, xpath, xmlValue)
12 > xpath <- "//div[@class='name']/following-sibling::div[@class='
      popupMenu']//span[@itemprop='postal-code']"
13 > zipcodes_vec <- xpathSApply(tb_parse, xpath, xmlValue)
14 > names_vec <- str_replace_all(names_vec, "(\\n|\\t|\\r| {2,})", "")
15 > zipcodes_vec <- as.numeric(zipcodes_vec)

1 > entries_df <- data.frame(plz = zipcodes_vec, name = names_vec)
2 > head(entries_df)
3   plz                                name
4 1 64625                Bertsch-Feuerstein Lilli
5 2 68549 Bierig-Feuerstein Brigitte u. Feuerstein Norbert
6 3 68526                Blatt Karl u. Feuerstein-Blatt Ursula
7 4 50733                                Feuerstein
8 5 69207                                Feuerstein
9 6 97769                                Feuerstein

```

Now, we need a dataset that links zip codes (*Postleitzahlen*, PLZ) and geographic coordinates. We can use datasets from the OpenGeoDB project (see

<http://opengeodb.org>)

```

1 > download.file("http://fa-technik.adfc.de/code/opengeodb/PLZ.tab",
2 + destfile = "geo_germany/plz_de.txt")
3 > plz_df <- read.delim("geo_germany/plz_de.txt", stringsAsFactors
4 + = FALSE, encoding = "UTF-8")
5 > plz_df[1:3, ]
6   X.loc_id  plz      lon      lat      Ort
7 1      5078 1067 13.72107 51.06003 Dresden
8 2      5079 1069 13.73891 51.03956 Dresden
9 3      5080 1097 13.74397 51.06675 Dresden

```

Now, if we merge the two

```

1 > places_geo <- merge(entries_df, plz_df, by = "plz", all.x = TRUE)
2 > places_geo[1:3, ]
3   plz      name X.loc_id      lon      lat      Ort
4 1 1159  Feuerstein Falk      5087 13.70069 51.04261  Dresden
5 2 1623  Feuerstein Regina      5122 13.29736 51.16516 Lommatzsch
6 3 2827  Feuerstein Wolfgang      5199 14.96443 51.13170  Grlitz

```

Now we simply need some shapefile (see slides on Spatial aspects),

```
1 > download.file("http://biogeo.ucdavis.edu/data/gadm2/shp/DEU_adm.zip",
2 + destfile = "geo_germany/ger_shape.zip")
3 > unzip("geo_germany/ger_shape.zip", exdir = "geo_germany")
4 > projection <- CRS("+proj=longlat +ellps=WGS84 +datum=WGS84")
5 > map_germany <- readShapePoly(str_c(getwd(), "/geo_germany/DEU_adm0.
6 + shp"),
7 + proj4string = projection)
8 > map_germany_laender <- readShapePoly(str_c(getwd(), "/geo_germany/
9 + DEU_adm1.shp"),
10 + proj4string=projection)
11 > coords <- SpatialPoints(cbind(places_geo$lon, places_geo$lat))
12 > proj4string(coords) <- CRS("+proj=longlat +ellps=WGS84 +datum=WGS84
13 + ")
14 > data("world.cities")
15 > cities_ger <- subset(world.cities,
16 + country.etc == "Germany" &
```

```

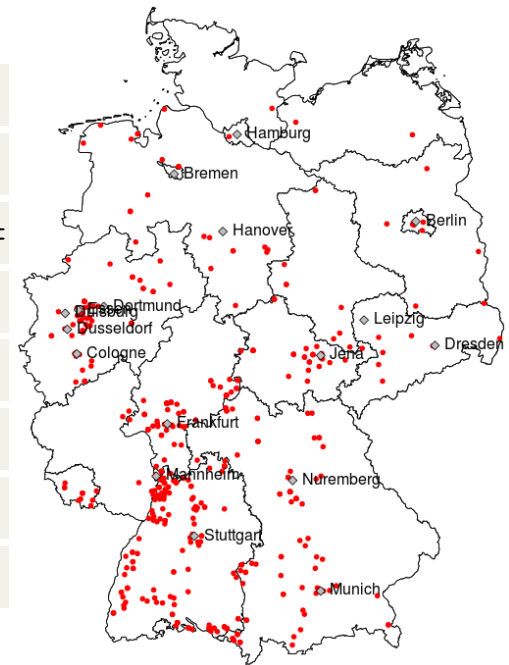
14 + (world.cities$pop > 450000 |
15 + world.cities$name %in%
16 + c("Mannheim", "Jena"))
17 > coords_cities <- SpatialPoints(cbind(cities_ger$long, cities_ger$lat
    ))

```

```

1 > plot(map_germany)
2 > plot(map_germany_laender, add = TRUE)
3 > points(coords$coords.x1, coords$coords.x2, pch =
    20, col = "red")
4 > points(coords_cities, col = "black", , bg = "
    grey", pch = 23)
5 > text(cities_ger$long, cities_ger$lat, labels =
    cities_ger$name, pos = 4)

```



Similarly, consider Petersen, Gruber and Schultze

