

*Mathematical Modeling  
with R*

# 数学建模

## 基于R

- 以R语言为载体介绍数学建模的常用方法，探索R扩展程序包强大的计算与求解能力。
- 非传统的数学建模教材，穿插丰富的案例，重点介绍实际应用广泛的统计模型和优化模型。

薛毅 编著



机械工业出版社  
China Machine Press

# 数学建模

基于R

*Mathematical Modeling*

薛毅 编著



机械工业出版社  
China Machine Press

## 图书在版编目 (CIP) 数据

---

数学建模：基于 R / 薛毅编著. —北京：机械工业出版社，2017.5  
(华章应用统计系列)

ISBN 978-7-111-57068-4

I. 数… II. 薛… III. 数学模型—高等学校—教材 IV. O141.4

中国版本图书馆 CIP 数据核字 (2017) 第 129150 号

---

本书以 R 语言为载体，介绍数学建模常用的统计方法，并着重介绍了如何从 CRAN 社区下载相关的 R 扩展程序包，以及如何使用这些程序包中的函数求解线性规划、最优化、图论与网络、数值分析方面的模型。

本书可作为“数学建模”课程的教材或数学建模竞赛的辅导教材，也可作为理工、经管、生物等专业的本科生、研究生或相关专业技术人员学习 R 软件的参考书。

出版发行：机械工业出版社（北京市西城区百万庄大街 22 号 邮政编码：100037）

责任编辑：和 静

责任校对：殷 虹

印 刷：北京市荣盛彩色印刷有限公司

版 次：2017 年 7 月第 1 版第 1 次印刷

开 本：186mm×240mm 1/16

印 张：21.5

书 号：ISBN 978-7-111-57068-4

定 价：69.00 元

---

凡购本书，如有缺页、倒页、脱页，由本社发行部调换

客服热线：(010) 88378991 88361066

投稿热线：(010) 88379604

购书热线：(010) 68326294 88379649 68995259

读者信箱：hzjsj@hzbook.com

版权所有·侵权必究

封底无防伪标均为盗版

本书法律顾问：北京大成律师事务所 韩光 / 邹晓东

# 前 言

R 是一款免费软件，主要用于统计分析、绘图和数据挖掘等。但随着 R 的广泛使用，R 软件的求解能力已不仅仅局限于统计计算的内容，特别是 R 扩展程序包的下载和安装，大大地增强了 R 软件的计算与求解能力，例如，能够完成优化、图论与网络、数值分析等方面的计算。

本书之所以命名为《数学建模：基于 R》，是因为除介绍数学建模常用的统计方法外，还着重介绍了如何从 CRAN(Comprehensive R Archive Network)社区下载相关的扩展程序包，如何使用这些程序包中的函数求解线性规划、最优化、图论与网络、数值分析方面的模型。

采用该命名的第二个原因，是在内容的编排和选取方面与传统的数学建模教材不同，基本上不再讲授传统数学建模课程的基本内容，而是将侧重点放在实际应用中使用较为广泛的两类模型——统计模型和优化模型，以及如何使用 R 软件求解这两类模型上。

本书共有 6 章。第 1 章“概率统计模型”和第 2 章“多元分析模型”属于统计模型的范畴，只需使用 R 基本库中的函数就可完成相应的求解工作。第 3 章“线性规划模型”、第 4 章“最优化模型”和第 5 章“图论与网络模型”属于运筹学的内容，使用 R 基本库中的函数无法完成此类模型的求解，需要在 CRAN 社区下载相关的扩展程序包，使用程序包中的函数完成运筹学模型的求解工作。第 6 章“数值分析”介绍数值代数和微分方程数值解等内容，这部分内容实际上是数值计算(包括统计计算)的基础，也可以看成前面内容的补充。

作为数学建模教材，本书的每一章都有一至两个数学建模案例分析，其目的有两个：一是让读者了解数学建模的整个过程；二是复习该章所讲授的知识及相关 R 函数，学会使用 R 软件求解问题。

本书介绍的模型完全可由其他软件完成求解工作，如 SPSS、LINGO 或 MATLAB 等，但这些都是商业软件，而且有的还价格昂贵。而 R 是一款免费的开源软件，从这一点来说，对读者更有意义：你不但能够享受到他人的工作成果，也能将你的成果放到网上，与他人分享。这正是 R 的魅力，也是 R 这些年来发展如此迅速的原因。

从严格意义上讲，本书不能算作数学建模的教材，也不是 R 软件使用手册，而是希望通过 R 对数学模型的求解，让读者了解并学会使用 R 求解统计或非统计模型，以及如何下载程序包来扩展 R 的计算能力。当然，在学习了这些内容之后，你可以下载其他的程序包<sup>⊖</sup>，帮助你完成工作或科研所需的计算工作。

本书所介绍的 R 函数均以 R-3.1.1 版本为基准，所有函数(包括下载程序包中的函数)均通过测试，读者如果需要书中例题的相关程序，以及例题和习题中的数据文件，可以发

⊖ 截至 2015 年 8 月 1 日，CRAN 网站共有 6 957 个 R 包，涵盖了不同领域的应用。

送电子邮件向作者索取，邮件地址：xueyi@bjut.edu.cn.

本书可作为“数学建模”课程的教材或教学参考书，也可作为数学建模竞赛的辅导教材，还可作为理工、经济、管理、生物等专业的本科生、研究生或者相关专业的技术人员学习 R 软件的参考书.

受编者水平所限，书中难免存在不足甚至错误之处，欢迎读者不吝指正.

在本书出版之际，谨向对本书提供帮助的各位老师和专家表示感谢，对北京工业大学研究生院对于数学建模课程的支持表示感谢，同时对机械工业出版社为本书出版所做的大量工作表示感谢.

# 目 录

前言	
<b>第 1 章 概率统计模型</b>	1
1.1 数据的描述性分析	1
1.1.1 数据的数字特征	1
1.1.2 随机变量的分布	5
1.1.3 常用的分布	6
1.1.4 数据的图形描述	9
1.2 参数的区间估计与假设检验	13
1.2.1 单个总体的区间估计与假设 检验	13
1.2.2 两个总体的区间估计与假设 检验	14
1.2.3 区间估计与假设检验的计算	16
1.2.4 两个正态总体方差比 $\sigma_1^2/\sigma_2^2$ 的 估计与检验	22
1.3 非参数检验	24
1.3.1 二项分布的检验	24
1.3.2 符号检验	28
1.3.3 符号秩检验与秩和检验	30
1.4 分布检验	33
1.4.1 Pearson 拟合优度 $\chi^2$ 检验	34
1.4.2 Kolmogorov-Smirnov 检验	37
1.4.3 正态性检验	39
1.5 列联表检验	39
1.5.1 Pearson $\chi^2$ 独立性检验	40
1.5.2 Fisher 精确独立性检验	42
1.6 相关性检验	44
1.6.1 Pearson 相关检验	44
1.6.2 Spearman 相关检验	45
1.6.3 Kendall 相关检验	45
1.6.4 cor. test 函数	46
1.7 数学建模案例分析——食品质量 安全抽检数据分析	49
1.7.1 问题的提出	49
1.7.2 问题 1: 三年各主要食品领域 安全情况的变化趋势	49
1.7.3 问题 2: 找出某些规律性的 东西	52
1.7.4 问题 3: 如何改进食品的抽检 办法	58
1.7.5 结论	59
习题 1	59
<b>第 2 章 多元分析模型</b>	64
2.1 回归分析	64
2.1.1 线性回归模型	64
2.1.2 回归诊断	69
2.1.3 逐步回归	77
2.2 方差分析	81
2.2.1 单因素方差分析	81
2.2.2 多重均值检验	85
2.2.3 进一步讨论	87
2.2.4 秩检验	89
2.2.5 双因素方差分析	90
2.3 判别分析	97
2.3.1 判别分析的基本原理	97
2.3.2 判别分析的计算	99
2.4 数学建模案例分析——气象观察 站的优化	102
2.4.1 问题的提出	102
2.4.2 假设	103
2.4.3 分析	103
2.4.4 问题的求解	104
2.4.5 结论	105

习题 2 .....	106	4.3.3 投资组合模型 .....	179
<b>第 3 章 线性规划模型</b> .....	110	4.3.4 选址问题 .....	181
3.1 线性规划的数学模型 .....	110	4.4 数学建模案例分析——飞行管理 问题 .....	183
3.1.1 数学模型 .....	110	4.4.1 问题的提出 .....	183
3.1.2 线性规划的图解法 .....	112	4.4.2 数学模型的建立 .....	185
3.2 线性规划问题求解 .....	114	4.4.3 问题的求解 .....	185
3.2.1 程序包的下载与安装 .....	114	4.4.4 结论 .....	188
3.2.2 lp()函数的使用 .....	115	习题 4 .....	188
3.2.3 灵敏度分析 .....	117	<b>第 5 章 图论与网络模型</b> .....	191
3.2.4 整数规划 .....	120	5.1 图的基本概念 .....	191
3.3 运输问题与最优指派问题 .....	123	5.1.1 柯尼斯堡七桥问题 .....	191
3.3.1 运输问题 .....	123	5.1.2 图的定义 .....	192
3.3.2 最优指派问题 .....	127	5.1.3 简单图与完全图 .....	195
3.4 线性规划模型的应用 .....	129	5.1.4 偶图 .....	196
3.4.1 城市规划 .....	130	5.1.5 邻接矩阵与赋权矩阵 .....	197
3.4.2 生产计划与库存控制 .....	131	5.1.6 子图与补图 .....	199
3.4.3 人力规划 .....	137	5.1.7 顶点度 .....	200
3.4.4 下料问题 .....	139	5.1.8 路和连通 .....	203
3.4.5 集合覆盖问题 .....	141	5.2 最短路问题 .....	205
3.5 数学建模案例分析 .....	142	5.2.1 计算固定两点间的最短路 .....	205
3.5.1 装货问题 .....	142	5.2.2 计算任意两点间的最短路 .....	209
3.5.2 DVD 在线租赁问题 .....	145	5.2.3 计算最短路 R 函数 .....	209
习题 3 .....	151	5.2.4 最短路问题的应用 .....	212
<b>第 4 章 最优化模型</b> .....	157	5.3 最优连线问题 .....	215
4.1 最优化问题的数学模型 .....	157	5.3.1 树 .....	215
4.1.1 无约束优化问题 .....	157	5.3.2 生成树 .....	217
4.1.2 约束优化问题 .....	159	5.3.3 最优树 .....	217
4.1.3 求解最优化问题的图解法 .....	162	5.4 图的连通度 .....	218
4.2 最优化问题的求解 .....	164	5.4.1 基本概念 .....	219
4.2.1 一元函数求极值 .....	164	5.4.2 计算图连通度的 R 函数 .....	220
4.2.2 多元无约束问题 .....	164	5.5 最大流问题 .....	222
4.2.3 多元约束问题 .....	168	5.5.1 最大流问题的基本概念 .....	222
4.2.4 求极值函数的扩展 .....	170	5.5.2 主要定理 .....	223
4.3 最优化模型的应用 .....	176	5.5.3 求解最大流问题的 R 函数 .....	224
4.3.1 曲线拟合 .....	176	5.6 中国邮递员问题 .....	225
4.3.2 路灯照明问题 .....	177		

5.6.1 Euler图 .....	225	6.3.3 案例分析——跟车安全距离 的讨论 .....	278
5.6.2 中国邮递员问题 .....	226	6.4 数值积分与数值微分 .....	280
5.7 旅行商问题 .....	228	6.4.1 数值积分 .....	280
5.7.1 Hamilton圈 .....	228	6.4.2 重积分的计算 .....	283
5.7.2 求解旅行商问题 .....	229	6.4.3 数值微分 .....	284
5.7.3 求解旅行商问题的R函数 .....	231	6.4.4 案例分析——估计水塔的 水流量 .....	288
5.7.4 旅行商问题的应用——印刷线路 板过孔问题 .....	233	6.5 求解非线性方程(组) .....	291
5.8 数学建模案例分析 .....	236	6.5.1 非线性方程求根 .....	291
5.8.1 通信网络最优连线问题 .....	236	6.5.2 求解非线性方程组 .....	293
5.8.2 灾情巡视路线 .....	240	6.5.3 案例分析——GPS定位 问题 .....	298
习题5 .....	245	6.6 非线性最小二乘问题 .....	299
<b>第6章 数值分析</b> .....	<b>250</b>	6.6.1 求解非线性最小二乘问题的 函数 .....	300
6.1 数值代数 .....	250	6.6.2 案例分析——GPS定位 问题(续) .....	302
6.1.1 矩阵运算 .....	250	6.7 常微分方程(组)的数值解 .....	303
6.1.2 矩阵分解 .....	251	6.7.1 常微分方程(组)初值问题 .....	303
6.1.3 求解线性方程组 .....	257	6.7.2 求解高阶微分方程 .....	308
6.1.4 线性方程组的应用——投入产出 模型 .....	259	6.7.3 常微分方程边值问题 .....	311
6.2 插值 .....	264	6.7.4 常微分方程建模 .....	312
6.2.1 多项式插值 .....	264	习题6 .....	317
6.2.2 分段线性插值 .....	266	<b>答案</b> .....	323
6.2.3 三次样条函数 .....	267	<b>索引</b> .....	330
6.2.4 二元插值函数 .....	271	<b>参考文献</b> .....	335
6.3 数据拟合 .....	276		
6.3.1 最小二乘原理 .....	276		
6.3.2 求解超定线性方程组的QR分解 方法 .....	277		



# 第 1 章 概率统计模型

数学模型是用变量及数学符号建立起来的一系列等式和不等式，是用来描述客观事物的特征、内在联系及其规律性的模型。

客观事物的某些特征的表现形式，往往具有某种不确定性，因此，代表其特征的变量的取值具有随机性。有时这些变量虽然不具有随机性，但由于观测条件的限制或随机因素的干扰，这些变量的观测值也具有随机性。

如果按变量的取值是否具有随机性来划分数学模型，数学模型就分成了确定性数学模型和不确定性数学模型。不确定性数学模型又称为随机模型。

本章和下一章将介绍随机模型的有关内容，以及如何使用 R 中的函数求解随机模型。后面的各章将介绍确定性模型，以及如何下载扩展程序包，完成确定性模型的求解工作。

## 1.1 数据的描述性分析

在建立随机模型之前，首先要分析数据的主要特征，也就是数据的数字特征。这些特征通常是均值、方差，或者是数据服从什么分布。只有在确定了这些特征之后，才能建立起符合实际的模型。

本节介绍数据的描述性分析的统计方法，以及完成此类分析的 R 函数。

### 1.1.1 数据的数字特征

已知一组试验(或观测)数据为  $X_1, X_2, \dots, X_n$ ，它们可以从所要研究的对象的全体——总体  $X$  中取出的，这  $n$  个观测值就构成一个样本。数据分析的任务就是要对这  $n$  个数据进行分析，提取数据中包含的有用信息。

#### 1. 样本均值

样本均值是数据的平均数，也就是样本的算术平均值，其计算公式为

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \quad (1.1)$$

它描述数据取值的平均位置。

在 R 中，`mean()` 函数用于计算样本均值，其使用格式为

```
mean(x, trim = 0, na.rm = FALSE, ...)
```

参数 `x` 为需要计算均值的对象(如向量、矩阵、数组或数据框)。

`trim` 为  $(0, 0.5)$  之间的数(默认值为 0)，表示在计算均值之前，去掉两端数据的百分比，即计算截尾均值。

`na.rm` 为逻辑变量，表示能否处理带有缺失数据(NA)的样本，默认值为 `FALSE`。

...为附加参数.

例如,

```
> x <- c(0:10, 50); xm <- mean(x)
> c(xm, mean(x, trim = 0.10))
[1] 8.75 5.50
```

程序中的 `c()` 为连接函数, 将数据连接成向量.

## 2. 样本方差

方差是描述数据取值分散性的一个度量. 样本方差是样本相对于均值的偏差平方和的平均, 记为  $S^2$ , 即

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 \quad (1.2)$$

其中  $\bar{X}$  为样本均值. 样本方差的开方称为样本标准差, 记为  $S$ , 即

$$S = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2} \quad (1.3)$$

在 R 中, 可用 `var()` 函数计算样本方差, 其使用格式为

```
var(x, y = NULL, na.rm = FALSE, use)
```

参数  $x$  为数值向量、矩阵或数据框. 当  $y$  为 `NULL`(默认值)时, 计算样本  $x$  的方差; 当  $y$  为数值向量、矩阵或数据框时, 计算样本  $x$  与  $y$  的协方差.

`na.rm` 为逻辑变量, 表示能否处理带有缺失数据(NA)的样本, 默认值为 `FALSE`.

计算样本标准差的函数是 `sd()`, 其使用格式为

```
sd(x, na.rm = FALSE)
```

参数的意义与 `var()` 函数相同.

例如,

```
> x <- c(12, 9, 11, 5, 1, 4, 8, 3, 2, 10, 6, 7)
> var(x)
[1] 13
> sd(x)
[1] 3.605551
```

## 3. 中位数

中位数就是数据排序位于中间位置的值. 如果数据按顺序排列, 对于奇数个数据, 中位数就是中间位置的数据; 对于偶数个数据, 中位数就是中间两个数据的平均值.

在 R 中, 可用 `median()` 函数计算样本的中位数, 其使用格式为

```
median(x, na.rm = FALSE)
```

参数  $x$  为需要计算中位数的数值型向量.

`na.rm` 为逻辑变量, 表示能否处理带有缺失数据(NA)的样本, 默认值为 `FALSE`.

例如,

```
> x <- c(12, 9, 11, 5, 1, 4, 8, 3, 2, 10, 6, 7, NA)
> median(x)
[1] NA
> median(x, na.rm = TRUE)
[1] 6.5
```

#### 4. 分位数

分位数可以看成中位数的推广, 如中位数就是 0.5 分位数. 常用的还有  $\frac{1}{4}$  分位数和  $\frac{3}{4}$  分位数, 分别称为下四分位数和上四分位数.  $p$  分位数又可称为第  $100p$  百分位数, 如中位数可称为第 50 百分位数, 下四分位数可称为第 25 百分位数, 上四分位数可称为第 75 百分位数.

在 R 中, `quantile()` 函数计算样本的分位数, 其使用格式为

```
quantile(x, probs = seq(0, 1, 0.25), na.rm = FALSE,
         names = TRUE, type = 7, ...)
```

参数  $x$  为样本构成的数值向量.

`probs` 为数值向量, 数值在 0 与 1 之间, 表示需要计算的分位数, 默认值为 0, 0.25, 0.50, 0.75 和 1.

`na.rm` 为逻辑变量, 表示能否处理带有缺失数据(NA)的样本, 默认值为 `FALSE`.

`names` 为逻辑变量, 表示是否将百分位数作为返回值中变量的属性, 默认值为 `TRUE`.

`type` 为 1~9 之间任何一个整数, 表示计算分位数的算法, 默认值为 7.

例如,

```
> quantile(1:10)
 0%   25%   50%   75%   100%
1.00  3.25  5.50  7.75  10.00
```

#### 5. 极差与四分位极差

样本极差(记为  $R$ )是一组数据的最大值与最小值之差, 其计算公式为

$$R = \max(X) - \min(X) \quad (1.4)$$

其中  $X$  是由样本构成的向量. 样本极差是描述样本分散性的数字特征. 数据越分散, 其极差越大. 由于极差是利用一组数据两端的信息, 因此容易受到极端值的影响.

在 R 中, 与极差有关的函数有 `max()`, `min()` 和 `range()`, 其使用格式为

```
max(..., na.rm = FALSE)
min(..., na.rm = FALSE)
range(..., na.rm = FALSE)
```

参数...为数据构成的向量.

`max()` 函数的返回值是最大值, `min()` 函数的返回值是最小值, `range()` 函数的返回

值是由最小值和最大值构成的二维向量。

因此，计算极差的程序为

```
R <- max(x) - min(x)
r <- range(x); R <- r[2] - r[1]
```

其中  $x$  是由样本构成的向量。

样本的上、下四分位数之差称为四分位极差(或半极差)，记为  $R_1$ ，即

$$R_1 = Q_3 - Q_1 \quad (1.5)$$

它也是度量样本分散性的重要数字特征，特别对于具有异常值的数据，它作为分散性度量具有稳健性。因此，它在稳健性数据分析中具有重要作用。

由半极差的定义(式(1.5))，可得到半极差的计算公式

```
R1 <- quantile(x, 3/4) - quantile(x, 1/4)
```

其中  $x$  是由样本构成的向量。

**例 1.1** 某班有 31 名学生，某门课的考试成绩如下：

```
25 45 50 54 55 61 64 68 72 75 75
78 79 81 83 84 84 84 85 86 86 86
87 89 89 89 90 91 91 92 100
```

求这门课程考试成绩的平均值、中位数、方差、极差和四分位极差。

**解** 编写一个统一计算平均值、中位数、方差、极差和四分位极差等统计量的函数(程序名: discript.R)。

```
discript <- function(x){
  R1 <- quantile(x, 3/4, names = FALSE)
    - quantile(x, 1/4, names = FALSE)
  data.frame(
    n = length(x), max = max(x), min = min(x),
    R = max(x) - min(x), R1 = R1,
    mean = mean(x), median = median(x),
    sd = sd(x)
  )
}
```

程序中的 `length()` 函数是计算向量的长度，即维数。

将数据写入数据文件(exam0101.data)，然后调用该函数计算：

```
X <- scan("exam0101.data")
source("discript.R"); discript(X)
```

在程序中，`scan()` 函数是从数据文件中读取数据，`source()` 是将编写好的函数调入内存，其计算结果如下：

```
n max min R R1 mean median sd
1 31 100 25 75 18 76.70968 84 16.69769
```

计算结果说明：共有 31 名学生的成绩，最高分为 100 分，最低分为 25 分，极差是 75 分，半极差是 18 分，平均分数为 76.7 分，中位数分数为 84 分，标准差为 16.7 分。

## 1.1.2 随机变量的分布

### 1. 随机变量

从一个总体中抽取不同的样本，分析各个样本所获得的点估计往往不尽相同，这种表现出变异性特征的量称为变量。

在进行统计试验以前，一般并不知道某一试验的确切结果，但是可以赋予试验结果以实际数量的一个函数。因此这一变量称为随机变量。随机变量常用大写字母表示，如  $X$ ,  $Y$ ,  $Z$ 。它们可能出现的具体结果或数值则可用小写字母表示，如  $x$ ,  $y$ ,  $z$ 。

最常见的随机变量有两类。一类是以计数形式表示的随机变量，称为离散型随机变量；另一类是取值在某个有限或无限区间的随机变量，称为连续型随机变量。

### 2. 分布函数

描述一个随机变量  $X$ ，不仅要说明它能够取哪些值，而且还要关心它取这些值的概率。对任意的实数  $x$ ，令

$$F(x) = P\{X \leq x\}, \quad x \in (-\infty, +\infty) \quad (1.6)$$

则称  $F(x)$  为随机变量  $X$  的分布函数，也称为累积分布函数。

从直观上看，分布函数  $F(x)$  是一个定义在  $(-\infty, +\infty)$  上的实值函数， $F(x)$  在点  $x$  处取值为随机变量  $X$  落在区间  $(-\infty, x]$  上的概率。

### 3. 概率函数与概率密度函数

如果随机变量  $X$  的全部可能取值只有有限多个或可列无穷多个，则称  $X$  为离散型随机变量。

对于离散型随机变量  $X$ ，可能取值为  $x_k$  的概率为

$$P\{X = x_k\} = p_k, \quad k = 1, 2, \dots \quad (1.7)$$

则称式(1.7)为离散型随机变量  $X$  的分布律。

离散型随机变量的分布函数为

$$F(x) = P\{X \leq x\} = \sum_{x_k \leq x} P\{X = x_k\} = \sum_{x_k \leq x} p_k \quad (1.8)$$

对于随机变量  $X$ ，如果存在一个定义在  $(-\infty, +\infty)$  上的非负函数  $f(x)$ ，使得对于任意实数  $x$ ，总有

$$F(x) = P\{X \leq x\} = \int_{-\infty}^x f(t) dt, \quad -\infty < x < +\infty \quad (1.9)$$

则称  $X$  为连续型随机变量， $f(x)$  为  $X$  的概率密度函数，简称概率密度。

### 4. 分位数

设  $X$  为随机变量，对任给的  $0 < \alpha < 1$ ，若存在  $x_\alpha$ ，使得

$$P\{X \leq x_\alpha\} \geq 1 - \alpha, \quad P\{X > x_\alpha\} \geq \alpha \quad (1.10)$$

则称  $x_\alpha$  为  $X$  的上  $\alpha$  分位数(或上  $\alpha$  分位点). 对任给的  $0 < p < 1$ , 若存在  $x_p$ , 使得

$$P\{X \leq x_p\} \geq p, \quad P\{X > x_p\} \geq 1 - p \quad (1.11)$$

则称点  $x_p$  为  $X$  的下  $p$  分位数(或下  $p$  分位点). 由式(1.10)和式(1.11)可以得到上、下分位数之间的关系: 上  $\alpha$  分位数就是下  $1 - \alpha$  分位数.

由分位数的关系式(1.10)(或式(1.11))可知, 分位数不是唯一的. 但对于连续型随机变量, 分位数确实是唯一的, 此时, 可将式(1.11)改写成

$$F(x_p) = p$$

其中  $F(x)$  为随机变量  $X$  的分布函数. 对于连续型随机变量, 分布函数是严格单调递增的, 所以可将下  $p$  分位数表示为逆分布函数, 即

$$x_p = F^{-1}(p) \quad (1.12)$$

由上、下分位数之间的关系式, 上  $\alpha$  分位数可表示为

$$x_\alpha = F^{-1}(1 - \alpha) \quad (1.13)$$

### 1.1.3 常用的分布

#### 1. 正态分布

若随机变量  $X$  的概率密度函数为

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}, \quad -\infty < x < +\infty \quad (1.14)$$

其中  $\mu$  和  $\sigma$  ( $\sigma > 0$ ) 为两个常数, 则称  $X$  服从参数为  $\mu$  和  $\sigma^2$  的正态分布, 也称为 Gauss(高斯)分布, 记作  $X \sim N(\mu, \sigma^2)$ . 若  $X \sim N(\mu, \sigma^2)$ , 则

$$F(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(t-\mu)^2}{2\sigma^2}} dt, \quad -\infty < x < +\infty \quad (1.15)$$

在 R 中, 正态分布的基本名称为 `norm`, 加上不同的前缀表示不同的函数, 如 `dnorm` 表示概率密度函数, `pnorm` 表示分布函数, `qnorm` 表示分位函数. 函数的使用格式为

```
dnorm(x, mean = 0, sd = 1, log = FALSE)
pnorm(q, mean = 0, sd = 1, lower.tail = TRUE, log.p = FALSE)
qnorm(p, mean = 0, sd = 1, lower.tail = TRUE, log.p = FALSE)
```

参数  $x$  或  $q$  为纯量或向量, 表示概率密度函数或分布函数的自变量.

$p$  为纯量或向量, 描述分位点的概率.

`mean` 为纯量或向量, 描述均值参数(即  $\mu$ )的取值, 默认值为 0.

`sd` 为纯量或向量, 描述标准差参数(即  $\sigma$ )的取值, 默认值为 1.

`log`, `log.p` 为逻辑变量, 表示概率值  $p$  是否由其对数值  $\log(p)$  给出, 默认值为 `FALSE`.

`lower.tail` 为逻辑变量, 表示是否作下尾运算. 取 `TRUE`(默认值)表示  $F(x) = P\{X \leq x\}$ , 对应的分位数是下分位数. 取 `FALSE` 表示  $F(x) = P\{X > x\}$ , 对应的分位数是上

分位数.

图 1.1 描绘的是不同参数的正态分布的概率密度函数图, 分别是:  $\mu=0, \sigma=0.5$ ;  $\mu=2, \sigma=0.5$ ;  $\mu=0, \sigma=1$ .

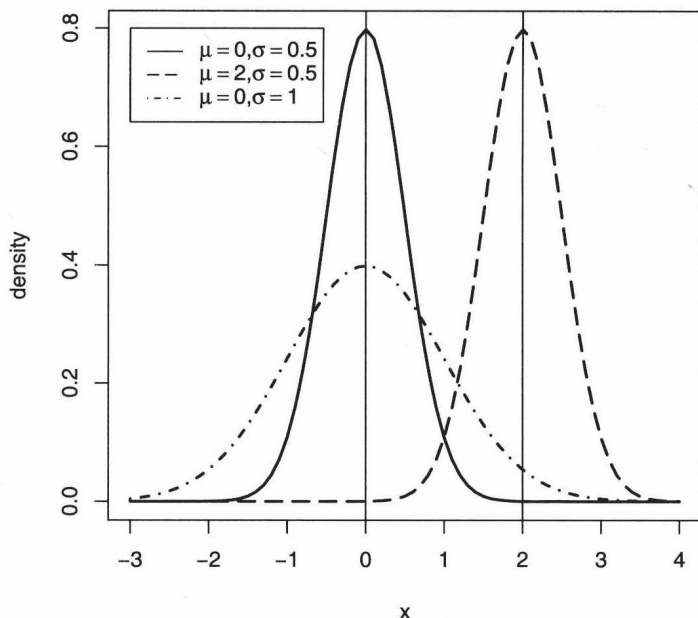


图 1.1 正态分布的概率密度函数

**例 1.2** 设  $X \sim N(\mu, \sigma^2)$ , 分别计算  $P\{|X-\mu| \leq \sigma\}$ ,  $P\{|X-\mu| \leq 2\sigma\}$  和  $P\{|X-\mu| \leq 3\sigma\}$ , 以及正态分布的上 0.025 分位点  $Z_{0.025}$ .

**解** 当  $X \sim N(\mu, \sigma^2)$  时,  $Z = \frac{X-\mu}{\sigma} \sim N(0, 1)$ , 所以用标准正态分布计算即可.

```
> x <- 1:3; p <- pnorm(x) - pnorm(-x); p
[1] 0.6826895 0.9544997 0.9973002
```

这就是通常所说的  $3\sigma$  原则, 即在  $1\sigma$ ,  $2\sigma$  和  $3\sigma$  区间内的概率分别为 68.3%, 95.4% 和 99.7%.

注意: `qnorm()` 函数提供的是下分位点, 计算上  $\alpha$  分位点等价于计算下  $1-\alpha$  分位点, 所以计算程序为

```
> alpha <- 0.025; z <- qnorm(1-alpha); z
[1] 1.959964
```

## 2. $\chi^2$ 分布

如果  $Z_i \sim N(0, 1) (i=1, 2, \dots, n)$ , 且  $Z_i$  是相互独立的, 则称

$$X = Z_1^2 + Z_2^2 + \dots + Z_n^2 \quad (1.16)$$

为自由度为  $n$  的  $\chi^2$  分布, 记为  $X \sim \chi^2(n)$ . 如果  $Z_i \sim N(\delta, 1)$ , 则称  $X$  为非中心化的  $\chi^2$  分

布，记  $X \sim \chi^2(n, \delta)$ ，称  $\delta$  为非中心化参数。

在 R 软件中，用 `chisq` 表示  $\chi^2$  分布，其调用格式如下：

```
dchisq(x, df, ncp = 0, log = FALSE)
pchisq(q, df, ncp = 0, lower.tail = TRUE, log.p = FALSE)
qchisq(p, df, ncp = 0, lower.tail = TRUE, log.p = FALSE)
```

参数 `df` 为自由度，`ncp` 为非中心化参数(即  $\delta$ )，其余参数的意义与正态分布函数相同。

图 1.2 描绘的是  $\chi^2$  分布的概率密度函数在不同参数下的图形。

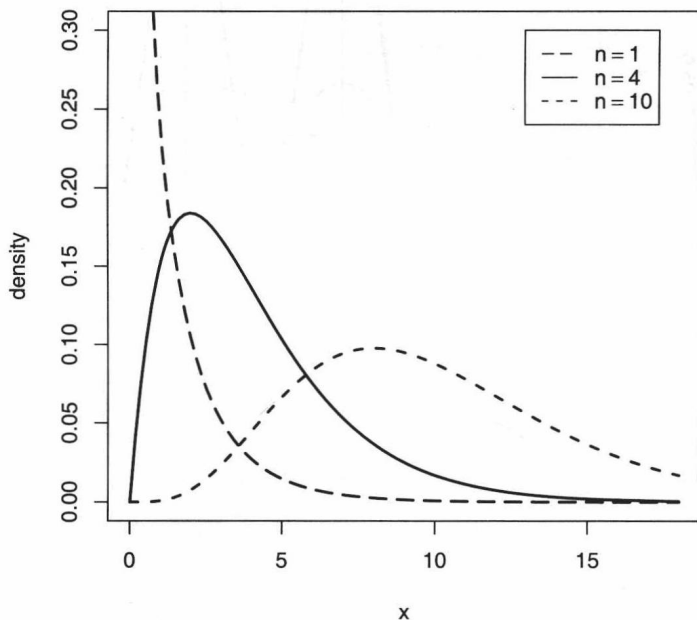


图 1.2  $\chi^2$  分布的概率密度函数

### 3. $t$ 分布

如果随机变量  $Z \sim N(0, 1)$ ， $X \sim \chi^2(n)$  且  $X$  与  $Z$  相互独立，则称

$$T = \frac{Z}{\sqrt{X/n}} \quad (1.17)$$

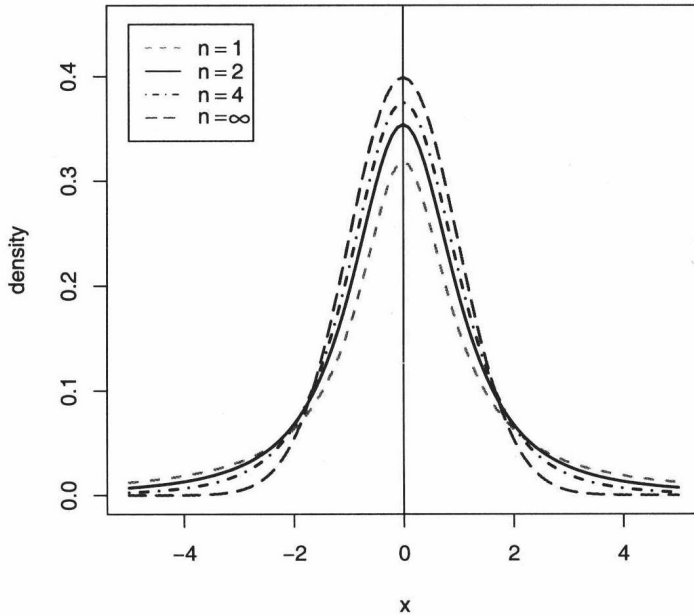
为自由度为  $n$  的  $t$  分布，记为  $T \sim t(n)$ 。如果  $Z \sim N(\delta, 1)$ ，则称  $T$  为非中心化  $t$  分布，记为  $T \sim t(n, \delta)$ ，称  $\delta$  为非中心化参数。

在 R 中， $t$  分布的使用格式是

```
dt(x, df, ncp = 0, log = FALSE)
pt(q, df, ncp = 0, lower.tail = TRUE, log.p = FALSE)
qt(p, df, ncp = 0, lower.tail = TRUE, log.p = FALSE)
```

参数 `df` 为自由度，`ncp` 为非中心化参数(即  $\delta$ )，其余参数的意义与正态分布函数相同。



图 1.3 描绘的是  $t$  分布的概率密度函数在不同参数下的图形.图 1.3  $t$  分布的概率密度函数

#### 4. $F$ 分布

如果随机变量  $X \sim \chi^2(n_1)$ ,  $Y \sim \chi^2(n_2)$  且相互独立, 则称

$$F = \frac{X/n_1}{Y/n_2} \quad (1.18)$$

为第 1 自由度为  $n_1$  和第 2 自由度为  $n_2$  的  $F$  分布, 记为  $F \sim F(n_1, n_2)$ . 如果  $X \sim \chi^2(n_1, \delta)$ , 则称  $F$  为非中心化  $F$  分布, 记为  $F \sim F(n_1, n_2; \delta)$ , 称  $\delta$  为非中心化参数.

在 R 中,  $F$  分布的使用格式为

```
df(x, df1, df2, ncp = 0, log = FALSE)
pf(q, df1, df2, ncp = 0, lower.tail = TRUE, log.p = FALSE)
qf(p, df1, df2, ncp = 0, lower.tail = TRUE, log.p = FALSE)
```

参数  $df1$  为第 1 自由度,  $df2$  为第 2 自由度,  $ncp$  为非中心化参数(即  $\delta$ ), 其余参数的意义与正态分布函数相同.

图 1.4 描绘的是  $F$  分布的概率密度函数在不同参数下的图形.

#### 1.1.4 数据的图形描述

可以通过数据的图形描述判断数据的分布情况, 比如是否来自于正态分布等.

##### 1. 直方图

直方图又称柱状图或质量分布图, 是一种统计报告图, 由一系列高度不等的纵条纹或