# 课程体系

📁 R语言数据分析

📁 上部：论道

- 📄 第1章 气象万千、数以等观
- 📄 第2章 所谓学习、归类而已
- 📄 第3章 格言联璧话学习
- 📄 第4章 源于数学、归于工程

📁 **中部：执具**

- 📄 第5章 工欲善其事必先利其器
- 📄 **第6章 基础编程**
- 📄 第7章 数据对象

📁

- 📄 第8章 人人都爱tidyverse
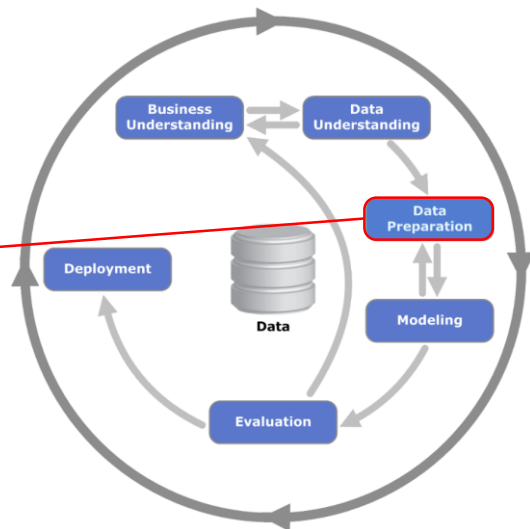- 📄 第9章 最美不过数据框

📁 下部 博术

- 📄 第10章 观数以形
- 📄 第11章 相随相伴、谓之关联
- 📄 第12章 既是世间法、自当有分别
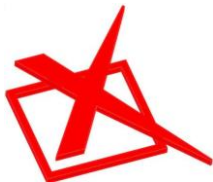- 📄 第13章 方以类聚、物以群分
- 📄 第14章 庐山烟雨浙江潮

# R语言数据分析Mini案例

```r
1  #读取数据
2  library(readxl)
3  cjb <- read_excel("data/cjb.xlsx")
4  View(cjb)
5  #对数据进行探索性分析
6  library(tidyverse)
7  cjb %>%
8    select(sx, wlfk) %>%
9    ggplot(aes(x = wlfk, y = sx, fill = wlfk)) +
10   geom_boxplot()
11 #数据预处理
12 as_five_grade_scores <- function(x) {
13   cut(x, breaks = c(0, seq(60, 100, by = 10)),
14       labels = c("不及格", "及格", "中", "良", "优"))
15 }
16 cjb <- cjb %>%
17   filter(zcj != 0) %>%#剔除脏数据
18   mutate_at(vars(xb, wlfk), factor) %>% #类型转换
19   mutate_at(vars(yw:sw), as_five_grade_scores)#数据分箱
20 View(cjb)
21 #建模
22 library(arulesViz)
23 my_model <- cjb %>%
24   select(xb:wlfk) %>%
25   apriori(parameter = list(supp = 0.06, conf = 0.8),
26           appearance = list(rhs = paste0("wlfk=", c("文科", "理科"))))
27 #模型评估
28 inspectDT(my_model)
29 #可视化
30 plot(my_model)
```

# 两种不同的写法

```r
12 ▾ as_five_grade_scores <- function(x) {
13     cut(x, breaks = c(0, seq(60, 100, by = 10)),
14         labels = c("不及格", "及格", "中", "良", "优"))
15   }
16   cjb <- cjb %>%
17     filter(zcj != 0) %>%#剔除脏数据
18     mutate_at(vars(xb, wlfk), factor) %>% #类型转换
19     mutate_at(vars(yw:sw), as_five_grade_scores)#数据分箱
```

```r
16   cjb <- cjb %>%
17     filter(zcj != 0) %>%#剔除脏数据
18     mutate(xb = factor(xb),
19            wlfk = factor(wlfk),
20            yw = cut(yw, breaks = c(0, seq(60, 100, by = 10)),
21                     labels = c("不及格", "及格", "中", "良", "优")),
22            sx = cut(sx, breaks = c(0, seq(60, 100, by = 10)),
23                     labels = c("不及格", "及格", "中", "良", "优")),
24            wy = cut(wy, breaks = c(0, seq(60, 100, by = 10)),
25                     labels = c("不及格", "及格", "中", "良", "优")),
26            zz = cut(zz, breaks = c(0, seq(60, 100, by = 10)),
27                     labels = c("不及格", "及格", "中", "良", "优")),
28            ls = cut(ls, breaks = c(0, seq(60, 100, by = 10)),
29                     labels = c("不及格", "及格", "中", "良", "优")),
30            dl = cut(ls, breaks = c(0, seq(60, 100, by = 10)),
31                     labels = c("不及格", "及格", "中", "良", "优")),
32            wl = cut(wl, breaks = c(0, seq(60, 100, by = 10)),
33                     labels = c("不及格", "及格", "中", "良", "优")),
34            hx = cut(hx, breaks = c(0, seq(60, 100, by = 10)),
35                     labels = c("不及格", "及格", "中", "良", "优")),
36            sw = cut(sw, breaks = c(0, seq(60, 100, by = 10)),
37                     labels = c("不及格", "及格", "中", "良", "优"))
38   )
```

copy-and-paste

You should consider writing a function whenever you've copied and pasted a block of code more than twice (i.e. you now have three copies of the same code).

—*Hadley Wickham, Chief Scientist at Rstudio*

图片来自：http://hadley.nz/

- **Rule of three:** the code can be copied once, but that when the same code is used three times, it should be extracted into a new procedure

- Do not Repeat Yourself (or DRY) principle

- Separation of Concerns

其它一些编码法则请参阅：https://en.wikipedia.org/wiki/Don%27t_repeat_yourself

## 编写函数

```
fun_name <- function(arg1, arg2 = default1, …) {
    #注释
    表达式(循环/判别/……)
    return(返回值)
}
```

- R里边一切都是对象，对象通过赋值来产生，函数也不例外
- 函数声明关键字是function， function返回值就是函数
- 参数列表是以逗号分割，函数主体可以是任何合法的R表达式，通常是一组由大括弧 ﹛ 和 ﹜ 括起来的表达式
- 若无return语句，最后一个表达式的值作为返回值
- 以fun_name(arg1,arg2,…)的形式调用函数

## 函数的定义

```r
#摄氏度（Celsius）到华氏度（Fahrenheit）的转换
ce2fa <- function(ce) { #参数ce为输入
  fa <- 1.8 * ce + 32 #对输入进行处理
  return(fa)  #输出相应的值
}
ce2fa(0)#0℃相当于32℉
#> [1] 32
ce2fa(0:10)#将0~10℃转换为相应的℉
#> [1] 32.0 33.8 35.6 37.4 39.2 41.0 42.8 44.6 46.4 48.2 50.0
```

## 函数的定义

```
ce2fa

#> function(ce) { #参数ce为输入

#>    fa <- 1.8 * ce + 32 #对输入进行处理

#>    return(fa) #输出相应的值

#> }

#> <bytecode: 0x00000000144b5d28>
```

## 位置参数和名义参数

```r
#位置参数和名义参数
frm <- function(name, frm = "BUPT") {
  cat(name," is frm ", frm)
}


frm()#出错
#> Error in cat(name, " is frm ", frm) :
#>    argument "name" is missing, with no default
frm("axb")#参数的缺省值
#> axb  is frm  BUPT
```

## 位置参数和名义参数

```r
#位置参数和名义参数
frm <- function(name, frm = "BUPT") {
  cat(name," is frm ", frm)
}


frm(name = "AXB", frm = "BJTU")
frm("AXB", "BJTU") #和上述语句等价
#> AXB  is frm  BJTU
frm(frm = "BJTU", name = "AXB")
#> AXB  is frm  BJTU
```

# 特殊的参数...

```r
#看几行我们比较熟悉的代码
xm <- c("周黎", "汤海明", "舒江辉", "翁柯", "祁强", "湛容")
yw <- c(94, 87, 92, 91, 85, 92)
xb <- c(FALSE, TRUE, TRUE)
#再看看sum函数
sum(94, 87, 92, 91, 85, 92)
sum(1, 3, 5, 7)
?c
?sum
```

Description

This is a generic function which comb

The default method combines its argu
names are removed.

Usage

```
## S3 Generic function
c(...)

## Default S3 method:
c(..., recursive = FALSE, use.names = TRUE)
```

Description

sum returns the sum of all the values p

Usage

```
sum(..., na.rm = FALSE)
```

Arguments

...
    numeric or complex or logical ve

na.rm

... 可以包含任意多的参数。它通常在参数个数未知或者某些参数会传递给其它函数的情况下使用

```r
my_func <- function(...) {

  cat("The second arg is ", ..2) #第二个参数

  dot_args <- list(...) #通过list来捕捉任意参数

  message("\nThe sum is ",sum(dot_args[[1]], dot_args[[5]]))

}

my_func(1, 'arg2', 3, 4, 5, 6, 7, 8)

#> The second arg is  arg2

#> The sum is 6
```

## 熟而不觉的函数：作为函数的二元操作符

```r
#+、-、*、/ 其实都是函数
1 + 2
"+"(1, 2)
#> [1] 3
'/'(2, 3)
#> [1] 0.6666667
'^'(10, 2)
#> [1] 100
">"(2, 1)
#[1] TRUE
```

```r
'<-'(new_var, 3)
new_var
#> [1] 3
#:本质上也是一个函数
':'(1, 10)
#> [1]  1  2  3  4  5  6
7  8  9 10
#下标也是函数
'['(1:10, 2)
#> [1] 2
```

## 熟而不觉的函数：作为函数的二元操作符

```r
#%in%运算符：左侧的每个元素是否在右侧的集合之中
c(1, 3, 9) %in% 1:3
'%in%'(c(1, 3, 9), 1:3)
#> [1]  TRUE  TRUE FALSE


#自己定义二元操作符函数：a、b为直角边，c为斜边
"%ab2c%" <- function(a, b) {
  sqrt(sum(a^2, b^2))
}
3 %ab2c% 4
# [1] 5
```

## 熟而不觉的函数：**作为函数的二元操作符**

```r
#看完%ab2c%之后，对下边的符号，也就觉得不过如此了
library(purrr)
x <- c(17, 28, 17, 12, 15, 12, 49)
x %>%
  unique() %>%
  sort()
#等价于下边的代码，不过是更加简洁优雅
x <- c(17, 28, 17, 12, 15, 12, 49)
x2 <- unique(x)
x3 <- sort(x2)
x3
#> [1] 12 15 17 28 49
```

## 熟而不觉的函数：作为函数的二元操作符

```r
5 + 2
#> [1] 7
#来点恶作剧
"+" <- function(x, y) {
  x * y
}

5 + 2
#[1] 10

rm("+")  #消除恶作剧的+运算
5 + 2
#> [1] 7
```

#特殊函数的帮助文档

?"+"  #双引号

?'+'  #单引号

?`+`  #反单引号

?'%in%'

?'round`

?round



```
Arithmetic {base}                                    R Documentation

                      Arithmetic Operators

Description

These unary and binary operators perform arithmetic on numeric or
complex vectors (or objects which can be coerced to them).

Usage

+ x
- x
x + y
x - y
x * y
x / y
x ^ y
x %% y
x %/% y

Arguments

x,
y   numeric or complex vectors or objects which can be coerced to such,
    or other objects for which methods have been written.
```

谢谢聆听
**Thank you**

艾新波

手机：**13641159546**

QQ：**23127789**

微信：**13641159546**

E-mail: **13641159546@126.com**

**axb@bupt.edu.cn**

地址：北京邮电大学科研楼**917**室

课程网址：**https://github.com/byaxb/RDataAnalytics**