

153 分钟学会 R

刘思喆

中国第二届 R 语言会议（上海）

2009 年 12 月 13 日

文档结构

1. 前言
2. 基础知识
3. 输入输出
4. 数据处理
5. 数学运算
6. 字符操作
7. 日期时间
8. 绘图相关
9. 统计模型
10. 其他

Getting Started:

R is a free software environment for statistical computing and graphics. It compiles and runs on a wide variety of UNIX platforms, Windows and MacOS. To download R, please choose your preferred CRAN mirror.

获取资料

- <http://cran.r-project.org/manuals.html>
- <http://cran.r-project.org/other-docs.html>
- vignette
- 统计之都 *bbs*

推荐阅读

1. R for beginner
2. An Introduction to R
3. 153
4. Modern Applied Statistics with S

安装 R?

在 R 的官方网址上，选择网站镜像 <http://cran.r-project.org/mirrors.html>，比如 UC Berkeley 下载软件副本。R 拥有在 Linux, MacOS X, Windows 平台下的各个版本，如果是 Windows 用户，进入镜像网站，选择 Windows (95 and later)，进入 base，下载 R-x.x.x-win32.exe。

Task Views?

CRAN Task Views

Bayesian	Bayesian Inference
ChemPhys	Chemometrics and Computational Physics
ClinicalTrials	Design, Monitoring, and Analysis of Clinical Trials
Cluster	Cluster Analysis & Finite Mixture Models
Distributions	Probability Distributions
Econometrics	Computational Econometrics
Environmetrics	Analysis of Ecological and Environmental Data
ExperimentalDesign	Design of Experiments (DoE) & Analysis of Experimental Data
Finance	Empirical Finance
Genetics	Statistical Genetics
Graphics	Graphic Displays & Dynamic Graphics & Graphic Devices & Visualization
gR	gRaphical Models in R
HighPerformanceComputing	High-Performance and Parallel Computing with R
MachineLearning	Machine Learning & Statistical Learning
MedicalImaging	Medical Image Analysis
Multivariate	Multivariate Statistics
NaturalLanguageProcessing	Natural Language Processing
Optimization	Optimization and Mathematical Programming
Pharmacokinetics	Analysis of Pharmacokinetic Data
Psychometrics	Psychometric Models and Methods
Robust	Robust Statistical Methods
SocialSciences	Statistics for the Social Sciences
Spatial	Analysis of Spatial Data
Survival	Survival Analysis
TimeSeries	Time Series Analysis

R 使用 cpu 不能超过 50%

这是 Windows 下任务管理器的误导，它将多个 CPU 看作是单个 CPU，同时计算使用比例。而 R 是单线程计算软件，它不能同时使用 2 个以上的 CPU。当你的计算机应用的是双核技术，你会发现 CPU 应用会定格在 50% 上。

获得在线帮助

- ?
- help.search()
- RSiteSearch()

在 2.5.0 版本以后，R 引入了命令自动补全功能，使用 Tab 键能自动补全 R 命令；或使用第二次 Tab 后，返回所有可能的补全命令列表。

R 需要编程么？

大多数时候不需要，因为 R 有很多函数和包，而且每天都在增加，你用的一般方法和函数都可以在 R 自带包中找到。

比如线性回归：

```
1 m <- lm( y ~ x , data = dd )  
summary(m)
```

R 的内存使用

R 的工作内存大小的设定值为 32Mb 到 3Gb 间的任意数值。但需要提示的是：32 位的 Windows 平台可用最大有效内存为 2Gb，也就是说，实际上 R 的工作内存区间为 32Mb 至 2Gb。

Windows 下升级 R ， 但不想重装 packages

```
update.packages()
```

R 初始加载的包

包	描述
stats	常用统计函数
graphics	基础绘图函数
grDevices	基础或 grid 图形设备
utils	R 工具函数
datasets	基础数据集
methods	用于 R 对象和编程工具的方法和类的定义
base	基础函数

获得 R 命令的源码

- 1.
2. `methods(foo)`
3. `*.tar.gz`

R 的数据类型？

常用数据类型

	类型	说明
1	字符 (character)	它们常常被引号包围
2	数字 (numeric)	实数向量
3	整数 (integer)	整数向量
4	逻辑 (logical)	逻辑向量 (TRUE=T、FALSE=F)
5	复数 (complex)	复数
6	列表 (list)	S 对象的向量
7	因子 (factor)	常用于标记样本

Everything in S is an object;

Every object in S has a class.

读取其他软件录入的文件

foreign 包，它可以读取 Minitab, S, SAS, SPSS, Stata, Systat, dBase 保存的数据

R 读取 Excel

- 拷贝至 clipboard, `read.table('clipboard')`;
- 另存为 csv 文件, `read.csv()` 读入;
- 加载 RODBC 包, 使用 `odbcConnectExcel()` 函数;
- `xlsReadWrite` 包中的 `read.xls` 函数。

R 输出 $\text{T}_{\text{E}}\text{X}$ 文本

1. Hmisc 包中的 `latex()`
2. xtable 包中的 `xtable()`
3. quantreg 包中的 `latex.table()`

R 处理缺失值

- `is.na()`
- `NA, TRUE, FALSE`

两个数据框是否相等

如果每个元素都相同，那么这两个数据框也相同

```
a1 <- data.frame(num = 1:8, lib = letters[1:8])
2 a2 <- a1
  a2[[3,1]] <- 2 -> a2[[8,2]]
4 any(a1 != a2) #      all(a1 == a2)
  identical(a1, a2)
6 which(a1 != a2, arr.ind = TRUE)
```

去除相同的行

```
x <- c(9:20, 1:5, 3:7, 0:8)
2 (xu <- x[!duplicated(x)])
unique(x) # is more efficient
```

如何对不规则数组进行统计分析？

```
1 attach(warpbreaks)
  tapply(breaks, list(wool, tension), mean)
3 aggregate(breaks, list(wool, tension), mean)
## from the help
5 aggregate(state.x77,
             list(Region = state.region,
                 Cold = state.x77[, "Frost"] > 130),
             mean)
```

随机抽取

<code>sample(n)</code>	随机组合 $1, \dots, n$
<code>sample(x)</code>	随机组合向量 $x, \text{length}(x) > 1$
<code>sample(x, replace = T)</code>	bootstrap
<code>sample(x,n)</code>	非放回的从 x 中抽取 n 项
<code>sample(x,n, replace = T)</code>	放回的从 x 中抽取 n 项
<code>sample(x,n, replace = T ,prob = p)</code>	以概率 p , 放回的从 x 中抽取 n 项

如何进行复数计算？

```
x <- 1 + 1i # x <- complex(1,1)  
2 Mod(x) ; Conj(x)
```


求矩阵的特征值和特征向量的函数是什么？

已知 $A = \begin{bmatrix} -1 & 2 & 2 \\ 2 & -1 & -2 \\ 2 & -2 & -1 \end{bmatrix}$ 试求 $B = \left(\frac{1}{2}A^{-1}\right) + E$ 的特征值。

```
A <- matrix(c(-1,2,2,2,-1,-2,2,-2,-1),3,3)
2 m <- solve(0.5*A) + diag(c(1,1,1))
eigen(m)
```

这里还使用了函数 `solve()`，这个函数用于运算

```
1 a%*%x = b
```

而得到 `x`，当然也可以用来求矩阵的逆。

求立方根如何运算？

$x^{1/3}$ 。在 R 里面 `sqrt()` 函数可以计算开平方，故新手容易推测开立方也有函数。事实上 R 里面使用 `^` 来作幂函数运算。`^` 不但是运算符号，还可以看作是函数：

```
1 "^(x , 1/3)
```

在 R 中的运算符号包括：

R 中的运算符号

数学运算	<code>+, -, *, /, ^, %%, %/%</code>	加、减、乘、除、乘方、余数、整除
逻辑运算	<code>>, <, >=, <=, ==, !=</code>	大于, 小于, 大于等于, 小于等于, 等于, 不等于

如何求矩阵各行 (列) 的均值?

- `apply()`
- `rowMeans()`, `colMeans()`

如何计算组合数或得到所有组合？

`choose()` 用于计算组合数 $\binom{n}{k}$ ，函数 `combn()` 可以得到所有元素的组合。使用 `factorial()` 计算阶乘。

如何模拟高斯（正态）分布数据？

使用 `rnorm(n, mean, sd)` 来产生 n 个来自于均值为 `mean`，标准差为 `sd` 的高斯（正态）分布的数据。在 R 里面通过分布前增加字母 '**d**' 表示概率密度函数，'**p**' 表示累积分布函数，'**q**' 表示分位数函数，'**r**' 表示产生该分布的随机数。这些分布具体可以参考第 30 页中“R 的分布函数”，或 R-intro 中的 **Probability distributions** 章节

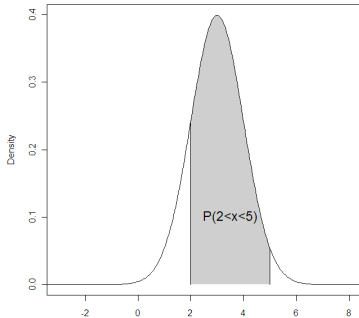
如何模拟高斯（正态）分布数据？ (2)

使用这些函数可以很轻松的进行相关的分布的概率计算，如已知 $X \sim N(3, 1)$ ，计算

$$P(2 \leq X \leq 5)$$

```
1 | pnorm(5, 3, 1) - pnorm(2, 3, 1)
```

计算结果为 0.8185946。



R 的分布函数

分布	R 函数	附加参数	默认参数
beta	beta	shape1(α),shape2(β)	
二项	binom	size(n),prob(p)	
χ^2	chisq	df	
均匀	unif	min(a),max(b)	<i>min = 0, max = 1</i>
指数	exp	rate	<i>rate = 1</i>
F	f	df1(r_1),df2(r_2)	
伽玛	gamma	shape(α),scale(θ)	<i>scale = 1</i>
超几何	hyper	$m = N_1, n = N_2, k = n$	
正态	norm	mean(μ),sd(σ)	<i>mean = 0, sd = 1</i>
泊松	pois	lamda(λ)	
t	t	df	
威布尔	weibull	shape(α),scale(θ)	<i>scale = 1</i>

对大小写敏感么？

R 中有很多基于 Unix 的包，故 R 对大小写是敏感的。可以使用 `tolower()`、`toupper()`、`casefold()` 这类的函数对字符进行转化。

```
1 x <- "MiXeD cAsE 123"  
  chartr("iXs", "why", x)  
3 chartr("a-cX", "D-Fw", x)  
  tolower(x)  
5 toupper(x)
```


日期可以做算术运算么？

一般我们需要使用 `as.Date()` , `as.POSIXct()` 函数将读取的日期（字符串）转化为“Date”类型数据，“Date”类型数据可以进行算术运算。

```
1 d1 <- c("06/29/07") ;    d2 <- c("07/02/07")
  D1 <- as.Date(d1, "%m/%d/%y")
3 D2 <- as.Date(d2, "%m/%d/%y")
  D1 + 2 ;    D1 - D2
5 difftime(D1, D2, units = "days")
```

如何将日期表示为“星期日, 22 七月 2007”?

使用 `format()` 函数。

```
1 format((Sys.Date(), format="%A, %d %B %Y"))
```

具体 `format` 参数可以参考 `help(strptime)` 的 `details` 部分。

如何在同一画面画出多张图？

- 绘图参数：`par(mfrow = c(2,2))` 或 `par(mfcol = c(2,2))`;
- 更为强大功能的 `layout` 函数，它可以设置图形绘制顺序和图形大小；
- `split.screen()` 函数。

```
1 layout(matrix(c(1, 1, 1,
```

```
2 3, 4,
```

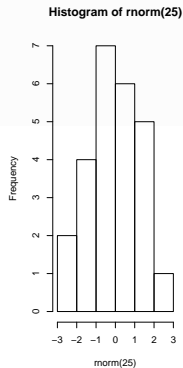
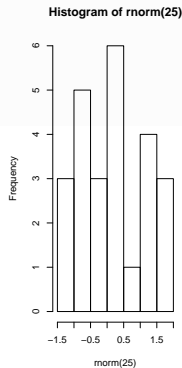
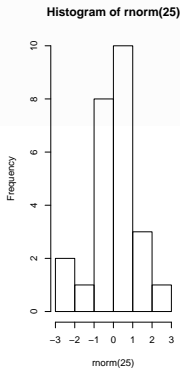
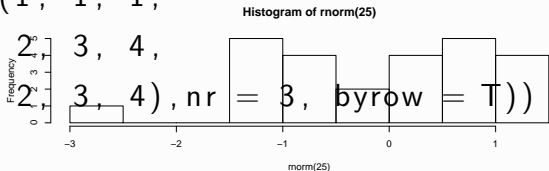
```
3 3, 4), nr = 3, byrow = T))
```

```
4 hist(rnorm(25))
```

```
5 hist(rnorm(25))
```

```
6 hist(rnorm(25))
```

```
7 hist(rnorm(25))
```



如何在已有图形上加一条水平线

使用低水平绘图命令 `abline()`，它可以作出水平线（ y 值 $h=$ ）、垂线（ x 值 $v=$ ）和斜线（截距 $a=$ ，斜率 $b=$ ）。

R 中的绘图命令可以分为“高水平”（`High_level`）、“低水平”（`Low_level`）和“交互式”（`Interactive`）三种绘图命令。

简要地说，“高水平”绘图命令可以在图形设备上绘制新图；“低水平”绘图命令将在已经存在图形上添加更多的绘图信息，如点、线、多边形等；使用“交互式”绘图命令创建的绘图，可以使用如鼠标这类的定点装置来添加或提取绘图信息。在已有图形上添加信息当然要使用“低水平”绘图命令。

常用的绘图设备都有哪些？

	名称	描述
屏幕 显示	x11	X 窗口
	windows	Windows 窗口
文件 设备	postscript	ps 格式文件
	pdf	pdf 格式文件
	pictex	供 L ^A T _E X 使用的文件
	png	png 格式文件
	jpeg	jpeg 格式文件
	bmp	bmp 格式文件
	xfig	供 XFIG 使用的图形格式
win.metafile	emf 格式的文件	


为什么 R 不能显示 8 种以上的颜色？

当绘图参数 `col` 使用数字来代替颜色名时会有这种情形，这是因为 R 内置调色板默认为 8 种颜色：

```
1 palette()  
   barplot(rnorm(15, 10, 3), col = 1:15)  
3 palette(rainbow(15))  
   barplot(rnorm(15, 10, 3), col = 1:15)  
5 palette("default")
```

如何用不同的颜色来代表数据？

高级绘图函数一般都有 `col` 参数可以设置。对于像 `barplot()` 这类图形，可以使用“颜色组”(color sets) 来设置颜色，颜色组包括如下几类：

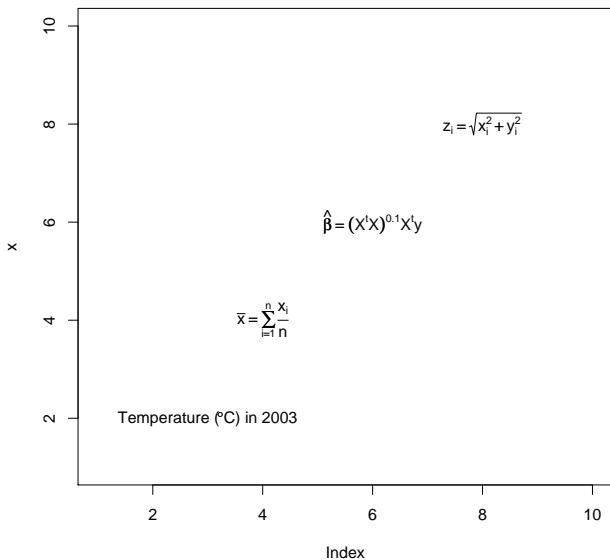
名称	描述
<code>rainbow()</code>	彩虹色 ()
<code>heat.colors()</code>	红色至黄色 ()
<code>terrain.colors()</code>	绿色、棕色至白色 ()
<code>topo.colors()</code>	深蓝色至浅棕色 ()
<code>cm.colors()</code>	浅蓝到白色，浅紫色 ()
<code>gay()</code> 、 <code>grey()</code>	灰色 ()

网格 (lattice) 绘图和普通绘图有什么区别？

网格 (lattice) 绘图实际上是 S-plus 中 Trellis 绘图在 R 中的实现，是多元数据可视化的方法。网格绘图相对于普通绘图来说，是一种拥有“固定格式”的绘图方式，当然它相对来说较难修改。适合对分属不同类数据绘图：

函数	说明
<code>xyplot(y~x)</code>	双变量散点图
<code>dotplot(y~x)</code>	Cleveland 点图 (逐行逐列累加图)
<code>barchart(y~x)</code>	y 对 x 的条形图
<code>stripplot(y~x)</code>	一维图，x 必须是数值型，y 可以是因子
<code>bwplot(y~x)</code>	箱线图
<code>histogram(~x)</code>	直方图

如何在 R 的绘图中加入数学公式或希腊字符？



在 word 里如何使用 R 生成的高质量绘图？

矢量绘图的效果是最好的，比如 eps、pdf，而不是位图（png、jpg、tiff 等）。在 word 里面，可以使用 eps，虽然在屏幕上显示不是很好，但打印效果却不错。

如何使用逐步回归？

在 R 里，可以使用计算逐步回归的 `step()` 函数。它以计算 AIC 信息统计量为准则，选取最小的 AIC 信息统计量来达到逐步回归的目的。

如何做聚类分析？

- `kmeans()`
- `hlust()`
- `cluster` 包

如何做主成分分析？

1. `princomp()`
2. `loadings()`
3. `screeplot()`

如何对样本数据进行正态检验？

比较常见的方法：`shapiro.test()`，`ks.test()`(Kolmogorov-Smirnov 检验)，`jarque.bera.test()` (需要 `tseries` 包)。或者参考专门用作正态检验的 `normtest` 包，`fBasics` 包中的相关函数。这几个包（包括基础包）大概提供了十几种检验函数。

假设检验？

bartlett.test	方差齐次性检验	binom.test	二项检验
chisq.test	χ^2 检验	cor.test	相关性检验
fisher.test	Fisher 精确检验	friedman.test	Friedman 秩和检验
kruskal.test	Kruskal-Wallis 秩和检验	mcnemar.test	McNemar 检验
pairwise.t.test	均值的多重比较	PP.test	Phillips-Perron 检验
var.test	方差比检验	wilcox.test	Wilcoxon 秩和检验

logistic 回归相关函数是？

logistic 回归是关于响应变量为 0-1 定性变量的广义线性回归问题，这里需要使用广义线性模型 `glm()` 函数，且广义线性模型的分布族为二项分布。

广义线性模型中的常用分布族

分布	函数	模型
高斯 (Gaussian) ¹	$E(y) = x^T \beta$	普通线性模型
二项 (Binomial)	$E(y) = \frac{\exp(x^T \beta)}{1 + \exp(x^T \beta)}$	Logistic 模型和概率单位 (probit) 模型
泊松 (Poisson)	$E(y) = \exp(x^T \beta)$	对数线性模型

如何求 Spearman 等级（或 kendall）相关系数

`cor()` 函数默认为求出 Pearson 相关系数，修改其 `method` 参数即可求得 Kendall τ 和 Spearman 秩相关系数。

```
1 cor(longley, method = "spearman")
```

名称	方法	用途（条件）
Pearson	线性	正态总体假定
Kendall τ	协同	非参数检验
Spearman	样本秩	非参数检验

如何使用时间序列相关模型？

- R 中使用 $arima(x, order = c(0, 0, 0), seasonal = list(order = c(0, 0, 0)))$ 对模型进行拟合；
- tseries 包: `adf.test`, `pp.test`, `runs.test`, `garch`...

如何做判别分析？

参考 MASS 包中的 `lda()` 函数（Fisher *Linear Discriminant Analysis*）和 `qda()` 函数。

R 支持的数据挖掘算法？

	主题	核心函数	对应扩展包
统计 模型	主成分分析	princomp	base
	因子分析	factanal	base
	回归模型	lm, nlm, rlm	base, stats, MASS
	Logistic	glm, polr, lrm	base, MASS, Design
	cox 比例模型	coxph, cph	survival, Design
	方差分析	aov, TukeyHSD	stats
时间序列	ar, arima, garch	stats, tseries	
数据 挖掘 算法	关联规则	apriori, Tertius	arules, RWeka
	K-methods	kmeans, clara	base, cluster
	朴素贝叶斯	LBR	Rweka
	判别分析	lda, qda, fda	MASS, mda
	决策树	rpart, J48, tree, ctree	rpart, RWeka, tree, party
	随机森林	randomForest	randomForest
	支持向量机	svm, SMO	e1071, RWeka
	层次聚类	hclust, agen	base, cluster
	k 近邻	knn, lBk, kknn	class, RWeka, kknn
神经网络	nnet	nnet	

R 有类似于 SPSS 的界面么？

安装包 Rcmdr ，加载包后，使用命令

```
1  Commander()
```

调出可供使用的图形使用界面。由于这个图形使用界面需要若干基础包外的其他函数，故还需要包 car 、 effects 、 abind 、 lmtest 、 multcomp 、 relimp 、 RODBC 、 rgl 的支持。

如何释放 R 运行后占用的内存？

因为 R 是在内存中运算，所以当 R 读入了体积比较大的数据后，即使删除了相关对象，内存空间仍不能释放。`gc()` 函数虽然主要用来报告内存使用情况，但是一个重要的用途便是释放内存。

用什么文本编辑器比较好？

比较常用的是 `Tinn-R`，`RWinEdt`²，`ESS`(Emacs Speaks Statistics)，甚至任意一款编辑器，如 `UltraEdit`³，这些都支持 R 语法的高亮显示。如果是 Windows 桌面环境下的用户，对这些不是很了解，记事本也不失为一种选择。

²下载、安装 WinEdt 后，在 R 中安装 RWinEdt 包即可使用

³需要下载、修改 wordfile

Thanks