

 新世纪高等学校教材

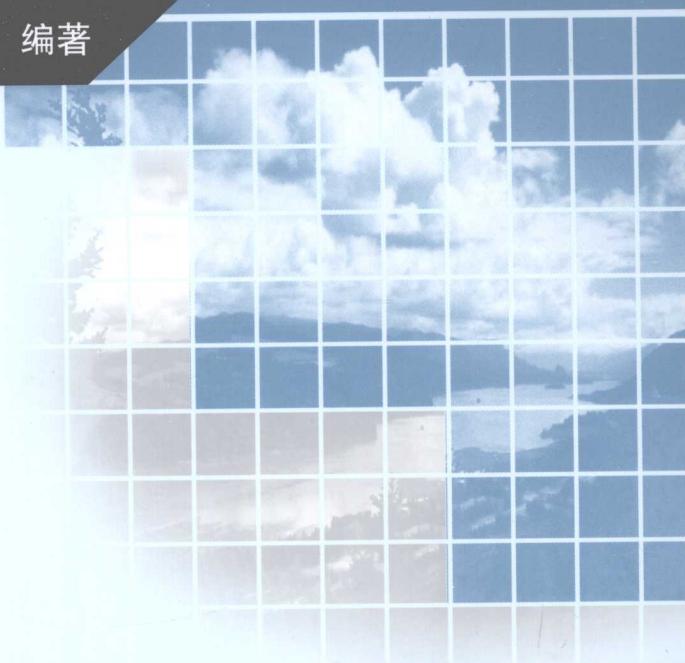
环境科学与工程系列教材

北京师范大学环境学院 组编

环境统计 分析

杨晓华 刘瑞民 曾 勇 编著

HUANJING TONGJI FENXI



北京师范大学出版集团
BEIJING NORMAL UNIVERSITY PUBLISHING GROUP
北京师范大学出版社

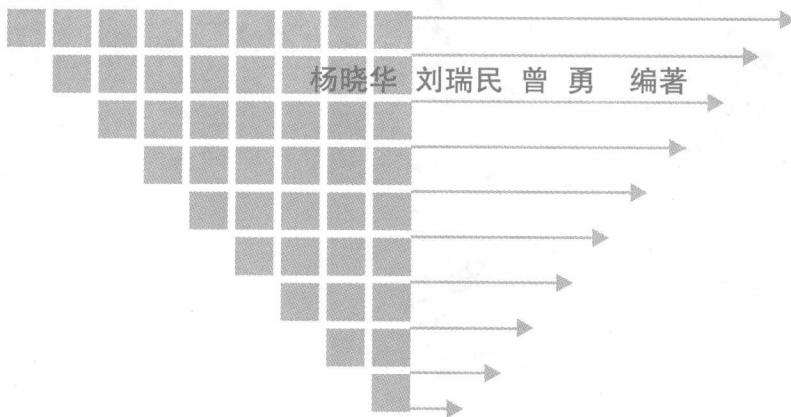
新世纪高等学校教材

环境科学与工程系列教材

北京师范大学环境学院 组编

环境统计分析

HUANJING TONGJI FENXI



北京师范大学出版集团

BEIJING NORMAL UNIVERSITY PUBLISHING GROUP

北京师范大学出版社

图书在版编目(CIP) 数据

环境统计分析 / 杨晓华, 刘瑞民, 曾勇编著. —北京: 北京师范大学出版社, 2008.8
(环境科学与工程系列教材)
ISBN 978-7-303-09502-5

I . 环… II . ①杨… ②刘… ③曾… III . 环境统计—统计分析—高等学校—教材 IV . X11

中国版本图书馆 CIP 数据核字(2008)第 113056 号

出版发行: 北京师范大学出版社 www.bnup.com.cn

北京新街口外大街 19 号

邮政编码: 100875

印 刷: 北京新丰印刷厂

经 销: 全国新华书店

开 本: 170mm × 230 mm

印 张: 20.25

字 数: 337 千字

印 数: 1 ~ 3 000 册

版 次: 2008 年 9 月第 1 版

印 次: 2008 年 9 月第 1 次印刷

定 价: 35.00 元

责任编辑: 毛 佳 装帧设计: 高 霞

责任校对: 李 菡 责任印制: 马鸿麟

版权所有 侵权必究

反盗版、侵权举报电话: 010 - 58800697

北京读者服务部电话: 010 - 58808104

外埠邮购电话: 010 - 58808083

本书如有印装质量问题, 请与印制管理部联系调换。

印制管理部电话: 010 - 58800825

前 言

环境统计分析是环境科学与环境工程的基础学科之一，是一门对环境系统不确定性问题进行数据处理、模型构建和分析的学科。环境系统，系指地球表面包括非生物、生物的各种环境因素及其相互关系的总和，是一个具有时、空、量、序变化的复杂巨系统。受人类活动、天文、气候和气象等众多因素的影响，环境系统中存在许多不确定性现象，并且有大量的数据需要进行统计分析和处理。环境的理论和实践对统计信息的需求急剧增加，对统计分析的理论和方法提出了更高的要求。在自然、社会与环境关系的基础上，用统计方法对环境问题予以量化描述和分析已成为环境研究的迫切需要。环境统计学的产生与发展使人们能够利用数理统计方法处理或解决环境中的不确定性问题，使其定量化，其中包括寻找变量之间的定量关系、从数据中发现环境趋势、探索环境系统变化规律。现代环境统计学一个很重要的标志就是模型技术的运用及量化分析。

全书分三大部分，共10章。其中，第1章属于基础篇，简要地介绍了环境统计分析的概率统计基础知识；第2~9章属于模型篇，阐述了环境一元线性回归分析、环境多元线性回归分析、环境系统聚类分析、环境模糊聚类分析、环境判别分析、环境主成分分析、环境因子分析、人工神经网络等方法、模型的原理，并给出了分析案例；第10章属于空间分析篇，介绍了环境空间统计分析的基本原理，

并给出了应用实例。全书的大多数例子都是用目前常用的统计分析语言 Matlab 编写实现的，是理论联系实际的经验总结，具有可操作性。本书适于做高等院校环境科学与环境工程专业的高年级本科生和研究生教材，对环境科学与环境工程、生态学、资源与管理、应用数学、地理科学等相关领域的学者和科研人员也有重要的参考价值。

本书第 1 章由杨晓华、曾勇执笔，第 2~9 章由杨晓华执笔，第 10 章由刘瑞民执笔，全书由杨晓华统稿。另外，尹心安参加了第 1 章的编写工作；王伟参加了第 3 章、第 4 章、第 10 章的编写工作；陈强、胡晓雪参加了第 5 章、第 6 章的编写工作；余敦先参加了第 1 章、第 3 章、第 6 章、第 8 章的编写工作。2004 级、2005 级的博士研究生、2005 级的硕士研究生也提供了部分例题和习题。另外，习题答案均是用 Matlab 语言计算完成。

在本书的编写和出版过程中，北京师范大学环境学院院长杨志峰教授，副院长沈珍瑶、刘静玲教授，还有牛军峰、孙涛副教授以及北京师范大学出版社的胡廷兰、毛佳等同志对本书提出了许多宝贵意见。书中若干例题选自所列参考文献，在此一并表示感谢。由于我们的水平有限，书中错误在所难免，欢迎读者批评指正。

衷心感谢北京师范大学出版社给予的大力支持！

本书的完成得到国家重点基础研究发展规划项目（G2003CB415204）的资助，在此表示衷心的感谢！

编著者

2007 年 7 月

内 容 提 要

本书阐述了常用的环境统计分析方法，并给出了分析案例。首先简明扼要地介绍了环境统计分析的概率统计基础知识，又重点阐述了环境一元线性回归分析、环境多元线性回归分析、环境系统聚类分析、环境模糊聚类分析、环境判别分析、环境主成分分析和环境因子分析这些常用的环境统计分析模型；另外还给出了现代环境数据处理常用的人工神经网络方法和空间统计分析方法。对每一种方法，本书除了讲明基本原理外，还给出了大量的计算分析例题和案例。本书的部分例子是用目前实用的统计分析语言 Matlab 编写实现的，是理论联系实际的经验总结，具有实用性。本书适于做高等院校环境科学与环境工程专业的高年级本科生和研究生教材，对环境科学与环境工程、生态学、资源与管理、应用数学、地理科学等相关领域的学者和科研人员也有重要的参考价值。

目 录

第1章 概率统计基础 (1)

1.1 四种重要的概率分布.....	(1)
1.1.1 正态分布.....	(1)
1.1.2 χ^2 分布	(4)
1.1.3 t 分布	(5)
1.1.4 F 分布	(6)
1.2 随机向量的数字特征.....	(7)
1.2.1 数学期望	(7)
1.2.2 方差和均方差	(10)
1.2.3 原点矩和中心矩	(11)
1.2.4 变异系数	(12)
1.2.5 协方差阵和自协方差阵	(12)
1.2.6 随机变量的相关系数	(13)
1.2.7 总体与样本	(15)
1.2.8 样本子样的一些数字特征	(16)
1.2.9 大数定律	(16)
1.2.10 中心极限定理	(18)
1.3 参数估计.....	(20)
1.3.1 点估计	(21)
1.3.2 区间估计	(21)
1.4 参数假设检验.....	(24)
1.4.1 假设检验的原理	(25)

1.4.2 假设检验的步骤	(26)
1.4.3 参数检验	(27)
1.5 方差分析与试验设计初步	(34)
1.5.1 方差分析概述	(34)
1.5.2 单因素方差分析	(35)
1.5.3 双因素方差分析	(39)
1.5.4 试验设计初步	(45)
思考题 1	(48)
参考文献	(49)

第 2 章 环境一元线性回归分析 (50)

2.1 一元线性回归模型	(50)
2.1.1 变量间的统计关系	(50)
2.1.2 一元线性回归模型	(52)
2.1.3 最小二乘法估计	(54)
2.2 线性回归方程的显著性检验	(55)
2.2.1 F 检验法	(56)
2.2.2 相关系数检验法	(58)
2.2.3 样本决定系数 r^2	(59)
2.3 线性回归式的误差估计	(60)
2.3.1 线性回归式的误差估计	(60)
2.3.2 线性回归的步骤	(61)
2.4 可化为一元线性回归的曲线回归	(62)
2.4.1 倒数变换	(62)
2.4.2 对数变换	(63)
2.4.3 混合变换	(64)
2.5 环境应用	(65)
思考题 2	(69)
参考文献	(70)

第 3 章 环境多元线性回归分析 (71)

3.1 多元线性回归模型	(71)
3.2 参数的最小二乘估计	(72)

3.3 回归方程的显著性检验.....	(74)
3.3.1 拟合优度检验.....	(75)
3.3.2 F 检验.....	(76)
3.4 回归系数的显著性检验.....	(77)
3.5 Matlab 语言在多元回归中的应用	(79)
3.6 环境应用.....	(81)
思考题 3	(84)
参考文献	(86)

第 4 章 环境系统聚类分析 (87)

4.1 聚类分析概述.....	(87)
4.2 聚类要素的数据处理.....	(88)
4.3 距离和相似系数的计算.....	(93)
4.3.1 距离的计算.....	(93)
4.3.2 相似系数的计算.....	(97)
4.3.3 距离和相似系数选择原则.....	(99)
4.4 系统聚类分析常用方法.....	(100)
4.4.1 最短距离系统聚类法原理.....	(102)
4.4.2 最远距离聚类法原理.....	(103)
4.4.3 系统聚类法公式的统一.....	(105)
4.5 环境应用.....	(107)
思考题 4	(112)
参考文献	(115)

第 5 章 环境模糊聚类分析 (116)

5.1 模糊集理论.....	(116)
5.1.1 模糊集的基本概念.....	(117)
5.1.2 模糊集的表示方法.....	(117)
5.1.3 模糊集的运算.....	(119)
5.1.4 模糊映射.....	(120)
5.2 模糊关系	(120)
5.3 模糊等价关系	(121)
5.4 模糊聚类分析步骤	(123)

5.4.1	数据标准化	(123)
5.4.2	模糊相似矩阵的建立	(124)
5.4.3	聚类分析	(126)
5.4.4	分类的 F 检验	(130)
5.5	环境应用	(132)
思考题 5	(137)	
参考文献	(139)	

第 6 章 环境判别分析 (140)

6.1	距离判别分析	(140)
6.1.1	两总体情况	(140)
6.1.2	多总体情况	(144)
6.2	Fisher 判别	(145)
6.3	Bayes 判别	(150)
6.4	环境应用	(153)
思考题 6	(161)	
参考文献	(163)	

第 7 章 环境主成分分析 (164)

7.1	主成分分析概述	(164)
7.2	主成分分析计算原理	(165)
7.3	主成分分析的性质	(169)
7.4	环境应用	(170)
思考题 7	(178)	
参考文献	(180)	

第 8 章 环境因子分析 (181)

8.1	因子分析概述	(181)
8.2	正交因子模型	(182)
8.3	正交因子模型的统计意义	(184)
8.4	正交因子模型的求解	(185)
8.5	因子旋转	(188)

8.6 因子得分.....	(191)
8.7 环境应用.....	(193)
思考题 8	(204)
参考文献	(205)

第 9 章 人工神经网络 (206)

9.1 人工神经网络概述.....	(206)
9.2 人工神经元模型.....	(209)
9.3 BP 神经网络	(212)
9.3.1 BP 神经网络原理	(212)
9.3.2 BP 算法	(213)
9.3.3 环境应用	(223)
9.4 RBF 神经网络	(225)
9.4.1 RBF 神经网络原理	(225)
9.4.2 RBF 神经网络模型	(226)
9.4.3 环境应用	(228)
思考题 9	(230)
参考文献	(230)

第 10 章 环境空间统计分析 (232)

10.1 环境空间信息概述	(232)
10.1.1 环境空间信息特征	(233)
10.1.2 环境空间信息种类	(234)
10.1.3 环境空间信息来源	(234)
10.2 环境空间统计分析	(236)
10.2.1 区域化变量	(237)
10.2.2 协方差函数	(238)
10.2.3 变差函数	(239)
10.2.4 普通克立格插值	(248)
10.2.5 环境应用	(252)
10.3 环境空间主成分分析	(261)
10.3.1 空间主成分分析步骤	(262)
10.3.2 环境应用	(263)

思考题 10	(268)
参考文献	(269)

部分思考题答案	(270)
---------------	-------

附录	(303)
----------	-------

附表 1 标准正态分布表	(303)
附表 2 相关系数检验表	(304)
附表 3 χ^2 分布临界值表	(305)
附表 4 t 分布临界值表	(306)
附表 5 F 分布临界值表	(307)

第1章 概率统计基础

环境的理论和实践对统计信息的需求急剧增加，对统计分析的理论和方法提出了更高的要求。在自然、社会与环境关系的基础上，用统计方法对环境问题予以量化分析已成为环境科学工作者的迫切需要。环境统计学的产生与发展使人们能够利用数理统计方法处理或解决环境中的不确定性问题，使其定量化，其中包括寻找变量之间的定量关系、从数据中发现环境趋势、探索环境系统变化规律。为了能深刻理解和分析环境数据的数量特征和内在关系，需要我们首先掌握数理统计的基础知识。本章重点阐述环境统计分析的概率统计基础。

本章的主要内容是：

- 四种重要的概率分布；
- 随机向量的数字特征；
- 参数估计；
- 参数假设检验；
- 方差分析与试验设计初步。

1.1 四种重要的概率分布

在环境科学中，弄清统计分析对象的理论分布是关键的一环。土壤中的某些污染物、重金属的分布，大气中若干种微粒的浓度分布、监测值的误差分布等均服从正态分布或接近正态分布或取对数后服从正态分布。 χ^2 分布、 t 分布、 F 分布是统计推断中经常碰到的另外三种分布。研究污染物在环境中的分布规律已是当前环境科学的研究中重要的课题之一。

1.1.1 正态分布

市场上的食品很多是 1 kg 袋装，袋上标有“净含量 1 kg”的字样。但当用稍微精确一些的天平称那些食品的重量时，会发现有些可能会重些，有些可能会轻些，但都在 1 kg 左右。其中，多数离 1 kg 不远，离 1 kg 越近就越可能出现，离 1 kg 越远就越不可能。一般认为这种重量分布近似地服从正态分布(normal distribution)。近似地服从正态分布的变量很常见，如实验误差、商品的重量或

尺寸、某年龄人群的身高和体重等。在一定条件下，许多不是正态分布的样本均值在样本量很大时，也可用正态分布来近似。

若随机变量 X 的分布密度为：

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad (-\infty < x < +\infty, \sigma > 0) \quad (1.1)$$

则称 X 服从正态分布 $N(\mu, \sigma^2)$ ，简记为 $X \sim N(\mu, \sigma^2)$ 。其中， μ 为均值， σ 为标准差， σ^2 为方差（标准差的平方）。

正态分布的密度曲线是一个对称的、呈钟形的曲线（最高点在均值处）（图 1-1）。正态分布是一族分布，各种正态分布根据它们的均值和标准差不同而有区别。标准差为 1 的正态分布 $N(0, 1)$ 称为标准正态分布（standard normal distribution）。标准正态分布的密度函数与分布函数记为：

$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} \quad (-\infty < x < +\infty) \quad (1.2)$$

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{y^2}{2}} dy \quad (-\infty < x < +\infty) \quad (1.3)$$

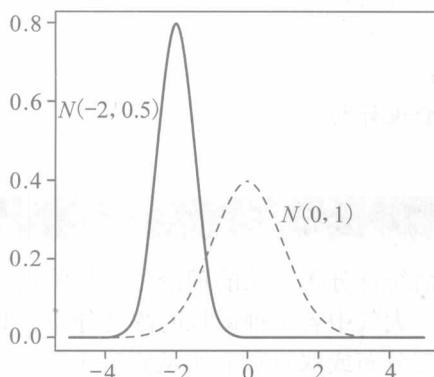


图 1-1 两条正态分布的密度曲线图
(左边是 $N(-2, 0.5)$ 分布，右边是 $N(0, 1)$ 分布)

在实际的生活中，我们经常会因为标准正态分布的优异特性而需要将一般的正态分布标准化，下面简单介绍一下正态分布的标准化过程。

设 $X \sim N(\mu, \sigma^2)$ ，作简单变换（减去其均值 μ ，再除以标准差 σ ），则很容易得到随机变量 $Y = \frac{X - \mu}{\sigma} \sim N(0, 1)$ 。

因为：

$$E(Y) = E\left(\frac{X - \mu}{\sigma}\right) = \frac{1}{\sigma}[E(X) - \mu] = 0$$

$$D(Y) = D\left(\frac{X-\mu}{\sigma}\right) = \frac{1}{\sigma^2} D(X) = 1$$

这样就将一个普通的正态分布变成了一个标准的正态分布。

标准正态分布中还有一个十分重要的概念就是分位点。为了便于今后应用，对于标准正态随机变量，本书引入上侧分位点的定义(盛聚，1998)。

设 $X \sim N(0, 1)$ ，若 z_α 满足条件

$$P(X > z_\alpha) = \alpha \quad (0 < \alpha < 1)$$

则称 z_α 为标准正态分布的上侧 α 分位点，如图 1-2 所示。

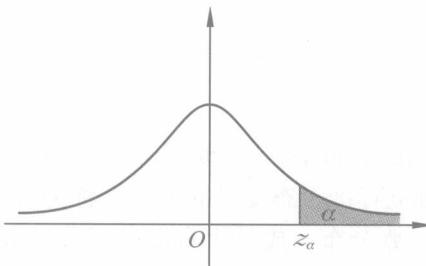


图 1-2 标准正态分布的上侧 α 分位点 z_α

例如，查附表 1 可知： $z_{0.01} = 2.326348$, $z_{0.05} = 1.644854$, $z_{0.10} = 1.281552$, $z_{0.154} = 1.019428$ 。

例 1.1 已知 $X \sim N(\mu, \sigma^2)$ ，求 X 在区间 $(\mu - k\sigma, \mu + k\sigma)$ 的概率，这里 $k = 1, 2, 3$ 。

解 $\forall a, b, 0 < a < b$, 有：

$$\begin{aligned} P(a < X < b) &= \int_a^b \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(x-\mu)^2} dx \\ &= \frac{x-\mu}{\sigma} \int_{\frac{a-\mu}{\sigma}}^{\frac{b-\mu}{\sigma}} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}t^2} dt = \Phi\left(\frac{b-\mu}{\sigma}\right) - \Phi\left(\frac{a-\mu}{\sigma}\right) \end{aligned}$$

这样在区间 $(\mu - k\sigma, \mu + k\sigma)$ 的概率 ($k = 1, 2, 3$) 为：

$$P(\mu - \sigma < X < \mu + \sigma) = \Phi(1) - \Phi(-1) = 0.6826$$

$$P(\mu - 2\sigma < X < \mu + 2\sigma) = \Phi(2) - \Phi(-2) = 0.9544$$

$$P(\mu - 3\sigma < X < \mu + 3\sigma) = \Phi(3) - \Phi(-3) = 0.9974$$

其中， $\Phi(-x) = 1 - \Phi(x)$ 。由此我们可以知道，属于正态分布的随机变量 X 之值，几乎都落在 $(\mu - 3\sigma, \mu + 3\sigma)$ 区间里，落在该区间外的机会极少。

例 1.2 某地水体 COD 浓度 $X \sim N(5, 2^2)$ ，求 COD 浓度落在区间 (4, 8) 的

概率。

解 $\mu=5, \sigma=2$

$$\begin{aligned} P(4 < X < 8) &= \Phi\left(\frac{8-\mu}{\sigma}\right) - \Phi\left(\frac{4-\mu}{\sigma}\right) \\ &= \Phi\left(\frac{8-5}{2}\right) - \Phi\left(\frac{4-5}{2}\right) \\ &= \Phi(1.5) - \Phi(-0.5) \\ &= 0.9332 - 0.3085 \\ &= 0.6247 \end{aligned}$$

1.1.2 χ^2 分布

一个由正态变量导出的分布是 χ^2 分布(chi-square distribution)。该分布在一些检验中会用到。 n 个独立标准正态变量的平方和称为有 n 个自由度的 χ^2 分布，记为 $\chi^2(n)$ 。 χ^2 分布为一族分布，成员由自由度区分。由于 χ^2 分布变量为正态变量的平方和，因此它不会取负值。

设 X_1, X_2, \dots, X_n 是取自标准正态总体 $N(0, 1)$ 的容量为 n 的样本，那么 $\chi^2 = \sum_{i=1}^n X_i^2$ 即为由正态分布导出的自由度为 n 的 $\chi^2(n)$ 分布。所谓自由度，就是指可以自由取值的数据的个数，或者指不受任何约束、可以自由变动的变量的个数。

对于任意一个 $\chi^2(n)$ 分布，它的概率密度函数为：

$$P(x) = \begin{cases} \frac{1}{2^{\frac{n}{2}} \Gamma\left(\frac{n}{2}\right)} x^{\frac{n}{2}-1} e^{-\frac{x}{2}} & (x > 0) \\ 0 & (x \leq 0) \end{cases} \quad (1.4)$$

记为 $\chi^2 \sim \chi^2(n)$ ，式中 n 为正整数， $\Gamma\left(\frac{n}{2}\right)$ 为 Γ 函数值， $\Gamma(z) = \int_0^\infty x^{z-1} e^{-x} dx$ ($z > 0$)。

同正态分布类似，对于 χ^2 分布也有上侧 α 分位点。如果 $P(\chi^2 > \chi_a^2(n)) = \alpha$ ，则称 $\chi_a^2(n)$ 为上侧 α 分位点。对于不同的 α, n ，上侧 α 分位点的值已制成表格（附表 3），可以查到。例如对于 $\alpha = 0.050, n = 9$ ，查得 $\chi_{0.050}^2(9) = 16.919$ 。但大部分书只给出到 $n=45$ 的上侧 α 分位点的值。费歇尔(R. A. Fisher)曾证明，当 n 充分大时，近似有

$$\chi_a^2(n) \approx \frac{1}{2} (z_\alpha + \sqrt{2n-1})^2 \quad (1.5)$$

其中, z_α 为标准正态分布的上侧 α 分位点。利用式(1.5)可以求当 $n > 45$ 时, $\chi^2_\alpha(n)$ 分布的上侧 α 分位点的近似值。

例如, 查附表 1 并计算, 可得:

$$\chi^2_{0.010}(100) \approx \frac{1}{2} (2.326348 + \sqrt{199})^2 \approx 135.0231$$

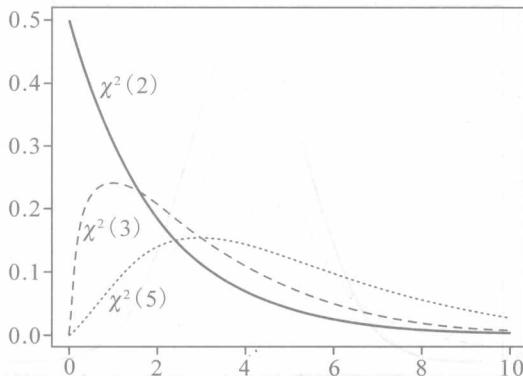


图 1-3 自由度分别为 2,3,5 的 χ^2 分布密度曲线图

1.1.3 t 分布

设 $X \sim N(0, 1)$, $Y \sim \chi^2(n)$, 并且 X, Y 独立, 则随机变量

$$t = \frac{X}{\sqrt{Y/n}}$$

服从自由度为 n 的 t 分布(t -distribution 或 student's t), 记为 $t \sim t(n)$ 。

对于任意一个 $t(n)$ 分布, 它的概率密度函数为:

$$P(x) = \frac{\Gamma(\frac{n+1}{2})}{\Gamma(\frac{n}{2}) \sqrt{n\pi}} \left(1 + \frac{x^2}{n}\right)^{-\frac{n+1}{2}} \quad (-\infty < x < +\infty) \quad (1.6)$$

式中, n 为正整数。

不同的样本量通过标准化所产生的 t 分布也不同, 这样就形成一族分布。 t 分布的分布曲线关于 $x=0$ 对称, 它的密度曲线看上去有些像标准正态分布, 但是中间瘦一些, 而且尾巴长一些。当自由度 k 无限增大时, t 分布将趋近于标准正态分布 $N(0, 1)$ 。

同样, 类似于前面的两个分布, t 分布也有上侧 α 分位点的概念。