

R In Action Notes

Author : Yunqiu Xu

R DataAnalysis

- 4 Parts, 16 chapters:
- Introduction-getting data into R-visualisation-data management
- Basic tutorials : 1-6,11,16
- Advanced tutorials : 7,10
- Some difficult topics : 8,9,12-15

Part I Getting Started

Chapter 1 Introduction to R

1.1 Managing the workspace

```
1. getwd() #工作区目录，相当于git中的"pwd"
2. setwd("C:/mydirectory")
3. #注意一点，R中"\"为转义，不可直接出现在路径内，需要改为"/"or "\\\"
4. dir() #展示wd中的文件
5. file.exists("directoryName")/dir.create("directoryName") #查询并创建工作目录
```

- 一个栗子：

```
1. if(!file.exists("xyq")){ #如果不存在则返回TRUE
2.     dir.create("xyq")} #创建工作目录
```

- 常用函数

```
1. source("filename") #从某个文件中读取R编码，见后
2. list.files("C:\data") #类似ls()
3. ls() #类似git "ls"
```

```

4. options() #view or set current options
5. history(#) #查询命令历史, 默认为前25个
6. savehistory("myfile") #将命令历史储存到myfile, 载入用loadhistory("myfile")
7. save.image("");
8. save(objectlist,file="myfile"); load("myfile")
9. q() #退出R, 退出前会询问是否保存
10. demo(graphics/Hershey/persp/image or none) #展示R自带的可视化类型

```

- 注意&为交集，返回向量，而&&则返回值

1.2 Getting help

```

1. help.start() #general help
2. ?foo or help("foo") #get help for a function
3. example("foo") #the examples of the function you want to know

```

1.3 Packages

- Install packages:install+library+search

```

1. install.packages("package_name")
2. library(package_name) #when the installation is ended, you need to load it
3. search() #you will find packages that have been loaded
4. help("package_name") #seekhelp

```

Chapter 2 Creating a datas

- R data structures
- data entry
- importing data
- annotating data
- Dataset：本书用observations/variables来称谓行或列
- 基本数据类型：

```

1. numeric#数字, 整数或小数
2. integer#整形

```

```
3. character#字母
4. logical#逻辑正误
5. complex#复数
```

2.1 数据结构

- scalars标量,vectors向量,matrices矩阵,arrays数组,data frames数据框架,lists

2.1.1 向量vectors

```
1. a <- c(1,3,5,7,9) #使用函数c()来形成向量
2. #向量元素必须为同一种类型
3. #如果添加不同元素也可以，但是你class一下这个向量，会发现会被强制转换成一种类型
4. b <- c("one","two","three") #字母向量
5. c <- c(TRUE,FALSE,TRUE,TRUE) #逻辑向量
6. a[1] #向量a中第一个元素（这个跟python从零开始不同）
7. a[c(2,4)] #列出向量a中第二个和第四个元素
8. a[2:6] #列出向量a中第二个到第六个元素，a <-c(2:6) 等价于 a <- c(2,3,4,5,6)
9. #两个向量x/y可以加减乘除，元素一一对应
```

2.1.2 矩阵Matrices

```
1. mymatrix <- matrix(vector, #矩阵中元素
2.                     nrow, #number of rows
3.                     ncol, #number of columns
4.                     byrow=logical_value, #矩阵排列形式，默认False按列排列，TRUE则为按行排
   列
5.                     dimnames=list(char_vector_rownames, char_vector_colnames) #名称
6.                     )
```

- 当然，也可以先把元素整出来，再给这些元素dim一下：

```
1. m<-1:10
2. dim(m)<-c(2,5)
```

- 再举个栗子：将1-20按列排列到5*4矩阵内，且规定矩阵的行名与列名

```
1. y<-matrix(1:20,
2.           nrow=5,ncol=4,
3.           dimnames=list(c("a","b","c","d","e"),
4.                         c("大傻","二傻","三傻","四傻")))
5. y[2,] #矩阵y第二行元素
```

```

6.   y[,2] #矩阵y第二列元素
7.   y[2,c(3,4)] #矩阵y第二行第3,4个元素（相当于第二行与第三列第四列的交叉点）

```

- 两个类型相同的矩阵可以加减乘除，元素一一对应
- 矩阵乘法：`x %*% y`

2.1.3. 阵列Arrays

- 矩阵是特殊的阵列（只有二维），但阵列可以有多个维度

```

1.   myarray <- array(vector,
2.                       dimensions, #每个维度的最大数值,如c(3,4,5)
3.                       dimnames=list()) #每个维度的名称,行名/列名...
4.   )

```

- 举个栗子：建立一个三维（3*4*5）阵列

```

1.   dim1=c("大傻子","二瘸子","三蹦子")
2.   dim2=c("001","002","003","004")
3.   dim3=c("d1","d2","d3","d4","d5")
4.   y <-array(1:60,c(3,4,5),dimnames=list(dim1,dim2,dim3))

```

2.1.4 数据框架DataFrames——比矩阵更适用于R

- 与矩阵的不同之处在于不同列可以放置不同类型的数据

```

1.   mydata <- data.frame(col1,col2,col3)

```

- 举个栗子：将姓名、年龄、学号、籍贯集成数据框架

```

1.   name <-c("大傻子","二瘸子","三蹦子")
2.   age <-c(18,14,22)
3.   studentID <-c(12306,57652,61365)
4.   birthplace <-c("天津","黑龙江","广东")
5.   mydata <- data.frame(name,age,studentID,birthplace)
6.   mydata[2:4] #从第二列到第四列的元素
7.   mydata[c("age","name")] #年龄与姓名信息
8.   mydata$name #这是一种新的查询子集方法

```

- 举个栗子：创建年龄与籍贯的交叉分析表

```

1.   table(mydata$age,mydata$birthplace)

```

```

1. attach/detach(mydata) #将数据框架引入/移出搜索路径
2. #后面学习时看到过，一般在最开始引入如attach(mtcars)
3. #attach/detach相当于创造了一个封闭空间，在这个空间里调用指定数据集子集不需要主数据集前缀
4. #attach后，之后再调用列不需要再标注mydata$
5. with(mydata,table(age,birthplace)) #另一种省略标注的办法，加with
6. #如果有多个表格，用{}括起来
7. #注意下with的限制，with里的生成任务只能在with括号里使用，如在外面使用age会显示错误
8. #解决措施，使用特殊的任务生成器<<-替代<-
9. with(mydata,xyq<<-table(age,birthplace))

```

- 建议使用with而不是attach
- 一些别的小tips

```

1. read.table()/read.csv() #这个不是创建，是从外部读取
2. data.matrix() #将数据框架转化为矩阵
3. nrow()/ncol() #行数与列数

```

2.1.5 系数Factors

- 一种特殊的整形向量，其中每个整形都带有一个标签
- factor要强于单单使用整型变量
- 例如使用标签"male/female"要强于整形1/2

```

1. bp=factor(birthplace,order=TRUE)

```

- 举个栗子：

```

1. #将逻辑分类的字符集组成一个向量
2. x<-factor(c("Y","Y","N","NG","Y","N","NG","Y","NG"))
3. table(x) #对不同的标签计数成表
4. unclass(x) #将不同标签转化为整形权重
5. #一般情况下，我们直接在factor函数里给出不同的level
6. #这在线型模拟中十分重要，因为第一个标签将作为基准物
7. levels=c("Y","N","NG")

```

2.1.6 Lists: 可添加不同类型的object，同数据框架比还能添加不同长度

```

1. mylist <- list(name1=object1,name2=object2,...)

```

- 举个栗子：

```
1. g<-"my first list"
2. h<-c(23,24,25,21)
3. j<-matrix(1:10,nrow=5,ncol=2)
4. k<-c("xyq1","xyq2","xyq3")
5. mylist=list(title=g, ages=h,j,k)
```

2.2 导入(读取)及写入数据

- 导入数据前需要考虑所需内存及硬盘的容量
- 所有的文件都要先存放到工作区目录中才能被找到,或者将你的工作区目录改到这个文件的所在地

2.2.1 键盘导入

```
1. fix()/edit() #编辑界面
```

2.2.2 从文本文件导入数据 read.table()/read.csv()

```
1. mydataframe <- read.table(
2.     file, #工作目录中的某个文件名, 如xxx.csv
3.     header=TRUE, #第一行是否要包含变量名, TRUE/FALSE
4.     sep="," , #分割标点
5.     nrows=10 #读取10行
6.     row.names="name" #可选参数, 选择一个或多个变量作为列标识符, 类似SQL中的
    主键, 例如"STUDENT'S ID"
7.     skip=8 #跳前8行再开始读取
8.     na.strings="sb", #使用该字符指代NA
9. )
```

2.2.3 从EXCEL导入数据

- xlsx package

```
1. library(xlsx)
2. read.xlsx("filename",
3.     sheetIndex=1, #which sheet do we want to store the data on
4.     header=TRUE,
5.     colIndex=2:4, #读取2-4列
6.     rowIndex=3:5, #读取3-5行
7. )
```

- RODBC包: 安装并使用RODBC包install+library

```
1. channel <- odbcConnectExcel("xxx.xls")
2. mydataframe <- sqlFetch(channel, "mysheet")
3. odbcClose(channel) #也可以导入Access数据
```

2.2.4 从XML导入数据 : XML package

```
1. library(XML)
2. doc<-xmlTreeParse(fileUrl,useInternal=TRUE) #连接url
3. rootNode<-xmlRoot(doc) #get root 并进入url
4. xmlName(rootNode)/names(rootNode)
5. rootNode[[1]]#the 1st component of rootNode
6. rootNode[[1]][[1]] #进一步查询子集
7. xmlSApply(rootNode,xmlValue) #从XML文件中提取信息, 类似爬虫!!
```

- 举个栗子 :

```
1. xpathSApply(rootNode,"//name",xmlValue) #获取商品名
2. xpathSApply(rootNode,"//price",xmlValue) #获取价格
```

2.2.5 从网络爬虫webscraping导入数据

- 注意: python爬取网页的能力强于R
- 简单爬取网页

```
1. con=url("uelname") #connect to url
2. htmlCode=readLines(con) #get the data
3. close(con)
```

- 结合之前学的XML: xmlTreeParse+xpathSApply

```
1. library(XML)
2. doc<-xmlTreeParse(fileUrl,useInternal=TRUE)
3. xpathSApply(html,"//title",xmlValue)
```

- httr package

```
1. library(httr)
2. html2<- GET(url)
3. content2<- content(html2,as="text")
```

```

4. parseHTML<- htmlParse(content2,asText=TRUE)
5. xpathSApply(html,"//title",xmlValue)

```

- 对于有password的：

```

1. GET(urlname,authenticate("user","passwd"))

```

2.3.6 其他统计语言SAS/SPSS

- foreign package #可以将经其他统计语言获取的数据转换为R

```

1. read.xport(SAS)
2. read.spss(SPSS)

```

2.3.7 从MySQL导入数据: #RMySQL package

- 连接数据库并进行查询dbConnect+dbGetQuary+dbDisconnect

```

1. ucscDb<-dbConnect(MySQL(),user="XXX",
2.                   host="xxx") #connect the database
3. result<-dbGetQuary(ucscDb,
4.                    "show databases;")#query command
5. dbDisconnect(ucscDb) #解除连接

```

- 打开某个数据库

```

1. ucscDb<-dbConnect(MySQL(),user="XXX",
2.                   db="hg19",host="xxx")
3. allTables<-dbListTables(hg19)
4. allTables[1:5]

```

- get dimensions of a table

```

1. dbListFields(hg19,"aircar") #给出aircar表中的所有列名
2. dbGetQuery(hg19,"select count(*) from airfly")

```

- 读取数据库中某个表dbReadTable

```

1. dbReadTable(hg19,"tablename")

```

- 选取某个子集


```

1. query<-dbSendQuery(hg19,"yourquery") #提交查询
2. ff<-fetch(query) #获取查询结果，滞后就可以对查询结果进行R操作
3. dbDisconnect(hg19)

```

• 其他DBMS:

```

1. RPostgreSQL package
2. RODBC package #PostgreSQL/MySQL/Access/SQLite
3. RMongo package #MongoDb

```

2.3.8 JSON 导入数据: jsonlite package

```

1. library(jsonlite)
2. jsonData<-fromJSON("urlname")
3. jsonData$subset1$subset2 #逐层次查询子集
4. myjson<-toJSON(data,
5.     pretty=TRUE #自动缩进，美化格式
6.     ) #将数据写入为JSON形式
7. cat(myjson) #连接并print，如果有多个objects相当于paste，这里仅仅为print
8. myjson2<-fromJSON(myjson) #使用刚刚的fromJson()函数可以再将转化成JSON的数据转化回来

```

2.3.9 导入图片及音频

- 读取image数据: jpeg/readbitmap/png/EBImage packages
- 读取GIS数据: rdgal/rgeos/raster packages
- 读取music数据: tuneR/seewave packages

2.3.10 结合 JHU Data Science Series总结一下读取和写入的办法

```

1. read.table()/read.csv() #常规读取表格类文本类数据的方法
2. write.table
3. readLines #读取文本文件的行，可以用来爬网页
4. writeLines
5. source/dget #读取文件中的R编码
6. dump/dput
7. load #读取已经保存的workplace
8. save
9. unserialize #二进制形式读取R对象
10. serialize

```

- 这里着重讲一下source-dump; dget-dput, dump/dput 最后生成的是textual format而

非一般的表格

- Textual Format的特点：
 - dumping/dputing的文本格式是可编辑的;
 - 同table/csv不同，dump/dput保存了元数据（牺牲部分可读性），不需要多次具体说明;
 - 文本格式与版本控制系统更兼容;
 - 文本格式可长期保存;
 - 文本格式的不足：对空间的利用效率低;
- 举一个dput+dget的栗子：适合写入单个对象的情况

```
1. y<- data.frame(a=1,b="a") #建立一个data frame
2. dput(y) #简单演示一下
3. dput(y,file="y.R") #将y写入y.R
4. new.y<-dget("y.R") #将读取y.R的结果赋值给new.y
```

- 举一个dumping+source的栗子：适用于多个对象

```
1. x<- "foo"
2. y<- data.frame(a=1,b="a")
3. dump(c("x","y"),file="data.R") #将对象x, y都写入data.R
4. rm(x,y) #删除x, y
5. source("data.R") #如果直接查询x, y会找不到因为已经被删掉了，但是source读取后又可以找到
```

2.3.12 连接文件或者网络connecting to file/gzip/bzip/url

- 连接网络：

```
1. file(description = "",
2.       open = "", #打开方式：r-reading/w-writing/a-appending;rb/wb/ab为
   #以二进制编码形式操作
3.       blocking = TRUE,
4.       encoding = getOption("encoding"),
5.       raw = FALSE)
6. url(description, #同理不加赘述
7.      open = "",
8.      blocking = TRUE,
9.      encoding = getOption("encoding"),
10.     method)
11. gzfile(description, #同理不加赘述
12.      open = "",
```

```
13.     encoding = getOption("encoding"),
14.     compression = 6)
```

- 举个栗子：

```
1.  connection<-file("foo.txt","r") #与foo.txt建立连接，操作方式为读取
2.  data<-read.csv(connection) #读取数据
3.  close(connection) #断开连接
4.  #等价于
5.  data<-read.csv("foo.txt") #在这个例子里，完全可以用一行代码来替代
```

- 再举个栗子，如果仅仅需要读取文件中的部分数据，使用file()更简单：

```
1.  x<-readlines(connection,10) #仅仅读取前10行数据就够了
```

- 再最后举个栗子，使用connection+readLines连接并读取网页内容

```
1.  con<-url("urlname") #连接网站
2.  x<-readLines(con) #读取数据
3.  head(x) #显示head数据
```

- 从网上下载数据：

```
1.  download.file(url, #copy link address
2.  destfile="C:\data\shabi.csv", #save as
3.  method="curl",
4.  #下载方式，对于Mac打开https时需要特别说明，win则是默认的
5.  )
```

2.4 标注数据集

2.4.1 给变量加标签

```
1.  names(patientdata)[2] <- "Age at hospitalization (in years)"
```

2.4.2 带值的标签: 使用 factor() function

2.5 将R数据以文档/pdf/表格等形式输出

2.5.1 以R编码文档输出

```
1. dput(y, file="y.R") #将y写入y.R, 使用dget()读取
2. dump(c("x", "y"), file="data.R") #将多个对象写入data.R, 使用source()读取
```

2.5.2 保存为文档

```
1. sink("123.txt") #创建一个文档
2. source("xxxx") #读取数据, 或者用别的方法导入数据
3. sink() #保存
4. #另一种形式: read.csv()
5. write.csv(data, file="D:/bearf2.csv")
```

2.5.3 保存为pdf

```
1. pdf("mygraph.pdf")
2. #写入需要保存的内容
3. dev.off() #关闭图像生成函数, 自动保存
```

Chapter 3 Getting started with graphs

- creating&saving graphs;
- customizing factors;
- annotating with text and titles;
- controlling a graph's dimensions;
- combining mutiple graphs into one

3.1 基本步骤：

```
1. pdf("mygraph.pdf") #保存为pdf图像
2. attach(mtcars) #连接Motor Trend Car Road Tests数据集
3. plot(dose, drugA, type="b") #散点图, typeb意指点+线状图, p仅仅是点, l仅仅是线,
   #别的可以参考help(plot)
4. abline(lm(mpg~wt)) #拟合
5. title("Regression of MPG on Weight") #add a title
6. detach(mtcars) #断开数据集
7. dev.off() #关闭图像生成函数
```

- 同时连接多个图表：

```
1. dev.new() #创建一个新图表
2. dev.new() #再创建一个新图表
```

- 一些其他用于操作图表的命令：

```
1. dev.next()
2. dev.prev()
3. dev.set()
4. dev.off()
5. dev.cur() #seek help
```

3.3 图表参数 par() function

```
1. opar <- par(no.readonly=TRUE)
```

- 使用par函数设置并查询图表参数

```
1. par(..., no.readonly = FALSE)
2.     #...表示所有类似于tag=value形式的参数
3.     #no.readonly默认为FALSE，如果为TRUE表示保存当前设置
4. par(lty=2,pch=17) #线型更改为曲线(lty=2)，并将点改为实心三角(pch=17)
5. #也可以分成两行写
6. plot(dose,drugA,type="b") #这里就和之前一样了
7. #也可以用另一种写法plot(dose, drugA, type="b", lty=2, pch=17)
8. par(opar)
```

3.3.1 symbols and lines

```
1. pch #点类型 P51图3.4 如之前的pch=17
2. cex #点大小，默认为1,1.5大50%，0.5小50%
3. lty #line_type
4. lwd #line_width，默认为1
```

3.3.2 colors col参数

```
1. col=c("red","blue") #点与线的颜色参数
2. #只有一条线，一种颜色，则整条线以及所有点都是红色的
3. #只有一条线，两种颜色，线线是一种颜色，点是两种交替
4. #两条线，两种颜色，一条线一种颜色
```

```

5.  #多条线, 两种颜色, 交替
6.  col=1
7.  col="red"
8.  col="#FFFFFF"
9.  col=rgb(1,1,1)
10. col=hsv(0,0,1)
11. col.axis #坐标轴颜色参数
12. col.lab #坐标轴标识颜色参数
13. col.main #title
14. col.sub #subtitle
15. fg #前景
16. bg #背景

```

3.3.3 Text characteristics

```

1.  cex #可以使用之前控制图标大小的参数来控制字体大小
2.  cex.axis
3.  cex.lab
4.  cex.main
5.  cex.sub
6.  font #字体类型, 默认为1,2加粗, 3斜体, 4粗斜体, 5Adobe默认编码字体
7.  #同理axis/lab/main/sub
8.
9.  ps #字体大小, 像素点表示, the text size=ps*cex
10. windowsFont(family) #字体集类型
11. windowsFont(
12.     A=windowsFont("字体A")
13.     B=windowsFont("字体B")
14. ) #将字体进行赋值, 调用更方便
15. #如果以pdf形式输出, 可以直接改成pdf里的字体
16. names(pdfFonts()) #查找pdf字体
17. pdf(file="xyq.pdf", family="fontname") #将xyq.pdf以某个字体输出

```

3.3.4 Graph and margin dimensions

```

1.  pin #控制维度(宽/高, 单位为英寸)
2.  mai #控制边距(上下左右, 单位为英寸)
3.  mar #和上面那个一样, 只是单位为1/12英寸

```

• 举个栗子:

```

1.  par(pin=c(4,3), mai=c(1,.5,1,.2)) #4*3, 上边距1以此类推

```

3.3.5 最后一个栗子，综合前面所有的图表设置参数

```
1.  opar <- par(no.readonly=TRUE) #保存当前设置
2.  par(pin=c(3,4)) #图表大小3*4
3.
4.  par(lwd=2,cex=1.5) #线宽2, 点大小1.5
5.
6.  par(cex.axis=0.75,font.axis=3) #坐标轴字体大小0.75, 斜体
7.
8.  plot(dose,drugA,type="b",pch=19,lty=2,col="red") #点+线, 点型19, 线型为2,
   #红色
9.
10. plot(dose,drugB,type="b",pch=23,lty=6,col="blue",bg="green") #点型23, 线
   #型6, 蓝色, 绿色背景
11. par(opar)
```

3.4 添加文本，个性化坐标轴

3.4.1 titles

```
1.  title(main="",sub="",xlab="",ylab="") #也可以添加颜色, 尺寸等参数
```

3.4.2 axes

```
1.  axis(side, #位置1底2左3顶4右
2.      at #加入刻度
3.      labels
4.      pos #交叉点
5.      lty #线型
6.      col #颜色
7.      las #labels平行(0)/垂直(2)于坐标轴
8.      tck #刻度长度, 正数为图表内负数为外默认-0.01
9.      )
```

● 举个栗子：

```
1.  x<-c(2:16)
2.  y<-x+1
3.  z<-10/y
4.  opar<-par(no.readonly=TRUE)
5.
6.  par(mar=c(5,4,4,8)+0.1) #确定边界, 单位为1/12英寸
7.
```

```

8. plot(x,y,type="b",pch=21,col="red", yaxt="n",lty=3,ann=FALSE) #x对y作图,
   点型, 线型
9.
10. lines(x,z,type="b",pch=22,col="blue",lty=2) #向已有的图表添加散点图而非重新
   做一个图表
11.
12. axis(2,at=x,labels=x,col.axis="blue",lty=2)
13.
14. axis(4,at=z,labels=round(z,digits=2),col.axis="blue",las=2,cex.axis=0.7
   ,tck=-.01) #平行, 字体大小0.7, 刻度长度及方向为-0.1
15.
16. mtext("y=x+1", side=4, line=3, cex.lab=1, las=2, col="blue")
17.
18. title("An Example of Creative Axes",xlab="X values",ylab="Y=X+1")
19.
20. par(opar)

```

3.4.3 reference lines基准线: `abline()` function

```

1. abline(a = NULL, #intercept截距
2.         b = NULL, #slope斜率
3.         h = yvalues, #水平数值, 只有这个则是一条y值固定的基准线
4.         v = xvalues, #竖直数值, 只有这个则是一条x值固定的基准线
5.         reg = NULL,
6.         coef = NULL,
7.         untf = FALSE,
8.         ...) #该函数用于添加线条

```

● 举个栗子：

```

1. abline(v=c(1,10,2),lty=2,col="blue") #x=1/x=10/x=2三条线
2. abline(h=seq(1,10,2),lty=4,col="yellow") #y=1,3,5,7,9

```

3.4.4 legend标注 `legend()` function

```

1. legend("topleft", #location
2.        inset=.05,
3.        title="Drug Type",
4.        c("A","B"),
5.        lty=c(1, 2),
6.        pch=c(15, 17),
7.        col=c("red", "blue"))

```


3.4.5 text annotations文本注释 `text()/mtext()` functions

```
1. text(location,  
2.     "text to place",  
3.     pos, #相对于location的位置1below, 2left, 3above, 4right  
4.     ...) #图表内部
```

- `text()`可以用于给图表里的点做标注，这里有个栗子：

```
1. attach(mtcars)  
2. plot(wt, mpg,  
3.     main="Mileage vs. Car Weight",  
4.     xlab="Weight", ylab="Mileage",  
5.     pch=18,  
6.     col="blue")  
7. text(wt, mpg, row.names(mtcars),  
8.     cex=0.6, pos=4, col="red")  
9. detach(mtcars)  
10. mtext("text to place",  
11.     side, #1为下边界, 2为左边界, 3为上边界, 4为右边界  
12.     line=n, #文本标注与边界线的相对位置, 0为最接近  
13.     ...) #四面边界
```

- 再举一个栗子：

```
1. opar <- par(no.readonly=TRUE)  
2. par(cex=1.5)  
3. plot(1:7, 1:7, type="n") #n for "no plotting"  
4. text(3, 3, "小婊砸")  
5. text(4, 4, family="mono", "我赵日天第一个不服")  
6. mtext(family="serif", "ASS WE CAN", side=1, line=1)  
7. par(opar)
```

3.5 组合图表 `par()` / `layout()`

- 举个栗子先：`par()`

```
1. attach(mtcars)  
2. opar <- par(no.readonly=TRUE)  
3. par(mfrow=c(2,2)) #组合图表2行2列  
4. plot(wt, mpg, main="Scatterplot of wt vs. mpg")
```

```

5. plot(wt, disp, main="Scatterplot of wt vs disp")
6. hist(wt, main="Histogram of wt") #hist是柱状图
7. boxplot(wt, main="Boxplot of wt") #箱图
8. par(opar)
9. detach(mtcars)

```

• 再举个栗子：`layout()`

```

1. attach(mtcars)
2. layout(matrix(c(1,1,2,3),2,2,byrow=TRUE),
3.         widths=c(3, 1), heights=c(1, 2))
4. #layout函数里通过矩阵指定位置，这里构造一个2*2按行排列的矩阵
5. #第一个图在第一行，剩下两个图在第二行
6. #指定相对高度与宽度
7. hist(wt)
8. hist(mpg)
9. hist(disp)
10. detach(mtcars)

```

3.5.1 精密地控制图表: `fig= graphical parameter`

```

1. opar <- par(no.readonly=TRUE)
2.
3. par(fig=c(0, 0.8, 0, 0.8)) #这里的数字都是相对值，把图表左下角原点当作(0,0)，
   #右上角当作(1,1)
4. #本图表的范围是x轴0-0.8, y轴0-0.8
5.
6. plot(mtcars$wt, mtcars$mpg, xlab="miles per gallon", ylab="car weight")
7. par(fig=c(0,0.8,0.55,1), new=TRUE) #在图表上方加箱图
8. #这个箱图在x轴上0-0.8, y轴上0.55-1, 相当于在散点图上方
9.
10. boxplot(mtcars$wt, horizontal=TRUE, axes=FALSE)
11. par(fig=c(0.65, 1, 0, 0.8), new=TRUE) #在图表右侧加箱图
12. #这个箱图在x轴上0.65-1, y轴上0-0.8, 相当于在散点图右侧
13. boxplot(mtcars$mpg, axes=FALSE)
14. mtext("Enhanced Scatterplot", side=3, outer=TRUE, line=-3)
15. par(opar)

```

Chapter 4 Basic data management

- Manipulating dates and missing values

- Understanding data type conversions
- Creating and recoding variables
- Sorting, merging, and subsetting datasets
- Selecting and dropping variables

4.1 Creating new variables

```
1. s1<-seq(1,10,by=2) #1-10,by=2
2. s2<-seq(1,10,length=3) #元素数量3个, 平均分配 (1,5.5,10)
```

- 创建二进制变量：

```
1. s3<-ifelse(data<0,TRUE,FALSE) #设置condition,如果data<0返回TRUE
2. table(s3,data<0)
```

- 分割数据

```
1. s4<-cut(data,breaks=quantile(data)) #获取分割数据, 此处分割点为默认分位点
2. table(s4)
3. #更简洁的分割方法: Hmisc package
4. s4S<-cut2(data,g=4) #使用包内函数cut2, 分割成4组数据
5. s5<-factor(data) #创建factor 变量
```

- 举个栗子：

```
1. yesno<--sample(c("Y","N"),size=10,replace = TRUE)
2. yesnofac<-factor(yesno,levels=c("Y","N")) #建立factor
3. relevel(yesnofac,ref="N") #向量顺序不变, 但是level顺序变成了N在前
4. as.numeric(yesnofac) #转换为数字形式, 也可以as.character
```

4.2 重新编码变量recoding

- 例如：
 - 将一系列连续变量导入新目录
 - 替换错误的值
 - 基于缺失的数据创建变量

- 可以使用一系列逻辑操作符返回 `TRUE/FALSE` , 如 `>`、`==` 等

```
1. x|y #或
2. x&y #与
```

- 举个栗子：将一系列年龄数据转换为老中青三个年龄层

```
1. total$age1[total$age==99]<-NA #将年龄为99的数据重编码为NA
2. total$age1[total$age>75]<-"老人"
3. total$age1[total$age>=55 & total$age <= 75] <- "中年人"
4. #青年以此类推，且将分组后的数据保存到新的列age1中
```

- 另一种更简化的写法，使用 `within` 函数：

```
1. total<-within(total,{age1<-NA #创建空的新目录
2.   age1[age>75]<- "老人"
3.   age1[age>=55 & age<=75]<-"中年人"
4.   age1[age<55]<-"青年人"
5. }) #within函数类似以前学的with，不同在于可以对数据框架进行修改
```

4.3 重命名变量：`fix()/rename()/names()`

```
1. fix(yourdataset) #方法1：使用fix调出交互编辑器
2. rename(dataframe,
3.   c(oldname="newname",
4.     oldname="newname",...))
5. ) #方法2：加入reshape包，使用其中的rename() function
6. names(leadership)[2] <- "testDate" #方法3，同样需要安装并导入reshape包
7. #如果后面不加赋值则是单纯显示该列列名
```

4.4 日期date values

- 三种基本日期格式：`Date class` 转化为日期, `POSIXct/POSIXlt Class` 转化为秒，`lt`指list

```
1. %d #日期
2. %a/%A #星期，前者为缩写
3. %m/%b/%B #数字月份，字母缩写月份，完整字母月份
4. %y/%Y #2位年份，4位年份
5. %H #HOUR
```

```

6. Sys.Date() #当前日期
7. date() #当前时间
8. format(yourtime,format="%B %d %Y %a") #格式转换
9. as.Date() #将数值转换为时间格式, as.POSIXct/POSIXlt()同理
10. #可以通过两个时间格式的相减得到时间差, 但必须为同一种时间格式
11. #unclass 这个object you will find the 时间差 between your object and
    the original time
12. as.character(dates) #将日期转为字符
13. strptime(datestring,"%B %d %Y %a") #将字符串向量结构转化为POSIXlt日期格式

```

- 可以看下帮助文档, 有关 `as.Date()/strptime` 字符日期转换/ `ISOdatetime` 数字日期转换

4.5 type conversions 数据类型转换

- is判断, as转换

```

1. is/as.numeric/character/vector/matrix/data.frame/factor/logical

```

4.6 sorting排序

- order():

```

1. newdata<-total[order(total$x,decreasing=FALSE),] #注意这个格式,按x排序, 如
    果是反序则decreasing为TRUE

```

- sort():

```

1. sort(X$var1,
2.       decreasing=TRUE, #倒序(默认为FALSE)
3.       na.last=NA/TRUE/FALSE,
4.       )

```

- plyr package:

```

1. library(plyr)
2. arrange(X,colname)
3. arrange(X,desc(colname)) #倒序
4. X$newvar<-newdata #增加一个变量(列)
5. # 另一种方法cbind(X,newdata)

```

4.7 merging datasets 合并数据集

4.7.1 新增列 `merge()`

```
1. cbind() #这个是简单的方法，把不同列合并成矩阵或数据集，适用于两个数据集没有相同列名的情况
2.
3. merge(dataframeA, dataframeB,
4.        by.x=c("ID", "Country"),
5.        by.y="solution_id", #by参数表示内键（相同行/列名），将两个合并的数据集相关联
6.        all=TRUE
7.        ) #如果某个变量在set A有而set B没有，合并后set B此处为NA
8.
9. merge(dataframeA, dataframeB, all=TRUE) #自动默认合并
```

4.7.2 新增行

```
1. rbind(dataframeA, dataframeB)
2. #注意，两个数据集必须有相同的变量（列名），但列的顺序不一定相同
3. #如果列数不同会发生：列多的数据集多出的列被删除or列少的数据集多出空列，用NA填充
```

4.8 缺失值missing values及移除缺失值的几种方法

- 缺失值的栗子：问卷调查中空缺或格式错误的填写

```
1. x<-c(1, 2, NA, NaN, 10, 3) #NA为缺失值，NaN为缺失数字
```

- 使用NA/NaN来标注，用于指代未定义的数学操作

```
1. is.na() #检测某个对象是否为NA
2. is.NaN() #检测某个对象是否为NaN
3. #NA 缺失值可以有不同的类型如integer NA
4. #NaN缺失值是NA的一种，但NA不一定为NaN
```

- 得到NA数量的方法，先用is.na()转换为逻辑格式，再sum(正确为1错误为0)
- 移除缺失值的几种方法：

```
1. #方法1：先查出缺失值并赋值给bad，再调用其反义
2. bad<-is.na(x)
```

```

3.      x[!bad]
4.
5.      #方法2, 适用于多组向量, 但无法移除NAN
6.      good<-complete.cases(x,y,z)
7.      x[good] #得到无缺失值的向量x
8.      #对于矩阵, 也可以使用complete.cases, 只是之后的调用方法有些不同
9.      x[good,][1:6,] #记住good后的逗号
10.     newdata<-na.omit(olddata) #方法3, na.omit(), 但不足是会把NAN变成"NAN"而不是
      消除

```

4.9 Subsetting查询子集

4.9.1 基本方法：[] 给出子集 / [[]]和 \$ 给出元素

- 查看子集名称并取名：可参考之前的重命名变量

```

1.     names(mydata) #查看名字
2.     names(mydata) <- c(list/matrix("name","name","name")) #取名

```

- 查询向量：

```

1.     a[1] #向量a中第一个元素（这个跟python从零开始不同）
2.     a[c(2,4)] #列出向量a中第二个和第四个元素, 同样适用于其他数据集
3.     a[2:6] #列出向量a中第二个到第六个元素, a <-c(2:6) 等价于 a <- c(2,3,4,5,6),
      同样适用于其他数据集
4.     a[a>5] #列出向量a中大于5的元素

```

- 查询矩阵：

```

1.     y[2,] #矩阵y第二行元素
2.     y[,2] #矩阵y第二列元素
3.     y[2,c(3,4)] #矩阵y第二行第3,4个元素（相当于第二行与第三列第四列的交叉点）
4.     x[1,2,drop=FALSE] #一般情况下查询返回的元素组成向量, 修改drop参数为FALSE, 返回
      一个矩阵而非向量

```

- 查询阵列及list：引入 [[]] 及 \$ 提取list中的某个子集，不同在于 \$ 后接的子集名称，而 [[]] 中的是数字或是"名称"

```

1.     table(mydata$age,mydata$birthplace)
2.     x[[c(1,3)]] #提取list中第一个子集的第三个元素, 等价于x[[1]][[3]]

```

- 部分匹配partial matching :

```
1. x<-list(aard=1:5)
2. x$a--1 2 3 4 5
3. x[["a"]]  
4. x[["a",exact=FALSE]] #将exact参数设为FALSE 开启partial matching
```

4.9.2 条件查询

- 栗子1：选取列查询（列即是变量）

```
1. newdata <- leadership[paste("q", 1:5, sep="")]
2. #等价于：
3. newdata <- leadership[c("q1", "q2", "q3", "q4", "q5")]
4. #paste() 创建相同的字母向量，在第五章还会见到
```

- 栗子2：排除列查询

- 注意，NULL为未定义，和NA空值不一样

```
1. myvars <- names(leadership) %in% c("q3", "q4") #转化为逻辑格式
2. newdata <- leadership[!myvars] #选取反义
3.
4. #等价于：
5. newdata <- leadership[c(-8,-9)] #排除第8、9列数据
6.
7. #还等价于：
8. leadership$q3<- leadership$q4<- NULL #将这两行赋值为NULL
```

- 栗子3：行查询（行即是观察值observations）

```
1. newdata <- leadership[which(leadership$gender=="M" & leadership$age > 30),] #大于30岁的男性
2. #行查询类似于SQL中where查询
3. #为了简化变量条件，可以attach/detach，这样就不用在变量前面加一长串修饰了
```

4.9.3 subset() 函数——比前面条件查询简单得多的方法

- 栗子1：不小于35或者小于24，显示q1,q2,q3,q4四列数据

```
1. newdata <- subset(leadership, age >= 35 | age < 24, select=c(q1, q2, q3, q4))
```


- 栗子2：大于35的男性，显示列gender与q4之间的所有数据

```
1. newdata <- subset(leadership, gender=="M" & age > 25,select=gender:q4)
```

4.9.4 随机样本 `sample()` function

```
1. sample(x, #the vector u wanna take samples
2.         size, #the number of samples you wanna take,不加的话则给出整条向量的随机排列
3.         replace = FALSE,
4.         prob = NULL)
```

- 栗子：

```
1. sample(1:10,4)
2. sample(letters,5) #选取5个字母
3. sample(1:10) #未设定样本数则为排序
4. sample(1:10,replace=TRUE) #数字可重复
```

4.10 使用SQL命令来操作数据框架：`sqldf` package

```
1. newdf <- sqldf(
2.     "select * from mtcars where carb=1 order by mpg",
3.     row.names=TRUE
4.     ) #直接使用SQL进行条件查询
```

可参考[帮助文档](#)

Chapter 5 Advanced data management

- mathematical and statistical functions
- character functions
- looping and conditional execution
- user-written functions
- ways to aggregate and reshape data

5.1 Numerical and character functions

5.1.1 Mathematical functions

```
1. abs(x)
2. sqrt(x)
3. ceiling(x) #不小于x的最小整数
4. floor(x) #不大于x的最大整数
5. trunc(x) #截取truncate x的整数部分
6. round(x,digits=n) #四舍五入保留n位小数（不包含整数数位）
7. signif(x,digits=n) #保留n位significant digits有效数字（相当于n包含了整数数位）
8. cos/sin/tan/acos/asin/atan(x)
9. cosh/sinh/tanh/acosh/asinh/atanh(x) #双曲线hyperbolic三角函数
10. log(x,base=n)
11. log(x)/log10(x)/log2(x) #默认base=e/10/2
12. exp(x)
```

5.1.2 Statistical functions

```
1. mean(mean(x,trim=0.05,na.rm=TRUE)) #trim[0,0.05]代表从首尾去除一定百分比的数据
2. median(x) #median中位数
3. sd(x) #standard deviation 标准差
4. var(x) #variance 方差
5.
6. mad(x) #Median absolute deviation 中位数绝对偏差
7.     #先求出中位数，然后测算每个元素的绝对值与中位数的差，得到一组非负向量
8.
9. quantile(x,probs=c(0.25,0.75)) #分位数，中位数是一种特殊的分位数
10.    (probs=0.5)，type参数有9种，默认为1
11.
12. range(x, na.rm = FALSE, finite = FALSE) #return c(minx,maxx), finite默
13.    认为FALSE，自动移除无限循环
14.
15. min/max(x)
16. sum(x)
17. diff(x, #lagged difference滞后差分
18.       lag=n1, #默认为1，对滞后n1阶的数据进行差分
19.       differences=n2)
    #就是说每个数据减去其前第n1个数据，形成一个新向量
    #差分阶数，默认为1即一阶差分，lag执行一次
```

```

20.         #differences=2 意味着差分执行一次后的新向量再按lag执行一次差分
21.
22.     scale(x, #数据的标准化与中心化—消除量纲对数据结构的影响
23.           #中心化即各项数据减去均值得到新向量
24.           #标准化指中心化后的新向量除以数据集标准差得到的新向量
25.           center=TRUE, #可以为逻辑变量或矩阵某列
26.           scale=TRUE) #可以为逻辑变量或矩阵某列

```

5.1.3 Probability functions

- 基本格式：`[dpqr]distribution_abbreviation()`
- 首字母d/p/q/r：
 - d-density密度
 - p-cumulative distribution 累积分布函数
 - q-quantile function分位函数
 - r-random generation/deviates 随机生成/模拟
- distribution_abbreviation:
 - Beta-beta
 - Binomial-binom
 - Cauchy-cauchy
 - Chi-squared (noncentral)-chisq
 - Exponential-exp
 - F-f
 - Gamma-gamma
 - Geometric-geom, 几何的
 - Hypergeometric-hyper, 超几何的
 - Lognormal-lnorm, 对数正态
 - Logistic-logis
 - Multinomial-multinom
 - Negative binomial-nbinom
 - Normal-norm 传说中的正态分布
 - Poisson-pois
 - Wilcoxon Signed Rank-signrank, 魏克森讯号等级检定
 - T-t
 - Uniform-unif

- Weibull-weibull
- Wilcoxon Rank Sum-wilcox, 魏克森秩和

- 栗子1：正态分布与泊松分布

```

1.  rnorm(n,mean=0,sd=1) #随机正态分布, n为输出元素数量, mean默认0, sd默认1
2.
3.  dnorm(d,mean=0,sd=1,log=FALSE) #正态概率密度, d为分位点, log参数为TRUE时还会
    给出结果的log值
4.
5.  pnorm(p,mean=0,sd=1,lower.tail=TRUE,log.p=FALSE) #正态累积分布, p为概率
6.  #lower.tail为TRUE意指计算分布左侧函数(即P(X<=p)), 若想计算
    upper.tail(P(X>=p)), 则为FALSE
7.
8.  qnorm(q,mean=0,sd=1,lower.tail=TRUE,log.p=FALSE) #正态分位分布(逆累积分布)
    , q为分位点
9.  #p/qnorm互为反函数: 如果F()为某个标准正态分布的累积分布函数:
10.     pnorm(q)=F(q)
11.     qnorm(p)=F-1(p)
12.
13.  rpois(n,rate) #随机泊松分布
14.  ppois(p,rate) #累积泊松分布, Pr(x<=p)

```

- 栗子2：从线性模型中产生随机数

- 线性模型 $y = b_0 + b_1x + e$
- 假定误差 $e = N(0, 2^2)$, $x = N(0, 1^2)$, $b_0 = 0.5$, $b_1 = 2$

```

1.  set.seed(20)
2.  x<-rnorm(100) #随机选取100个x点
3.  e<-rnorm(100,0,2) #随机选取100个e点, 根据假设, mean=0, sd=2
4.  y<-0.5+2*x+e
5.  summary(y) #大致看一下结果
6.  plot(x,y) #绘制散点图
7.  #如果期望x为二进制数, 将x<-rbinom(100,1,0.5)即可, 其余不变

```

- 栗子3：从泊松模型中产生随机数(广义generalized线性模型)

- 模拟泊松模型 $Y \sim \text{Poisson}(\mu)$
- $\log \mu = b_0 + b_1x$; $b_0 = 0.5$, $b_1 = 0.3$

```

1.  set.seed(1)
2.  x<-rnorm(100) #随机选取100个x点
3.  logμ<-0.5+0.3*x

```

```

4. y<-rpois(100,exp(logμ)) #随机泊松模型, rate=exp(logμ)=μ
5. summary(y)
6. plot(x,y)

```

- 栗子4：均匀取样：`set.seed+runif`

- `set.seed`的作用：初始化随机种子，如果以后还想产生相同的一组随机样本的话，调用相同的种子即可
- `set.seed`不仅适用于随机取样，别的需要随机数的场合，也可以引入`setseed`进行重复

```

1. set.seed(1234) #初始化随机种子1234
2. runif(5) #在[0,1]之间产生5个随机数

```

- 注意几种调用的情况：

```

1. #情况1：
2. set.seed(1234)
3. a<-runif(5)
4. b<-runif(5) #a不等于b
5. #情况2：
6. set.seed(1234)
7. a<-runif(5)
8. set.seed(1234)
9. b<-runif(5) #a等于b

```

5.1.4 Character functions

```

1. nchar(x) #返回x中每个元素的字符数量
2. substr(x,start,stop) #返回x中从start到stop的字符值，如果赋值的话则为替换
3. grep/grepl(pattern, #打算进行匹配的内容，如"A"
4. x,
5. ignore.case = FALSE, #是否区分大小写
6. fixed = FALSE #FALSE-pattern为正则表达式格式；TRUE-pattern为单纯的
   字符串
7. ) #字符匹配，grep返回匹配位置，grepl返回TRUE/FALSE的数量
8.
9. sub/gsub(pattern,replacement,x,
10. ignore.case=FALSE,fixed=FALSE
11. ) #在x中查找pattern并替换，
12. #sub和gsub的区别是前者只做一次替换（不管有几次匹配）
13. #而gsub把满足条件的匹配都做替换
14.
15. strsplit(x,split,fixed=FALSE) #将x中字符用split(如" ")分割开来

```

```

16. #注意, 分隔符需要是x中有的字符
17.
18. paste (... , #多个字符串
19.         sep = " " #分隔符, 可以不加
20.         ) #将多个字符串连接起来
21.
22. toupper/tolower(x) #全部大写/小写

```

5.1.5 其他常用函数

```

1. length(object) #元素数量
2. seq(from,to,by) #创建一个大小为[from,to]的数列, 间隔数为by
3. seq_along(vector_name) #相当于seq_len(length(vector_name))
4. rep(x,n) #repeat x for n times
5. cut(x,n) #相当于在n中加n个点, 将x分为n+1个区间
6. pretty(x,n) #在x这个范围里(如c(1:5))均匀选取n个数字

```

```

1. cat(... , #这里同paste, 可以添加不同元素及分隔符, 但注意结尾要加一个换行字符"\n"
2.         file="", #可以加文件地址, 将paste写入文件
3.         append = FALSE, #覆盖原文, 如果TRUE是不覆盖
4.         ) #连接并打印, 也可以写入文件, 如果只有一个object的话等同于print
5.
6. #举个栗子; cat("i = ", 1, "\n", file="c:/work/result.txt")将i=1写入结果

```

```

1. dim(object) #维数(秩)
2. str(object) #利用str()可以对函数结构、运算结果及矩阵等object进行了解, 比
   head()/summary()更直观
3. fivenum() #5-number
4. summary(min,lower-hinge,median,upper-hinge,max)
5. head(object, n=6) #前6行, 等价于object[c(1,2),]
6. tail(object, n=6) #后6行
7. colSums()/rowSums() #行列sums
8. table(data %in% c("1212","1234")) #查看表中数据包含"1212"或"1234"的数量
9. restData(data %in% c("1212","1234")) #得到表中包含这两个或数据的元素
10. xtabs(Xcol~Ycol+Zcol,data=mtcars) #形成交叉表cross tables

```

```

1. class(object) #对象类型
2. object.size(data) #查看data的大小, 默认用bytes表示
3. mode(object) #存储方式
4. names(object) #元素名称
5. c(object1,object2...) #将多个对象组成一个向量
6. object/print(object) #打印对象

```

```

7.  cbind(object1,object2...) #将多个对象组成一列，可以形成矩阵等
8.  rbind(object1,object2...) #将多个对象组成一行，可以形成矩阵等
9.  ls() #列出当前的所有对象
10. rm(object1,object2...) #移除某个或多个对象
11. newobject <- edit(object) #编辑对象并另存为新对象
12. fix(object) #在原位编辑对象

```

```

1.  lm(y~x) #Fitting Linear Model拟合线性模型，会给出intercept截距与slope斜率
2.  gl(n,k,length) #分级，n为级数，k为每级重复次数，length为向量长度默认为n*k
3.
4.  gl(3,5) #栗子1：3levels，每个level重复5次，15个元素
5.  gl(3,1,10) #栗子2：3levels，每个level重复1次，10个元素 (1231231231)

```

5.1.6 apply系列函数，对matrix/array/dataframe/list进行操作

- lapply: Loop over a [list] and evaluate a function on each element
 - 对list的每一个子集做计算
 - lapply 可加入辅助函数[split]
- sapply : Same as lapply but try to [simplify the result]
- apply: Apply a function over the margins of an [array] 阵列，保留某一个维度做计算
- tapply: Apply a function over [subsets of a vector]
- mapply: Multivariate version of lapply，多元的lapply

5.1.6.1 lapply

```

1.  lapply(X, #should be a list.If not, use as.list() to change its type
2.         FUN, #这个函数(e.g. sum,mean)会对X[1]计算一次，再计算X[2]...
3.         ...)

```

- 这里有个例子再深入了解下：

```

1.  x<-1:4
2.  lapply(x,runif) #最后第一行得到一个结果，第二行两个...第四行四个
3.  lapply(x,runif,max=10,min=0) #默认为[0,1]，这里改成[0,10]

```

- lapply可以用于匿名函数，举个构建匿名函数的栗子：

```

1.  lapply(x,
2.      function(elt) elt[,1] #注意不加逗号！
3.      #匿名函数提取list中每一个子集中第一行元素，对x中每列子集都代入这个函数
4.      )

```

5.1.6.2 sapply

- sapply是简化的lapply
 - 如果结果是1*n list(每个元素长度均为1)，返回一个向量
 - 如果结果是m*n list(每个元素长度相等)，返回一个矩阵

5.1.6.3 apply

- apply 对一个数组按行/列计算
 - most often used to apply a function to the rows/columns of a matrix;
 - can be used with general arrays;
 - works in one line

```

1.  array(X, #array
2.      MARGIN, #integer vector, indicate which margins should be retained
3.      #就是保留一个维度，然后打散其他的维度进行函数运算
4.      #Margin也可以为向量，如c(1,2)表示保留第一维度和第二维度
5.      FUN,
6.      ...)

```

- 举个栗子：

```

1.  x<- matrix(c(1:60),10,6) #10行6列矩阵，维度为2
2.
3.  apply(x,2,mean) #保留6列，打散10行进行平均值计算
4.  >5.5 15.5 25.5 35.5 45.5 55.5 #最后得到6个数据为每列的平均值
5.
6.  apply(x,1,mean) #保留10行，打散6列进行平均值计算
7.  >26:35 #最后得到10个数据为每行的平均值

```

- Apply的应用，几个简化的计算（这几个函数可以直接使用）：


```

1. rowSums=apply(x,1,sum) #保留行，做每行的sum，最后得到的数据数量=行数
2. rowMeans=apply(x,1,mean)
3. colSums=apply(x,2,sum)
4. colMeans=apply(x,2,mean)

```

- 再举个栗子：array

```

1. x<-array(c(1:60),c(3,4,5)) #3*4*5阵列
2. apply(x,1,mean) #行平均值，3行总共3个数据(把第三维度接起来)
3. #列平均值同理
4. apply(x,3,mean) #第三维为每个3*4matrix的平均值，总共5个数据

```

5.1.6.3 mapply

- mapply 需要引入多个参数时

```

1. mapply(FUN,
2.       ...,
3.       MoreArgs = NULL, #a list of other arguments in Function
4.       SIMPLIFY = TRUE, #whether the results should be simplified
5.       USE.NAMES = TRUE #use names if the first ... argument has
   names
6.       )

```

- 举个栗子：

```

1. mapply(rep,1:4,4:1)#rep函数的第一个参数执行1:4，第二个参数执行4:1
2. #等价于list(rep(1,4),rep(2,3),rep(3,2),rep(4,1))

```

5.1.6.4 tapply

- tapply 对不规则阵列使用函数（分组统计）

```

1. tapply(X, #vector
2.       INDEX, #该list中的每一个元素都是与x有同样长度的因子，可通过gl()实现
3.       FUN = NULL,
4.       ...,
5.       simplify = TRUE) #可见前面的sapply

```

- 举个栗子：

```
1. index<-gl(3,3) #111222333分成3个level
2. tapply(c(1:9),index,mean)
3. >2 5 8
```

5.1.6.5 split

- split: lapply的辅助函数,可用于将一个object进行分组

```
1. function (x, #vector/list/dataframe
2.           f, #a factor or a list of factors
3.           drop = FALSE, #whether empty factors levels should be dropped
4.           ...)
```

- 举个栗子：将1:30均分为3组并取平均值

```
1. x<-c(1:30)
2. f<-gl(3,1,30)
3. y<-split(x,f) #可以得到一个有三个子集的list, 每个子集10个元素
4. #[ (1,4,...28), (2,5,...29), (3,6,...30) ]
5. lapply(y,mean)
```

- 再举个复杂一点的栗子：将空气质量按月分组并求取均值

```
1. library(datasets) #得到空气质量总数据集
2. aqbymonth<-split(airquality,airquality$Month) #将空气质量按月分组
3. lapply(aqbymonth,function(x) colMeans(x[,c("Ozone","Solar.R","Wind")],n
4.      a.rm=TRUE))
#匿名函数colMeans计算其中3列的平均值
```

- 最后一个栗子：多个levels 1.1,1.2,2.1,...

```
1. split(x,list(f1,f2),drop=TRUE) #and drop the empty levels
```

5.2 Control Flow

- statement/cond/expr/seq

- if/else/for/while/repeat/break/next/return

5.2.1 if-else

```
1.  if(x<10) {  
2.    10  
3.  }else if(x<15) {  
4.    15  
5.  }else {  
6.    0  
7.  }
```

5.2.2 for loop

```
1.  for(i in 1:10) {  
2.    print(i)  
3.  }
```

5.2.3 while loop

```
1.  count<-0  
2.  while(count<10) {  
3.    print(count)  
4.    count<-count+1  
5.  }
```

5.2.4 repeat: 初始化无限循环，只有break才能退出，用得比较少(因为需要尽可能避免无限循环)

```
1.  x<-0  
2.  repeat{  
3.    x<-x+1  
4.    if(x>100) {  
5.      print(x)  
6.      break}  
7.    }
```

5.2.5 next/return

```
1.  next : skip  
2.  return : exit and return a given value
```

5.3 User-written functions

5.3.1 Create an easy function

- create a new R script, and your function should be written there
- 举个栗子:

```
1. summy<-function(x,y){x+y}
2. highlight these statements
3. then press run to run your function in console windows
4. #注意一定要选中，否则无法实现
```

5.3.2 Take a matrix and calculate the mean of each column

- for loop

```
1. columnmean<- function(y,removeNA=TRUE){ #除去NA
2.   nc<-ncol(y) #nc为y的列数
3.   means<- numeric(nc)
4.   for(i in 1:nc){
5.     means[i]<-mean(y[,i],na.rm=removeNA) #求y第i列的平均
   值,na.rm=TRUE
6.     print(means[i])
7.   }
8. }
```

5.3.3 参数

```
1. f<-function(a,b=1,c=2,d=NULL,...){} #b,c,d均为默认参数,...为可变参数，这个可以参见python
```

- 惰性编译lazy evaluation:

```
1. f<-function(a,b){a+2}
2. f(2)->4
3. #但函数式中涉及到关于b的操作而输入没有b会返回错误
```

- 可变参数的栗子：

```
1. args(paste)
2. function(..., sep = " ", collapse = NULL)
3. paste("a","b",sep=":")—"a:b" #传入字符a,b并使用：分隔
4. paste("a","b",se=":")—"a b :" #这里传入的是可变参数se，而分隔符则采用默认参数
```

5.3.4 scoping rules作用域规则

- 自由变量的值首先取决于函数被定义时它们被赋予的值，而不是函数被调用时的值。
- 只要不是局部参数的都是自由变量
- 感觉R与python类似，先查看局部，在查看全局
- 两种环境：global environment及calling environment
 - global environment: 即为workspace，指代某函数被创建的环境，调用词汇/静态作用域
 - calling environment: 指某函数内部框架的环境，调用动态作用域
- 举个栗子来说明lexical scoping与动态作用域:

```

1.  y = 10 #全局变量，这里是函数f/g被创建的环境
2.  f = function(x) { #不是x的都是自由变量
3.      y = 2 #x为形参，y为自由变量
4.      y ^ 2 + g(x) #g同样为自由变量
5.  }
6.  g = function(x) {
7.      x * y #这里的y也是自由变量，调用g被创建的环境中的y，即y=10
8.  }
9.
10. f(3) #f调用局部y=2，这里为动态作用域，而g则调用全局y=10，最后返回34
    
```

- 在一个Function中，对一个变量的lexical scoping rule的检索次序是：
 - 先检查本Function中的calling environment，如果能找到该变量，则返回该变量;
 - 如果不能，则检索这个Function被创建(不是被调用)的Environment;
 - 对于"< <-"函数变量的赋值符，检索时跳过函数内部环境直接全局;
 - 如果一个函数被定义且调用于全局环境，那么calling/global环境没什么区别，相当于动态作用域
- 查询某个函数所处的环境：

```

1.  ls(environment(function_name)) #返回一系列参数
2.  get("para_name", environment(function_name)) #返回该参数的值
    
```

- 静态作用域造成的结果：
 - All objects must be stored in memory

- All functions must carry a pointer to their environments
- In S-PLUS, free variables are always looked up in the global workspace, so everything can be stored on the disk because the “defining environment” of all functions is the same.

5.4 Aggregation聚类 and reshaping重构

- Aggregation : replace groups of observations with summary statistics based on those observations
- Reshaping: alter the structure (rows and columns) determining how the data is organized

5.4.1 Transpose 转置t(x)

```
1. t(x) #行列互换
```

5.4.2 Aggregating data

```
1. aggregate(x, #the data object to be collapsed
2.           by, #[a list] of variables that will be crossed to form new obs
           ervations
3.           FUN)
```

5.4.3 reshaping data:melting+casting

```
library(reshape2)
```

- melting: 保留一些变量，将指定数据变为新变量

```
1. carM<-melt(mtcars,id=c("cyl","gear"), #当做维度的变量，每个变量占一列
2.           measure.vars =c("mpg","hp")) #被当做观测值的变量，列变量名称和值组成v
           ariable和value两列
```

- casting: recast data set into a particular shape
- 举个栗子：得到cyl关于mpg, hp的数据数量

```
1. clyData<-dcast(carM, #刚melt的数据
2.               cyl~variable)
3. #得到结果为"4轮车的mpg数据有11个hp数据有11个，6轮车***，8轮车数据***"
```

- 再举个栗子：得到平均值

```
1. clyData<-dcast(carM,cyl~variable,mean)
```

Part II Basic Method

Chapter6 Basic graphs

- bar,box,dot plots
- pie,fan charts
- Histograms直方图,kernel density plots

6.1 Bar plots柱状图

```
1. barplot(height,#height is a vector or matrix
2.      main="the Title of the plot",
3.      sub="subtitle",
4.      xlab=" ",ylab=" ",
5.      names.arg=c("a","b","c"), #定义各柱子的名称
6.      col=c("red","yellow","green"), #三列的颜色
7.      legend=rownames(counts), #添加标注
8.      )
```

- 其他参数:

```
1. width=1 #柱宽度
2. space=c(0,1) #柱间隔,这个例子是双柱矩阵,间隔0+间隔1
3. beside=FALSE #对于矩阵向量, FALSE为堆叠的柱状图, 而TRUE为并列的
4. horiz=FALSE #FALSE通常的竖直柱状图, TRUE为水平的柱状图
5. density=NULL #暗线密度
6. angle=45 #暗线角度
7. border="red" #边界颜色
8.      border=NA #无边界线
9.      border=TRUE #边线与暗影线颜色一致
10. x/ylim #limits for the axis
11. xpd=TRUE #是否允许柱形图超越图表区域
12. axes=TRUE #是否添加axis
```

```

13. axisnames=TRUE
14. cex.axis/cex.names #点与字符大小, 之前学过, 1.5为1.5倍
15. inside=TRUE #当space=0时是否添加柱与柱的分界线
16. plot=TRUE
17. axis.lty=0 #line type
18. offset=0 #a vector indicating how much the bars should be shifted relative to the x axis
19. add=FALSE #是否将该柱形图加入已存在的散点图默认为否

```

- spinegrams: 一种特殊的柱状图 `spine()` function 柱高为1, 堆叠柱状图, 通过比例来展示不同

6.2 Pie charts

```

1. pie(x, labels = names(x), #或者可以为c(),
2.      col = rainbow(length(x)) #使用颜色中的rainbow来填充
3.      main = "TITLE",
4.      edges = 200, #边数, 默认200边近似圆形, 而edge=length(x)则给出多边形
5.      radius = 0.8, #调整大小, pie图画在一个边长为[-1,1]的正方形内
6.      #如果标签过长, 需要让直径小一些
7.      clockwise = FALSE, #默认逆时针
8.      init.angle = if(clockwise) 90 else 0,
9.      #起始角, 默认为3点钟方向, 如果顺时针为TRUE则默认为12点钟方向
10.     density = NULL, angle = 45, , border = NULL, #跟前面一样有关暗线的
    参数
11.     lty = NULL , ...)

```

- `pie3D()/fanplot()` 3D pie/扇形图, 需要先 `library(plotrix)`

6.3 Histograms直方图

```

1. hist(x, breaks = "Sturges", #可以指定直方柱数量, breakpoints基于pretty函数
2.      freq = NULL, #FALSE纵坐标密度, TRUE纵坐标频率
3.      #仅有一种情况freq默认为TRUE: breaks为等距且下面的probability未指明
4.      probability = !freq, #!freq的别名
5.      include.lowest = TRUE,
6.      right = TRUE, # the histogram cells are right-closed (left open)
    intervals
7.      density = NULL, angle = 45, col = NULL, border = NULL,
8.      main = paste("Histogram of" , xname), xlab , ylab,

```



```

9.         xlim = range(breaks), ylim = NULL, #坐标系范围
10.        axes = TRUE, plot = TRUE, labels = FALSE,
11.        nclass = NULL, warn.unused = TRUE, ...)

```

- 举个栗子：

- 给出直方图
- 添加密度曲线density curve
- 添加rug plots

```

1. hist(mtcars$mpg, freq=FALSE, breaks=12, col="yellow",
2.      xlab="Miles Per Gallon", main="FUCK YOUR ASSHOLE")
3.
4. rug(jitter(mtcars$mpg)) #这个函数是下一张密度点的一种
5. #一维函数rugplots可以给出在不同区间有多少个x, 如果样本太多则jitter大致数量
6.
7. lines(density(mtcars$mpg), col="red", lwd=2) #密度曲线

```

6.4 Kernel density plots

```
`plot/line(density(x))`
```

- 举个栗子：

```

1. d<-density(mtcars$mpg)
2. plot(d, main="son of bitch")
3. polygon(d, col="red", border="blue") #填充红色, 边界线蓝色
4. rug(mtcars$mpg, col="brown") #棕色密度曲线

```

- 再举个栗子：把不同直方图的密度点叠加在一张图上

```

1. #需要用到sm.density.compare(x, factor) 新的package sm
2. par(lwd=2) #将线宽加倍, 更易于观察
3. library(sm)
4. attach(mtcars)
5. cyl.f<-factor(cyl, levels=c(4, 6, 8),
6.              labels=c("4 cylinder, 6 cylinder, 8 cylinder")
7.              ) #将cyl转换为factor
8.
9. sm.density.compare(mpg, cyl, xlab="MPG") #将4轮、6轮、8轮车的MPG按密度点展示出来

```

```

10. title(main="FUCK YOU")
11.
12. colfill<-c(2:(1+length(levels(cyl.f)))) #这里colfill=c(2:4)
13. legend(locator(1), levels(cyl.f), fill=colfill) #添加标签
14. detach(mtcars)

```

6.5 Box plots 箱图

- 箱图可以给出minimum , lower quartile(25%分度点),median,upper quartile(75%分度点),maximum
- 一个简单的箱图：

```

1. boxplot(mtcars$mpg, main="", ylab="")

```

6.5.1 使用平行箱图来比较组间数据parallel box

```

1. boxplot(mpg~cyl, #给出一个formula,
2.         data=mtcars,
3.         main,xlab="cyl",ylab="mpg",col="red"
4.         varwidth=TRUE, #箱图宽度与样本大小的平方根成比例
5.         horizontal=TRUE #水平坐标系
6.         notch=TRUE #有凹口的箱图
7.         ) #更细化的栗子可见P136 Figure 6.14

```

6.5.2 violin plots `vioplot()` function

```

1. vioplot(x1, x2, ... , names=, col=)

```

6.6 Dot plots `dotchart()` function

- 适用于：将大量带标签的值展示在水平坐标轴图中

```

1. dotchart(mtcars$mpg, labels=row.names(mtcars),...)

```

- 举个栗子：将标签点排序输出，且对不同的范围给予不同的颜色

```

1. x<- mtcars[order(mtcars$mpg)] #按mpg进行排序
2. x$cyl<-factor(x$cyl) #创建factor

```

```

3. x$color[x$cyl==4] <- "red"
4. x$color[x$cyl==6] <- "blue"
5. x$color[x$cyl==8] <- "darkgreen"
6. dotchart(x$mpg,
7.          labels=row.names(x),
8.          cex=.7, pch=19 #之前学过的点大小,点类型
9.          groups=x$cyl,
10.         color=x$color,
11.         main="", xlab="")

```

Chapter 7 Basic statics

- Descriptive statistics
- Frequency and contingency tables
- Correlations and covariances
- t-tests
- Nonparametric statistics

7.1 Descriptive statistics描述统计

- 基本任务：
 - mtcars dataset ;
 - examine descriptive statistics by transmission type (am) and number of cylinders (cyl) ;
 - 传动方式am : [0:automatic,1:manual];cyl:4,5,6

7.1.1 A menagerie of methods

`summary()/apply()/sapply()/fivenum`

- you can get mean/sd/var/min/max/median/length/range/quantile
 - skew/偏度：用于衡量分布的不对称程度或偏斜程度的指标
 - α =变量的三阶中心动差除以标准差的三次方
- ```
skew<-sum((x-m)^3/s^3)/n
```
- m为均值，s为标准差
- $\alpha$ 图像为以x均值为基准线向左边或右边偏斜的曲线
  - $\alpha > 0$ 正偏斜；

- $\alpha < 0$  负偏斜；
- $\alpha = 0$ ，对称的正态分布

○ kurtosis峰度：衡量分布的集中程度或分布曲线的尖锐程度的指标

- $\beta$  = 四阶中心动差除以标准差的四次方(即方差平方)的差-3
- `kurt<-sum((x-m)^4/s^4)/n-3` m为均值，s为标准差
- $\beta$  图像为以x均值为中心的对称曲线，区别为尖锐/平缓
- $\beta > 0$ , 分布比正态分布更集中在平均数周围，更尖锐
- $\beta = 0$ ，正态分布
- $\beta < 0$ , 分布比正态分布更分散，更平缓

● 举个栗子：建立函数来求取上述统计量

```
1. mystats<-function(x, na.omit=FALSE) {
2. if (na.omit)
3. x<-x[!is.na(x)]
4. m<-mean(x)
5. n<-length(x)
6. s<-sd(x)
7. skew<-sum((x-m)^3/s^3)/n
8. kurt<-sum((x-m)^4/s^4)/n-3
9. return(c(n=n, mean=m, stdev=s, skew=skew, kurtosis=kurt))
10. }
11.
12. sapply(mtcars[vars], mystats)
```

● 这里推荐了几个外部的packages:Hmisc/pastecs/psych

```
1. describe() #Hmisc, 返回变量数/样本数/缺失值数/均值/分位/5个最大最小值
2. #psych package也有一个同名函数
3.
4. stat.desc() #pastecs, 同样可以给出 descriptive statistics
```

## 7.1.2 Descriptive statistics by group

● 使用aggregate():

```
1. aggregate(mtcars[vars], by=list(am=mtcars$am), mean)
2. aggregate(mtcars[vars], by=list(am=mtcars$am), sd)
3. #不足：aggregate() 仅能使用单值函数如mean，无法一次返回多个统计学结果
```

- alternative : `by(data, INDICES, FUN) #INDICES向量用于定义groups`

```
1. dstats <- function(x) (c(mean=mean(x), sd=sd(x)))
2. by(mtcars[vars], mtcars$am, dstats)
```

- other alternatives:

```
1. summaryBy() in doBy package
2. describe.by() in psych package
3. melt+cast in reshape package #可参见5.4.3
```

### 7.1.3 Visualizing results 见chapter 6

## 7.2 Frequency and contingency tables

### 7.2.1 fenerating frequency tables

```
1. table()/xtabs() #生成频数表, 交叉表
2. prop.table(table, margins) #将频数表转换为比例
3. margin.table(table, margins) #得到表中某个变量的频数
4. addmargins(table, margins) #增加一行 summary margins (sums by default)
5. ftable(table) #将多维表格转化为 compact and attractive manner
```

- 举个栗子说明以上函数 : ONE-WAY TABLES

```
1. mytable <- with(Arthritis, table(Improved)) #生成简单的频数表
2. prop.table(mytable) #将频数按比例显示(0.500, 0.167...)
3. prop.table(mytable)*100 #比例*100, 便于观察
```

- 再举个栗子 : TWO-WAY TABLES

```
1. mytable <- table(mpg, cyl)
2. mytable <- xtabs(~ mpg + cyl, data=mydata) #需要交叉的数据在~右边, 这里mpg
 为行名, cyl为列名
3. margin.table(mytable, 1) #得到某个变量的频数, 1代表表中第一个变量, 输出不同mpg的
 counts
4. prop.table(mytable, 1)
5.
6. addmargins(mytable) #在交叉表里给所有行列增加总数值
7. addmargins(prop.table(mytable, 1), 2) #仅仅给出行变量(1)的总值, 输出为一新列
8. addmargins(prop.table(mytable, 2), 1) #仅仅给出列变量(2)的总值, 输出为一新行
```

```

9.
10. #alternatives:
11. CrossTable() in gmodels package

```

- 最后一个栗子：多维表格

```

1. mytable <- table(mpg,cyl,hp)
2. ftable(mytable) #三维表扁平化
3. margin.table(mytable,c(1,2)) #给出第1,2个变量的counts
4. ftable(addmargins(prop.table(mytable, c(1, 2)), 3)) #添加第1,2个变量的sum

```

## 7.2.2 Test of independence in 3 methods

- chi-square test of independence

```
chisq.test(mytable) #给出X-squared, df, p-value
```

- Fisher exact test

```
fisher.test(), 给出p+alternative hypothesis
```

- Cochran-Mantel-Haenszel test

```
mantelhaen.test(), 给出Cochran-Mantel-Haenszel M^2, df, p
```

## 7.2.3 Measures of association

- assocstats() in vcd package

- phi coefficient+contingency coefficient+Cramer' s V for a two-way table
- In general, larger magnitudes indicated stronger associations

- kappa() in vcd package

- Cohen' s kappa and weighted kappa for a a confusion matrix

## 7.2.4 visualisation (chapter 6+mosaic and association plots in chapter 11)

## 7.2.5 converting tables to flat files(自己写一个函数)

```

1. table2flat <- function(mytable) {
2. df <- as.data.frame(mytable)
3. rows <- dim(df)[1]
4. cols <- dim(df)[2]
5. x <- NULL
6. for (i in 1:rows){
7. for (j in 1:df$Freq[i]){
8. row <- df[i,c(1:(cols-1))]
9. x <- rbind(x,row)
10. }

```

```

11. }
12. row.names(x) <- c(1:dim(x)[1])
13. return(x)
14. }

```

- 举个栗子来实现这个函数：

```

1. treatment <- rep(c("Placebo", "Treated"), times=3)
2. improved <- rep(c("None", "Some", "Marked"), each=2)
3. Freq <- c(29,13,7,17,7,21)
4. mytable <- as.data.frame(cbind(treatment, improved, Freq))
5. mydata <- table2flat(mytable)

```

## 7.3 Correlations

### 7.3.1 type of correlation `spearson/spearman/kendall`

```

1. cor(x, #your dataframe, or simply cor(a,b)
2. use= , #对缺失数据的操作, 包括all.obs (assumes no missing data)
3. #everything (any correlation involving a case with missing
values will be set to missing);
4. #complete.obs (listwise deletion);
5. #pairwise.complete.obs (pairwise deletion)
6. method= #the type of correlation: spearson/spearman/kendall
7.) #default: everything+pearson
8. cov() #covariances 协方差
9. #notice: 显著性需要参考7.3.2

```

- Other types of correlations:
  - partial correlations 偏相关 `pcor()` in `ggm` package
  - correlation with 2 quantitative variables, controlling for one or more other quantitative variables  
即将其他变量都看做常数，就找两个变量的相关性

```

1. pcor(u, # numeric vector, 前两个数字是打算找相关性的, 后面为控制变量
2. S) #covariance matrix among the variables
3. hetcor() in polycor package

```

- Pearson product-moment correlations #numeric variables

- polyserial correlations 多项相关 between numeric and ordinal variables(有序变量)
- polychoric correlations 多序频相关,ordinals
- tetrachoric correlations 四分相关between two dichotomous variables  
假定有序变量和二分变量从正态分布中被移除

### 7.3.2 testing correlations for [significance]

- 对于null hypothesis：假定变量是彼此独立无关的 `cor.test()`

```
1. cor.test(x, y,
2. alternative = "two.side/less/greater", #two-tailed/one-tailed test
3. #hepothesis:correlation<0时, "less";同理, >0时用"greater"
4. #two.side is the default 默认双侧检验
5. method="pearson/kendall/spearman"
6.)
7. #不足: only one correlation per time
```

- `corr.test()` in psych package
  - correlations and significance levels for [matrices] of P/S/K

```
1. library(psych)
2. corr.test(data,
3. use="complete/pairwise", #整列或成对删除缺失值
4. method) #默认为pearson
```

- `r.test()` in psych package
  - significance of a correlation coefficient;
  - difference between 2 independent correlations ;
  - difference between 2 dependent correlations sharing 1 single variable ;
  - difference between 2 dependent correlations based on completely different variables ;

### 7.3.3 visualisation correlations (scatter plots/correlograms in cpter 11)

## 7.4 t tests

- 回归模型个体系数的显著性进行假设检验[样品与均值的差异]



- $|t|$ 越大-样品偏离均值越远,样品与均值的差别越显著
- 举几个应用的栗子 :
  - 接受新药相比接受原有治疗是否显著提高疗效 ?
  - 某个制造过程是否比其他过程缺陷率更低 ?

#### 7.4.1. Independent t-test独立t检验

- Are you more likely to be imprisoned if you commit a crime in the South?
  - 比较南方/非南方州被监禁的probability-two-group independent t-test
- hypothesis : 人口mean相等 , two groups are independent

```
1. t.test(y ~ x, #y:numeric,x: dichotomous二分变量;或者t.test(y1,y2)
2. alternative #同前
3. mu=0,#a number indicating the true value of the mean
4. paired = FALSE, #是否需要paired t-test对偶t检验
5. var.equal = FALSE, #默认var不等, 且使用 Welch degrees of freedom modification
6. #if TRUE,方差相等and specify a pooled var estimate联合方差估计法
7. conf.level = 0.95,
8. data)
```

- 实现实例 :

```
1. head(UScrime)
2. t.test(Prob ~ So, data=UScrime) #入狱概率~南方州
3. #结论:t大p小, 不接受原假设, 南方/非南方州入狱概率不等
```

#### 7.4.2 Dependent t-test

- if unemployment rate for younger males is greater than for older males the two groups aren' t independent
- hypothesis : 不同组的差异是正态分布的

```
t.test(y1, y2, paired=TRUE)
```

- 实现实例 :

```
1. sapply(UScrime[c("U1","U2")], function(x) (c(mean=mean(x),sd=sd(x))))
2. with(UScrime, t.test(U1, U2, paired=TRUE))
3. #结论: 差异显著, 不接受原假设, 年轻人就业率更高
```

### 7.5 Nonparametric tests of group differences

- 在总体方差未知或知道甚少的情况下，利用样本数据对总体分布形态等进行推断
- 不涉及有关总体分布的参数

### 7.5.1 Comparing two groups#note that these 2 groups are independent

- Wilcoxon rank sum test(Mann–Whitney U test): 用于检验观察值是否取自同一概率分布(某一组取得更好结果的可能性是否大于另一组)

```
wilcox.test(y ~ x, data)/wilcox.test(y1, y2) #其余参数跟前面类似
```

- 这里举个栗子：之前南方州入狱问题：

```
1. with(UScrime, by(Prob, So, median)) #给出median
2. wilcox.test(Prob ~ So, data=UScrime) #结果p很小, 不接受原假设
3. #若t-tests的假定很reasonable, 参数检验更易找到差异
4. #但当假定unreasonable时, 非参数检验更适合(如有序排列的数据)
```

### 7.5.2 Comparing more than 2 groups

- 比较四个地区的文盲率——需要进行one-way design单项设计
- Kruskal–Wallis test : non-paramater, independent groups

```
kruskal.test(y ~ A, data) #A:grouping variable with >= 2 levels
```

- Friedman test: 感觉跟上面那货差不多

```
friedman.test(y ~ A | B, data) #B:blocking variable限制条件
```

- 实现实例：结论，文盲率并不一样(p很小)

```
1. states <- as.data.frame(cbind(state.region, state.x77))
2. kruskal.test(Illiteracy ~ state.region, data=states)
3. #以上两个test的不足：无法了解which regions differ signifly froms other
```

- Mann–Whitney U test: `npmc()` in `npmc` package

```
1. class <- state.region
2. var <- state.x77[,c("Illiteracy")]
3. mydata <- as.data.frame(cbind(class, var))
4. rm(class, var)
5. A <- summary(npmc(mydata), type="BF") #给出所有组两两对比的结果C(n, 2)
6. B <- aggregate(mydata, by=list(mydata$class), median)
7.
8. #A：南方-其他的p明显要小, 而其他地区的两两对比差异不大
9. #B：南方中位文盲率更高
```

- `anova()` 比较几个模型哪个更合适

## Part III Intermediate methods