

第六章 S语言介绍

S语言是由AT&T贝尔实验室开发的一种用来进行数据探索、统计分析、作图的解释型语言。它的丰富的数据类型(向量、数组、列表、对象等)特别有利于实现新的统计算法,其交互式运行方式及强大的图形及交互图形功能使得我们可以方便地探索数据。

目前S语言的实现版本主要是S-PLUS。它基于S语言,并由Insightful公司进一步扩充和完善。作为统计学家及一般研究人员的通用方法工具箱, S强调数据可视化、数据探索、统计建模、统计方法的编程,以及可扩展性。

S-PLUS可以直接用来进行标准的统计分析得到所需结果,但是它的主要的特点是它可以交互地从各个方面去发现数据中的信息,并可以很容易地实现一个新的统计方法。

S-PLUS有微机版本和 workstation 版本,它是一个商业软件。Auckland大学的Robert Gentleman 和Ross Ihaka 及其他志愿人员开发了一个R系统,其语法形式与S语言基本相同,但实现不同,两种语言的程序有一定的兼容性。R是一个GPL自由软件,现在的版本是1.8.1版,在几乎所有操作系统都可以得到。R比S-PLUS还少许多功能,但已经具有了很强的实用性。R编程计算速度经常比S-PLUS快。

我们在这里尽量介绍S-PLUS和R都能使用的功能,且以R为主。下面我们统称S-PLUS和R。

6.1 S快速入门

S的基本界面是一个交互式命令窗口,命令提示符是一个大于号,命令的结果马上显示在命令下面。S命令主要有两种形式:表达式或赋值运算(用<-表示)。在命令提示符后键入一个表达式表示计算此表达式并显示结果。赋值运算把赋值号右边的值计算出来赋给左边的变量。可以用向上光标键来找回以前运行的命令再次运行或修改后再运行。

S是区分大小写的,所以x和X是不同的名字。

我们用一些例子来看S的特点。假设我们已经进入了R(或S-PLUS)的交互式窗口。输入以下语句(大于号是系统提示,不是自己输入的):

```
> x1 <- 0:100
> x2 <- x1*2*pi/100
> y <- sin(x2)
> plot(x,y, type='l')
```

这些语句可以绘制正弦曲线图。其中,“<-”是赋值运算符。0:100表示一个从0到100的等差数列向量。从第二个语句可以看出,我们可以对向量直接进行四则运算,计算得到的x2是向量x1的所有元素乘以常数 $2\pi/100$ 的结果。从第

三个语句可以看到函数可以以向量为输入, 并可以输出一个向量, 结果向量y的每一个分量是自变量x2的每一个分量的正弦函数值。从最后一个语句可以看出函数的调用也很自由, 可以按位置给出自变量, 也可以用“自变量名=”的形式指定自变量值, 这样可以使用缺省值。

下面我们看一看S的统计功能。

```
> marks <- c(10, 6, 4, 7, 8)
> mean(marks)
[1] 7
> sd(marks)
[1] 2.236068
> median(marks)
[1] 7
> min(marks)
[1] 4
> max(marks)
[1] 10
> boxplot(marks)
>
```

第一个语句输入若干数据到一个向量, 函数c()用来把数据组合为一个向量。后面用了几个函数来计算数据的均值、标准差、中位数、最小值、最大值。最后的函数绘制数据的盒形图。例中sd()是R中才有的函数, 在S-PLUS中要用sqrt(var())来计算。在S命令方式中要显示一个表达式的值只要键入它。

为了演示S的回归计算, 我们把SAS中的例子数据SASUSER.CLASS 输出到了一个文本文件CLASS.TXT, SAS程序如下:

```
data _null_;
  set sasuser.class;
  file 'class.txt';
  put name sex age height weight;
run;
```

我们在R中把这个文件读入为一个数据框(data frame, 相当于SAS中的数据集), 并进行回归, 绘制数据散点图和回归直线。程序及结果如下, 注意第二行的加号是系统给出的续行提示。

```
> cl <- read.table("c:/work/class.txt", as.is=c(1),
+ col.names=c("Name", "Sex", "Age", "Height", "Weight"))
> cl
```

	Name	Sex	Age	Height	Weight
1	Alice	F	13	56.5	84.0
2	Becka	F	13	65.3	98.0
3	Gail	F	14	64.3	90.0
4	Karen	F	12	56.3	77.0
5	Kathy	F	12	59.8	84.5
6	Mary	F	15	66.5	112.0
7	Sandy	F	11	51.3	50.5
8	Sharon	F	15	62.5	112.5

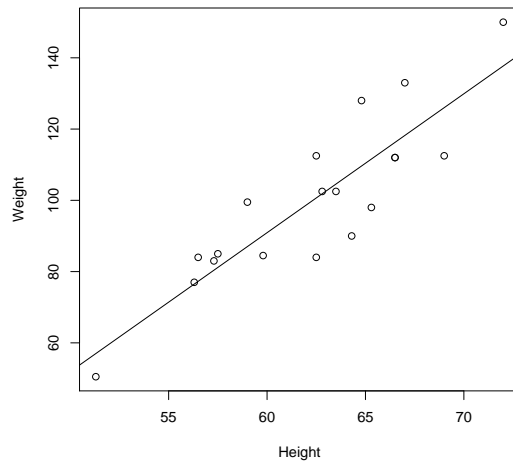


图 6.1: 用R产生的回归图形

```

9   Tammy   F   14   62.8  102.5
10  Alfred  M   14   69.0  112.5
11   Duke   M   14   63.5  102.5
12  Guido   M   15   67.0  133.0
13  James   M   12   57.3   83.0
14 Jeffrey  M   13   62.5   84.0
15   John   M   12   59.0   99.5
16  Philip  M   16   72.0  150.0
17  Robert  M   12   64.8  128.0
18  Thomas  M   11   57.5   85.0
19 William  M   15   66.5  112.0

```

```

> attach(cl)
> plot(Height, Weight)
> fit1 <- lm(Weight ~ Height)
> fit1

```

```

Call:
lm(formula = Weight ~ Height)

```

```

Coefficients:
(Intercept)      Height
    -143.027      3.899

```

```

> abline(fit1)
>

```

结果可以产生图6.1。我们从这个例子可以看到，分析由函数完成，结果也是对象，可以作为进一步分析的输入。

S-PLUS也提供了一般的计算功能。比如, 求一个矩阵的逆:

```
> A <- matrix(c(1,2, 7,3), ncol=2, byrow=T)
> A
      [,1] [,2]
[1,]    1    2
[2,]    7    3
> Ai <- solve(A)
> Ai
      [,1] [,2]
[1,] -0.2727273  0.18181818
[2,]  0.6363636 -0.09090909
```

还可以进行矩阵运算, 如:

```
> b <- c(2,3)
> x <- Ai %*% b
> x
      [,1]
[1,] -2.220446e-16
[2,]  1.000000e+00
```

这实际上是解了一个线性方程组。“%*%”表示矩阵乘法。

可以把若干行命令保存在一个文本文件(比如C:\WORK\MYPROG.R)中, 然后用source函数来运行整个文件:

```
> source("C:\\WORK\\MYPROG.R")
```

注意字符串中的反斜杠要写成两个。用sink()函数可以把以后的输出从屏幕窗口转向到一个外部文本文件, 例如:

```
> sink("C:\\WORK\\MYPROG.OUT")
```

要恢复输出到屏幕窗口, 使用:

```
> sink()
```

在S中可以用“?”号后面跟要查询的函数来显示帮助信息, 例如

```
> ?c
> ?"=="
```

(查询特殊符号要用撇号括起来)。另外, 在MS Windows版的S-PLUS中可以用Windows帮助来得到帮助, 而在R中可以用HTML浏览器来查询帮助(用“help.start()函数或Help - R language (html)”启动)。用好S语言重要的一点就是要用熟系统的帮助功能, 因为S的教材或手册都不能列出所有的细节, 如果列出所有细节则过于繁琐了。

要退出S-PLUS或R, 可以用q()函数, 也可以用菜单命令。R在退出时提问是否保存当前工作空间, 它可以把当前定义的所有对象(有名字的向量、矩阵、列表、函数等)保存到一个文件。S-PLUS可以选择保存哪些对象。

6.2 S向量

S语言是基于对象的语言, 不过它的最基本的数据还是一些预先定义好的数据类型, 如向量、矩阵、列表(list)等。更复杂的数据用对象表示, 比如, 数据框对象, 时间序列对象, 模型对象, 图形对象, 等等。

S语言表达式可以使用常量和变量。变量名的规则是: 由字母、数字、句点组成, 第一个字符必须是字母, 长度没有限制。大小写是不同的。特别要注意句点可以作为名字的合法部分, 而在其它面向对象语言中句点经常用来分隔对象与成员名。另外, 下划线不能用在名字中, 因为它是赋值符号“<-”的缩写。

6.2.1 常量

常量可以笼统地分为逻辑型、数值型和字符型三种, 实际上数值型数据又可以分为整型、单精度、双精度等, 非特殊需要不必太关心其具体类型。例如, 123, 123.45, 1.2345e30是数值型常量, “Weight”, “李明”是字符型(用两个双撇号或两个单撇号包围)。逻辑真值写为TRUE(注意区分大小写, 写true没意义), 逻辑假值写为FALSE。复数常量就用3.5 - 2.1i这样的写法表示。

S中的数据可以取缺失值, 用符号NA代表缺失值。函数is.na(x)返回x是否缺失值(真还是假)。

6.2.2 向量(Vector)与赋值

向量是具有相同基本类型的元素序列, 大体相当于其他语言中的一维数组。实际上在S中标量也被看作是长度为1的向量。

定义向量的最常用办法是使用函数c(), 它把若干个数值或字符串组合为一个向量, 比如:

```
> marks <- c(10, 6, 4, 7, 8)
> x <- c(1:3, 10:13)
> x
[1] 1 2 3 10 11 12 13
> x1 <- c(1, 2)
> x2 <- c(3, 4)
> x <- c(x1, x2)
> x
[1] 1 2 3 4
```

在显示向量值时我们注意到左边总出现一个“[1]”, 这是代表该显示行的第一个数的下标, 例如:

```
> 1234501:1234520
[1] 1234501 1234502 1234503 1234504 1234505 1234506 1234507 1234508 1234509
[10] 1234510 1234511 1234512 1234513 1234514 1234515 1234516 1234517 1234518
[19] 1234519 1234520
```

第二行输出从第10个数开始, 所以在行左边显示“[10]”。

S中用符号“<-”(这是小于号紧接一个减号)来为变量赋值。另一种赋值的办法是用assign函数, 比如

```
> x1 <- c(1, 2)
```

和

```
> assign("x1", c(1, 2))
```

效果相同。

函数length(x)可以计算向量x的长度。

6.2.3 向量运算

可以对向量进行加(+)减(-)乘(*)除(/)、乘方(^)运算, 其含意是对向量的每一个元素进行运算。例如:

```
> x <- c(1, 4, 6.25)
> y <- x*2+1
> y
[1] 3.0 9.0 13.5
```

另外, %/%表示整数除法(比如5 %/% 3为1), %%表示求余数(如5 %% 3为2)。

也可以用向量作为函数自变量, sqrt、log、exp、sin、cos、tan等函数都可以用向量作自变量, 结果是对向量的每一个元素取相应的函数值, 如:

```
> sqrt(x)
[1] 1.0 2.0 2.5
```

函数min和max分别取自变量向量的最小值和最大值, 函数sum计算自变量向量的元素和, 函数mean计算均值, 函数var计算样本方差, 函数sd计算标准差(在Splus中用sqrt(var())计算), 函数range返回包含两个值的向量, 第一个值是最小值, 第二个值是最大值。例如:

```
> max(x)
[1] 6.25
```

如果求var(x)而x是 $n \times p$ 矩阵, 则结果为样本协方差阵。

sort(x)返回x的元素从小到大排序的结果向量。order(x)返回使得x从小到大排列的元素下标向量(x[order(x)]等效于sort(x))。

任何数与缺失值的运算结果仍为缺失值。例如,

```
> 2*c(1, NA, 2)
[1] 2 NA 4
> sum(c(1, NA, 2))
[1] NA
```

6.2.4 产生有规律的数列

在S中很容易产生一个等差数列。例如, 1:n产生从1到n的整数列, -2:3产生从-2到3的整数列, 5:2产生反向的数列:

```
> n <- 5
> 1:n
[1] 1 2 3 4 5
> -2:3
[1] -2 -1 0 1 2 3
> 5:2
[1] 5 4 3 2
```

要注意1:n-1不是代表1到n-1而是向量1:n减去1, 这是一个常犯的错误:

```
> 1:n-1
[1] 0 1 2 3 4
> 1:(n-1)
[1] 1 2 3 4
```

seq函数是更一般的等差数列函数。如果只指定一个自变量 $n > 0$, 则seq(n)相当于1:n。指定两个自变量时, 第一自变量是开始值, 第二自变量是结束值, 如seq(-2,3)是从-2到3。S函数调用的一个很好的特点是它可以使用不同个数的自变量, 函数可以对不同类型的自变量给出不同结果, 自变量可以用“自变量名=自变量值”的形式指定。例如, seq(-2,3)可以写成seq(from=-2, to=3)。可以用一个by参数指定等差数列的增加值, 例如:

```
> seq(0, 2, 0.7)
[1] 0.0 0.7 1.4
```

也可以写成seq(from=0, to=2, by=0.7)。有参数名的参数的次序任意, 如:

```
> seq(0, by=0.7, to=2)
[1] 0.0 0.7 1.4
```

可以用length参数指定数列长度, 如seq(from=10, length=5)产生10到14。seq函数还可以用一种seq(along=向量名)的格式, 这时只能用这一个参数, 产生该向量的下标序列, 如:

```
> x
[1] 1.00 4.00 6.25
> seq(along=x)
[1] 1 2 3
```

另一个类似的函数是rep, 它可以重复第一个自变量若干次, 例如:

```
> rep(x,3)
[1] 1.00 4.00 6.25 1.00 4.00 6.25 1.00 4.00 6.25
> rep(c(100, 200), c(3, 2))
[1] 100 100 100 200 200
```

第一个参数名为x, 第二个参数名为times(重复次数)。

6.2.5 逻辑向量

向量可以取逻辑值, 如:

```
> l <- c(TRUE, TRUE, FALSE)
> l
[1] TRUE TRUE FALSE
```

当然, 逻辑向量是一个比较的结果, 如:

```
> x
[1] 1.00 4.00 6.25
> l <- x > 3
> l
[1] FALSE TRUE TRUE
```

一个向量与常量比较大小, 结果还是一个向量, 元素为每一对比较的结果逻辑值。

两个向量也可以比较, 如:

```
> log(10*x)
[1] 2.302585 3.688879 4.135167
> log(10*x) > x
[1] TRUE FALSE FALSE
```

比较运算符包括<, <=, >, >=, ==(相等), !=(不等)。

两个逻辑向量可以进行与(&)、或(|)运算, 结果是对应元素运算的结果。对逻辑向量x计算!x表示取每个元素的非。例如:

```
> (x > 1.5) & log(10*x)>1.5
[1] FALSE FALSE TRUE
```

判断一个逻辑向量是否都为真值的函数是all, 如:

```
> all(log(10*x) > x)
[1] FALSE
```

判断是否其中有真值的函数是any, 如:

```
> any(log(10*x) > x)
[1] TRUE
```

函数is.na(x)用来判断x的每一个元素是否缺失。如

```
> is.na(c(1, NA, 3))
[1] FALSE TRUE FALSE
```

逻辑值可以强制转换为整数值, TRUE变成1, FALSE变成0。例如, 我们以age>65为老年人, 否则为年轻人, 可以用c("young", "old")[(age>65)+1]这样的表示。当年龄大于65时age>65等于TRUE, 加1则把TRUE转换为数值型的1, 结果得2, 于是返回第二个下标处的“old”。否则等于0+1下标处的“young”。

6.2.6 字符型向量

向量元素可以取字符串值。例如：

```
> c1 <- c("x", "sin(x)")
> c1
[1] "x" "sin(x)"
> ns <- c("Weight", "Height", "年龄")
> ns
[1] "Weight" "Height" "年龄"
```

paste函数用来把它的自变量连成一个字符串, 中间用空格分开, 例如：

```
> paste("My", "Job")
[1] "My Job"
```

连接的自变量可以是向量, 这时各对应元素连接起来, 长度不相同且较短的向量被重复使用。自变量可以是数值向量, 连接时自动转换成适当的字符串表示, 例如：

```
> paste(c("X", "Y"), "=", 1:4)
[1] "X = 1" "Y = 2" "X = 3" "Y = 4"
```

分隔用的字符可以用sep参数指定, 例如下例产生若干个文件名：

```
> paste("result.", 1:5, sep=".")
[1] "result.1" "result.2" "result.3" "result.4" "result.5"
```

如果给paste()函数指定了collapse参数, 则paste()可以把一个字符串向量的各个元素连接成一个字符串, 中间用collapse指定的值分隔。比如paste(c('a', 'b'), collapse='.')得到'a.b'。

6.2.7 复数向量

S支持复数运算。复数常量只要用象3.5+2.1i这样的格式即可。复向量的每一个元素都是复数。可以用complex()函数生成复向量, 如给定实部和虚部向量定义复数向量：

```
> x <- (0:100)/100*2*pi
> y <- sin(x)
> z <- complex(re=x, im=y)
> plot(z)
```

也可以给定模和辐角定义复数向量：

```
> zz <- complex(mod=rep(1, 12), arg=(0:11)/12*2*pi)
> plot(zz)
```

Re()计算实部, Im()计算虚部, Mod()计算复数模, Arg()计算复数幅角, Conj()计算共轭。基本的数学函数也支持复数运算, 比如:

```
> sqrt(-2+0i)
[1] 0+1.414214i
```

注意自变量是复数时sqrt等函数才进行复数运算。

6.2.8 向量下标运算

S提供了十分灵活的访问向量元素和向量子集的功能, 有效运用这些功能才能发挥S语言的优势。

某一个元素只要用x[i]的格式访问, 其中x是一个向量名, 或一个取向量值的表达式, i是一个下标, 如:

```
> x
[1] 1.00 4.00 6.25
> x[2]
[1] 4
> (c(1, 3, 5) + 5)[2]
[1] 8
```

可以单独改变一个元素的值, 例如:

```
> x[2] <- 125
> x
[1] 1.00 125.00 6.25
```

事实上, S提供了四种方法来访问向量的一部分, 格式为x[v], x为向量或向量值的表达式, v是如下的一种下标向量:

一、取正整数值的下标向量

在x[v]中, v为一个向量, 元素取值在1到length(x)之间, 取值允许重复, 例如,

```
> x[c(1,3)]
[1] 1.00 6.25
> x[1:2]
[1] 1 125
> x[c(1,3,2,1)]
[1] 1.00 6.25 125.00 1.00
> c("a", "b", "c")[rep(c(2,1,3), 3)]
[1] "b" "a" "c" "b" "a" "c" "b" "a" "c"
```

二、取负整数值的下标向量

在 $x[v]$ 中, v 为一个向量, 元素取值在 $-\text{length}(x)$ 到 -1 之间, 表示扣除相应位置的元素。例如:

```
> x[-(1:2)]  
[1] 6.25
```

三、取逻辑值的下标向量

在 $x[v]$ 中, v 为和 x 等长的逻辑向量, $x[v]$ 表示取出所有 v 为真值的元素, 如:

```
> x  
[1] 1.00 125.00 6.25  
> x<10  
[1] TRUE FALSE TRUE  
> x[x<10]  
[1] 1.00 6.25  
> x[x<0]  
numeric(0)
```

可见 $x[x<10]$ 取出所有小于10的元素组成的子集。这种逻辑值下标是一种强有力的检索工具, 例如 $x[\sin(x)>0]$ 可以取出 x 中所有正弦函数值为正的元素组成的向量。

如果下标都是假值则结果是一个零长度的向量, 显示为`numeric(0)`。

四、取字符型值的下标向量

在定义向量时可以给元素加上名字, 例如:

```
> ages <- c(Li=33, Zhang=29, Liu=18)  
> ages  
  Li Zhang  Liu  
  33    29   18
```

这样定义的向量可以用通常的办法访问, 另外还可以用元素名字来访问元素或元素子集, 例如:

```
> ages["Zhang"]  
Zhang  
    29  
> ages[c("Li", "Liu")]  
Li Liu  
33 18
```

向量元素名可以后加, 例如:

```
> ages1 <- c(33, 29, 18)
> names(ages1) <- c("Li", "Zhang", "Liu")
> ages1
  Li Zhang  Liu
33   29   18
```

上面我们看到了如何访问向量的部分元素。在S中还可以改变一部分元素的值, 例如:

```
> x
[1] 1.00 125.00 6.25
> x[c(1,3)] <- c(144, 169)
> x
[1] 144 125 169
```

注意赋值的长度必须相同, 例外是可以把部分元素赋为一个统一值:

```
> x[c(1,3)] <- 0
> x
[1] 0 125 0
```

要把向量所有元素赋为一个相同的值而又不想改变其长度, 可以用`x[]`的写法:

```
> x[] <- 0
```

注意这与“`x <- 0`”是不同的, 前者赋值后向量长度不变, 后者使向量变为标量0。

改变部分元素值的技术与逻辑值下标方法结合可以定义向量的分段函数, 例如, 要定义 $y=f(x)$ 为当 $x<0$ 时取 $1-x$, 否则取 $1+x$, 可以用:

```
> y <- numeric(length(x))
> y[x<0] <- 1 - x[x<0]
> y[x>=0] <- 1 + x[x>=0]
```

6.3 多维数组和矩阵

6.3.1 数组(array)和矩阵(matrix)

数组(array)可以看成是带多个下标的类型相同的元素的集合, 常用的是数值型的数组如矩阵, 也可以有其它类型(如字符型、逻辑型、复型数组)。S可以很容易地生成和处理数组, 特别是矩阵(二维数组)。

数组有一个特征属性叫做维数向量(dim属性), 维数向量是一个元素取正整数值的向量, 其长度是数组的维数, 比如维数向量有两个元素时数组为二维数组(矩阵)。维数向量的每一个元素指定了该下标的上界, 下标的下界总为1。

一组值只有定义了维数向量(dim属性)后才能被看作是数组。比如:

```
> z <- 1:24
> dim(z) <- c(2, 3, 4)
```

这时z已经成为了一个维数向量为c(2,3,4)的三维数组。也可以把向量定义为一维数组, 例如:

```
> dim(z) <- 24
```

数组元素的排列次序缺省情况下是采用FORTRAN的数组元素次序(按列次序), 即第一下标变化最快, 最后下标变化最慢, 对于矩阵(二维数组)则是按行存放。例如, 数组z各元素次序为z[1,1,1], z[2,1,1], z[1,2,1], z[2,2,1], z[1,3,1], z[2,3,1], z[1,2,2], ..., a[2,3,4]。

用函数array()或matrix()可以直接定义数组。array()函数的完全使用为array(x, dim=length(x), dimnames=NULL), 其中x是第一自变量, 应该是一个向量, 表示数组的元素值组成的向量。dim参数可省, 省略时定义一维数组(但不同于向量)。dimnames属性可以省略, 不省略时是一个长度与维数相同的列表(list, 见后面), 列表的每个成员为一维的名字。例如上面的z可以这样定义:

```
> z <- array(1:24, dim=c(2,3,4))
```

函数matrix()用来定义最常用的一种数组: 二维数组, 即矩阵。其完全格式为

```
matrix(data = NA, nrow = 1, ncol = 1,
       byrow = FALSE, dimnames = NULL)
```

其中第一自变量data为数组的数据向量(缺省值为缺失值NA), nrow为行数, ncol为列数, byrow表示数据填入矩阵时按行次序还是列次序, 一定注意缺省情况下按列次序, 这与我们写矩阵的习惯是不同的。dimnames缺省是空值, 否则是一个长度为2的列表, 列表第一个成员是长度与行数相等的字符型向量, 表示每行的标签, 列表第二个成员是长度与列数相同的字符型向量, 表示每列的标签。例如, 定义一个3行4列, 由1:12按行次序排列的矩阵, 可以用:

```
> b <- matrix(1:12, ncol=4, byrow=TRUE)
> b
      [,1] [,2] [,3] [,4]
[1,]    1    2    3    4
[2,]    5    6    7    8
[3,]    9   10   11   12
```

注意在有数据的情况下只需指定行数或列数之一。指定的数据个数允许少于所需的数据个数, 这时循环使用提供的数据。例如:

```
> b <- matrix(0, nrow=3, ncol=4)
```

生成3行4列的元素都为0的矩阵。

6.3.2 数组下标

要访问数组的某个元素, 只要象上面那样写出数组名和方括号内用逗号分开的下标即可, 如`z[2,1,2]`。

更进一步我们还可以在每一个下标位置写一个下标向量, 表示对这一维取出所有指定下标的元素, 如`z[1, 2:3, 2:3]`取出所有第一下标为1, 第二下标为2或3, 第三下标为2或3的元素。注意因为第一维只有一个下标所以退化了, 得到的是一个维数向量为`c(2,2)`的数组。

另外, 如果略写某一维的下标, 则表示该维全选。例如, `z[1, ,]`取出所有第一下标为1的元素, 得到一个形状为`c(3,4)`的数组。`z[, 2,]`取出所有第二下标为2的元素得到一个形状为`c(2,4)`的数组。`z[1,1,]`则只能得到一个长度为4的向量, 不再是数组(`dim(z[1,1,])`值为NULL)。`z[, ,]`或`z[]`都表示整个数组。比如

```
z[] <- 0
```

可以在不改变数组维数的条件下把元素都赋成0。

还有一种特殊下标是对于数组只用一个下标向量(是向量, 不是数组), 比如`z[3:4]`, 这时忽略数组的维数信息, 把下标表达式看作是对数组的数据向量取子集。

6.3.3 不规则数组下标

在S中甚至可以把数组中的任意位置的若干个元素作为一个子集访问, 其方法是用一个矩阵作为数组的下标, 矩阵的每一行是一个元素的下标, 数组有几维下标矩阵的每一行就有几列。例如, 我们要把上面的形状为`c(2,3,4)`的数组`z`的第`[1,1,1]`, `[2,2,3]`, `[1,3,4]`, `[2,1,4]`号共四个元素作为一个整体访问, 先定义一个包含这些下标作为行的二维数组:

```
> b <- matrix(c(1,1,1, 2,2,3, 1,3,4, 2,1,4),
               ncol=3, byrow=T)
> b
      [,1] [,2] [,3]
[1,]    1    1    1
[2,]    2    2    3
[3,]    1    3    4
[4,]    2    1    4
> z[b]
[1] 1 16 23 20
```

注意取出的是一个向量。我们还可以对这几个元素赋值, 如:

```
> z[b] <- c(101, 102, 103, 104)
> z
```

或

```
> z[b] <- 0
> z
```

6.3.4 数组四则运算

数组可以进行四则运算(+, -, *, /, ^), 解释为数组对应元素的四则运算, 参加运算的数组一般应该是相同形状的(dim属性完全相同)。例如, 假设A, B, C是三个形状相同的数组, 则

```
> D <- C + 2*A/B
```

计算得到的结果是A的每一个元素除以B的对应元素再乘以2然后加上C的对应元素得到相同形状的数组。四则运算遵循通常的优先级规则。

形状不一致的向量和数组在少数情况下也可以进行四则运算, 一般的规则是数组的数据向量对应元素进行运算, 把短的循环使用来与长的匹配, 并尽可能保留共同的数组属性。例如:

```
> x1 <- c(100, 200)
> x2 <- 1:6
> x1+x2
[1] 101 202 103 204 105 206
> x3 <- matrix(1:6, nrow=3)
> x3
      [,1] [,2]
[1,]     1     4
[2,]     2     5
[3,]     3     6
> x1+x3
      [,1] [,2]
[1,]  101  204
[2,]  202  105
[3,]  103  206
```

除非你清楚地知道规则应避免使用这样的办法(标量与数组或向量的四则运算除外)。

6.3.5 矩阵运算

矩阵是二维数组, 但因为其应用广泛所以对它定义了一些特殊的运算和操作。

函数t(A)返回矩阵A的转置。nrow(A)为矩阵A的行数, ncol(A)为矩阵A的列数。

矩阵之间进行普通的加减乘除四则运算仍遵从一般的数组四则运算规则, 即数组的对应元素之间进行运算, 所以注意A*B不是矩阵乘法而是矩阵对应元素相乘。

要进行矩阵乘法, 使用运算符%*%, A %*% B表示矩阵A乘以矩阵B(当然要求A的列数等于B的行数)。例如:

```
> A <- matrix(1:12, nrow=4, ncol=3, byrow=T)
> B <- matrix(c(1,0), nrow=3, ncol=2, byrow=T)
> A
      [,1] [,2] [,3]
```

```

[1,] 1 2 3
[2,] 4 5 6
[3,] 7 8 9
[4,] 10 11 12
> B
      [,1] [,2]
[1,] 1 0
[2,] 1 0
[3,] 1 0
> A %*% B
      [,1] [,2]
[1,] 6 0
[2,] 15 0
[3,] 24 0
[4,] 33 0
>

```

另外, 向量用在矩阵乘法中可以作为行向量看待也可以作为列向量看待, 这要看哪一种观点能够进行矩阵乘法运算。例如, 设 x 是一个长度为 n 的向量, A 是一个 $n \times n$ 矩阵, 则“ $x \%*\% A \%*\% x$ ”表示二次型 $x'Ax$ 。但是, 有时向量在矩阵乘法中的地位并不清楚, 比如“ $x \%*\% x$ ”就既可能表示内积 $x'x$ 也可能表示 $n \times n$ 阵 xx' 。因为前者较常用, 所以S选择表示前者, 但内积最好还是用`crossprod(x)`来计算。要表示 xx' , 可以用“`cbind(x) \%*\% x`”或“`x \%*\% rbind(x)`”。

函数`crossprod(X, Y)`表示一般的交叉乘积(内积) $X'Y$, 即 X 的每一列与 Y 的每一列的内积组成的矩阵。如果 X 和 Y 都是向量则是一般的内积。只写一个参数 X 的`crossprod(X)`计算 X 自身的内积 $X'X$ 。注意即使结果是标量 S 仍会返回一个 1×1 矩阵。

其它矩阵运算还有`solve(A,b)`解线性方程组 $Ax = b$, `solve(A)`求方阵 A 的逆矩阵, `svd()`计算奇异值分解, `qr()`计算QR分解, `eigen()`计算特征向量和特征值。详见随机帮助, 例如:

```
> ?qr
```

函数`diag()`的作用依赖于其自变量。`diag(vector)`返回以自变量(向量)为主对角元素的对角矩阵。`diag(matrix)`返回由矩阵的主对角元素组成的向量。`diag(k)`(k 为标量)返回 k 阶单位阵。

6.3.6 矩阵合并与拉直

函数`cbind()`把其自变量横向拼成一个大矩阵, `rbind()`把其自变量纵向拼成一个大矩阵。`cbind()`的自变量是矩阵或者看作列向量的向量, 自变量的高度应该相等(对于向量, 高度即长度, 对于矩阵, 高度即行数)。`rbind`的自变量是矩阵或看作行向量的向量, 自变量的宽度应该相等(对于向量, 宽度即长度, 对于矩阵, 宽度即列数)。如果参与合并的自变量比其它自变量短则循环补足后合并。例如:

```

> x1 <- rbind(c(1,2), c(3,4))
> x1

```



```

      [,1] [,2]
[1,]    1    2
[2,]    3    4
> x2 <- 10+x1
> x3 <- cbind(x1, x2)
> x3
      [,1] [,2] [,3] [,4]
[1,]    1    2   11   12
[2,]    3    4   13   14
> x4 <- rbind(x1, x2)
> x4
      [,1] [,2]
[1,]    1    2
[2,]    3    4
[3,]   11   12
[4,]   13   14
> cbind(1, x1)
      [,1] [,2] [,3]
[1,]    1    1    2
[2,]    1    3    4

```

因为`cbind()`和`rbind()`的结果总是矩阵类型(有`dim`属性且为二维), 所以可以用它们把向量表示为 $n \times 1$ 矩阵(用`cbind(x)`)或 $1 \times n$ 矩阵(用`rbind(x)`)。

设 z 是一个数组, 要把它转化为向量(去掉`dim`和`dimnames`属性), 只要用函数`as.vector(z)`返回值就可以了(注意函数只能通过函数值返回结果而不允许修改它的自变量, 比如`t(X)`返回 X 的转置矩阵而 X 本身并未改变)。另一种由数组得到其数据向量的简单办法是使用函数`c()`, 例如`c(z)`的结果是 z 的数据向量。这样的另一个好处是`c()`允许多个自变量, 它可以把多个自变量都看成数据向量连接起来。例如, 设 A 和 B 是两个矩阵, 则`c(A,B)`表示把 A 按列次序拉直为向量并与把 B 按列次序拉直为向量的结果连接起来。一定注意拉直时是按列次序拉直的。

6.3.7 数组的维名字

数组可以有一个属性`dimnames`保存各维的各个下标的名字, 缺省时为`NULL`(即无此属性)。我们以矩阵为例, 它有两个维: 行维和列维。比如, 设 x 为2行3列矩阵, 它的行维可以定义一个长度为2的字符向量作为各行的名字, 它的列维可以定义一个长度为3的向量作为各列的名字, 属性`dimnames`是一个列表, 列表的每个成员是一个维名字的字符向量或`NULL`。例如:

```

> x <- matrix(1:6, ncol=2,
+ dimnames=list(c("one", "two", "three"),
+               c("First", "Second")),
+ byrow=T)
> x
      First Second
one      1      2
two      3      4
three    5      6

```

我们也可以先定义矩阵x然后再为dimnames(x)赋值。

对于矩阵, 我们还可以使用属性rownames和colnames来访问行名和列名。
如:

```
> x <- matrix(1:6, ncol=2, byrow=T)
> colnames(x) <- c("First", "Second")
> rownames(x) <- c("one", "two", "three")
```

在定义了数组的维名后我们对这一维的下标就可以用它的名字来访问, 例如:

```
> x[c("one", "three"),]
      First Second
one       1      2
three     5      6
```

6.3.8 数组的外积

两个数组a和b的外积是由a的每一个元素与b的每一个元素搭配在一起相乘得到一个新元素, 这样得到一个维数向量等于a的维数向量与b的维数向量连起来的向量, 即若d为a和b的外积, 则 $\text{dim}(d)=c(\text{dim}(a), \text{dim}(b))$ 。

a和b的外积记作a %o% b。如

```
> d <- a %o% b
```

也可以写成一个函数调用的形式:

```
> d <- outer(a, b, '*')
```

注意outer(a, b, f)是一个一般性的外积函数, 它可以把a的任一个元素与b的任一个元素搭配起来作为f的自变量计算得到新的元素值, 外积是两个元素相乘的情况。函数当然也可以是加、减、除, 或其它一般函数。当函数为乘积时可以省略不写。

例如, 我们希望计算函数 $z = \cos(y)/(1+x^2)$ 在一个x和y的网格上的值用来绘制三维曲面图, 可以用如下方法生成网格及函数值:

```
> x <- seq(-2, 2, length=20)
> y <- seq(-pi, pi, length=20)
> f <- function(x, y) cos(y)/(1+x^2)
> z <- outer(x, y, f)
> persp(x, y, z)
```

用这个一般函数可以很容易地把两个数组的所有元素都两两组合一遍进行指定的运算。

下面考虑一个有意思的例子。我们考虑简单的 2×2 矩阵 $\begin{pmatrix} a & b \\ c & d \end{pmatrix}$, 其元素均在0,1,...,9中取值。假设四个元素a, b, c, d都是相互独立的服从离散均匀分布的随机变量, 我们设法求矩阵行列式 $ad - bc$ 的分布。首先, 随机变量ad和bc同分布, 它的取值由以下外积矩阵给出, 每一个取值的概率均为1/100:

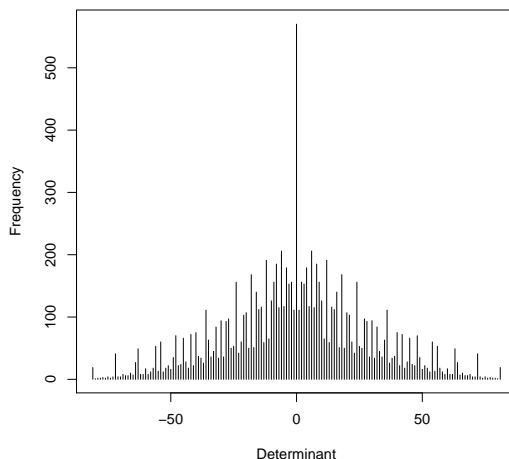


图 6.2: 随机矩阵行列式分布

```
> d <- outer(0:9, 0:9)
```

这个语句产生一个 10×10 的外积矩阵。

为了计算 ad 的100个值(有重复)与 bc 的100个值相减得到的10000个结果, 可以使用如下外积函数:

```
> d2 <- outer(d, d, "-")
```

这样得到一个维数向量为 $c(10,10,10,10)$ 的四维数组, 每一个元素为行列式的一个可能取值, 概率为万分之一。因为这些取值中有很多重复, 我们可以用一个`table()`函数来计算每一个值的出现次数(频数):

```
> fr <- table(d2)
```

得到的结果是一个带有元素名的向量`fr`, `fr`的元素名为行列式的一个取值, `fr`的元素值为该取值出现的频数。比如`fr[1]`的元素名为`-81`, 值为`19`, 表示值`-81`在数组`d2`中出现了19次。通过计算`length(fr)`可以知道共有163个不同值。还可以把这些值绘制一个频数分布图(除以10000则为实际概率), 见图6.2:

```
> plot(as.numeric(names(fr)), fr, type="h",
+ xlab="Determinant", ylab="Frequency")
```

其中`as.numeric()`把向量`fr`中的元素名又转换成了数值型, 用来作为作图的横轴坐标, `fr`中的元素值即频数作为纵轴, `type="h"`表示是画垂线型图。

6.3.9 数组的广义转置

可以用`aperm(a, perm)`函数把数组a的各维按perm中指定的新次序重新排列。例如：

```
> a <- array(1:24, dim=c(2,3,4))
> b <- aperm(a, c(2, 3, 1))
```

结果a的第2维变成了b的第1维，a的第3维变成了b的第2维，a的第1维变成了b的第3维。这时有 $a[i_1, i_2, i_3] \equiv b[i_2, i_3, i_1]$ 。注意 $c(i_1, i_2, i_3)[2,3,1] = c(i_2, i_3, i_1)$ 。一般地，若 $b \leftarrow \text{aperm}(a, p)$ ，i是数组a的一个下标向量，则 $a[\text{rbind}(i)] \equiv b[\text{rbind}(i[p])]$ ，即a的一个元素下标经过p变换后成为b的对应元素的下标。

对于矩阵a，`aperm(a, c(2,1))`恰好是矩阵转置。对于矩阵转置可以简单地用`t(a)`表示。

6.3.10 apply函数

对于向量，我们有`sum`、`mean`等函数对其进行计算。对于数组，如果我们想对其中一维(或若干维)进行某种计算，可以用`apply`函数。其一般形式为：

```
apply(X, MARGIN, FUN, ...)
```

其中X为一个数组，MARGIN是固定哪些维不变，FUN是用来计算的函数。例如，设a是 3×4 矩阵，则`apply(a, 1, sum)`的意义是对a的各行求和(保留第一维即第一个下标不变)，结果是一个长度为3的向量(与第一维长度相同)，而`apply(a, 2, sum)`意义是对a的各列求和，结果是一个长度为4的向量(与第二维长度相同)。

如果函数FUN的结果是一个标量，MARGIN只有一个元素，则apply的结果是一个向量，其长度等于MARGIN指定维的长度，相当于固定MARGIN指定的那一维的每一个值而把其它维取出作为子数组或向量送入FUN中进行运算。如果MARGIN指定了多个维，则结果是一个维数向量等于 $\text{dim}(X)[\text{MARGIN}]$ 的数组。如果函数FUN的结果是一个长度为N的向量，则结果是一个维数向量等于 $c(N, \text{dim}(X)[\text{MARGIN}])$ 的数组，注意这时不论是对哪一维计算，结果都放在了第一维。所以，比如我们要把 3×2 矩阵a的两列分别排序，只要用`apply(a, 2, sort)`，这样对每一列排序得到一个长度为3的向量，用第一维来引用，结果的维向量为 $c(N, \text{dim}(a)[2]) = c(3, 2)$ ，保留了列维，恰好得到所需结果，运行如下例：

```
> a <- cbind(c(4,9,1), c(3,7,2))
> a
      [,1] [,2]
[1,]    4    3
[2,]    9    7
[3,]    1    2
> apply(a, 2, sort)
      [,1] [,2]
[1,]    1    2
[2,]    4    3
[3,]    9    7
>
```

但是, 如果要对行排序, 则`apply(a, 1, sort)`把`a`的每一行2个元素排序后的结果用第一维来引用, 结果的维向量为`c(N, dim(a)[1])=c(2, 3)`, 把原来的列变成了行, 所以`t(apply(a,1,sort))`才是对`a`的每一行排序的结果。如:

```
> apply(a, 1, sort)
      [,1] [,2] [,3]
[1,]     3     7     1
[2,]     4     9     2
> t(apply(a,1,sort))
      [,1] [,2]
[1,]     3     4
[2,]     7     9
[3,]     1     2
```

上面我们只用了矩阵(二维数组)作为例子讲解`apply`的用法。实际上, `apply`可以用于任意维数的数组, 函数`FUN`也可以是任意可以接收一个向量或数组作为第一自变量的函数。比如, 设`x`是一个维数向量为`c(2,3,4,5)`的数组, 则`apply(x, c(1,3), sum)`可以产生一个2行4列的矩阵, 其每一元素是`x`中固定第1维和第3维下标取出子数组求和的结果。

6.4 因子

概率论中随机变量可以分为连续型和离散型, 相应地, 统计中的变量可以分为区间变量和名义变量、有序变量。区间变量取连续的数值, 可以进行求和、平均等运算。名义变量和有序变量取离散值, 可以用数值代表也可以是字符型值, 其具体数值没有加减乘除的意义, 不能用来计算而只能用来分类或者计数。名义变量例子有性别、省份、职业, 有序变量例子有班级名次、年龄组别。

因为离散变量有各种不同表示方法, 在S中为统一起见使用因子(`factor`)来表示这种分类变量。还提供了有序因子(`ordered factor`)来表示有序变量。因子是一种特殊的向量, 其中每一个元素取一组离散值中的一个, 因子对象有一个特殊属性`levels`表示这组离散值。例如:

```
> x <- c("男", "女", "男", "男", "女")
> y <- factor(x)
> y
[1] 男 女 男 男 女
Levels: 男 女
```

函数`factor()`用来把一个向量编码成为一个因子。其一般形式为:

```
factor(x, levels = sort(unique(x), na.last = TRUE),
      labels, exclude = NA, ordered = FALSE)
```

可以自行指定各离散取值(水平`levels`), 不指定时由`x`的不同值来求得。`labels`可以用来指定各水平的标签, 不指定时用各离散取值的对应字符串。`exclude`参数用来指定要转换为缺失值(`NA`)的元素值集合。如果指定

了levels, 则因子的第i个元素当它等于水平中第j个时元素值取“j”, 如果它的值没有出现在levels中则对应因子元素值取NA。ordered取真值时表示因子水平是有次序的(按编码次序)。例如:

```
> f <- factor(c(1,0,1,1,0),
              levels=c(1,0), labels=c("男","女"))
```

结果和刚才一样。

可以用is.factor()检验对象是否因子, 用as.factor()把一个向量转换成一个因子。

因子的基本统计是频数统计, 用函数table()来计数。例如,

```
> sex <- factor(c("男", "女", "男", "男", "女"))
> res.tab <- table(sex)
> res.tab
sex
男 女
3  2
```

表示男性3人, 女性2人。table()的结果是一个带元素名的向量, 元素名为因子水平, 元素值为该水平的出现频数。

从这个简单例子我们也可以看出S与SAS的一个区别: SAS的结果一般只显示出来, 保存只能存为数据集一种类型, 不便于对中间结果进一步编程处理; 而S的结果除了可以显示外本身都是S对象(如这里的向量结果), 可以很方便地进一步处理。

可以用两个或多个因子进行交叉分类。比如, 性别(sex)和职业(job)交叉分组可以用table(sex, job)来统计每一交叉类的频数, 结果为一个矩阵, 矩阵带有行名和列名, 分别为两个因子的各水平名。

因子可以用来作为另外的同长度变量的分类变量。比如, 假设上面的sex是5个学生的性别, 而

```
> h <- c(165, 170, 168, 172, 159)
```

是这5个学生的身高, 则

```
> tapply(h, sex, mean)
```

可以求按性别分类的身高平均值。这样用一个等长的因子向量对一个数值向量分组的办法叫做不规则数组(ragged array)。后面我们还可以看到更多的因子的应用。

6.5 列表(list)

6.5.1 列表定义

列表是一种特别的对象集合, 它的元素也由序号(下标)区分, 但是各元素的类型可以是任意对象, 不同元素不必是同一类型。元素本身允许是其它复杂数据类型, 比如, 列表的一个元素也允许是列表。这样的数据类型称为是递归(recursive)数据类型。例如:

```
> rec <- list(name="李明", age=30, scores=c(85, 76, 90))
> rec
$name
[1] "李明"

$age
[1] 30

$scores
[1] 85 76 90
```

列表元素总可以用“列表名[[下标]]”的格式引用。例如：

```
> rec[[2]]
[1] 30
> rec[[3]][2]
[1] 76
```

但是，列表不同于向量，我们每次只能引用一个元素，如rec[[1:2]]的用法是不允许的。

注意：“列表名[下标]”或“列表名[下标范围]”的用法也是合法的，但其意义与用两重括号的记法完全不同，两重记号取出列表的一个元素，结果与该元素类型相同，如果使用一重括号，则结果是列表的一个子列表(结果类型仍为列表)。

在定义列表时如果指定了元素的名字(如rec中的name, age, scores)，则引用列表元素还可以用它的名字作为下标，格式为“列表名[“元素名”]”，如：

```
> rec[["age"]]
[1] 30
```

另一种格式是“列表名\$元素名”，如：

```
> rec$age
[1] 30
```

其中“元素名”可以简写到与其它元素名能够区分的最短程度，比如“rec\$s”可以代表“rec\$score”。这种写法方便了交互运行，编写程序时一般不用简写以免降低程序的可读性。

使用元素名的引用方法可以让我们不必记住某一个下标代表那一个元素，而直接用易记的元素名来引用元素。事实上，向量和矩阵也可以指定元素名、行名、列名。

定义列表使用list()函数，每一个自变量变成列表的一个元素，自变量可以用“名字=值”的方式给出，即给出列表元素名。自变量的值被复制到列表元素中，自变量如果是变量并不会与该列表元素建立关系(改变该列表元素不会改变自变量的值)。

6.5.2 修改列表

列表的元素可以修改, 只要把元素引用赋值即可。如:

```
> rec$age <- 45
```

甚至

```
> rec$age <- list(19, 29, 31)
```

(可以任意修改一个列表元素)。如果被赋值的元素原来不存在, 则列表延伸以包含该新元素。例如, rec现在共有三个元素, 我们定义一个新的命名元素, 则列表长度变为4, 再定义第六号元素则列表长度变为6:

```
> rec$sex <- "男"
> rec[[6]] <- 161
> rec
$name
[1] "李明"

$age
[1] 30

$scores
[1] 85 76 90

$sex
[1] "男"

[[5]]
NULL

[[6]]
[1] 161
```

第五号元素因为没有定义所以其值是“NULL”, 这是空对象的记号。如果rec是一个向量, 则其空元素为“NA”, 这是缺失值的记号。从这里我们也可以体会“NULL”与“NA”的区别。

几个列表可以用连接函数c()连接起来, 结果仍为一个列表, 其元素为各自变量的列表元素。如:

```
> list.ABC <- c(list.A, list.B, list.C)
```

注意在S中句点是名字的合法部分, 一般没有特殊意义。

6.5.3 几个返回列表的函数

列表的重要作用是把相关的若干数据保存在一个数据对象中, 这样在编写函数时我们就可以返回这样一个包含多项输出的列表。因为函数的返回结果可以完整地存放在一个列表中, 我们可以继续对得到的结果进行分析, 这是S比SAS灵活的一个地方。下面给出几个返回列表的例子。

一、特征值和特征向量

函数eigen(x)对对称矩阵x计算其特征值和特征向量, 返回结果为一个列表, 列表的两个成员(元素)为values和vectors。例如:

```
> ev <- eigen((1:3) %o% (1:3))
> ev
$values
[1] 1.400000e+01 0.000000e+00 -8.881784e-16

$vectors
      [,1]      [,2]      [,3]
[1,] 0.2672612 0.8944272 -0.3585686
[2,] 0.5345225 -0.4472136 -0.7171372
[3,] 0.8017837 0.0000000 0.5976143
```

可见三个特征值只有第一个不为零(由于数值计算精度所限, 第三个特征值应为零但结果只是近似为零)。特征向量按矩阵存放, 每一列为一个特征向量。

二、奇异值分解及行列式

函数svd()进行奇异值分解 $X = UDV'$, 其中 X 是任意 $n \times m$ 阵, U 为 $n \times n$ 正交阵, V 为 $m \times m$ 正交阵, D 为 $n \times m$ 对角阵(只有主对角线元素不为零且非负)。svd(x)返回有三个成员d, u, v的列表, d为包含奇异值的向量(即 D 的主对角线元素), u, v分别为上面的两个正交阵。易见如果矩阵x是方阵则x的行列式的绝对值等于奇异值的乘积, 所以:

```
> absdetx <- prod(svd(x)$d)
```

或者我们可以为此定义一个函数:

```
> absdet <- function(x) prod(svd(x)$d)
```

三、最小二乘拟合与QR分解

函数lsfit(x,y)返回最小二乘拟合的结果。最小二乘的模型为线性模型

$$y = X\beta + \varepsilon$$

lsfit(x,y)的第一个参数x为模型中的设计阵 X , 第二个参数y为模型中的因变量y(可以是一个向量也可以是一个矩阵), 返回一个列表, 成员coefficients为上面模型的 β (最小二乘系数), 成员residuals为拟合残差, 成员intercept用来指示是否有截距项, 成员qr为设计阵 X 的QR分解, 它本身也是一个列表。模型拟合缺省情况有截距项, 可以用intercept=FALSE选项指定无截距项。关于最小二乘拟合还可参见ls.diag()函数(查看帮助)。

函数qr(x)返回x的QR分解结果。矩阵 X 的QR分解为 $X = QR$, Q 为正交阵, R 为上三角阵。函数结果为一个列表, 成员qr是压缩了 Q 和 R 在内的一个矩阵, 包括对角线在内的上三角部分为 R , Q 的信息压缩存放在下三角部分。结果的其它成员为一些辅助信息。用qr.Q()和qr.R()可以从qr()的结果中提取 Q 和 R 。

6.6 数据框(data.frame)

数据框是S中类似SAS数据集的一种数据结构。它通常是矩阵形式的数据,但矩阵各列可以是不同类型的。数据框每列是一个变量,每行是一个观测。

但是,数据框有更一般的定义。它是一种特殊的列表对象,有一个值为“data.frame”的class属性,各列表成员必须是向量(数值型、字符型、逻辑型)、因子、数值型矩阵、列表,或其它数据框。向量、因子成员为数据框提供一个变量,如果向量非数值型则会被强制转换为因子,而矩阵、列表、数据框这样的成员为新数据框提供了和其列数、成员数、变量数相同个数的变量。作为数据框变量的向量、因子或矩阵必须具有相同的长度(行数)。

尽管如此,我们一般还是可以把数据框看作是一种推广了的矩阵,它可以用矩阵形式显示,可以用对矩阵的下标引用方法来引用其元素或子集。

6.6.1 数据框生成

数据框可以用data.frame()函数生成,其用法与list()函数相同,各自变量变成数据框的成分,自变量可以命名,成为变量名。例如:

```
> d <- data.frame(name=c("李明", "张聪", "王建"),  
+ age=c(30, 35, 28), height=c(180, 162, 175))  
> d  
  name age height  
1 李明  30   180  
2 张聪  35   162  
3 王建  28   175
```

如果一个列表的各个成分满足数据框成分的要求,它可以用as.data.frame()函数强制转换为数据框。比如,上面的d如果先用list()函数定义成了一个列表,就可以强制为一个数据框。

一个矩阵可以用data.frame()转换为一个数据框,如果它原来有列名则其列名被作为数据框的变量名,否则系统自动为矩阵的各列起一个变量名(如X1, X2)。

6.6.2 数据框引用

引用数据框元素的方法与引用矩阵元素的方法相同,可以使用下标或下标向量,也可以使用名字或名字向量。如d[1:2, 2:3]。数据框的各变量也可以用按列表引用(即用双括号[]或\$符号引用)。

数据框的变量名由属性names定义,此属性一定是非空的。数据框的各行也可以定义名字,可以用rownames属性定义。如:

```
> names(d)  
[1] "name" "age" "height"  
> rownames(d)  
[1] "1" "2" "3"
```

6.6.3 attach()函数

数据框的主要用途是保存统计建模需要的数据。S的统计建模功能都需要以数据框为输入数据。我们也可以把数据框当成一种矩阵来处理。

在使用数据框的变量时可以用“数据框名\$变量名”的记法。但是, 这样使用较麻烦, S提供了attach()函数可以把数据框“连接”入当前的名字空间。例如,

```
> attach(d)
> r <- height / age
```

后一语句将在当前工作空间建立一个新变量r, 它不会自动进入数据框d, 要把新变量赋值到数据框中, 可以用

```
> d$r <- height / age
```

这样的格式。

为了取消连接, 只要调用detach()(无参数即可)。

注意: S-PLUS和R中名字空间的管理是比较独特的。它在运行时保持一个变量搜索路径表, 在读取某个变量时到这个变量搜索路径表中由前向后查找, 找到最前的一个; 在赋值时总是在位置1赋值(除非特别指定在其它位置赋值)。attach()的缺省位置是在变量搜索路径表的位置2, detach()缺省也是去掉位置2。所以, S编程的一个常见问题是当你误用了一个自己并没有赋值的变量时有可能不出错, 因为这个变量已在搜索路径中某个位置有定义, 这样不利于程序的调试, 需要留心这样的问题。

除了可以连接数据框, 也可以连接列表。

6.7 输入输出

6.7.1 输出

在S交互运行时要显示某一个对象的值只要键入其名字即可, 如:

```
> x <- 1:10
> x
[1] 1 2 3 4 5 6 7 8 9 10
```

这实际上是调用了print()函数, 即print(x)。在非交互运行(程序)中应用print()来输出。print()函数可以带一个digits=参数指定每个数输出的有效数字位数, 可以带一个quote=参数指定字符串输出时是否带两边的撇号, 可以带一个print.gap=参数指定矩阵或数组输出时列之间的间距。

print()函数是一个通用函数, 即它对不同的自变量有不同的反应。对各种特殊对象如数组、模型结果等都可以规定print的输出格式。

cat()函数也用来输出, 但它可以把多个参数连接起来再输出(具有paste()的功能)。例如:

```
> cat("i = ", i, "\n")
```

注意使用cat()时要自己加上换行符“\n”。它把各项转换成字符串, 中间隔以空格连接起来, 然后显示。如果要使用自定义的分隔符, 可以用sep=参数, 例如:

```
> cat(c("AB", "C"), c("E", "F"), "\n", sep="")
ABCDEF
```

cat()还可以指定一个参数“file=文件名”，可以把结果写到指定的文件中，如：

```
> cat("i = ", 1, "\n", file="c:/work/result.txt")
```

如果指定的文件已经存在则原来内容被覆盖。加上一个append=TRUE参数可以不覆盖原文件而是在文件末尾附加，这很适用于运行中的结果记录。

cat()函数和print()都不具有很强的自定义格式功能，为此可以使用cat()与format()函数配合实现。format()函数为一个数值向量找到一种共同的显示格式然后把向量转换为字符型。例如：

```
> format(c(1, 100, 10000))
[1] "      1" "    100" " 10000"
```

S-PLUS中的format()函数功能较强，具有较多的控制参数，请参见帮助。R中目前format()函数功能仍较弱，但R有一个formatC函数可以提供类似C语言的printf格式功能。formatC对输入向量的每一个元素单独进行格式转换而不生成统一格式，例如：

```
> formatC(c(1, 10000))
[1] "1" "1e+004"
```

在formatC()函数中可以用format=参数指定C格式类型，如“d”(整数)，“f”(定点实数)，“e”，“E”(科学记数法)，“g”，“G”(选择位数较少的输出格式)，“fg”(定点实数但用参数digits指定有效位数而不是总宽度)，“s”(字符串)。可以用width指定输出宽度，用digits指定有效位数(格式为e,E,g,G,fg时)或小数点后位数(格式为f)时。可以用flag参数指定一个输出选项字符串，字符串中有“-”表示输出左对齐，有“0”表示左空白用0填充，有“+”表示要输出正负号，等等。例如，我们有一个矩阵da中保存了三个日期的年、月、日：

```
> da
      [,1] [,2] [,3]
[1,]   99    1    3
[2,]   96   11    9
[3,]   65    5   18
```

为了输出这三个日期，可以用apply函数指定对每一行作用一个输出函数，此输出函数利用cat()和formatC来控制：

```
> apply(da, 1, function(r)
+ cat(formatC(r[1], format='d', width=2, flag='0'), '- ',
+ formatC(r[2], format='d', width=2, flag='0'), '- ',
+ formatC(r[3], format='d', width=2, flag='0'), '\n', sep=''))
```

```
99-01-03
96-11-09
65-05-18
NULL
```

这里我们知道`apply`函数第一个参数指定了一个矩阵, 第二个参数说明对行操作还是对列操作, 第三个参数是一个函数, 这里我们使用了直接定义一个函数作为参数的办法。输出结果中多了一个`NULL`函数, 这是因为我们在交互运行, `apply`的结果作为一个表达式的值(`NULL`)会被显示出来。为避免显示, 可以把结果赋给一个临时变量名, 或者把整个表达式作为`invisible()`函数的参数, 这时不显示表达式值。

`S`的输出缺省显示在交互窗口。可以用`sink()`函数指定一个文件以把后续的输出转向到这个文件, 并可用`append`参数指定是否要在文件末尾附加:

```
> sink("c:/work/result.txt", append=TRUE)
> ls()
> d
> sink()
```

调用无参数的`sink()`把输出恢复到交互窗口。

如果要把一个矩阵`x`输出到文件中, 可以用`write(t(x), file="文件名", ncol=ncol(x))`的形式。这里之所以要把`x`转置后再输出是因为`R`中矩阵是列优先的, 如果不转置则输出是按列输出的。如果不指定列数则缺省使用5列, 即使`x`的列数多于或少于5。文件名缺省使用"data"。

要把一个数据框`d`输出到文件中, 可以用`write.table(x, file="文件名")`, 输出文件中包含变量名表头和行名。

6.7.2 输入

为了从外部文件读入一个数值型向量, `S`提供了`scan()`函数。如果指定了`file`参数(也是第一参数), 则从指定文件读入, 缺省情况下读入一个数值向量, 文件中各数据以空白分隔, 读到文件尾为止。例如:

```
> cat(1:12, '\n', file='c:/work/result.txt')
> x <- scan('c:/work/result.txt')
```

如果文件中是一个用空白分隔的矩阵(或数组), 我们可以先用`scan()`把它读入到一个向量然后用`matrix()`函数(或`array()`函数)转换。如:

```
> y <- matrix(scan('c:/work/result.txt'),
              ncol=3, byrow=T)
```

实际上, `scan()`也能够读入一个多列的表格, 只要用`what`参数指定一个列表, 则列表每项的类型为需要读取的类型。用`skip`参数可以跳过文件的开始若干行不读。用`sep`参数可以指定数据间的分隔符。详见帮助。

`scan()`不指定读取文件名时是交互读入, 读入时用一个空行结束。

如果要读取一个数据框, `S`提供了一个`read.table()`函数。它只要给出一个文件名, 就可以把文件中用空白分隔的表格数据每行读入为数据框的一行。比如, 文件`c:\work\d.txt`中内容如下:

```
Zhou      15 3
"Li Ming" 9 李明
Zhang 10.2 Wang
```

用read.table读入:

```
> x <- read.table('c:/work/d.txt',
+ colClasses=c('character', 'numeric', 'character'))
> x
      V1    V2    V3
1   Zhou 15.0     3
2 Li Ming  9.0 李明
3   Zhang 10.2 Wang
```

读入结果为数据框。函数自动为数据框变量指定“V1”、“V2”这样的变量名, 指定“1”、“2”这样的行名。可以用col.names参数指定一个字符型向量作为数据框的变量名, 用row.names参数指定一个字符型向量作为数据框的行名。函数可以自动识别表列是数值型还是字符型, 并在缺省情况下把字符型数据转换为因子, 这里的第一列和第三列我们并不想转换, 所以人为地加上了一个colClasses参数说明每一列的数据类型。

read.table()可以读入带有表头的文件, 只要加上header=TRUE参数即可。可以用sep参数指定表行各项的分隔符。例如, 为了读入如下带有表头的逗号分隔文件c:\work\d.csv:

```
Name,score, cn
Zhou,15,3
Li Ming, 9, 李明
Zhang, 10.2, Wang
```

使用如下语句:

```
> x <- read.table('c:/work/d.csv', header=T, sep=',',
+ colClasses=c('character', 'numeric', 'character'))
> x
      Name score    cn
1     Zhou  15.0     3
2  Li Ming   9.0 李明
3    Zhang 10.2  Wang
```

其它一些用法见帮助。

6.8 程序控制结构

S是一个表达式语言, 其任何一个语句都可以看成是一个表达式。表达式之间以分号分隔或用换行分隔。表达式可以续行, 只要前一行不是完整表达式(比如末尾是加减乘除等运算符, 或有未配对的括号)则下一行为上一行的继续。

若干个表达式可以放在一起组成一个复合表达式, 作为一个表达式使用。组合用大括号表示, 如:

```
> {  
> x <- 15  
> x  
> }
```

S语言也提供了其它高级程序语言共有的分支、循环等程序控制结构。

6.8.1 分支结构

分支结构包括if结构：

if (条件) 表达式1

或

if (条件) 表达式1 else 表达式2

其中的“条件”为一个标量的真或假值，表达式可以用大括号包围的复合表达式。有else子句时一般写成：

```
if(条件) {  
  表达式组.....  
} else {  
  表达式组.....  
}
```

这样的写法可以使else不至于脱离前面的if。

例如，如果变量lambda为缺失值就给它赋一个缺省值，可用：

```
if(is.na(lambda)) lambda <- 0.5;
```

又比如要计算向量x的重对数，这只有在元素都为正且对数都为正时才能做到，因此需要先检查：

```
if(all(x>0) && all(log(x))>0) {  
  y <- log(log(x));  
  print(cbind(x,y));  
} else {  
  cat('Unable to comply\n');  
}
```

注意“&&”表示“与”，它是一个短路运算符，即第一个条件为假时就不计算第二个条件，如果不这样此例中计算对数就可以有无效值。在条件中也可以用“||”（两个连续的竖线符号）表示“或”，它也是短路运算符，当第一个条件为真时就不再计算第二个条件。

注：对于计算重对数的问题只要判断all(x>1)也可以，例子是为了说明短路与。

在用S编程时一定要时刻牢记S是一个向量语言，几乎所有操作都是对向量进行的。而S中的if语句却是一个少见的例外，它的判断条件是标量的真值或假值。比如，我们要定义一个分段函数f(x)，当x为正时返回1，否则返回0，马上可以想到用if语句实现如下：

```
if(x>0) 1 else 0
```

当 x 是标量时这个定义是有效的,但是当自变量 x 是一个向量时,比较的结果也是一个向量,这时条件无法使用。所以,这个分段函数应该这样编程:

```
y <- numeric(length(x))
y[x>0] <- 1
y[x<=0] <- 0
y
```

有多个if语句时else与最近的一个配对。可以使用“if ... else if ... else if ... else ...”的多重判断结构表示多分支。多分支也可以使用switch()函数。

6.8.2 循环结构

循环结构中常用的是for循环,是对一个向量或列表的逐次处理,格式为“for(name in values) 表达式”,如:

```
s <- 0
for(i in seq(along=x)){
  cat('x(', i, ') = ', x[i], '\n', sep='');
  s <- s+x[i];
}
```

这个例子我们需要使用下标的值,所以用seq(along=x)生成了x的下标向量。如果不需要下标的值,可以直接如此使用:

```
s <- 0
for(xi in x){
  cat(xi, '\n')
  s <- s + xi
}
```

当然,如果只是要求各元素的和,只要调用sum(x)即可。从这里我们也可以看出,显式的循环经常是可以避免的,利用函数对每个元素计算值、使用sum等统计函数及apply、lapply、sapply、tapply等函数往往可以代替循环。因为循环在S中是很慢的(S-PLUS和R都是解释语言),所以应尽可能避免使用显式循环。

我们再举一个例子。比如,我们要计算同生日的概率。假设一共有365个生日(只考虑月、日),而且各生日的概率是相等的(这里忽略了闰年的情况以及可能存在的出生日期分布的不均匀)。设一个班有 n 个人,当 n 大于等于365时{至少两个人有生日相同}是必然事件。当 n 小于365时,我们可以计算 $\Pr\{\text{至少有两人生日}\} = 1 - \Pr\{n\text{个人生日彼此不同}\}$,这时, n 个人的生日可取值数为 365^n ,而 n 个人彼此不同的可能数为365中取 n 个的排列数,彼此不同的概率为 $365 \times 364 \times \cdots \times (365 - (n - 1)) = 365! / (365 - n)!$ 。因此,为了计算 $n=1, 2, \dots, 364$ 的情况下的同生日概率,可以用如下循环实现:


```
x <- numeric(364)
for(n in 1:364){
  x[n] <- 1
  for(j in 0:(n-1)){
    x[n] <- x[n] * (365-j)/365
  }
  x[n] <- 1 - x[n]
}
```

这段程序运行了0.980秒。我们可以尽量用向量运算来实现, 速度要快得多:

```
x <- numeric(364)
for(n in 1:364){
  x[n] <- 1 - prod((365:(365-n+1))/365)
}
```

这段程序只用了0.030秒, 比第一个程序快32倍。注意不能直接去计算365!, 这会超出数值表示范围。

这个问题还可以用S的cumprod函数进一步简化:

```
x <- 1 - cumprod((365:2)/365)
```

这段程序没有任何循环, 运行时间已经无法测量。多次重复运行可以发现它比第二段程序快165倍。由此可见使用向量运算以及利用系统提供的函数对提高S程序效率的意义。

另外要注意使用for(i in 1:n)格式的计数循环时要避免一个常见错误, 即当n为零或负数时1:n是一个从大到小的循环, 而我们经常需要的是当n为零或负数时就不进入循环。为达到这一点, 可以在循环外层判断循环结束值是否小于开始值。

while循环是在开始处判断循环条件的当型循环, 如:

```
while(b-a>eps){
  cc <- (a+b)/2;
  if(f(cc)>0) b <- cc else a <- cc
}
```

是一段二分法解方程的程序, 要求函数f(x)是单调上升的。可以很容易地把这个程序修改成允许f(x)单调下降。

还可以使用

repeat 表达式

循环, 在循环体内用break跳出。

在一个循环体内用next表达式可以进入下一轮循环。

分支和循环结构主要用于定义函数。

6.9 S程序设计

对于复杂一些的计算问题我们应该编写成函数。这样做的好处是编写一次可以重复使用,并且可以很容易地修改,另外的好处是函数内的变量名是局部的,运行函数不会使函数内的局部变量被保存到当前的工作空间,可以避免在交互状态下直接赋值定义很多变量使得工作空间杂乱无章。

6.9.1 工作空间管理

前面我们已经提到, S在运行时保持一个变量搜索路径表, 要读取某变量时依次在此路径表中查找, 返回找到的第一个; 给变量赋值时在搜索路径的第一个位置赋值。在函数内部, 搜索路径表第一个位置是局部变量名空间, 所以变量赋值是局部赋值, 被赋值的变量只在函数运行期间有效。在函数外部, 搜索路径表的第一个位置是当前工作空间, 这个工作空间将保存所有在函数外部定义的变量以及函数定义, 而且可以在退出系统时保存工作空间, 下一次还可以看到以前定义的变量和函数。

这样, 因为工作空间里对象越来越多, 出错的机会增大了。尽量把工作都用函数实现可以避免这个问题, 函数内定义的变量是局部的, 不会进入当前工作空间。

可以直接管理工作空间中的对象。用ls()函数可以查看当前工作空间保存的变量和函数, 用rm()函数可以剔除不想要的对象。如:

```
> ls()
[1] "A"      "Ai"      "b"       "cl"      "cl.f"    "fit1"    "g1"      "marks"   "ns"
[10] "p1"     "rec"     "tmp.x"   "x"       "x1"      "x2"      "x3"      "y"
> rm(x, x1, x2, x3)
> ls()
[1] "A"      "Ai"      "b"       "cl"      "cl.f"    "fit1"    "g1"      "marks"   "ns"
[10] "p1"     "rec"     "tmp.x"   "y"
```

ls()可以指定一个pattern参数, 此参数定义一个匹配模式, 只返回符合模式的对象名。模式格式是UNIX中grep的格式。比如, ls(pattern="tmp[.]")可以返回所有以“tmp.”开头的对象名。

rm()可以指定一个名为list的参数给出要删除的对象名, 所以rm(list=ls(pattern="tmp[.]"))可以删除所有以“tmp.”开头的对象名。

6.9.2 函数定义

S中函数定义的一般格式为“函数名 <- function(参数表) 表达式”。定义函数可以在命令行进行, 例如:

```
> hello <- function(){
+   cat("Hello, world!\n")
+   cat("\n")
+ }
> hello
function(){
cat("Hello, world!\n")
cat("\n")
}
```

函数体为一个复合表达式, 各表达式的之间用换行或分号分开。不带括号调用函数时显示函数定义, 而不是调用函数。

在命令行输入函数程序很不方便修改, 所以我们一般是打开一个其他的编辑程序(如Windows的记事本), 输入以上函数定义, 保存文件, 比如保存到了C:\work\hello.r, 我们就可以用

```
> source("c:\\work\\hello.r")
```

运行文件中的程序。实际上, 用source()运行的程序不限于函数定义, 任何S程序都可以用这种方式编好再运行, 效果与在命令行直接输入基本相同, 只不过不自动显示表达式值。

对于一个已有定义的函数, 可以用fix()函数来修改, 如:

```
> fix(hello)
```

将打开一个编辑窗口显示函数的定义, 修改后关闭窗口函数就被修改了。fix()调用的编辑程序缺省为记事本, 可以用“options(editor=”编辑程序名”)”来指定自己喜欢的编辑程序。

函数可以带参数, 可以返回值, 例如:

```
larger <- function(x, y){  
  y.is.bigger <- (y>x);  
  x[y.is.bigger] <- y[y.is.bigger]  
  x  
}
```

这个函数输入两个向量(相同长度)x和y, 然后把x中比y对应元素小的元素替换为y中对应元素, 返回y的值。S返回值为函数体的最后一个表达式的值, 不需要使用return()函数。不过, 也可以使用“return(对象)”函数从函数体返回调用者。S只能返回一个对象, 所以要返回多个结果时需要用一个列表把这些结果包装起来返回。

6.9.3 参数(自变量)

函数可以带虚参数(形式自变量)。S函数调用方式很灵活, 例如, 如下函数:

```
fsub <- function(x, y) x-y
```

有两个虚参数x和y, 我们用它计算 $100 - 45$, 可以调用fsub(100,45), 或fsub(x=100,y=45), 或fsub(y=45, x=100), 或fsub(y=45, 100)。即调用时实参与虚参可以按次序结合, 也可以直接指定虚参名结合。实参先与指定了名字的虚参结合, 没有指定名字的按次序与剩下的虚参结合。

函数在调用时可以不给出所有的实参, 这需要在定义时为虚参指定缺省值。例如上面的函数改为:

```
fsub <- function(x, y=0) x-y
```

则调用时除了可以用以上的方式调用外还可以用`fsub(100)`, `fsub(x=100)`等方式调用, 只给出没有缺省值的实参。

即使没有给虚参指定缺省值也可以在调用时省略某个虚参, 然后函数体内可以用`missing()`函数判断此虚参是否有对应实参, 如:

```
trans <- function(x, scale){
  if(!missing(scale)) x <- scale*x
  . . . . .
}
```

此函数当给了`scale`的值时对自变量`x`乘以此值, 否则保持原值。这种用法在其它语言中是极其少见的, S可以实现这一点是因为S的函数调用在用到参数的值时才去计算这个参数的值(称为“懒惰求值”), 所以可以在调用时缺少某些参数而不会被拒绝。

S函数还可以有一个特殊的“...”虚参, 表示所有不能匹配的实参, 调用时如果有需要与其它虚参结合的实参必须用“虚参名=”的格式引入。例如:

```
> f <- function(...){
+   for(x in list(...)){
+     cat(min(x), '\n')
+   }
+ }
> f(c(5,1,2), c(9, 4, 7))
1
4
```

6.9.4 作用域

函数的虚参完全是按值传递的, 改变虚参的值不能改变对应实参的值。例如:

```
> x <- list(1, "abc")
> x
[[1]]
[1] 1

[[2]]
[1] "abc"

> f <- function(x) x[[2]] <- "!!"
> f(x)
> x
[[1]]
[1] 1

[[2]]
[1] "abc"
```

函数体内的变量也是局部的, 对函数体内的变量赋值当函数结束运行后变量值就删除了, 不影响原来同名变量的值。例如:

```
> x <- 2
> f <- function(){
+ print(x)
+ x <- 20
+ }
> f()
[1] 2
> x
[1] 2
```

这个例子中原来有一个变量x值为2, 函数中为变量x赋值20, 但函数运行完后原来的x值并未变化。但是也要注意, 函数中的显示函数调用时局部变量x还没有赋值, 显示的是全局变量x的值。这是S编程比较容易出问题的地方: 你用到了一个局部变量的值, 你没有意识到这个局部变量还没有赋值, 而程序却没有出错, 因为这个变量已有全局定义。

6.9.5 程序调试

S-PLUS和R目前还不象其它主流程序设计语言那样具有图形界面单步跟踪、设置断点、观察表达式等强劲的调试功能。调试复杂的S程序, 可以用一些通用的程序调试方法, 另外S也提供了一些调试用函数。

对任何程序语言, 最基本的调试手段当然是在需要的地方显示变量的值。可以用print()或cat()显示。例如, 我们为了调试前面定义的larger()函数, 可以显示两个自变量的值及中间变量的值:

```
larger <- function(x, y){
  cat('x = ', x, '\n')
  cat('y = ', y, '\n')
  y.is.bigger <- (y>x);
  cat('y.is.bigger = ', y.is.bigger, '\n')
  x[y.is.bigger] <- y[y.is.bigger]
  x
}
```

S提供了一个browser()函数, 当调用时程序暂停, 可以用ls()列出局部变量, 用户可以查看变量或表达式的值, 还可以修改变量。例如:

```
larger <- function(x, y){
  y.is.bigger <- (y>x);
  browser()
  x[y.is.bigger] <- y[y.is.bigger]
  x
}
```

我们运行此程序:

```

> larger(c(1,5), c(2, 4, 9))
Called from: larger(c(1, 5), c(2, 4, 9))
Browse[1]> x
[1] 1 5
Browse[1]> y
[1] 2 4 9
Browse[1]> y>x
[1] TRUE FALSE TRUE
Browse[1]> x[y.is.bigger]
[1] 1 NA
Browse[1]> c
[1] 2 5 9

```

退出R的browser()菜单可用c(在S中用return())。在R的browser()状态下用n命令可以进入单步执行状态, 用n或者回车可以继续, 用c可以退出browser, 用Q可以返回到命令行状态。

R提供了一个debug()函数, debug(f)可以打开对函数f()的调试, 执行到函数f时自动进入单步执行的browser()菜单, 回车就可以单步执行。用undebug(f)关闭调试。

6.9.6 程序设计举例

设计S程序是很容易的, 在初学时我们只要使用我们从一般程序设计中学来的知识并充分利用S中现成的各种算法及绘图函数就可以了。但是, 如果要用S编制计算量较大的程序, 或者程序需要发表, 就需要注意一些S程序设计的技巧。

用S语言开发算法, 最重要的一点是要记住S是一个向量语言, 计算应该尽量通过向量、矩阵运算来进行, 并充分利用S提供的现成的函数, 避免使用显式循环。显式循环会大大降低S的运算速度, 因为S是解释执行的。

例1. 比如, 考虑核回归问题。核回归是非参数回归的一种, 假设变量Y与变量X之间的关系为:

$$Y = f(X) + \varepsilon$$

其中函数f未知。观测到X和Y的一组样本 $X_i, Y_i, i=1, \dots, n$ 后, 对f的一种估计为:

$$f(x) = \frac{\sum_{i=1}^n K\left(\frac{x-X_i}{h}\right) Y_i}{\sum_{i=1}^n K\left(\frac{x-X_i}{h}\right)}$$

其中K叫做核函数, 一般是一个非负的偶函数, 原点处的函数值最大, 在两侧迅速趋于零。例如正态密度函数, 或所谓双三次函数核:

$$K(x) = \begin{cases} \left(1 - |x|^3\right)^3 & |x| \leq 1 \\ 0 & |x| > 1 \end{cases}$$

R函数density()可以进行核回归估计, 这里作为举例我们对这个算法进行编程。先来编制双三次核函数的程序:

```
kernel.dcube <- function(u){
  y <- numeric(length(u))
  y[abs(u)<1] <- (1 - abs(u[abs(u)<1])^3)^3
  y[abs(u)>=1] <- 0
  y
}
```

注意上面分段函数不用if语句而是采用逻辑向量作下标的办法, 这样定义出的函数允许以向量和数组作自变量。

假设我们要画出核估计曲线, 一般我们取一个范围与各 X_i 范围相同的等间距向量 x 然后计算估计出的 $f(x)$ 的值。设观测数据自变量保存在向量 X 中, 因变量保存在向量 Y 中, 则按我们一般的程序设计思路, 可以写成下面的S程序:

```
kernel.smooth1 <- function(X, Y, kernel=kernel.dcube,
                           h=1, m=100, plot.it=T){
  x <- seq(min(X), max(X), length=m)
  fx <- numeric(m)
  for(j in 1:m){
    # 计算第j个等间距点的回归函数估计值
    fx[j] <- sum(kernel((x[j]-X)/h)*Y) /
              sum(kernel((x[j]-X)/h))
  }
  if(plot.it){
    plot(X, Y, type="p")
    lines(x, fx)
  }
  cbind(x=x, fx=fx)
}
```

注意上面程序中用了sum函数来避免一重对 i 的循环。但是, 上面的写法中仍有一重对 j 的循环, 使得程序运行较慢。如何改写程序把这个循环也取消呢? 办法是把计算看成是矩阵运算。首先, 如果 x 是一个标量, 则 $f(x)$ 可以写成:

$$f(x) = (K \cdot Y) / (K \cdot \mathbf{1})$$

其中 $K = K\left(\frac{x-X}{h}\right)$ 是一个与 X 长度相等(即长度为 n)的行向量, $\mathbf{1}$ 是一列1。现在, x 实际是一个长度为 m 的向量, 对 x 的每一个元素可以计算一个长度为 n 的行向量 $K\left(\frac{x[j]-X}{h}\right)$, 把这些行向量上下合并为一个矩阵 $K(m \times n)$, 则 KY 为长度为 m 的向量, 每一个元素对应于一个 $x[j]$ 。合并的矩阵可以用S的outer()函数来计算:

```
f <- function(u, v, kernel, h) kernel((u-v)/h)
K <- outer(x, X, f, kernel=kernel, h=h)
# outer()的第一个自变量对应于结果矩阵的行,
# 第二个自变量对应于列
```

分母为了算每一行的和只要对 K 右乘 $\mathbf{1}$, 于是结果的估计值向量为:

```
fx <- c( (K %*% Y) /
         (K %*% matrix(1,ncol=1,nrow=length(Y)) ) )
```

这样修改kernel.smooth1可以得到更精简的函数kernel.smooth2。

核回归中窗宽 h 的选择是比较难的,我们写的核回归函数应该允许用户输入 h 为一个向量,对向量中每一个窗宽计算一条拟合曲线并画在图中,结果作为函数返回值的一列。读者可以作为练习实现这个函数,并模拟 X 和 Y 的观测然后画核估计曲线, $f(x)$ 用 $\sin(x)$ 。对 h 我们可能必须用循环来处理了。

例2. Daubechies小波函数计算

小波是近年来得到广泛注意的一个数学工具,可以应用于各种数学分析问题。小波中一个重要的函数叫做尺度函数(scale function),它满足所谓双尺度方程:

$$\phi(x) = \sqrt{2} \sum_k h_k \phi(2x - k)$$

一种特殊的尺度函数是只在有限区间上非零的,叫做紧支集的。紧支集尺度函数可以在给定 $\{h_k\}$ 后用以下迭代公式生成:

$$\eta_0(x) = I_{[-0.5, 0.5]}(x)$$

$$\eta_{n+1}(x) = \sqrt{2} \sum_{k=0}^{2N-1} h_k \eta_n(2x - k)$$

其中 N 是正整数, $N=2$ 时 $h_0=0.482962913145$, $h_1=0.836516303738$, $h_2=0.224143868042$, $h_3 = -0.129409522551$ 。已知 $\phi(x)$ 的支集(不为零的区间)为 $[0, 2N-1]$, $\eta_n(x)$ 的支集包含于 $[-0.5, 2N-1]$ 中。作为例子,我们来编写计算 $\phi(x)$ 的S程序。

(1)因为 $\phi(x)$ 是函数,所以我们只能得到 x 在 $\phi(x)$ 的支集中的一些等间隔点上的函数值。为了用迭代计算这些值,我们先写出算法的S伪代码(基本程序结构使用S语法但某些复杂作用用自然语言描述):

```
Daubechies.phi <- function(nsamples=256, nrep=20){
  # nsamples --- 需要计算的点的个数
  # nrep --- 迭代次数
  x <- 需要计算的x坐标向量
  N <- 2
  h <- c(0.482962913145, 0.836516303738,
         0.224143868042, -0.129409522551)
  y0 <- 向量, 保存eta(x)在n=0时的初值,
        当x在[-0.5, 0.5]内时为1,其它为0
  for(r in 1:nrep){ # 迭代即循环
    y <- 与x长度相同的0向量
    for(k in 1:(2*N-1)){ # 计算公式中的求和
      yk <- eta(.)在自变量取2x-k处的值*sqrt(2),
            用上一步的y0来表达
      y <- y + yk
    }
  }
}
```



```

    } # for k
    y0 <- y # 把新的eta(.)放进y0中
            # 作为下一步迭代的开始值
  } # for r
  把x和y只保留x在[0, 2*N-1]中的部分
  返回x和y作为(x[i], phi(x[i]))
}

```

(2) 下面我们逐步细化这个算法, 并根据需要对算法进行必要的修改。设 $x[i] = a + (i-1)/(n_{\text{sample}}-1) \cdot (b-a)$, 我们发现, 这个算法中关键问题是如何把 $2x - k$ 处的函数值从 $y0$ 中查到。理想的情况是 $2x[i]-k$ 恰好等于某一个 $x[j]$, 这时 $\eta(2x[i]-k) = \eta(x[j]) = y0[j]$ 。如果 $2x[i]-k$ 在两个 $x[j]$ 之间, 可以取最近的一个 $x[j]$ 然后使用 $y0[j]$ 。经试验, 这样的做法不能保证迭代收敛。所以, 应该调整算法使得每一个 $2x[i]-k$ 都恰好落在某个 $x[j]$ 上面。这要求

$$2a + \frac{2(i-1)}{n-1}(b-a) = a + \frac{j-1}{n-1}(b-a)$$

其中 n 为采样点数。推导得

$$j = 1 + 2(i-1) - (k+1) \frac{n-1}{b-a}$$

所以只要 $n-1$ 能被 $b-a$ 整除即可。为此, 取 $a = -1, b = 2N-1, n = 1 + (2N)m$, m 为整数。但是, 这样取的 n 个点包含了区间 $[-1, 0)$, 在最后的結果中是要丢弃的, 所以如果采样至少要 n_{sample} 个点的话应该保证 $n \cdot (2N-1)/(2N)$ 大于等于 n_{sample} 。

(3) 如果必须要求按指定的采样点数给结果的话可以用线性插值得到所需的结果。线性插值的函数为 `approx(x, y, n)`, 其中 (x, y) 为要插值的散点(两个向量), n 是要求的采样点数(从 x 的最小值到最大值均匀采样)。

(4) 最后的程序如下:

```

# S-examp2.r
# Daubechies 迭代有限支集正交小波基构造

Daubechies.phi <- function(nsample=256, nrep=20,
                             plot.it=T, debugging=F,
                             nsample.exact=F){

  # nsample: 采样点数
  # plot.it: 是否每次绘图
  # debuggin: 是否输出调试信息
  # nsample.exact: 是否严格要求采样点数.
  #              取FALSE时可以适当放大.

  # 注意: 所有eta(x)的支撑的并集为[-0.5, 2N-1],
  # 但为了迭代时计算的点都落在采样点上所以要把
  # 左端点设为-1, 把采样个数适当放大.

  N <- 2
  a <- -1

```

```

b <- 2*N-1 # 采样区间[-1, 2N-1]
ns <- ceiling((nsample*(2*N)/(2*N-1)-1)/(2*N))*(2*N) + 1
# ns: 新的采样点数
m <- (ns-1)/(2*N)
# m: phi(x)的采样密度(x轴每单位有多少个点)
h <- c(0.482962913145, 0.836516303738,
       0.224143868042, -0.129409522551) * sqrt(2)

x <- seq(from=-1, to=2*N-1, length=ns)
# x: phi(x)的采样位置
y0 <- numeric(ns) # 初始函数
cond1 <- (x>=-0.5 & x<=0.5)
y0[cond1] <- 1.0
y0[!cond1] <- 0.0

for(r in 1:nrep){ # 迭代次数
  if(debugging && plot.it){
    plot(x, y0, type="l", main=paste("Iteration", r))
    #locator(n=1)
  }
  y <- numeric(ns) # 迭代结果
  for(k in 0:(2*N-1)){
    # 首先计算eta(2x-k),
    # 而x对应于a, a+1/m, a+2/m, ..., b=a+(ns-1)/m
    # 共ns个点, 我们只要找出2x-k对应的下标即可.
    # 这些下标里面有些是出界的,
    # 只要找在1到ns内的即可.
    yk <- numeric(ns) # 保存eta(2x-k)
    i <- seq(along=x)
    j <- 2*i - 1 - (k+1)*m
    cond2 <- (j>=1) & (j<=ns)
    yk[cond2] <- y0[j[cond2]]*h[k+1]
    y <- y + yk
  } # k
  y0 <- y # 新的迭代初值
} # r

xphi <- x[(1+m):ns]
yphi <- y[(1+m):ns]
ns <- ns - m
# 当前phi(x)从0到2N-1采样的个数.
# 去掉了[-1, 0)的m个点

if(plot.it){
  plot(xphi, yphi, type="l",
       main=expression(phi(x)), xlab="x", ylab="")
  abline(h=0)
}

if(nsample.exact){

```

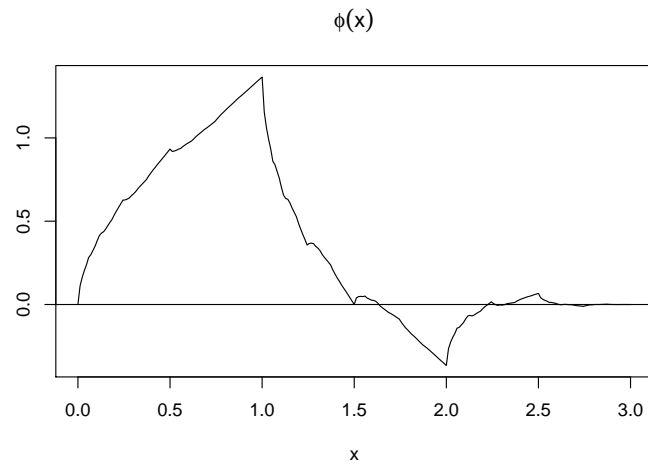


图 6.3: N=2的Daubechies小波尺度函数

```
# 必须要求nsample个点.
# 使用线性插值(approx, approxfun)
# 或样条插值(spline, splinefun)
xy1 <- approx(xphi, yphi, n=nsample)
xphi <- xy1$x
yphi <- xy1$y
}

invisible(list(xphi=xphi, yphi=yphi))
}
```

结果见图6.3。

6.10 图形

6.10.1 常用图形

S有很强的图形功能, 它可以用简单的函数调用迅速作出数据的各种图形, 当你熟悉了S图形的技术之后也可以指定许多图形选项按自己的要求定制图形。它的另一个特色是同一个绘图函数对不同的数据对象可以作出不同的图形。例如, 用6.1.2读入的cl数据框:

```
> attach(cl)
> plot(Height)
> plot(Sex)
```

第一个plot()绘制身高的散点图(图6.4), 第二个plot()绘制性别频数条形图(图6.5)。

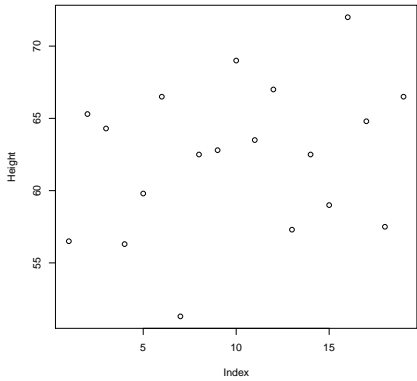


图 6.4: 身高散点图

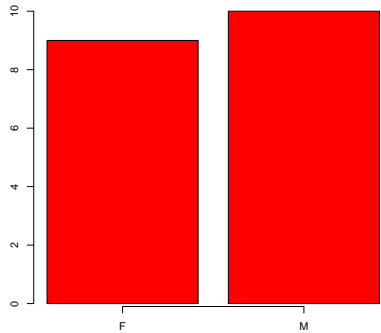


图 6.5: 性别频数条形图

最常用的绘图函数为plot(), 用plot()作两个变量x与y的散点图, 使用如下例的方法:

```
> plot(Height, Weight, main="体重对身高的回归",  
+ xlab="身高", ylab="体重")
```

见图6.6。上例也演示了S中如何输入较长的语句: 只要语句明显地未完成(比如, 缺右括号), 系统将给出一个加号作为续行提示。如果输入“x <- 1+2”时要拆行, 可以在赋值号后拆, 可以在加号与2之间拆, 但是如果在x后拆则只能显示x的当前值, 如果在1与加号之间拆只能把1赋给x。

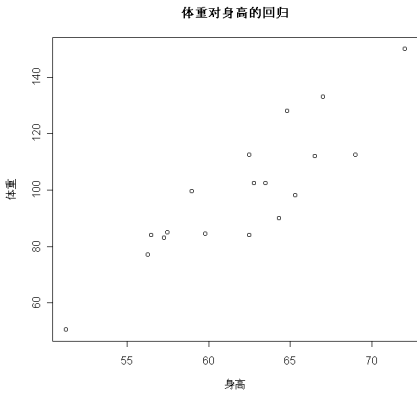


图 6.6: 体重对身高散点图

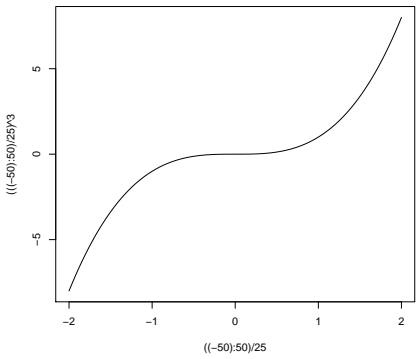


图 6.7: 三次多项式曲线

为了绘制连线图, 只要在plot()函数中加type="l"选项, 如:

```
> plot(((-50):50)/25, (((-50):50)/25)^3, type="l")
```

见图6.7。可以绘制变量的茎叶图, 如:

```
> stem(Weight)

The decimal point is 1 digit(s)
to the right of the |

 5 | 1
 5 | 67789
 6 | 033344
 6 | 557779
 7 | 2
```

绘制一个变量的盒形图, 如:

```
> boxplot(Weight)
```

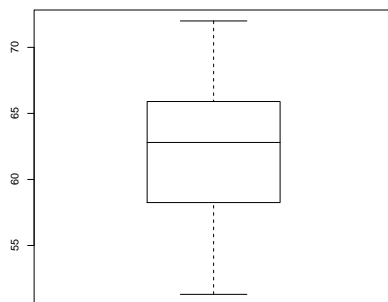


图 6.8: 体重的盒形图

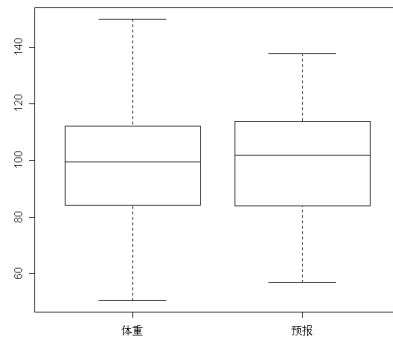


图 6.9: 并排盒形图

结果见图6.8。可以绘制几个变量并排的盒形图, 比如先计算Weight对Height的回归把拟合结果存入p1, 然后绘制并排盒形图:

```
> fit1 <- lm(Weight ~ Height)
> p1 <- predict(fit1, cl)
> boxplot(list("体重"=Weight, "预报"=p1))
```

见图6.9。用hist()函数可以绘制直方图。例如:

```
> hist(Weight)
```

见图6.11。用qqnorm()函数绘制正态概率图, 如:

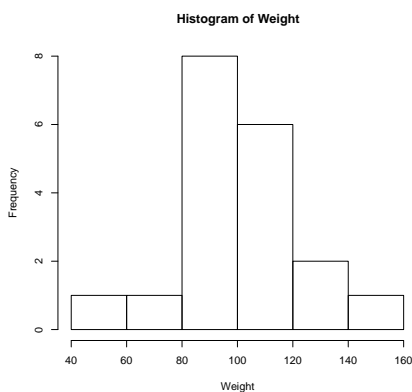


图 6.10: 体重的直方图

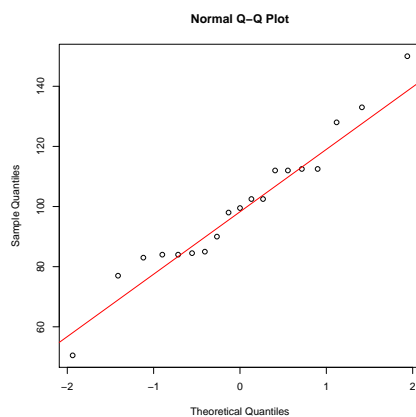


图 6.11: 体重的正态概率图

```
> qqnorm(Weight)
> qqline(Weight, col="red")
```

见图6.11。

6.10.2 高级图形函数

S的图形函数分为两类：高级图形函数—直接绘制图形并可自动生成坐标轴等附属图形元素；低级图形函数—可以修改已有的图形或者为绘图规定一些选择项。高级图形函数总是开始一个新图。下面我们介绍常用的高级图形函数，以及用来修饰这些高级图形函数的常用可选参数。

最常用的是plot()函数。比如，plot(x,y)(其中x, y是向量)对两个变量画散点图。如果z是一个定义了x变量和y变量的列表，或者一个两列的矩阵则用plot(z)也可以画散点图。如果x是一个时间序列对象(时间序列对象用ts()函数生成)，plot(x)绘制时间序列曲线图。如果x是一个普通向量，则绘制x的值对其下标的散点图。如果x是复数向量则绘制虚部对实部的散点图。如果f是一个因子，则plot(f)绘制f的条形图(每个因子水平的个数)。如果f是因子，y是同长度的数值向量，则plot(f,y)对f的每一因子水平绘制y中相应数值的盒形图。如果d是一个只有数值型数据的数据框，则plot(d)对d的每两个变量之间作图(散点图等)。

如果X是一个数值型矩阵或数据框，用pairs(X)可以绘制每两列之间的散点图矩阵。这在变量个数不太多时可以同时看到多个变量的两两关系，变量太多时则每个图太小，意义不大。

协同图(coplot)是一种多变量的探索性分析图形。其形式为coplot(y ~ x | z)，其中x和y是数值型向量，z是同长度的因子。对z的每一水平，绘制相应组的x和y的散点图。如：

```
> attach(cl)
> coplot(Weight ~ Height | Sex)
```

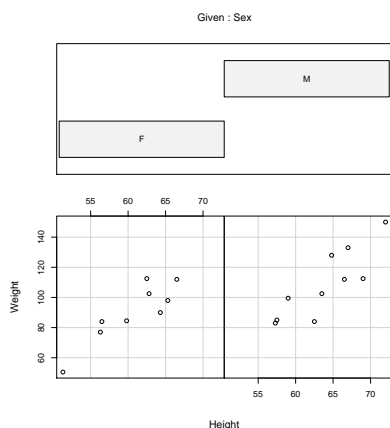


图 6.12: 以性别为条件的体重对身高的散点图

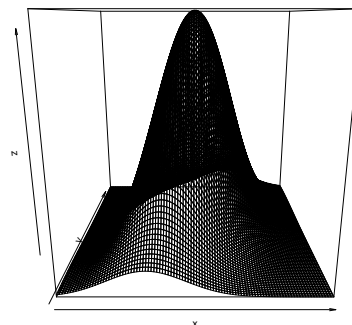


图 6.13: 二维正态密度曲面

产生图6.12, 对不同性别分别绘制了体重对身高的散点图。如果 z 是一个数值型变量, 则`coplot()`先对 z 的取值分组, 然后对 z 的每一组取值分别绘图。甚至可以用如`coplot(y ~ x | x1+x2)`表示对 x_1 和 x_2 的每一水平组合绘图。`coplot()`和`pairs()`函数缺省绘制散点图, 但可以用一个`panel=`参数指定其它的低级绘图函数, 如`lines`, `panel.smooth`等。

`qqnorm(x)`, `qqline(x)`, `qqplot(x,y)`作分位数-分位数图。`qqnorm(x)`对向量 x 作正态概率(纵轴为次序统计量值, 横轴为对应该次序统计量的标准正态分布分位数值)。`qqline(x)`在`qqnorm(x)`图上面画一条拟合曲线。`qqplot(x,y)`把 x 和 y 的次序统计量分别画在 x 轴和 y 轴以比较两个变量的分布。

`hist(x)`作向量 x 的直方图。缺省时自动确定分组, 也可以用`nclass=`参数指定分组个数, 或者用`breaks=`参数指定一个分组点向量。如果指定了`prob=T`则纵轴显示密度估计。

S 也可以作三维图、等值线图 and 等值色图, 函数为`persp()`, `contour()`和`image()`, 如:

```
x <- seq(-3,3, length=100)
y <- x
f <- function(x,y,ssq1=1, ssq2=2, rho=0.5){
  det1 <- ssq1*ssq2*(1 - rho^2)
  s1 <- sqrt(ssq1)
  s2 <- sqrt(ssq2)
  1/(2*pi*det1) * exp(-0.5 / det1 * (
    ssq2*x^2 + ssq1*y^2 - 2*rho*s1*s2*x*y))
}
z <- outer(x, y, f)
persp(x, y, z)
contour(x, y, z)
image(x, y, z)
```


表 6.1: 高级图形选项

add=T	使函数向低级图形函数那样不是开始一个新图形而是在原图基础上添加。
axes=F	暂不画坐标轴, 随后可以用axis()函数更精确地规定坐标轴的画法。缺省值是axes=T, 即有坐标轴。
log="x" log="y" log="xy"	把x轴, y轴或两个坐标轴用对数刻度绘制。
type= type="p" type="l" type="b" type="o" type="h" type="s" type="S" type="n"	规定绘图方式: 绘点 画线 绘点并在中间用线连接 绘点并画线穿过各点 从点到横轴画垂线 阶梯函数; 左连续 阶梯函数; 右连续 不画任何点、线, 但仍画坐标轴并建立坐标系, 适用于后面用低级图形函数作图。
xlab="字符串" ylab="字符串" main="字符串" sub="字符串"	定义x轴和y轴的标签。缺省时使用对象名。 图形的标题。 图形的小标题, 用较小字体画在x轴下方。

垂线, 等等)。locator()的返回值是一个列表, 有两个变量(元素)x和y, 分别保存点击位置的横坐标和纵坐标。

例如, 为了在已经绘制的曲线图中找一个空地方标上一行文本, 只要使用如下程序:

```
> text(locator(1), "Normal density", adj=0)
```

text()函数的adj参数用一个数字表示文本串相对于给定的坐标的画法, adj=0表示给定坐标为文本串左侧的坐标, adj=1表示给定坐标为文本串右侧的坐标, adj=0.5表示给定坐标为文本串中间的坐标。

函数identify(x, y, labels)在运行时也会停下来等待用户点击, 直到按了鼠标右键, 然后返回用户在图形中用鼠标点击的点的序号, 点击时对点击的点加标签。参数x和y给出要识别的各个点的坐标。labels参数指定点击某个点时要在旁边绘制的文本标签, 缺省时标出此点的序号, 如果只需要返回值而不想画任何标记则可以在调用此函数时加一个plot=F参数。注意identify()与locator()不同, locator()返回图中任意点击位置的坐标, 而identify()只返回离点击位置最近的点的序号。

例如, 我们在向量x和y中有若干个点的坐标, 运行如下程序:

```
> plot(x, y)
> identify(x, y)
```

这时显示转移到图形窗口, 进入等待状态, 用户可以点击图中特别的点, 该点的序号就会在旁边标出。为了结束, 只要单击右键并选择停止。返回结果为你点

表 6.2: 低级图形函数

points(x,y) lines(x,y)	在当前图形上叠加一组点或线。可以使用plot()的type=参数来指定绘制方法, 缺省时points() 画点,lines()画线。
text(x,y, labels, ...)	在由坐标x和y给出的位置标出由labels指定的字符串。labels可以是数值型或字符型的向量,labels[i]在x[i],y[i]处标出。
abline(a, b) abline(h=y) abline(v=x) abline(lm.obj)	在当前图形上画一条直线。两个参数a, b分布给出截距和斜率。指定h=参数时绘制水平线,指定v=参数时绘制垂直线。以一个最小二乘拟合结果lm.obj作为参数时由lm.obj的\$coefficients 成员给出直线的截距和斜率。
polygon(x, y, ...)	按向量x给出的横坐标和向量y给出的纵坐标确定顶点绘制多边形。可以用col=参数指定一个颜色填充多边形内部。
legend(x, y, legend, ...) legend(, angle=v) legend(, density=v) legend(, fill=v) legend(, col=v) legend(, lty=v) legend(, pch=v) legend(, marks=v)	legend函数用来在当前图形的指定坐标位置绘制图例。图例的说明文字由向量legend提供。至少下面的v值要给出以确定要对什么图例进行说明, v是长度与legend相同的向量。 angle参数指定几种阴影斜角。 density参数指定几种阴影密度。 fill参数指定几种填充颜色。 col参数指定几种颜色。 lty参数指定几种线型。 pch参数指定几种散点符号。为字符型向量。 marks参数也指定几种散点符号,但使用散点符号数值代号,为数值型向量。
title(main, sub)	绘制由main指定的标题和由sub指定的小标题。
axis(side, ...)	绘制一条坐标轴。这之前的绘图函数必须已经用axes=F选项抑制了自动的坐标轴。参数side 指定在哪一边绘制坐标轴,取值为1到4,1为下边,然后逆时针数。可以用at=参数指定刻度位置,用labels参数指定刻度处的标签。

击的各个点的序号。

6.10.6 图形参数的使用

前面我们已经看到了如何用main=, xlab=等参数来规定高级图形函数的一些设置。在实际绘图,特别是绘制用于演示或出版的图形时, S用缺省设置绘制的图形往往不能满足我们的要求。S提供了一系列所谓图形参数, 通过使用图形参数可以修改图形显示的所有各方面的设置。图形参数包括关于线型、颜色、图形排列、文本对齐方式等各种设置。每个图形参数有一个名字, 比如col代表颜色, col="red"表示红色。每个图形设备(作图的目的地)有一套单独的图形参数。

设置图形参数分为两种：永久设置与临时设置。永久设置使用par()函数进行设置, 设置后在退出前一直保持有效；临时设置则是在图形函数中加入图形参数, 如上面的例子：

```
> text(locator(1), "Normal density", adj=0)
```

中的adj参数。

par()函数用来访问或修改当前图形设备的图形参数。如果不带参数调用, 如:

```
> par()
$adj
[1] 0.5

$ann
[1] TRUE

. . . . .

$ylog
[1] FALSE
```

结果为一个列表, 列表的各元素名为图形参数的名字, 元素值为相应图形参数的取值。

如果调用时指定一个图形参数名的向量作为参数, 则只返回被指定的图形参数的列表:

```
> par(c("col", "lty"))
$col
[1] "black"

$lty
[1] "solid"
```

调用时指定名字为图形参数名的有名参数, 则修改指定的图形参数, 并返回原值的列表:

```
> oldpar <- par(col=4, lty=2)
> oldpar

$col
[1] "black"

$lty
[1] "solid"
```

因为用par()修改图形参数是保持到退出以前都有效的, 而且即使是在函数内此修改仍是全局的, 所以我们可以利用如下的惯用法, 在完成任务后恢复原来的图形参数:

```
> oldpar <- par(col=4, lty=2)
. . . . . (需要修改图形参数的绘图任务)
> par(oldpar) # 恢复原始的图形参数
```

除了象上面那样用`par()`函数永久修改图形参数, 我们还可以在几乎任何图形函数中指定图形参数作为有名参数, 这样的修改是临时的, 只对此函数起作用。例如:

```
> plot(x, y, pch="+")
```

就用图形参数`pch`指定了绘散点的符号为加号。这个设定只对这一张图有效, 对以后的图形没有影响。

6.10.7 图形参数详解

鉴于绘制有特殊需要的图形是S的一个强项, 而使用图形参数是完成此类任务的重要手段, 我们在这里较详细地介绍S的各种图形参数。这些图形参数可以大体上分为以下的几个大类, 我们将分别介绍:

- 图形元素控制
- 坐标轴与坐标刻度
- 图形边空
- 一页多图

一、图形元素

图形由点、线、文本、多边形等元素构成。下列的图形参数用来控制图形元素的绘制细节:

- `pch="+"` 指定用于绘制散点的符号。绘制的点往往略高于或低于指定的坐标位置, 只有`pch="."`没有这个问题。
- `pch=4` 如果`pch`的值为从0到18之间的一个数字, 将使用特殊的绘点符号。下例可以显示所有特殊绘点符号:

```
> plot(c(0, 100), c(0, 100), type="n",
+      axes=F, xlab='', ylab='')
> legend(10,90, as.character(0:9), pch=0:9)
> legend(50,90, as.character(10:18), pch=10:18)
```

- `lty=2` 指定画线用的线型。缺省值`lty=1`是实线。从2开始是各种虚线。
- `lwd=2` 指定线粗细, 以标准线粗细为单位。这个参数影响数据曲线的线宽以及坐标轴的线宽。下例绘制正弦曲线图:

```
> oldpar <- par(lwd=2)
> x <- (0:100)/100*2*pi
> plot(x, sin(x), type="l", axes=F)
> abline(h=0)
> abline(v=0)
> par(oldpar)
```

- `col=2` 指定颜色,可应用于绘点、线、文本、填充区域、图象。颜色值也可以用象"red","blue"这样的颜色名指定。
- `font=2` 用来指定字体的整数。一般`font=1`是正体,2是**黑体**, 3是斜体,4是**黑斜体**。
- `font.axis`, `font.lab`, `font.main`, `font.sub` 分别用来指定坐标刻度、坐标轴标签、标题、小标题所用的字体。
- `adj=-0.1` 指定文本相对于给定坐标的对齐方式。取0表示左对齐,取1表示右对齐,取0.5表示居中。此参数的值实际代表的是出现在给定坐标左边的文本的比例, 所以`adj=-0.1`的效果是文本出现在给定坐标位置的右边并空出相当于文本10%长度的距离。
- `cex=1.5` 指定字符放大倍数。

二、坐标轴与坐标刻度

许多高级图形带有坐标轴, 还可以先不画坐标轴然后用`axis()`单独加。函数`box()`用来画坐标区域四周的框线。

坐标轴包括三个部件: 轴线(用`lty`可以控制线型), 刻度线, 刻度标签。它们可以用如下的图形参数来控制:

- `lab=c(5, 7, 12)` 第一个数为x轴希望画几个刻度线,第二个数为y轴希望画几个刻度线,这两个数是建议性的; 第三个数是坐标刻度标签的宽度为多少个字符,包括小数点,这个数太小会使刻度标签四舍五入成一样的值。
- `las=1` 坐标刻度标签的方向。0表示总是平行于坐标轴,1表示总是水平,2表示总是垂直于坐标轴。
- `mgp=c(3,1,0)` 坐标轴各部件的位置。第一个元素为坐标轴位置到坐标轴标签的距离,以文本行高为单位。第二个元素为坐标轴位置到坐标刻度标签的距离。第三个元素为坐标轴位置到实际画的坐标轴的距离,通常是0。
- `tck=0.01` 坐标轴刻度线长度,单位是绘图区域大小,值为占绘图区域的比例。`tck`小于0.5时x轴和y轴的刻度线将统一到相同的长度。取1时即画格子线。取负值时刻度线画在绘图区域的外面。
- `xaxs="s", yaxs="d"` 控制x轴和y轴的画轴方法。

取值为"s" (即standard) 或"e" (即extended) 的时候数据范围控制在最小刻度和最大刻度之间。取"e"时如果有数据点十分靠近边缘轴的范围会略微扩大。这种画轴方式有时会在轴的一边留下太大的空白。

取值为"i" (即internal) 或"r" (此为缺省) 使得刻度线都落在数据范围内部,而"r"方式所留的边空较小。

取值设为"d"时会锁定此坐标轴,后续的图形都使用与它完全相同的坐标轴, 这在要生成一系列可比较的图形的时候是有用的。要解除锁定需要把这个图形参数设为其它值。

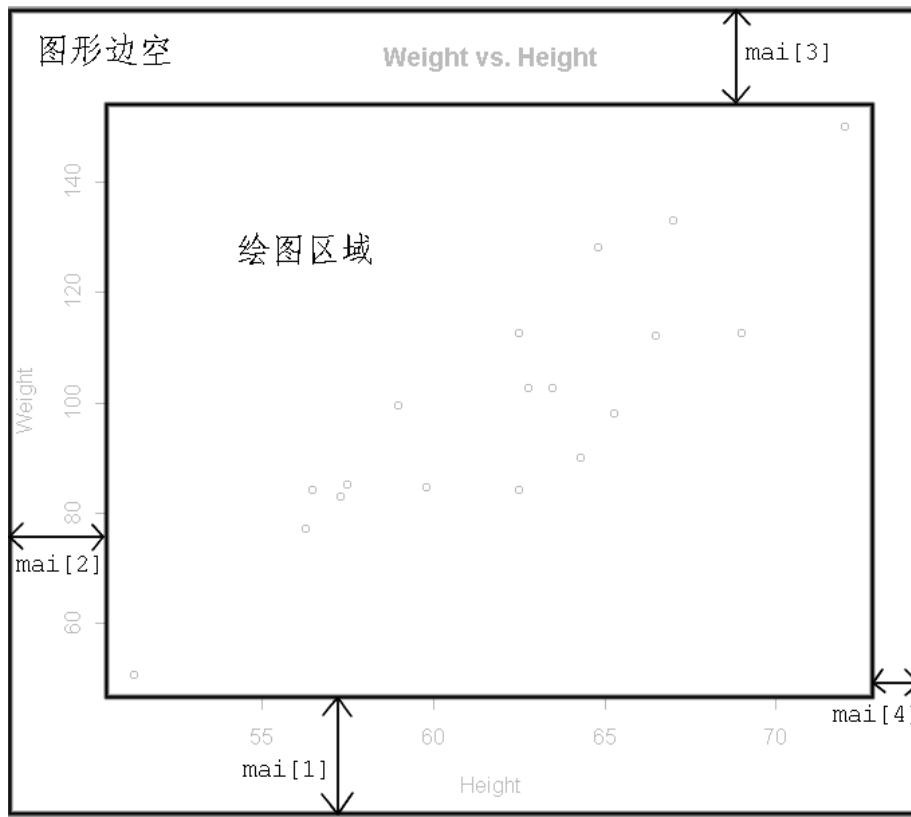


图 6.16: 图形边空

三、图形边空

S中一个单独的图由绘图区域(绘图的点、线等画在这个区域中)和包围绘图区域的边空组成, 边空中可以包含坐标轴标签、坐标轴刻度标签、标题、小标题等, 绘图区域一般被坐标轴包围。见图6.16。

边空的大小由mai参数或mar参数控制, 它们都是四个元素的向量, 分别规定下方、左方、上方、右方的边空大小, 其中mai取值的单位是英寸, 而mar的取值单位是文本行高度。例如:

```
> par(mai=c(1, 0.5, 0.5, 0))
> par(mar=c(4, 2, 2, 1))
```

这两个图形参数不是独立的, 设定一个会影响另一个。S缺省的图形边空常常太大, 以至于有时图形窗口较小时边空占了整个图形的很大一部分。通常我们可以取消右边空, 并且在不用标题时可以大大缩小上边空。例如下例可以生成十分紧凑的图形:

```
> oldpar <- par(mar=c(2,2,1,0.2))
> plot(x,y)
```

在一个页面上画多个图时边空自动减半, 但我们往往还需要进一步减小边空才能使多个图有意义。

四、一页多图

R可以在同一页面开若干个按行、列排列的窗格, 在每个窗格中可以作一幅图。每个图有自己的边空, 而所有图的外面可以包一个“外边空”, 见图6.17。

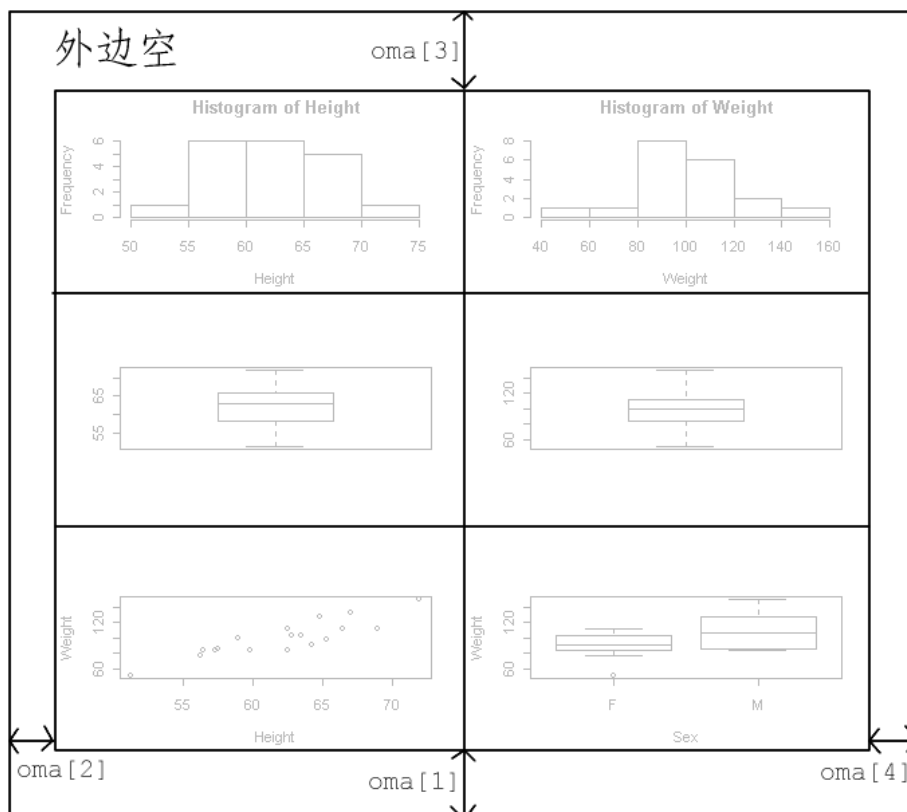


图 6.17: 一页多图

一页多图用mfrow参数或mfccl参数规定, 如:

```
> par(mfrow=c(3,2))
```

表示同一页有三行两列共六个图, 而且次序为按行填放。类似地,

```
> par(mfcol=c(3,2))
```

规定相同的窗格结构, 但是次序为按列填放, 即先填满第一列的三个再填第二列。要取消一页多图只要再运行

```
> par(mfrow=c(1,1))
```

即可。

缺省时无外边空。为了规定外边空大小, 可以在`par()`中用`omi`参数或`oma`参数。`omi`参数使用英寸为单位, `oma`参数以文本行高为单位, 两个参数均为四个元素的向量, 分别给出下、左、上、右方的边空大小。如:

```
> par(oma=c(2,0,3,0))
```

函数`mtext`用来在外边空加文字标注。其用法为

```
mtext(text, side = 3, line = 0, outer = FALSE)
```

其中`text`为要加的文本内容, `side`表示在哪一边写(1为下, 2为左, 3为上, 4为右), `line`表示边空从里向外数的第几行, 最里面的一行是第0号, `outer=TRUE`时使用外边空, 否则会使用当前图的边空。例如:

```
> par(mfrow=c(2,2), oma=c(0,0,3,0), mar=c(2,1,1,0.1))
> plot(x);plot(y);boxplot(list(x=x,y=y));plot(x,y)
> mtext("Simulation Data", outer=T, cex=1.5)
```

在多图环境中还可以用`mfg`参数来直接跳到某一个窗格, 比如

```
> par(mfg=c(2,2,3,2))
```

表示在三行两列的多图环境中直接跳到第二行第二列位置。`mfg`参数的后两个表示多图环境的行、列数, 前两个表示要跳到的位置。

可以不使用多图环境而直接在页面中的任意位置产生一个窗格来绘图, 参数为`fig`, 如:

```
> par(fig=c(4,9,1,4)/10)
```

此参数为一个向量, 分别给出窗格的左、右、下、上边缘的位置, 取值为占全页面的比例, 比如上面的例子在页面的右下方开一个窗格作图。

6.10.8 图形设备

S作图支持各种图形设备, 其中常用的是显示器和PostScript打印机。在一个S运行期间可以有多个图形设备同时存在。在R中, 用

```
> x11()
```

打开图形窗口绘图, 在S-PLUS中, 用

```
> win.graph()
```

打开图形窗口绘图。再次调用这样的函数将打开第二个图形窗口。用


```
> dev.list()
```

可显示以打开的图形设备的列表。

要关闭一个图形设备, 用

```
> dev.off()
```

这可以使得图形得以完成, 例如对于postscript设备关闭设备时可完成打印或存盘。用graphics.off()函数可以关闭所有打开的图形设备。

MS Windows下的R可以把显示窗口中的图形复制到剪贴板或存为各种格式的图形文件, 包括WMF, PostScript, PDF, PNG, BMP, JPEG, 这样我们可以用R生成所需图形然后存为需要的格式。MS Windows下的S-PLUS也具有类似功能。

各版本的R和S-PLUS都支持生成PostScript图形的功能, 生成的图形可以直接用于 \LaTeX 排版。如果用MS Word排版则可以把屏幕图形存为WMF等格式。生成PostScript文件的设备可以用如下函数打开:

```
> postscript(file="result1.ps",  
+ horizontal=FALSE, width=5, height=3)
```

这时用图形命令生成一个页面的图形, 然后用dev.off()关闭设备, 则可生成文件result1.ps。postscript()函数中horizontal参数指定是否将图旋转90度使得x轴平行于纸的长边, width和height规定图的宽和高, 单位是英寸(1英寸=2.54厘米)。

在打开了多个图形设备后可以用dev.set()函数来选择当前设备, dev.next()和dev.prev()分别返回下一个和上一个图形设备。比如dev.set(dev.prev())选择上一个图形设备。

6.11 S的对象

S是一种面向对象的语言。一般说来, S的对象包含了若干个元素作为其数据, 这些元素的个数叫做此对象的长度(length), 这些元素的共同的类型叫做此对象的模式(mode)。另外, 对象还可以包含一些特殊数据, 称为属性(attribute), 如列表的每一个成员(元素)都可以有变量名, 这些变量名组成的字符型向量为此列表的names属性。

S的面向对象能力依赖于对类(class)属性的使用。S的对象有的是简单对象, 这样的对象没有类(class)属性, 可以认为其类为缺省类(default); 有的是属于某一类的对象, 这些对象有一个类(class)属性, 同类的对象具有相同的特征, 可以为同类的对象定义针对这一类的特殊操作(如显示、绘图), 在面向对象术语中叫做方法。比如, 向量是简单对象, 它没有类属性; 数据框也是对象, 但是数据框有一个类属性class=data.frame。

6.11.1 固有属性: mode和length

S对象都有两个基本的属性(attribute): 类型(mode)属性和长度(length)属性。

S对象可分为单纯的(atomic)和复合的(recursive)两种, 单纯对象的所有元素都是同一种基本类型(如数值、字符串), 元素不再是对象, 这样的对象的类型(mode)有logical(逻辑型)、numeric(数值型)、complex(复数型)、character(字符型)等等; 复合对象的元素可以是不同类型, 每一个元素是一个对象, 这样的对象最常用的是列表。例如, 向量(vector)是单纯对象, 它的所有元素都必须是相同类型, 数值型向量的所有元素必须为数值型, 字符型向量的所有元素必须为字符型; 列表(list)是复合对象, 类型(mode)为列表(list), 列表的每一个元素(变量)都可以是一个S对象, 比如列表元素可以为一个数, 一个字符串, 一个向量, 甚至一个列表。

S对象有一种特别的空值类型, 只有一个特殊的NULL值为这种类型, 表示没有值(不同于NA, NA是一种特殊值, 而NULL根本没有对象值)。

为了判断对象的类型, S定义了许多个类似于is.numeric()这样的函数。比如, is.numeric(x)用来检验对象x是否数值型, 返回一个逻辑型标量结果; is.character()检验对象是否字符型, 等等。

长度属性表示S对象元素的个数, 比如length(2:4)等于3。注意向量允许长度为0, 数值型向量长度为零表示为numeric()或numeric(0), 字符型向量长度为零表示为character()或character(0)。

S可以强制进行类型转换, 例如

```
> z <- 0:9
> digits <- as.character(z)
> digits
[1] "0" "1" "2" "3" "4" "5" "6" "7" "8" "9"
> d <- as.numeric(digits)
> d
[1] 0 1 2 3 4 5 6 7 8 9
```

第二个赋值把数值型的z转换为字符型的digits。第三个赋值把digits又转换为了数值型的d, 这时d和z是一样的了。S还有许多这样的以as.开头的类型转换函数。

S允许对超出对象长度的下标赋值, 这时对象长度自动伸长以包括此下标, 未赋值的元素取缺失值(NA), 例如:

```
> x <- numeric()
> x[3] <- 100
> x
[1] NA NA 100
```

要缩短对象的长度又怎么办呢? 只要给它赋一个子集就可以了。例如:

```
> x <- 1:4
> x <- x[1:2]
> x
[1] 1 2
```

6.11.2 访问对象属性

对象属性是对象包含的数据中除元素以外的特殊数据, 每个属性有一个属性名, 有一个属性值。S定义了两个函数attributes和attr来访问对象的属

性。attributes(object)返回对象object的各特殊属性组成的列表, 其中不包括固有属性mode和length。例如:

```
> x <- c(apple=2.5, orange=2.1)
> attributes(x)
$names
[1] "apple" "orange"
```

可以用attr(object, name)的形式存取对象object的名为name的属性。例如:

```
> attr(x, "names")
[1] "apple" "orange"
```

也可以把attr()函数写在赋值的左边以改变属性值或定义新的属性, 例如:

```
> attr(x, "names") <- c("apple", "grapes")
> x
  apple grapes
    2.5    2.1
> attr(x, "type") <- "fruit"
> x
  apple grapes
    2.5    2.1
attr(,"type")
[1] "fruit"
> attributes(x)
$names
[1] "apple" "grapes"

$type
[1] "fruit"
```

这种对一个函数赋值的语法是在其它语言中极为少见的, 而S中则经常使用这样的写法。实际上, attr(x, "names")在这里不应该看成是一个函数值, 而应该看成是用来保存对象x的names属性的变量名。

6.11.3 对象的类

S用类(class)属性来支持面向对象的编程风格。对象的类属性区分对象的类, 对于同一类的对象可以定义一组特殊操作, 这一点和其它面向对象语言类似。面向对象风格的最重要的特点就是数据抽象与封装。所谓数据抽象与封装是指对象的用户要访问或修改对象只能通过对象提供的服务来进行, 用户不能看到对象内部的实现细节。这样用户不会直接修改对象的数据从而保护了数据的完整性, 而且用户只需要知道对象提供了哪些服务, 即使对象内部的实现改变了, 只要接口不变则用户程序不必改变。这样的做法可以提高程序的安全性和可重用性。

常见的面向对象语言一般先定义一个类, 这个类定义了一些数据结构, 然后有一些函数叫做“方法”可以操作这些数据。所谓对象, 是由某个类生成的实例, 其数据结构由所属的类定义, 而实际存储的数据则是属于对象本身的。对象拥有其所属类的所有方法, 方法在调用时操作的是属于这个对象的数据。

S也支持面向对象编程, 但是做法与常见的面向对象语言有很大差别。S对象的类由其类(class)属性指定, 每一个类都可以定义本类的服务, 服务以函数形式定义, 调用格式为“函数名(对象, 其它自变量)”。可见S的类机制是比较松散的, 它不像常见的面向对象语言那样必须先定义类的所有数据结构与方法, 而是可以随时定义函数作为类对象的服务。另外, S还定义了一系列的所谓“通用函数(general functions)”, 通用函数也是对象提供的服务, 但不同类的对象都可以使用相同的通用函数名字调用, 同一个通用函数可以针对不同类的对象起到相似的作用。用户只需要记忆很少的几个通用函数的名字, 就可以对几乎所有对象调用这些函数。比如, 通用print()函数用来显示对象, 它可以显示向量和矩阵, 但显示方法不同; 通用函数plot()函数用来画对象的图形, 对一个向量画图plot()画散点图, 纵轴为各元素值, 横轴为元素下标; 对一个时间序列对象画图plot()将画一条时间序列曲线, 并用年月等标记时间轴。

S的每一个通用函数实际是一组函数, 有一个共同的名字, 在调用时根据自变量的类(class)的不同决定调用一组中的哪一个函数。例如, 对向量x调用print(x)实际调用的是print.default(x), 对数据框x调用print(x)则实际调用的是print.data.frame(x)。如果自变量没有类属性, 或者此通用函数没有为此类自变量设计特殊的操作, 通用函数总有一个缺省方法可以调用(如print.default)。通用函数针对某一类的对象的特殊函数的命名为“通用函数名.类名()”。

对某一种类的对象有特殊操作的通用函数可以有很多个, 比如, 对data.frame类的对象定义了特殊操作的通用函数就有:

```
[, [[<-, any, as.matrix,
[<-, model, plot, summary,
```

等等。如果对data.frame类的对象d调用plot(d), 实际调用的函数是plot.data.frame(d)。要列出所有对某类有特殊操作的通用函数, 可以用

```
> methods(class="data.frame")
[1] "Math.data.frame" "Ops.data.frame"
. . . . .
[27] "summary.data.frame" "t.data.frame"
[29] "transform.data.frame" "xpdrows.data.frame"
```

其中t.data.frame就是调用t(d)时实际调用的函数。

也可以列出某通用函数对各类的特殊定义, 例如:

```
> methods(plot)
[1] "plot.data.frame" "plot.default" "plot.density" "plot.factor"
[5] "plot.formula" "plot.function" "plot.histogram" "plot.lm"
[9] "plot.mlm" "plot.mts" "plot.new" "plot.POSIXct"
[13] "plot.POSIXlt" "plot.table" "plot.ts" "plot.window"
[17] "plot.xy"
```

比如, plot.factor是对因子对象调用plot()函数是实际调用的函数。

为了暂时去掉一个有类的对象的class属性, 可以使用unclass(object)函数。

6.12 S初等统计

S-PLUS和R中已经集成了常用的统计功能,一些比较新的统计方法也得到了实现。

调用S函数可以完成基本的统计。这里我们介绍如何用S进行描述统计、探索性数据分析、分布研究、常见的假设检验。

6.12.1 单变量数据分析

对放在向量或数据框中的变量我们可以用summary()函数来计算几个描述统计量。例如:

```
> summary(c1)
      Name      Sex      Age      Height      Weight
Length:19      F: 9   Min.   :11.00   Min.   :51.30   Min.   : 50.50
Mode   :character M:10  1st Qu.:12.00  1st Qu.:58.25  1st Qu.: 84.25
                        Median :13.00  Median :62.80  Median : 99.50
                        Mean   :13.32  Mean   :62.34  Mean   :100.03
                        3rd Qu.:14.50  3rd Qu.:65.90  3rd Qu.:112.25
                        Max.   :16.00  Max.   :72.00  Max.   :150.00
```

对数值型变量计算了平均值,最大、最小值,中位数,四分之一和四分之三分位数;对因子计算了频数统计。

对数值型变量x也可以用mean, sd, var, median, min, max等函数计算其简单统计量。可以用stem(x)绘制x的茎叶图,用hist(x)画直方图,用boxplot(x)画盒形图,用qqnorm(x)和qqline(x)画正态QQ图。

为了估计数值型变量的分布密度,除了用hist()画直方图外还可以用density()函数进行非参数密度估计。例如:

```
plot(density(Height), main="Height density")
```

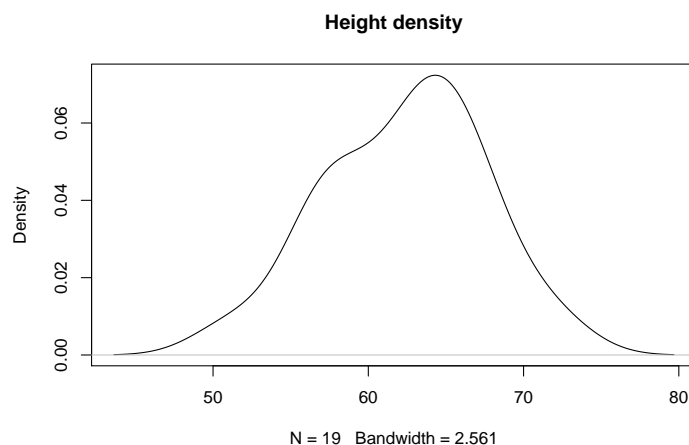


图 6.18: 身高的非参数密度估计

见图6.18。还可以设法把密度曲线加到直方图上,如:

```
h1 <- hist(Height, prob=T, plot=T)$density
h2 <- density(Height)
hist(Height, prob=T, ylim=range(h1, h2$y))
lines(h2)
```

见图6.19。

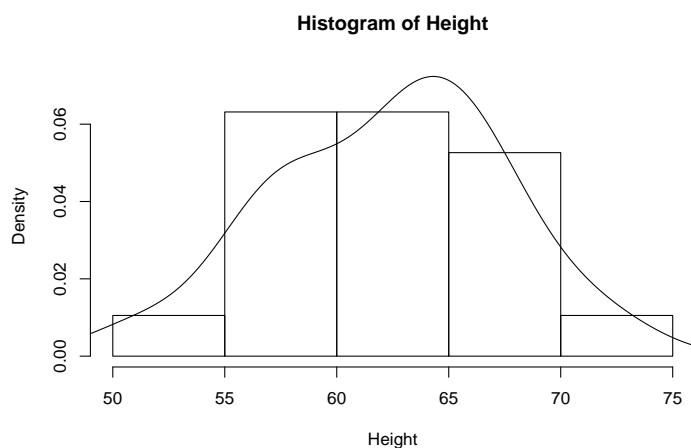


图 6.19: 身高的直方图及密度估计

对因子(分类变量) f 可以用`table(f)`统计其频数分布, 用`plot(f)`画其频数分布图。如

```
> table(Sex)
Sex
 F  M
 9 10
```

我们可以定义函数来画出常见的探索性数据分析图形。对连续型变量可以画直方图、盒形图、分布密度估计图和正态概率图。对离散型变量只要画其分布频数条形图即可, 分布频数用`table`函数计算。研究序列相关性可以作时间序列图和自相关函数图。因为这些图经常重复使用, 我们把它定义为函数, 在同一页面画出:

```
eda.shape <- function(x) {
  oldpar <- par(mfcol = c(2, 2),
    mar=c(2,2,0.2,0.2), mgp=c(1.2,0.2,0))
  hist(x, main="", xlab="", ylab="")
  plot(density(x), xlab = "x",
    ylab = "", main="")
  boxplot(x)
  qqnorm(x, main="", xlab="", ylab="")
}
```

```

    qqline(x)
    par(oldpar)
    invisible()
  }

eda.ts <- function(x) {
  oldpar <- par(mfrow=c(2,1),
    mar=c(2,2,1,0.2), mgp=c(1.2, 0.2, 0))
  plot.ts(x, main="", xlab="")
  acf(x, main="", xlab="")
  par(oldpar)
  invisible()
}

```

函数中最后的invisible()表示在命令行调用此函数时不要显示任何返回值。函数density用来作核密度曲线估计。

6.12.2 假设检验

R的ctest程序包(package)提供了常见的假设检验功能。为了调用程序包中的函数和数据, 需要先用library(包名)把它调入:

```
> library(ctest)
```

为了查看一个包中的内容可以在它提供的HTML帮助中查看, 也可以用R命令help(package=包名)查看。例如

```
> help(package=ctest)
```

为了检验正态性, 只要调用shapiro.test()函数就可以Shapiro-Wilk检验:

```

> shapiro.test(Height)

      Shapiro-Wilk normality test

data:  Height
W = 0.9791, p-value = 0.9312

```

可见身高数据的正态性很好。

函数t.test()可以进行单总体的t检验、两独立总体的t检验和成对t检验。例如,下面程序检验身高均值是否为 $\mu = 60$:

```

> t.test(Height, mu=60)

      One Sample t-test

data:  Height

```

```
t = 1.9867, df = 18, p-value = 0.06239
alternative hypothesis: true mean is not equal to 60
95 percent confidence interval:
 59.86567 64.80801
sample estimates:
mean of x
 62.33684
```

缺省进行的是双侧检验。为了进行单侧检验可以指定`alternative="greater"`或`"less"`。如:

```
> t.test(Height, mu=60, alternative="less")

      One Sample t-test

data:  Height
t = 1.9867, df = 18, p-value = 0.9688
alternative hypothesis: true mean is less than 60
95 percent confidence interval:
 -Inf 64.3765
sample estimates:
mean of x
 62.33684
```

对于用分组变量表示分组的两个独立组, 可以用`t.test(分析变量~ 分组变量)`的调用形式进行独立两组的t检验。如:

```
> t.test(Height ~ Sex)

      Welch Two Sample t-test

data:  Height by Sex
t = -1.4513, df = 16.727, p-value = 0.1652
alternative hypothesis:
true difference in means is not equal to 0
95 percent confidence interval:
 -8.155098  1.512875
sample estimates:
mean in group F mean in group M
 60.58889      63.91000
```

注意上面用的是不需要假定两组方差相等的Welch检验。如果假定两组方差相等可以在函数调用中加参数`var.equal=TRUE`。如果两个独立组放在了两个变量X1和X2中则可以直接用`t.test(X1, X2)`进行比较。

用`wilcox.test(分析变量~ 分组变量)`或`wilcox.test(X1, X2)`可以进行独立两组的Wilcoxon 秩和检验。

对于成对观测的变量X, Y, 要比较只要用`t.test(X, Y, paired=TRUE)`。用`wilcox.test(X, Y, paired=TRUE)`可以进行Wilcoxon符号秩检验。

6.13 S统计模型简介

这一节我们简单介绍S的统计模型。S中实现了几乎所有常见的统计模型,而且多种模型可以用一种统一的观点表示和处理。这方面S-PLUS较全面,它实现了许多最新的统计研究成果, R因为是自愿无偿工作所以统计模型部分还相对较欠缺。事实上,许多统计学家的研究出的统计算法都以S-PLUS程序发表,因为S语言是一种特别有利于统计计算编程的语言。

学习这一节需要我们具备线型模型、线型回归、方差分析的基本知识。

6.13.1 统计模型的表示

很多统计模型可以用一个线型模型来表示:

$$y_i = \sum_{j=0}^p \beta_j x_{ij} + e_i, \quad e_i \sim \text{iid}N(0, \sigma^2), \quad i = 1, 2, \dots, n$$

写成矩阵形式为

$$y = \mathbf{X}\beta + e$$

其中 y 为因变量, \mathbf{X} 为模型矩阵或称设计阵, 各列为 x_0, x_1, \dots, x_p 等各自变量, x_0 常常是一列1, 定义一个模型截距项。

在S中模型是一种对象, 其表达形式叫做一个公式(formula), 我们先举几个例子来看一看。假定 $y, x, x_0, x_1, x_2, \dots$ 是数值型变量, \mathbf{X} 是矩阵, A, B, C, \dots 是因子。

- $y \sim x, y \sim 1+x$ 两个式子都表示 y 对 x 的简单一元线型回归。第一个式子带有隐含的截距项,而第二个式子把截距项显式地写了出来。
- $y \sim -1+x, y \sim x-1$ 都表示 y 对 x 的通过原点的回归,即不带截距项的回归。
- $\log(y) \sim x_1 + x_2$ 表示 $\log(y)$ 对 x_1 和 x_2 的二元回归,带有隐含的截距项。
- $y \sim \text{poly}(x, 2), y \sim 1 + x + I(x^2)$ 表示 y 对 x 的一元二次多项式回归。第一种形式使用正交多项式,第二种形式直接使用 x 的各幂次。
- $y \sim \mathbf{X} + \text{poly}(x, 2)$ 因变量为 y 的多元回归,模型矩阵包括矩阵 \mathbf{X} ,以及 x 的二次多项式的各项。
- $y \sim A$ 一种方式分组的方差分析,指标为 y ,分组因素为 A 。
- $y \sim A + x$ 一种方式分组的协方差分析,指标为 y ,分组因素为 A ,带有协变量 x 。
- $y \sim A*B, y \sim A + B + A:B, y \sim B \%in\% A, y \sim A/B$ 非可加两因素方差分析模型,指标为 y , A, B 是两个因素。前两个公式表示相同的交叉分类设计,后两个公式表示相同的嵌套分类设计。
- $y \sim (A + B + C)^2, y \sim A*B*C - A:B:C$ 表示三因素试验,只考虑两两交互作用而不考虑三个因素间的交互作用。两个公式是等价的。
- $y \sim A * x, y \sim A / x, y \sim A / (1 + x) - 1$ 都表示对因子 A 的每一水平拟合 y 对 x 的线型回归,但三个公式的编码方式不同。最后一种形式对 A 的每一水平都分别估计截距项和斜率项。

- $y \sim A*B + \text{Error}(C)$ 表示有两个处理因素A和B,误差分层由因素C确定的设计。

在S中 \sim 运算符用来定义模型公式。一般的线型模型的公式形式为

因变量 \sim 第一项 $[\pm]$ 第二项 $[\pm]$ 第三项 $[\pm]$...

其中因变量可以是向量或矩阵,或者结果为向量或矩阵的表达式。 $[\pm]$ 是加号+或者减号-,表示在模型中加入一项或去掉一项,第一项前面如果是加号可以省略。

公式中的各项可以取为:

- 一个值为向量或矩阵的表达式,或1。
- 一个因子
- 一个“公式表达式”,由“公式运算符”把因子、向量、矩阵连接而成。

每一项定义了要加入模型矩阵或从模型矩阵中删除的若干列。一个1表示一个截距项列,除非显式地删除总是隐含地包括在模型公式中。

“公式运算符”的定义和Glim、Genstat软件中的定义类似,不过那里的“.”运算符这里改成了“:”,因为在S中句点是名字的合法字符。下表列出了各运算符的简要说明。

- $Y \sim M$ Y作为因变量由M解释。
- $M_1 + M_2$ 加入 M_1 和 M_2 。
- $M_1 - M_2$ 加入 M_1 但去掉 M_2 指定的项。
- $M_1 : M_2$ M_1 和 M_2 的张量积。如果两项都是因子,则为因子的交互作用。
- $M_1 \%in\% M_2$ 与 $M_1 : M_2$ 类似但模型矩阵编码方式不同。
- $M_1 * M_2$ 等于 $M_1 + M_2 + M_1 : M_2$ 。
- M^n M中所有项以及所有到n阶为止的交互作用项。
- $I(M)$ 将M隔离,使得M中的运算符按照原来的算术运算符解释而不是按公式运算符解释。表达式M的结果作为公式的一项。

注意在函数调用的括号内的表达式按普通四则运算解释。函数I()可以把一个计算表达式封装起来作为模型的一项使用。

注意S的模型表示只给出了因变量和自变量及自变量间的关系,这样只确定了线型模型的模型矩阵,而模型参数向量是隐含的,并没有在模型公式中体现出来。这种做法适用于线性模型,但不具有普遍性,例如非线性模型就不能这样表示。

6.13.2 线性回归模型

拟合普通的线性模型的函数为lm(),其简单的用法为:

```
fitted.model <- lm(formula, data=data.frame)
```

其中data.frame为各变量所在的数据框,formula为模型公式,fitted.model是线性模型拟合结果对象(其class属性为lm)。例如:

```
mod1 <- lm(Weight ~ Height + Age, data=c1)
```

可以拟合一个y对x1和x2的二元回归(带有隐含的截距项), 数据来自数据框production。拟合的结果存入了对象mod1中。注意不论数据框production是否以用attach()连接入当前运行环境都可被lm()使用。lm()的基本显示十分简练:

```
> mod1
Call:
lm(formula = Weight ~ Height + Age, data = c1)

Coefficients:
(Intercept)      Height          Age
    -141.224       3.597       1.278
```

只显示了调用的公式和参数估计结果。

6.13.3 提取信息的通用函数

lm()函数的返回值叫做模型拟合结果对象, 本质上是一个具有类属性值lm的列表, 有model、coefficients、residuals等成员。lm()的结果显示十分简单, 为了获得更多的拟合信息, 可以使用对lm类对象有特殊操作的通用函数, 这些函数包括:

add1 coef effects kappa predict residuals
alias deviance family labels print summary
anova drop1 formula plot proj

下面列出了lm类(拟合模型类)常用的通用函数的简单说明。

- anova(对象1,对象2) 把一个子模型与原模型比较,生成方差分析表。
- coefficients(对象) 返回回归系数(矩阵)。可简写为coef(对象)。
- deviance(对象) 返回残差平方和,如有权重则加权。
- formula(对象) 返回模型公式。
- plot(对象) 绘制模型诊断的几种图, 如残差对预测值图。
- predict(对象,newdata=数据框), predict.gam (对象,newdata=数据框) 有了模型拟合结果后对新数据进行预报。指定的新数据必须与建模时用的数据具有相同的变量结构。函数结果为对数据框中每一观测的因变量预报结果(为向量或矩阵)。
predict.gam()与predict()作用相同但适用性更广, 可应用于lm、glm和gam的拟合结果。比如,当多项式基函数用了正交多项式时, 加入了新数据导致正交多项式基函数改变, 用predict.gam() 函数可以避免由此引起的偏差。
- print(对象) 简单显示模型拟合结果。一般不用print()而直接键入对象名来显示。

- `residuals`(对象) 返回模型残差(矩阵),若有权重则适当加权。可简写为`resid`(对象)。
- `summary`(对象) 可显示较详细的模型拟合结果。

6.13.4 方差分析

方差分析是研究取离散值的因素对一个数值型指标的影响的经典工具。S进行方差分析的函数是`aov()`, 格式为`aov(公式, data=数据框)`, 用法与`lm()`类似, 提取信息的各通用函数仍有效。

我们以前面用过的不同牌子木板磨损比较的数据为例。假设`veneer`数据框保存了该数据:

```
> veneer
  Brand Wear
1  ACME  2.3
2  ACME  2.1
3  ACME  2.4
4  ACME  2.5
5 CHAMP  2.2
6 CHAMP  2.3
7 CHAMP  2.4
8 CHAMP  2.6
9  AJAX  2.2
10 AJAX  2.0
11 AJAX  1.9
12 AJAX  2.1
13 TUFFY  2.4
14 TUFFY  2.7
15 TUFFY  2.6
16 TUFFY  2.7
17 XTRA  2.3
18 XTRA  2.5
19 XTRA  2.3
20 XTRA  2.4
```

首先我们把每个牌子的木板的磨损情况画盒形图并且放在同一页面中, 作图如下:

```
plot(Wear ~ Brand, data=veneer)
```

见图6.20。这种图可以直观地比较一个变量在多个组的分布, 或者比较几个类似的变量。从图中可以看出, AJAX牌子较好, TUFFY较差, 其它三个牌子差别不明显。

为了检验牌子这个因素对指标磨损量有无显著影响, 只要用`aov()`函数:

```
> aov.veneer <- aov(Wear ~ Brand, data=veneer)
> summary(aov.veneer)
          Df Sum Sq Mean Sq F value    Pr(>F)
```

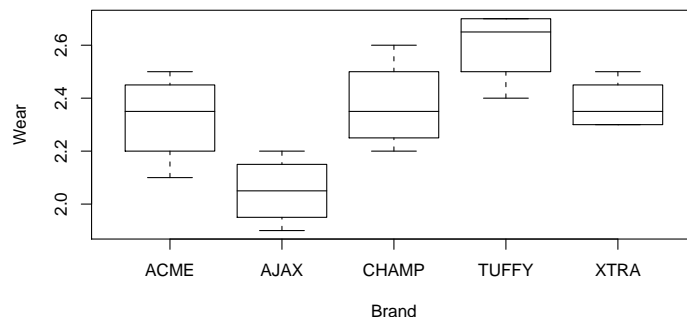


图 6.20: 不同品牌木板磨损量比较

```
Brand      4 0.61700 0.15425  7.404 0.001683 **
Residuals 15 0.31250 0.02083
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

可见因素是显著的。

6.14 统计分析实例

下面我们以6.1中的那个学生班的情况为例进行一些分析。我们希望了解体重、身高、年龄、性别等变量的基本情况及互相之间的关系。

6.14.1 数据输入

假设数据放在了文本文件c:\work\class.txt中, 没有列标题, 各变量上下对齐。我们先把数据读入一个S数据框对象中:

```
> cl <- read.table("c:/work/class.txt",
+ col.names=c("Name", "Sex", "Age", "Height", "Weight"),
+ row.names="Name")
> cl
      Sex Age Height Weight
Alice   F  13   56.5   84.0
Becka   F  13   65.3   98.0
Gail    F  14   64.3   90.0
. . . . .
William M  15   66.5  112.0
```

6.14.2 探索性数据分析(EDA)

首先我们先研究各变量的分布情况, 看分布是否接近正态, 有无明显的异常值, 有没有明显的序列相关, 等等。

先用6.12.1中定义的eda.shape()函数来研究各数值型变量的分布情况。在调用前先把数据框cl连接入当前的搜索路径中以直接使用cl中的变量名：

```
> attach(cl)
> summary(cl)
> eda.shape(Age)
> eda.shape(Height)
> eda.shape(Weight)
> tab.sex <- table(Sex)
> barplot(tab.sex)
```

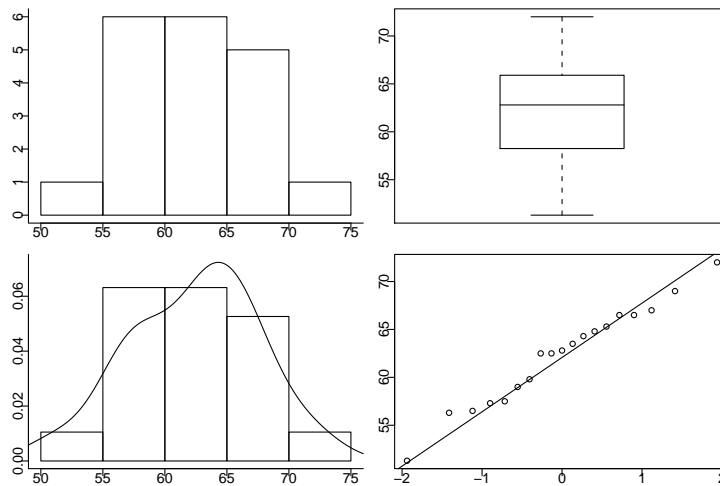


图 6.21: 身高分布探索性数据分析

因为数据是不同个体的观测所以不可能有序列相关, 未画时间序列图。这里给出了身高的分布图(图6.21)。可以看出, 身高和体重都相当接近正态且无明显的异常点, 体重因为取离散值所以直方图不接近正态, 但从核密度估计曲线看仍可作为正态处理。

为了研究数值型变量Weight、Height、Age间的关系, 我们画它们的散点图矩阵:

```
> pairs(cbind(Height, Weight, Age))
```

从散点图矩阵(图6.22)可以看出三个变量之间都可能线性相关关系。

为了研究因子Sex对其它变量的影响, 可以画Sex不同水平上各变量的盒形图, 如:

```
> oldpar <- par(mfcol=c(1,3))
> boxplot(Weight ~ Sex, ylab="Weight")
> boxplot(Height ~ Sex, ylab="Height")
> boxplot(Age ~ Sex, ylab="Age")
> par(oldpar)
```

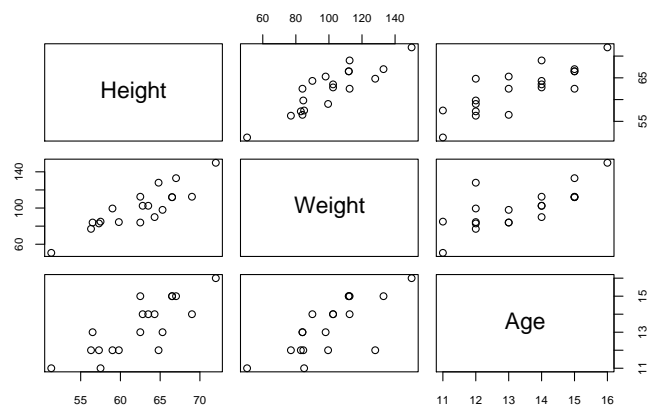


图 6.22: 身高、体重、年龄的散点图矩阵

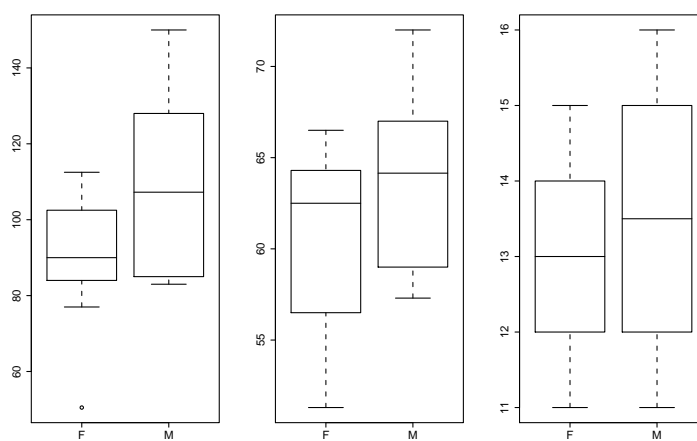


图 6.23: 性别对其他变量的影响

从图6.23可以看出, 男女的体重、身高有明显的差别, 而年龄则差别不明显。我们也可以分不同性别对某一变量分别作图或计算, 这里只要使用象`Weight[Sex=="F"]`, `Weight[Sex=="M"]`这样的取子集的办法就可以把观测分组。更进一步还可以用函数`tapply`直接按一个因子对观测分组然后作用某个函数:

```
> tapply(Weight, Sex, hist)
```

为了研究因子Sex的不同水平对其它变量间的相关关系的影响, 可以作协同图:

```
> coplot(Weight ~ Height | Sex)
```

结果图6.12没有反映明显的差别。

6.14.3 组间比较

我们来分析男女的身高有无显著差异, 这是两组比较的问题。上面EDA部分的并排盒形图已经提示男女身高有明显差异, 这里我们用统计假设检验给出统计结论。

男女两组可以认为是独立的, 而且每组内的观测也可以认为是相互独立的。根据EDA结果可以认为两组都来自正态总体。这样, 我们可以使用两样本t检验。因为方差是否相等为止, 我们干脆用不要求方差相等的近似两样本t检验:

```
> t.test(Height ~ Sex)
```

结果p值为0.1652, 按我们一般采用的0.05水平是不显著的。所以从这组样本看男女的身高没有发现显著差异。

类似可以进行男女体重的比较, p值为0.06799, 也不显著。

6.14.4 回归分析

下面我们研究对体重的预报。从散点图矩阵看, 体重与身高之间有明显的线性相关, 所以我们先拟合一个体重对身高的一元线性回归模型:

```
> lm.fit1 <- lm(Weight ~ Height, data=c1)
> lm.fit1

Call:
lm(formula = Weight ~ Height, data = c1)

Coefficients:
(Intercept)      Height
    -143.027         3.899

> summary(lm.fit1)

Call:
lm(formula = Weight ~ Height, data = c1)
```



```

Residuals:
      Min       1Q   Median       3Q      Max
-17.6807  -6.0642   0.5115   9.2846  18.3698

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept) -143.0269    32.2746  -4.432 0.000366 ***
Height       3.8990     0.5161   7.555 7.89e-07 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 11.23 on 17 degrees of freedom
Multiple R-Squared: 0.7705,    Adjusted R-squared: 0.757
F-statistic: 57.08 on 1 and 17 DF,  p-value: 7.887e-007

```

拟合的模型方程为 $\text{Weight} = -143.0269 + 3.8990 \times \text{Height}$, 复相关系数平方为 0.7705, 检验模型的斜率为 0 的 p 值为 $7.887e-007$, 可见模型是显著的。对于一元回归, 我们可以在因变量对自变量的散点图上叠加回归直线来看回归拟合的效果:

```

> plot(Weight ~ Height)
> abline(lm.fit1)

```

一般地, `lm` 拟合结果对象的 `plot()` 函数可以作出若干张检查拟合效果的图形, `R` 可以作四个图: 残差对拟合值图、残差的正态概率图、标准化残差绝对值平方根对拟合值图、Cook 距离图。

```

> oldpar <- par(mfrow=c(2,2),
+   mar=c(2.5,2,1.5,0.2), mgp=c(1.2, 0.2, 0))
> plot(lm.fit1)
> par(oldpar)

```

见图 6.24。如果想每个图都用全窗口来看则不要设置图形参数。四个图中, 残差对拟合值图可以反映残差中残留的结构, 如果模型充分的话残差应该是随机变化没有任何模式的。残差的正态概率图可以检验线性回归假设检验的重要假定—误差项服从正态分布是否合理, 可以看出残差的分布重尾、轻尾、左偏、右偏等情况。标准化残差绝对值平方根对拟合值图可以发现残差的异常值点, 即拟合最差的点。Cook 距离衡量每一观测对拟合结果的影响大小, 数值大的为强影响点。图中自动标出了最突出的点。

从 `lm.fit1` 的回归诊断图看残差没有明显的模式, 但残差分布有轻尾倾向。没有明显的异常值点。

下面我们看加入其它变量能否进一步改善模型的预报能力。用 `add1` 函数可以判断加入新的变量可以改善模型:

```

> add1(lm.fit1, ~ . + Age + Sex)
Single term additions

Model:

```

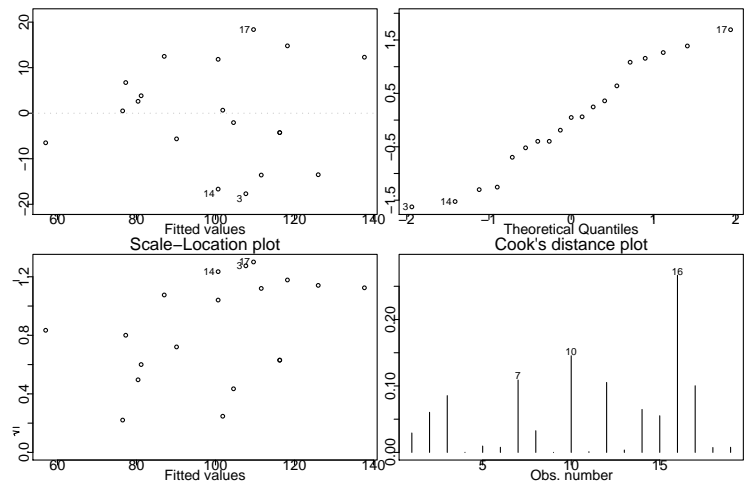


图 6.24: 回归诊断图

Weight ~ Height				
	Df	Sum of Sq	RSS	AIC
<none>			2142.49	93.78
Age	1	22.39	2120.10	95.58
Sex	1	184.71	1957.77	94.07

add1的结果显示一个方差分析表, 各行中<none>一行为不加变量的情况, Age一行为加入一个变量Age后的情况, Sex一行为加入一个变量Sex的情况。各列中DF为此变量的自由度, Sum of Sq为该变量对应的平方和, RSS为加入该变量后的残差平方和, AIC为加入该变量后的AIC统计量值。AIC较小的模型为较好的, 所以如果加入某个变量后的AIC减小就可以加入此变量。这里加入Age和加入Sex都使AIC变大, 所以不应加入这两个变量。

如果一开始就加入了所有变量, 可以用drop1()函数考察去掉一个变量后AIC是否可以变小:

```
> lm.fit2 <- lm(Weight ~ Height + Age + Sex, data=c1)
> summary(lm.fit2)

Call:
lm(formula = Weight ~ Height + Age + Sex, data = c1)

Residuals:
    Min       1Q   Median       3Q      Max
-19.6540  -6.5737   0.4602   7.6708  20.8515

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) -125.1151    33.9038  -3.690  0.00218 **
Height        2.8729     0.9971   2.881  0.01142 *
Age           3.1131     3.2362   0.962  0.35132
SexM          8.7443     5.8350   1.499  0.15472
---

```

```

Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 11.09 on 15 degrees of freedom
Multiple R-Squared:  0.8025,    Adjusted R-squared:  0.763 
F-statistic: 20.31 on 3 and 15 DF,  p-value: 1.536e-005

> drop1(lm.fit2)
Single term deletions

Model:
Weight ~ Height + Age + Sex
      Df Sum of Sq  RSS   AIC
<none>                 1844.01  94.93
Height  1   1020.61 2864.62 101.30
Age      1    113.76 1957.77  94.07
Sex      1    276.09 2120.10  95.58

```

从summary()的结果看Age和Sex就不显著。用drop1()发现去掉Age可以使AIC从94.93变小为94.07, 所以应该去掉Age。对去掉Age后的模型再用drop1()发现Sex也应该去掉。所以我们最后得到的模型是lm.fit1。

在修改模型或数据改变后重新拟合时还可以使用update函数, 比如要从lm.fit2中去掉Age就可以用:

```
> lm.fit3 <- update(lm.fit2, . ~ . - Age)
```

在自变量个数较多时S提供了一个step()函数用来进行逐步回归, 它从一个初始模型开始自动判断增加或去掉变量, 最后得到较好的模型:

```

> lm.fit0 <- lm(Weight ~ 1, data=c1)
> lm.step <- step(lm.fit0, ~ . + Height + Age + Sex)
Start:  AIC= 119.75
Weight ~ 1

      Df Sum of Sq  RSS   AIC
+ Height  1   7193.2 2142.5   93.8
+ Age      1   5124.5 4211.2  106.6
+ Sex      1   1681.1 7654.6  118.0
<none>                 9335.7  119.7

Step:  AIC= 93.78
Weight ~ Height

      Df Sum of Sq  RSS   AIC
<none>                 2142.5   93.8
+ Sex      1    184.7 1957.8   94.1
+ Age      1     22.4 2120.1   95.6
- Height  1   7193.2 9335.7  119.7

```

得到适当的模型以后, 我们可以用模型进行拟合或预报。拟合只要对模型拟合结果用predict()函数, 如:

```
> predict(lm.fit1)
```

对新数据如果想作预报, 只要在predict函数调用时加入新数据的数据框作为参数:

```
> new.data <- data.frame(Height=c(50, 51.2, 68, 69.7))
> predict(lm.fit1, new.data)
      1      2      3      4
51.92459 56.60343 122.10714 128.73549
```

6.15 用S作随机模拟计算

作为统计工作者, 我们除了可以用S迅速实现新的统计方法, 还可以用S进行随机模拟。随机模拟可以验证我们的算法、比较不同算法的优缺点、发现改进统计方法的方向, 是进行统计研究的最有力的计算工具之一。

随机模拟最基本的需要是产生伪随机数, S中已提供了大多数常用分布的伪随机数函数, 可以返回一个伪随机数序列向量。这些伪随机数函数以字母r开头, 比如rnorm()是正态伪随机数函数, runif()是均匀分布伪随机数函数, 其第一个自变量是伪随机数序列长度n。关于这些函数可以参见6.16以及系统帮助文件。下例产生1000个标准正态伪随机数:

```
> y <- rnorm(1000)
```

这些伪随机数函数也可以指定与分布有关的参数, 比如下例产生1000个均值为150、标准差为100的正态伪随机数:

```
> y <- rnorm(1000, mean=150, sd=100)
```

产生伪随机数序列是不重复的, 实际上, S在产生伪随机数时从一个种子出发, 不断迭代更新种子, 所以产生若干随机数后内部的随机数种子就已经改变了。有时我们需要模拟结果是可重复的, 为此我们可以在模拟开始之前调用setseed(正整数) 就可以保证每次产生的随机数是相同的序列。

作为例子, 我们来产生服从一个简单的线性回归的数据。

```
# 简单线性回归的模拟
lm.simu <- function(n){
  # 先生成自变量。假设自变量x的取值范围
  # 在150到180之间,大致服从正态分布。
  x <- rnorm(n, mean=165, sd=7.5)
  # 再生成模型误差。假设服从N(0, 1.2)分布
  eps <- rnorm(n, 0, 1.2)
  # 用模型生成因变量
  y <- 0.8 * x + eps
  return(data.frame(y,x))
}
```

S没有提供多元随机变量的模拟程序, 这里给出一个进行三元正态随机变量模拟的例子。假设要产生三元正态随机向量 $X \sim N(\mu, \Sigma)$ 的 n 个独立观测, 可以先产生 n 个服从三元标准正态分布的观测, 放在一个 n 行3列的矩阵中:

```
U <- matrix(rnorm(3*n), ncol=3, byrow=T)
```

可以认为矩阵U的每一行是一个标准的三元正态分布的观测。设矩阵 Σ 的Choleski分解为 $\Sigma = A'A$, A 为上三角矩阵, 若随机向量 $\xi \sim N(0, I)$, 则 $\mu + A'\xi \sim N(\mu, \Sigma)$ 。因此, $\mu + A'U'$ 作为一个三行 n 列的矩阵每一列都是服从 $N(\mu, \Sigma)$ 分布的, 且各列之间独立。经过转置, 产生的 X

```
X <- matrix(rep(mu,n), ncol=3, byrow=T) + U %*% A
```

是一个 n 行三列的矩阵。

有时模拟需要的计算量很大, 多的时候甚至要计算几天的时间。对于这种问题我们要善于把问题拆分成可以单独计算的小问题, 然后单独计算每个小问题, 把结果保存在S对象中或文本文件中, 最后综合保存的结果得到最终结果。

如果某一个问题的计算时间比较长, 我们在编程时可以采用以下的技巧: 每隔一定时间就显示一下任务的进度, 以免计算已经出错或进入死循环还不知道; 应该把中间结果每隔一段时间就记录到一个文本文件中(`cat()`函数可以带一个`file`参数和`append`参数, 对这种记录方法提供了支持), 如果需要中断程序, 中间结果可能是有用的, 有些情况下还可以根据记录的中间结果从程序中断的地方继续执行。

因为S是一个解释型语言, 所以有些需要大量迭代的任务纯用S编程效率很低。为此我们可以把主要部分用S编程, 把需要迭代的部分用C、FORTRAN或C++编程, 在S中调用由其他语言编译得到的DLL库。这里只简单介绍如何在MS Windows 版本的R中调用C或C++编译得到的DLL库。注意Borland公司免费提供了一个MS Windows下的C++编译程序, 有一个自由软件的DevC++集成编程环境也可以进行C和C++编辑和编译。

假设我们写了一个计算卷积的C程序:

```
void convolve(double *a, int *na,
              double *b, int *nb, double *ab){
    int i, j, nab = *na + *nb - 1;
    for(i = 0; i < nab; i++)
        ab[i] = 0.0;
    for(i = 0; i < *na; i++)
        for(j = 0; j < *nb; j++)
            ab[i + j] += a[i] * b[j];
}
```

注意C函数的所有自变量都是指针, 这是因为S的基本数据类型是一维数组。假设此C程序编译成了一个DLL库testdll.dll放在当前工作目录。注意编译时要生成Windows的控制台DLL而不是窗口环境的DLL。在资源管理器中查看可以发现DLL中输出了函数‘convolve’。为了在R中调用这个C函数, 需要如下步骤:

```
conv <- function(a, b){
    dll.filename <-
```

```

        file.path(paste("testdll",
                        .Platform$dynlib.ext, sep=""))
    if(!is.loaded(symbol.C("convolve"),
                        PACKAGE="testdll")){
        dyn.load(dll.filename)
    }
    .C("convolve",
        as.double(a), as.integer(length(a)),
        as.double(b), as.integer(length(b)),
        ab=double(length(a)+length(b)-1))$ab
}

```

这个S函数把C函数“convolve”进行了包装。其中，函数is.loaded(symbol.C(函数名))检查DLL中的C函数是否已经调入了R中，如果没有调入就要用dyn.load(DLL文件名)调入。因为在Windows环境下和Unix环境下DLL库的扩展名不同所以对文件名做了兼容性处理。为了调用C函数，使用“.C”函数。“C”函数的第一个参数是DLL中的C函数名(注意可能带下划线)，其他参数是C函数的自变量。所有从R传给C的自变量都要用as.double和as.integer包裹起来。“C”函数返回一个列表，列表元素为传入的自变量经过函数运行后的值，如果调用“.C”时给某些自变量加了名字(如上面的ab)则返回的列表中将使用这些名字。

对于反复调用的函数可以设法预先调入动态链接库而不是每次在函数内判断。

6.16 S常用函数参考

这一节分类列出常用的函数，需要时可以参看帮助。

6.16.1 基本

一、数据管理

vector:向量 numeric:数值型向量 logical:逻辑型向量 character:字符型向量
 list:列表 data.frame:数据框 c:连接为向量或列表
 length:求长度 subset:求子集 seq, from:to, sequence:等差序列 rep:重复
 NA:缺失值 NULL:空对象 sort, order, unique, rev:排序 unlist:展平列表
 attr, attributes:对象属性 mode, typeof:对象存储模式与类型
 names:对象的名字属性

二、字符串处理

character:字符型向量 nchar:字符数 substr:取子串 format, formatC:把对象用格式转换为字符串
 paste, strsplit:连接或拆分 charmatch, pmatch:字符串匹配 grep, sub, gsub:模式匹配与替换

三、复数

complex, Re, Im, Mod, Arg, Conj:复数函数

四、因子

factor:因子 codes:因子的编码 levels:因子的各水平的名字
nlevels:因子的水平个数 cut:把数值型对象分区间转换为因子 table:交叉频数表
split:按因子分组 aggregate:计算各数据子集的概括统计量
tapply:对“不规则”数组应用函数

6.16.2 数学

一、计算

+, -, *, /, %, %%, %/%:四则运算 ceiling, floor, round, signif, trunc,
zapsmall:舍入 max, min, pmax, pmin:最大最小值 range:最大值和最小值
sum, prod:向量元素和, 积 cumsum, cumprod, cummax, cummin:累加、累乘
sort:排序 approx和approx fun:插值 diff:差分 sign:符号函数

二、数学函数

abs, sqrt:绝对值, 平方根 log, exp, log10, log2:对数与指数函数 sin,
cos, tan, asin, acos, atan, atan2:三角函数 sinh, cosh, tanh, asinh,
acosh, atanh:双曲函数 beta, lbeta, gamma, lgamma, digamma, trigamma,
tetragamma, pentagamma, choose, lchoose:与贝塔函数、伽玛函数、组合数有关的特殊函数
fft, mvfft, convolve:富利叶变换及卷积 polyroot:多项式求根
poly:正交多项式 spline, splinefun:样条差值 bessell, bessellK,
besselJ, besselY, gammaCody:Bessel函数 deriv:简单表达式的符号微分或算法微分

三、数组

array:建立数组matrix:生成矩阵 data.matrix:把数据框转换为数值型矩阵
lower.tri:矩阵的下三角部分 mat.or.vec:生成矩阵或向量 t:矩阵转置
cbind:把列合并为矩阵 rbind:把行合并为矩阵 diag:矩阵对角元素
向量或生成对角矩阵 aperm:数组转置 nrow, ncol:计算数组的行数和列数
dim:对象的维向量 dimnames:对象的维名 row/colnames:行名或列名
%*%:矩阵乘法 crossprod:矩阵交叉乘积(内积) outer:数组外积
kronecker:数组的Kronecker积 apply:对数组的某些维应用函数
tapply:对“不规则”数组应用函数 sweep:计算数组的概括统计量
aggregate:计算数据子集的概括统计量 scale:矩阵标准化 matplot:对矩阵各列绘图
cor:相关阵或协差阵 Contrast:对照矩阵 row:矩阵的行下标集
col:求列下标集

四、线性代数

solve:解线性方程组或求逆 eigen:矩阵的特征值分解 svd:矩阵的奇异值分解
backsolve:解上三角或下三角方程组 chol:Choleski分解 qr:矩阵的QR分解
chol2inv:由Choleski分解求逆

五、逻辑运算

<, >, <=, >=, ==, !=:比较运算符 !, &, &&, |, ||, xor():逻辑运算符
logical:生成逻辑向量 all, any:逻辑向量都为真或存在真 ifelse():二者择

一 match, %in%:查找 unique:找出互不相同的元素 which:找到真值
下标集合 duplicated:找到重复元素

六、优化及求根

optimize, uniroot, polyroot:一维优化与求根

6.16.3 程序设计

一、控制结构

if, else, ifelse, switch:分支 for, while, repeat, break, next:循环 apply,
lapply, sapply, tapply, sweep:替代循环的函数。

二、函数

function:函数定义 source:调用文件 call:函数调用 .C, .Fortran:调用C或者Fortran子程序的动态链接库。 Recall:递归调用
browser, debug, trace, traceback:程序调试 options:指定系统参数 miss-
ing:判断虚参是否有对应实参 nargs:参数个数 stop:终止函数执行
on.exit:指定退出时执行 eval, expression:表达式计算 system.time:表达式
计算计时 invisible:使变量不显示 menu:选择菜单(字符列表菜单)

其它与函数有关的还有:delay, delete.response, deparse, do.call, dput, en-
vironment, , formals, format.info, interactive, is.finite, is.function, is.language,
is.recursive, match.arg, match.call, match.fun, model.extract, name, parse, sub-
stitute, sys.parent, warning, machine。

三、输入输出

cat, print:显示对象 sink:输出转向到指定文件 dump, save, dput,
write, write.table:输出对象 scan, read.table, load, dget:读入

四、工作环境

ls, objects:显示对象列表 rm, remove:删除对象 q, quit:退出系统
.First, .Last:初始运行函数与退出运行函数。 options:系统选项 ?, help,
help.start, apropos:帮助功能 data:列出数据集

6.16.4 统计计算

一、统计分布

每一种分布有四个函数:d—density(密度函数), p—分布函数, q—分位数函
数, r—随机数函数。比如, 正态分布的这四个函数为dnorm, pnorm, qnorm,
rnorm。下面我们列出各分布后缀, 前面加前缀d、p、q或r就构成函数名:

norm:正态, t:t分布, f:F分布, chisq:卡方(包括非中心),
unif:均匀, exp:指数, weibull:威布尔, gamma:伽玛, beta:贝
塔, lnorm:对数正态, logis:逻辑分布, cauchy:柯西, binom:二项
分布, geom:几何分布, hyper:超几何, nbinom:负二项, pois:泊松,
signrank:符号秩, wilcox:秩和, tukey:学生化极差

二、简单统计量

sum, mean, var, sd, min, max, range, median, IQR(四分位间距)等为统计量, sort, order, rank与排序有关, 其它还有ave, fivenum, mad, quantile, stem等。

三、统计检验

R中在ctest包中已实现的有shapiro.test, t.test, wilcox.test, chisq.test, prop.test, 等等。

四、多元分析

cor, cov.wt, var:协方差阵及相关阵计算 biplot, biplot.princomp:多元数据biplot图
cancor:典则相关 princomp:主成分分析 hclust:谱系聚类
kmeans:k-均值聚类 cmdscale:经典多维标度
其它有dist, mahalanobis, cov.rob。

五、时间序列

ts:时间序列对象 diff:计算差分 time:时间序列的采样时间 window:时间窗
在ts包中实现了acf等函数。

六、统计模型

lm, glm, aov:线性模型、广义线性模型、方差分析

练习

- 写出元素为3, -1.5, 3E-10的向量。
 - 写出从3开始每次增加3, 长度为100的向量。
 - 写出(0, 2)重复10次的向量。
 - 对向量x, 写出其元素大于等于0小于1的条件。
 - 对向量x, 写出其元素都等于0的条件。
 - 写出包含12个月份名称的向量。
 - 生成一个包含文件名tab1.txt到tab18.txt的字符串向量。
 - 写出包含方程 $z^6 = 1$ 的根的向量, 并写出其幅角的余弦和正弦值。
- 设x为一个长100的整数向量。比如, $x \leftarrow \text{floor}(10 * \text{runif}(100))$ 。
 - 显示x第21到30号元素。
 - 把x第31, 35, 39号元素赋值为0。
 - 显示x中除了第1号和第50号的元素之外的子集。
 - 列出x中个位数等于3的元素。
 - 列出x中个位数等于3的元素的下标位置。
 - 给x的每一个元素加上名字, 为x1到x100。
 - 求x的平均值并求每一个元素减去平均值后的离差, 计算x的离差平方和及元素的平方和。

- (h) 把x从大到小排序。计算x的10%分位数到90%分位数之间的距离。
3. 定义一个维数为(3, 4, 2)的数组其第一层(第三下标为1)取从1开始的奇数, 第二层取从2开始的偶数。显示每一层的第2行元素。把第(1,1,1),(2,2,2), (2, 2, 1)号元素赋值为零。把第一层加上100,把第二层加上200。分别计算第一层和第二层的平均值。
 4. 对线性模型 $Y = X\beta + \varepsilon$,写出当X满秩时计算 β 的S表达式。写出估计 ε 的方差的S表达式。
 5. 把SASUSER.GPA数据中的SEX, SATM, SATV分别输入到S中。计算不同性别的人数,并计算每一组的平均SATM分。把这些变量组合成一个列表。把SASUSER.GPA数据输入为S的数据框。
 6. 把语句 $x \leftarrow \text{floor}(\text{runif}(100))$ 所生成的向量保存到一个文本文件中,数据项用空格和换行分隔。从此文件中读入数据到向量y中。
 7. 设x是一个长度为n的向量,写一段程序,计算x的长度为s的滑动和:

$$S_x(t) = \sum_{i=0}^{s-1} x_{t-i}, \quad t = s, s+1, \dots, n$$

8. 写一个AR(1)的模拟函数:

$$x_t = a + bx_{t-1} + \varepsilon_t, t = 1, 2, \dots, n, \text{Var}(\varepsilon_t) = \sigma^2$$

函数的参数为n、a、b、 x_0 和 σ , 缺省时 $n=100$, $a=0$, $b=1$, $x_0=0$, $\sigma=1$ 。