

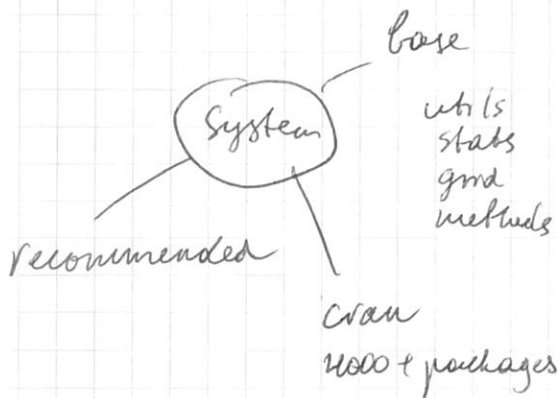
Coursera Computing for Data Analysis

Installing - just ~~use~~ follow
the wizard

About R

Not ideal for everything

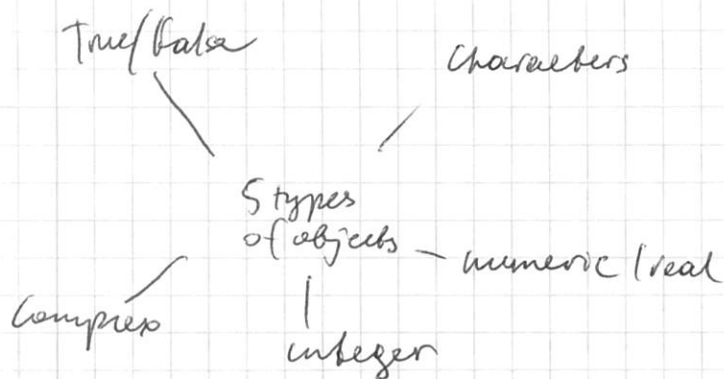
- needs to load a whole object into memory \Rightarrow
impossible to work with huge objects



Books: Software for DA
Use R!



Data types

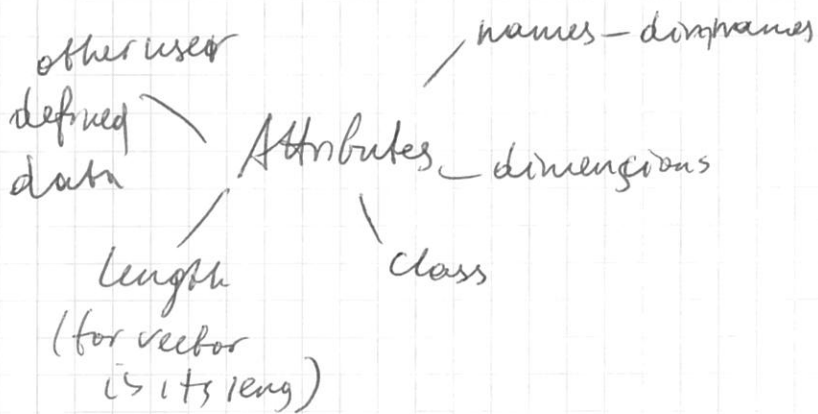
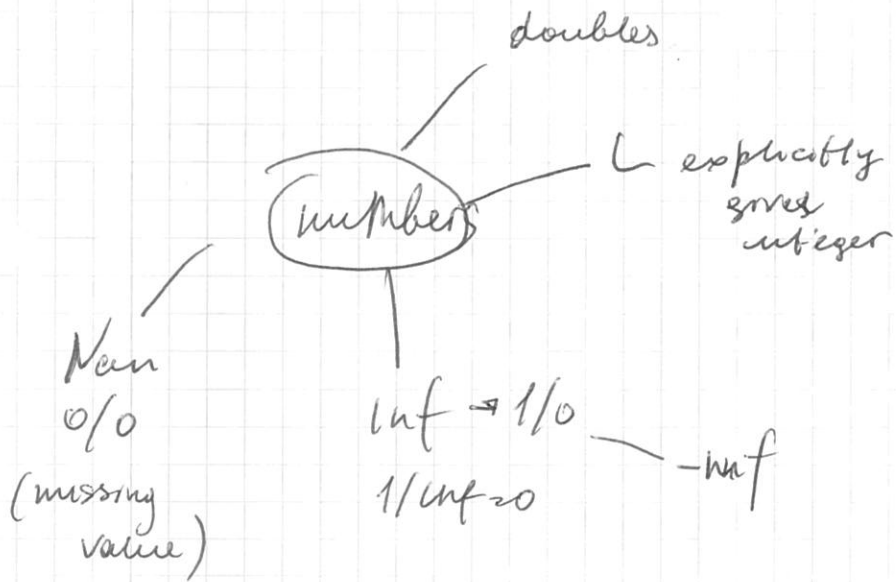


the most basic object - is a vector

↓
everything is of the same class

but there's a list - vector
that can contain different types

`vector()` creates an empty vector



attributes() - general function

At the prompt ~~the~~ type expression
we

> x ← 1 assignment

> print(x)

[1] 1 ← Numeric vector
with one expression

> x - just prints

> msg ← "hello"

- comment

When a complete expression is typed,
it's evaluated and the result of
evaluation is returned

✗

[1] 5

↑
it's vector

↖ first element

$x \leftarrow 1:20$
↑
sequence

> print x
[1] 1... 15
[16] 16... 20

c. other function for creating
vectors
(stands for concatenation)

$x \leftarrow c(0.5, 0.6)$ Numeric
 $\leftarrow c(TRUE, FALSE)$ Logical
 $c(T, F)$
 $c("a", "b")$ character
 ~~$c(0, 9:29)$~~ integer
 $c(1+0i, 2+4i)$ complex

$x \leftarrow \text{vector}("numeric", \text{length} = 10)$

$y \leftarrow c(1.7, "a") \leftarrow \text{character}$

different types

no error!

R will try to COERCE the elements
(tries to find "least common denominator")

$1.7 \rightarrow "1.7"$

$c(1, 2) \Rightarrow [1, 2]$

$c("a", TRUE) \Rightarrow ["a", "TRUE"]$

can be explicitly coerced

as. * function

$x \leftarrow 0.6$

as.numeric(x)

as.logical(x)

as.character(x) \Rightarrow gets converted

as.complex(x)

if no conversion is possible,
NA is returned

Matrices

vectors with dimension attribute

↓
(nrow, ncol)

$m \leftarrow \text{matrix}(\text{nrow}=2, \text{ncol}=3) \leftarrow$

$\text{dim}(m) \Rightarrow 2 \ 3$

$\text{attributes}(m)$

$\#dim - 2, 3$

initialized
with
NA

M. constructed column-wise

$m \leftarrow \text{matrix}(1:6, \text{nrow}=2, \text{ncol}=3)$



from "upper left"
con

1 3 5
2 4 6

another way

$$m \leftarrow 1:10$$

changing dimension attribute

$$? \dim(m) \leftarrow c(2,5)$$

does it work?...

column-binding (row-binding)

$cbind()$

$rbind()$

$$x \leftarrow 1:3 \quad y \leftarrow 10:12$$

$cbind(x, y)$

\Downarrow

1	10
2	11
3	12

$rbind(x, y)$

\Downarrow

1	2	3
10	11	12

Lists

vector with different data

$x \leftarrow \text{list}(1, "a", T, 1+4i)$

$[[1]] \neq \text{double-bracketed}$

$[1] 1$

...

Factor

categorical data vector

can be

ordered

unordered

factor — integer vector with each integer having a label

$\text{lm}()$

$\text{glm}()$

self-describing \leftarrow better
than just numbers

`x <- factor(c("yes", "yes", "no",
"yes", "no"))`

`table(x) => frequency`

`no`
2
`yes`
3

↓
levels
sorted
alphabetically

`unclass(x)` ← strips out the class

↓
2 2 1 2 1
↓ ↓
yes no

`attr("levels") => "no" "yes"`

`x <- factor(c(...),
levels = c("yes", "no"))`

↑
explicit ← auto (not always good)

↓
`levels = yes no`
(note the order!)

Missing Values

NaN - not a number

NA - missing value, any types

is.na() - if NA

is.nan

`x <- c(1, 2, NA, 10, 3)`

`is.na(x)`

returns logical vector

↓ ↓
FALSE TRUE

NaN is NA

is.na ↓
true

↓

`is.na(NA) => True`

`is.nan(NA) => false`

`is.nan(NaN) => true`

`is.na(NaN) => True`

Data Frames

for storing tabular data

Special type of list

with every el of the same len

special attributes:

row.names

`read.table()`

`read.csv()`

`data.matrix` converts to a matrix

↓
will coerce!

`x <- data.frame(foo = 1:4`

`bar = c(T, T, F, F))`

`nrow(x)`

`ncol(x)`

returns number of rows
number of cols

Names

good for describing
data

$x \leftarrow 1:3$

$\text{names}(x) \Rightarrow \text{NULL}$

$\text{names}(x) \leftarrow c("foo", "bar", "n")$

$x \Rightarrow$ will be printed with the names

$\text{names}(x)$ returns the vector

list can also have names

$x \leftarrow \text{list}(a=1, b=2, c=3)$

matrices can also

$m \leftarrow \text{matrix}(1:4, \text{row}=2, \text{col}=2)$

$\text{dimnames}(m) \leftarrow \text{list}(c("a", "b"),$
 $c("c", "d"))$

$m \Rightarrow$

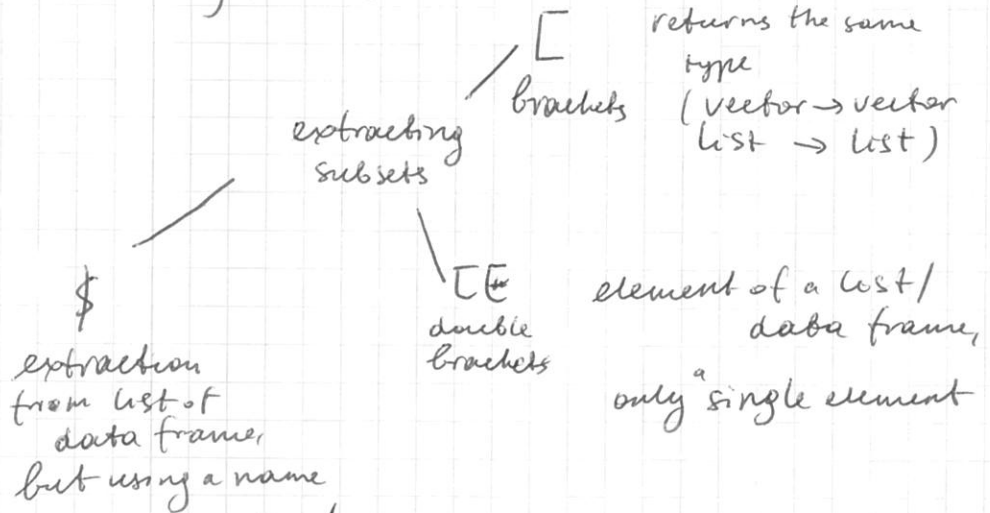
	c	d
a	1	3
b	2	4

Summary of Data Types

- atomic classes
numeric, logical, char, int, complex
- vector, list
- factors
- missing values
- data frames
- names

Subsetting

extracting subsets



starts from 1

```
x <- c("a", "b", "c", "c", "d", "a")
x[1] => "a"
```

```
x[1:4] => a b c c
```

```
x[x > "a"] b c c d
```

```
u <- x > "a"
```

\Downarrow

```
F T T T T F
```

```
x[u] b c c d
```

(only trues)

(+1819)

$X \leftarrow \text{matrix}(1:6, 2, 3)$
 $[1, 1]$
 \downarrow

$X[1, 2] \Rightarrow 3$
 $X[2, 1] \Rightarrow 2$

1	3	5
2	4	6

$X[1,] \leftarrow$ first row
 \swarrow missing

$X[, 2] \leftarrow$ second col

\downarrow
3, 4
}

returns a VECTOR!

$X[1, 2, \text{drop} = F]$

will return a matrix of 1 el
preserves dimensions

$X[1, , \text{drop} = F]$

returns a matrix

List

$x \leftarrow \text{list}(\text{foo} = 1:4, \text{bar} = 0.6)$

$x[1]$

\Downarrow

foo (returns)
[1] 1 2 3 4 {list}

$x[[1]] \Rightarrow [1] \ 1 \ 2 \ 3 \ 4$ (vector)

$x\$bar \Leftrightarrow x[["bar"]] \Rightarrow 0.6$ (element)

$x["bar"]$

\Downarrow

bar (returns)
[1] 0.6 {list}

$x \leftarrow \text{list}(\text{foo}, \dots, \text{bar}, \text{bar})$

$x[c(1, 3)]$

$\text{name} \leftarrow "bar"$

$x[\text{name}] ; x\$name$ ("name" doesn't exist)
 \Downarrow \Downarrow
 bar null!

X [[c(1, 3)]]

↓

1st 3th element

||

X [[1]] [[3]] ← the same

partial matching - for command line

x ← list(aard... ~ 1:5)

x\$a

← partially matches!

x[["a"]]

⇒ N4CC

↑

doesn't work!

x[["a", exact = F]]

↓

→ the same as

removing NAs

bad ← is.na(x)
x[!bad]
↑
inversion!

x ← with.naas
y ← with.naas

good ← complete.cases(x, y)
x[good], y[good]

~~new~~ r ← matrix

good ← complete.cases(r)
r[good,] [1:6,]

Vectorised operators

no loops!

$x \leftarrow 1:4$ $y \leftarrow 6:9$

$x + y$

$x > 2$

$x >= 2 \quad \Leftarrow$ each with each

$y == 8$

$x * y$

x / y

$x \leftarrow \text{matrix}(1:4, 2, 2)$

$y \leftarrow \text{matrix}(\text{rep}(10, 4), 2, 2)$

$x * y$ # element-wise mult!

x / y - element-wise

$x \%*\% y$ - matrix multiplication

Read/Write data

read.table read.csv

readLines - from text file
(char vector)

source - R code

dget - R object

load

unserialize

↓ analogues

write.table

writeln

dump

dput

save

serialize

read.table

args:

- file
- header - logical (var. names in the 1st line)
- sep - separator (",", "
with this
- colClass - char vector types
(not required)
- n rows
- comment.char - ~~comment~~ to be ignored after it
- skip - n of lines to skip
- stringsAsFactors

should character vars be coded as factors?

data <- read.table("foo.txt")

↑ only
file is required others
set to default
default values: # - comment, R figures out types

for small data

read.csv - simulator,
def. separator is comma

For larger sets

read.table - ~~read~~ RTFM!

make a rough calculation
if RAM is enough

another opt:

comment.char = "#" - faster

much
colClasses - faster if explicitly defined

How to figure out the classes?

```
initial <- read.table("db.txt", nrow=100)
# reads the classes
classes <- sapply(initial, class)
head(initial) # just shows a head of data
tabAll <- read.table("db.txt",
# colClasses = classes)
```

rows - doesn't make R faster, but
better for memory

str(.Platform)

Rough Calc for RAM

15 mln rows \times 120 cols | numeric |
 \approx 8 bytes

$1500000 \times 120 \times 8 \text{ bytes} = 1,346 \text{ B}$

(then multiply by 2)

Textual Format

dump and dput include the metadata
source \rightarrow dget

it's editable, good for ~~CSV~~ Version Control!
(not very space efficient)


```
y <- data.frame(a = 1, b = "a")
```

```
dput(y) => to the console  
(some R code)
```

```
dput(y, file = "y.R")  
new.y <- dget("y.R")
```

dump can be used for multiple objects

```
dump(c("x", "y"), file = "data.R")  
      names of  
      objects
```

```
rm(x, y) # removes vars  
source("data.R")
```

Interfaces to ~~our~~ outer world

file - connects to a file

gzip gzfile (gzip)

bzfile (bzip2)

url - to a web page

File Connections

str(file)

description - name of the file

open - code indicator

"r"	"w"	"a"	"rb"
read	write	append	"wb"
	(and create a new file)		"ab"
			binary

```
con ← file("foo.txt", "r")
```

```
data ← read.csv(con)
```

```
close(con)
```

⇒ read.csv(~~file~~
"data.txt")

```
con <- gzfile("w.gz")
```

```
x <- readLines(con, 10)
```

↑

first 10 lines

writes

```
con <- curl("http://...", "r")
```

```
x <- readLines(con)
```

Working directory and Editing R files

getwd() ← should be there
read.csv("file")

File → change dir...

dir() → ~~ls~~

↳ - list of functions

File → new script...

```
myfunc ← function() {  
  x ← rnorm(100)  
  mean(x)  
}
```

save as ~ mycode.R

source("mycode.R")

STR

structure of an R object
for diagnosis

`str(str)` ← brief summary of function,
data

`summary(x)` min, max
str on vector - mean, etc

PA 1

PA 1

7. mean of Ozone layer without NA values

`ozone <- data$Ozone`

`mean(ozone[!is.na(ozone)])`

8. The mean of Temp when Month = 6.

`mean(data[data$Month == 6,]
$Temp)`

9. Extract the subset of rows of the data
frame where Ozone values are above 31
and Temp > 90

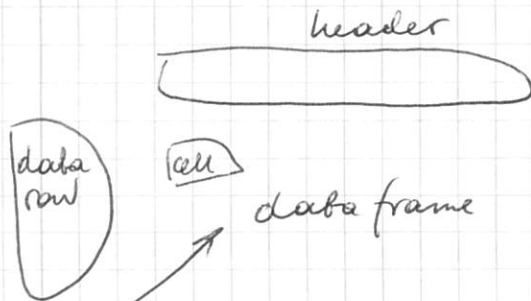
What is the mean of Solar.R
in this subset?

`good <- data$Ozone > 31 & data$Temp > 90`

`good[is.na(good)] <- FALSE`

`mean(data[good,]$Solar.R)`

Data Frames (Again)



`df[1]` - returns whole first column
`df["name"]` name col

`df[c("name", "name2")]`

↓
returns 2 cols

`df[24,]` - 24th row

`L <- df$name == 0`

`df[L,]`

`df[L,]$name`