

① 高等学校统计学类系列教材

实用回归分析

○ 何晓群 编著



高等教育出版社
HIGHER EDUCATION PRESS

高等学校统计学类系列教材

主要内容

本书是“高等学校统计学类系列教材”中的一本，主要介绍回归分析的基本原理、方法和应用。全书共分五章，第一章介绍回归分析的基本概念和简单线性回归分析；第二章介绍多元线性回归分析；第三章介绍非线性回归分析；第四章介绍 Logistic 回归分析；第五章介绍判别分析。本书可作为高等院校统计学专业及相关专业的教材，也可供从事统计工作的工程技术人员参考。

实用回归分析

何晓群 编著

责任编辑：王桂英

北京：高等教育出版社，2008.2

ISBN 978-7-04-023894-2

定价：25.00元

1. 统计学类 - 回归分析 - 教材 - 高等学校 - 2008.2

IV. O212.1

中国图书馆分类号：O212.1 (2008) 第 030214 号

责任编辑：王桂英
封面设计：王桂英
版式设计：王桂英
责任印制：王桂英
印刷：北京印刷厂

出版发行	高等教育出版社	地址	北京市西城区德胜大街4号	邮政编码	100120	电话	010-28281000	网址	http://www.hep.com.cn	电子邮箱	http://www.widernet.com
总发行	蓝色畅想图书发行有限公司	地址	北京市东城区德胜门内大街1号	邮政编码	100120	电话	010-28281000	网址	http://www.landao.com.cn	电子邮箱	http://www.landao.com.cn
印刷	北京印刷厂	地址	北京市东城区德胜门内大街1号	邮政编码	100120	电话	010-28281000	网址	http://www.landao.com.cn	电子邮箱	http://www.landao.com.cn
开本	787×960 1/16	印张	17.25	字数	330 000	版次	2008年2月第1版	印次	2008年2月第1次印刷	定价	25.00元

高等教育出版社

本书如有缺页、错页、倒页、漏页等质量问题，请向出版单位联系调换。

版权所有 侵权必究

ISBN 978-7-04-023894-2

内容提要

回归分析是现代统计学中应用较为活跃的模型分析技术。本书旨在提高社会、经济、管理类本科生的量化分析水平,选择众多的回归分析方法中最为实用的基本模型分析技术,结合社会经济与管理中的实际问题,利用 SPSS 统计软件对回归建模分析方法作了系统介绍。

本书既可作为统计学类专业教材,也可作为人文社会科学、财经管理类专业工作者的参考书。

图书在版编目(CIP)数据

实用回归分析/何晓群编著. —北京:高等教育出版社, 2008.5

ISBN 978 - 7 - 04 - 023894 - 5

I. 实… II. 何… III. 回归分析 - 高等学校 - 教材 IV. O212.1

中国版本图书馆 CIP 数据核字(2008)第 030214 号

策划编辑 李蕊 责任编辑 蒋青 封面设计 王凌波
责任绘图 杜晓丹 版式设计 王艳红 责任校对 俞声佳
责任印制 宋克学

出版发行 高等教育出版社
社 址 北京市西城区德外大街 4 号
邮政编码 100120
总 机 010 - 58581000

经 销 蓝色畅想图书发行有限公司
印 刷 北京地质印刷厂

开 本 787 × 960 1/16
印 张 17.75
字 数 330 000

购书热线 010 - 58581118
免费咨询 800 - 810 - 0598
网 址 <http://www.hep.edu.cn>
<http://www.hep.com.cn>
网上订购 <http://www.landaco.com>
<http://www.landaco.com.cn>
畅想教育 <http://www.widedu.com>

版 次 2008 年 5 月第 1 版
印 次 2008 年 5 月第 1 次印刷
定 价 24.20 元

本书如有缺页、倒页、脱页等质量问题,请到所购图书销售部门联系调换。

版权所有 侵权必究

物料号 23894 - 00

II 前言

本书也可作为现代统计分析方法课的教材。此书还可作为从事社会、经济、管理等研究和实际工作的同志进行量化研究的参考书。

由于本人学识有限,书中谬误之处在所难免,恳请读者批评指正。

中国人民大学应用统计科学研究中心
中国人民大学六西格玛质量管理研究中心

何晓群

2007年10月1日

于北京九鼎山庄长安仁居

目 录

89		
90		
93		
101		
108		
110		
110		
111		
116	第 1 章 统计学基础	1
133	§ 1.1 统计数据的整理与描述	1
150	§ 1.2 几种重要的概率分布	5
150	§ 1.3 参数估计	13
150	§ 1.4 假设检验	14
153	思考与练习	18
130	第 2 章 回归分析概述	19
141	§ 2.1 变量间的相关关系	19
141	§ 2.2 回归方程与回归名称的由来	22
144	§ 2.3 回归分析的主要内容及其一般模型	23
140	§ 2.4 建立实际问题回归模型的过程	26
131	§ 2.5 回归分析应用与发展述评	32
123	思考与练习	34
120	第 3 章 一元线性回归	35
120	§ 3.1 一元线性回归模型	35
120	§ 3.2 回归参数 β_0, β_1 的估计	39
100	§ 3.3 最小二乘估计的性质	42
134	§ 3.4 回归方程的显著性检验	45
180	§ 3.5 预测和控制	54
185	§ 3.6 建模总结和应注意的问题	58
187	思考与练习	63
180	第 4 章 多元线性回归	65
193	§ 4.1 多元线性回归模型	65
200	§ 4.2 多元回归参数的估计	68
200	§ 4.3 参数估计量的性质	72
200	§ 4.4 回归方程的显著性检验	76
215	§ 4.5 中心化和标准化	80
215	§ 4.6 相关矩阵与偏相关系数	83
215	§ 4.7 建模总结与评注	88
215	思考与练习	94

第 5 章 残差分析	96
§ 5.1 残差与残差图	96
§ 5.2 有关残差的性质	98
§ 5.3 异常值与强影响值	101
思考与练习	108
第 6 章 关于异方差性问题	110
§ 6.1 异方差产生的背景	110
§ 6.2 异方差性的诊断	111
§ 6.3 异方差问题的建模处理	116
思考与练习	122
第 7 章 关于自相关性问题	126
§ 7.1 自相关产生的背景	126
§ 7.2 自相关性的诊断	128
§ 7.3 自相关问题的建模处理	132
思考与练习	139
第 8 章 关于多重共线性问题	141
§ 8.1 多重共线性的产生和原因	141
§ 8.2 多重共线性的诊断	144
§ 8.3 消除多重共线性的方法	149
§ 8.4 本章补充	151
思考与练习	153
第 9 章 自变量选择与逐步回归	156
§ 9.1 自变量选择对估计和预测的影响	156
§ 9.2 所有子集回归	159
§ 9.3 逐步回归	166
§ 9.4 实例与评注	174
思考与练习	180
第 10 章 非线性回归	182
§ 10.1 可化为线性回归的曲线回归	182
§ 10.2 多项式回归	189
§ 10.3 非线性模型	195
§ 10.4 小结与评注	201
思考与练习	206
第 11 章 含定性变量的回归模型	209
§ 11.1 自变量中含有定性变量的回归模型	209
§ 11.2 含有定性变量的回归模型及应用	212
§ 11.3 因变量是定性变量的回归模型	218

§ 11.4 Logistic 回归基本理论和方法	219
§ 11.5 小结与评注	230
思考与练习	233
附录	236
附表 1 简单相关系数的临界值表	236
附表 2 t 分布表	237
附表 3 F 分布表	238
附表 4 D. W 检验上下界表	244
思考与练习参考答案	246
参考文献	272

第1章 统计学基础

为了更顺利地学习本课程的内容,本章将对统计学中的一些基本概念和术语作一简要回顾.

§ 1.1 统计数据的整理与描述

统计学是研究数据规律的方法论学科,统计数据是统计学研究的主要内容.借助统计学方法研究任何实际问题,首先要做的工作就是收集数据,收集数据是一项很重要的基础工作.收集数据的一般方法是查阅各种统计年鉴和报表,再就是运用某种调查方法获取欲研究问题的有关数据.抽样调查获取数据的方式在我国方兴未艾,抽样调查的方法很多,专业性很强,现在已有不少抽样技术的专著.需要利用抽样方法获取数据的研究者,还需很好地学习有关抽样技术的知识.

一、总体与样本

在一个统计问题中,通常把所要调查研究的事物或现象的全体称为总体,把组成总体的每个元素(成员)称为个体,一个总体中所含的个体的数量称为总体的容量.例如要研究某城市居民的家庭收入状况,那么这个城市所有家庭的收入状况就是我们研究的总体,而每个家庭的收入状况就是个体.

为了推断总体的某些特征,需要从总体中按一定的抽样技术抽取若干个体,将这一抽取过程称为抽样.所抽取的部分个体称为样本,样本中所含个体的数量称为样本容量.如在研究居民家庭收入时,随机抽取 1 000 户来进行调查,这 1 000 户就是一个样本,样本容量就是 1 000.

二、统计量

通过抽样或查统计年鉴得到的原始数据,一般是杂乱无章的,很难从中直接看出有价值的东西.因此,对获取的原始数据一般需要加以整理,以便把人们感兴趣的信息提取出来,并用简明醒目的方式加以表述.画原始数据的散点图、饼图、直方图等方法直观表达数据的常见方式.统计学中最主要的提取信息方式就是对原始数据进行一定的运算,以算出某些代表性的数字,足以反映出数据某

些方面的特征,这种数字被称为统计量.用统计学语言表述就是:统计量是样本的函数——它不依赖于任何未知参数.

例如样本均值和样本方差就是最重要的常用统计量.

均值是对数据集中特征的描述,方差是对数据波动特征的描述.

设 x_1, x_2, \dots, x_n 是一组独立的随机样本,则样本均值为

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i,$$

样本方差为

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2,$$

样本标准差为

$$S = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}.$$

例如有两组数据:

(4, 6, 8, 10, 12),

(6, 7, 8, 9, 10),

它们的均值 \bar{x} 都是 8, 这说明两组数据都以 8 为中心.读者可计算它们的方差, 第一组数据的方差比第二组的要大, 说明第一组数据相对均值 8 来说比较分散, 第二组数据相对均值 8 来说比较集中. 由这两组数据可以很直观地看出均值及方差的意义.

需要注意的是, 方差带单位是没有意义的, 只有标准差带上单位才有实际意义.

三、变异系数

如果两组数据的计量单位相同, 且均值一样, 可以利用标准差来比较两组数据的离散程度. 但当两组数据的计量单位不同或均值不同时, 就不能直接利用比较两组数据的标准差来分析两组数据的离散程度. 由此引入变异系数 V :

$$V = \frac{S}{\bar{x}}$$

例如 (4, 5, 6, 7, 8) 与 (40, 50, 60, 70, 80) 两组数据的标准差分别是 1.58 和 15.8, 如果仅从标准差来看显然第二组数据分散程度较大. 但是由于两组数据的均值不同, 分别为 6 和 60, 单纯由标准差来判断数据的分散程度就不合适. 实际上, 当我们算出两组数据的变异系数时, 得到 V 都是 0.26. 比较而言, 两组数据的分散程度就是相同的了.

四、偏度与峰度

偏度和峰度是描述统计数据分布偏斜和陡峭程度的统计量。

偏度用偏度系数 V_1 来描述：

$$V_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})^3}{S^3(n-1)},$$

其中 S 为样本标准差。

偏度系数 V_1 的意义可由图 1.1 表示。

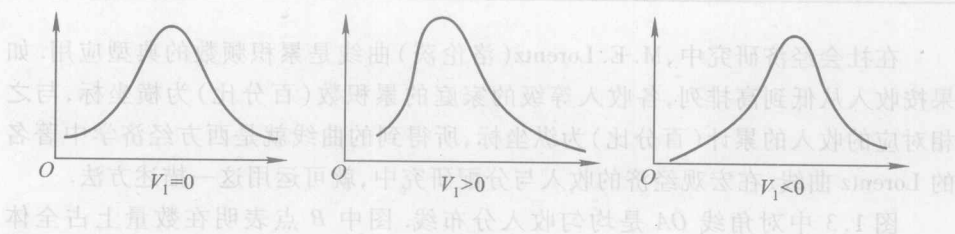


图 1.1

峰度用峰度系数 V_2 表示：

$$V_2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^4}{S^4(n-1)}.$$

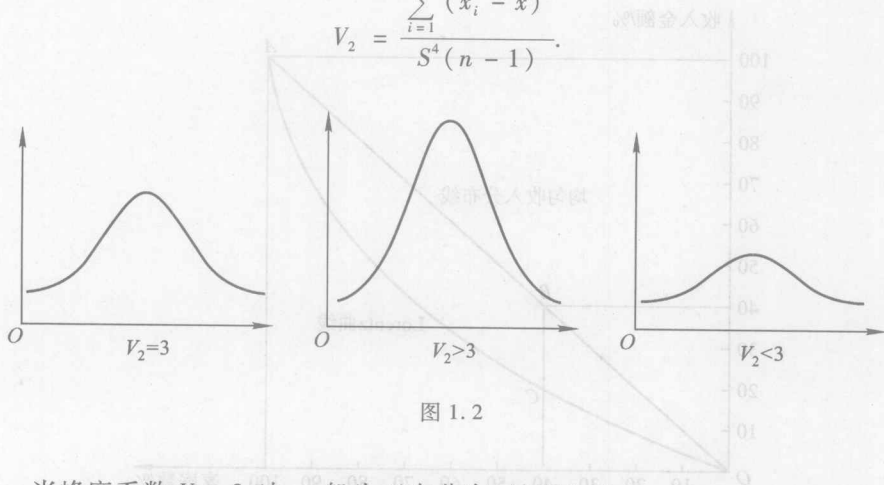


图 1.2

当峰度系数 $V_2 = 3$ 时，一般为正态分布，见图 1.2。

五、累积频数分布

在社会经济调查中，经常会得到频数数据。例如家庭月收入按等级划分时，就会得到每个等级的家庭数，通常将这些数据列在表中或画成直方图。

读者可依收入等级从低到高画出累积频数的直方图。

表 1.1 累积频数分布表

收入等级/元	家 庭 数	
	频 数	累积频数
5 000 ~ 6 000	800	800
6 001 ~ 7 000	700	1 500
7 001 ~ 8 000	500	2 000
8 001 ~ 9 000	300	2 300

在社会经济研究中, M. E. Lorentz(洛伦茨)曲线是累积频数的典型应用. 如果按收入从低到高排列, 各收入等级的家庭的累积数(百分比)为横坐标, 与之相对应的收入的累计(百分比)为纵坐标, 所得到的曲线就是西方经济学中著名的 Lorentz 曲线. 在宏观经济的收入与分配研究中, 就可运用这一描述方法.

图 1.3 中对角线 OA 是均匀收入分布线. 图中 B 点表明在数量上占全体 40% 的家庭在收入上也占 40%. 收入分布不大可能绝对平均, 所以 Lorentz 曲线一般并不是一条直线. 图中 C 点表示从最低收入开始的 40% 的家庭收入的合计还占不到总收入的 20%.

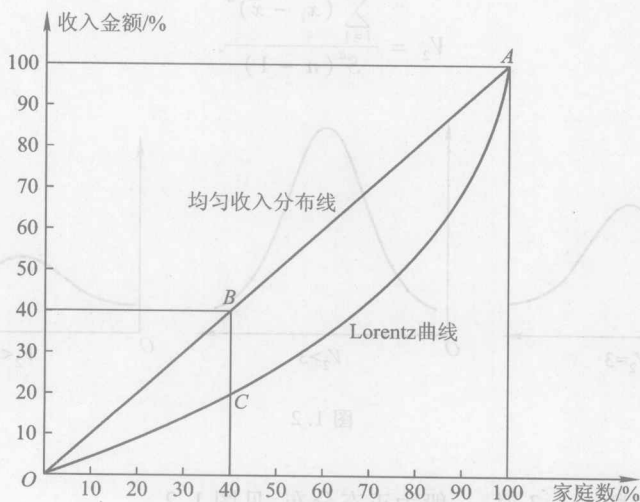


图 1.3

关于累积频数的百分比曲线可拓宽到衡量贫富差距的 Gini(基尼)系数. Gini 系数理论在中国当今的宏观经济研究中非常有用.

§ 1.2 几种重要的概率分布

一、正态分布

在经济研究和工商管理中,有许多随机变量的概率分布都可用正态分布来描述.例如一个城市居民的家庭收入、消费支出,某种股票月收益的百分比,某种产品的某质量特性指标都可近似用正态分布来描述.在实际问题的研究中,可以通过该随机变量的抽样数据的频数直方图与正态概率分布的钟形曲线相比较,来判断该随机变量是否为正态随机变量.

正态随机变量 X 的概率密度函数的形式如下:

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

其中, μ 为正态随机变量 X 的均值, σ^2 为正态随机变量 X 的方差.

通常对具有均值为 μ , 方差为 σ^2 的正态概率分布, 记为 $N(\mu, \sigma^2)$. 于是有正态随机变量 $X \sim N(\mu, \sigma^2)$.

一般来说, 正态分布的密度曲线是以 μ 为中心, 在 μ 的两侧呈对称的形状, 曲线的形状像一个钟的剖面, 故称为钟形曲线. σ 越大, 密度曲线的峰度越低; σ 越小, 密度曲线的峰度越高. 无论参数 μ 和 σ 取何值, 密度曲线下所覆盖的面积均等于 1. 正态分布的密度曲线见图 1.4.

正态分布曲线下, 位于 $\mu \pm \sigma$ 、 $\mu \pm 2\sigma$ 、 $\mu \pm 3\sigma$ 之间的面积分别约占总面积的 68.26%、95.45%、99.73%, 如图 1.5 所示.

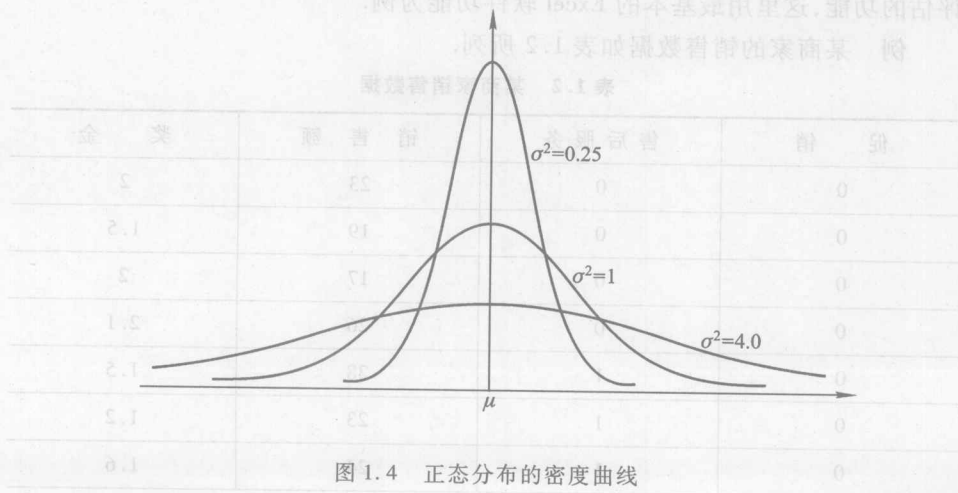


图 1.4 正态分布的密度曲线

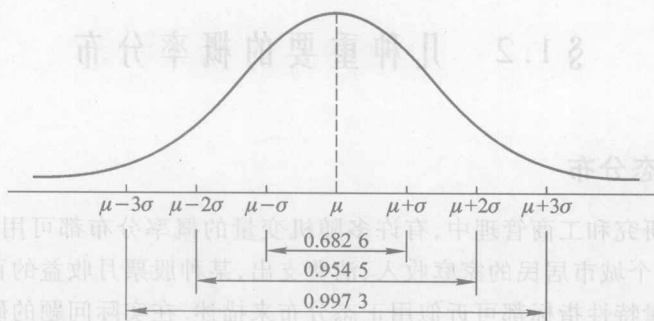


图 1.5

在正态分布的概率密度中,当 $\mu=0, \sigma=1$ 时,称随机变量 X 遵从标准正态分布,记为 $X \sim N(0,1)$ 。

关于正态分布的理论已很完善,数学上也易于处理。此外,当一个经济问题的模型误差是由很多因素构成的时候,总体的分布与正态分布近似。所以,在计量经济学和一些经济问题的建模研究中常假定为正态分布。另外,当总体概率分布为正态分布时,作为从中抽出的样本,其统计量的样本概率分布有 χ^2 分布、 t 分布、 F 分布等。因此,正态分布成为计量经济学乃至统计学中最重要的概念之一。

二、用 Excel 绘制正态概率图进行数据正态性的评估

在实际应用中,经常需要验证数据的正态性,许多统计软件都有数据正态性评估的功能,这里用最基本的 Excel 软件功能为例。

例 某商家的销售数据如表 1.2 所列。

表 1.2 某商家销售数据

促 销	售 后 服 务	销 售 额	奖 金
0	0	23	2
0	0	19	1.5
0	0	17	2
0	0	26	2.1
0	1	28	1.5
0	1	23	1.2
0	1	24	1.6

续表

促 销	售后服务	销 售 额	奖 金
0	1	30	1.8
1	0	26	1.8
1	0	22	1.1
1	0	20	0.9
1	0	30	2.1
1	1	36	2.1
1	1	28	1.21
1	1	30	1.91
1	1	32	2.15
2	0	30	1.8
2	0	23	1.2
2	0	25	1.3
2	0	32	1.92
2	1	48	1.7
2	1	40	1.3
2	1	41	1.2
2	1	46	1.81

在进行方差分析之前,必须检验数据的正态性.选取销售额这一列数据进行正态性评估.其评估步骤为:

第1步:对销售额按照升序排列,记为 $X(j)$.

第2步:计算 $(j-0.5)/24$.

第3步:根据公式 $(j-0.5)/24 = P(Z \leq z_i) = \Phi(z_i)$, 求出正态分位数 z_i .

单击 D2 单元格,选择“插入”→“函数”选项,在出现对话框中:

从“选择类别”窗口中选择“统计”.

从“选择函数”窗口中选择“NORMSINV”选项,选择“确定”,其界面如图 1.6.

再单击 D2 单元格,鼠标指向单元格右下角填充控点,按住鼠标左键往下拖至 D25 单元格,这样,计算出 C2:C24 区域中概率值对应的标准正态分位数如图 1.7 所示.

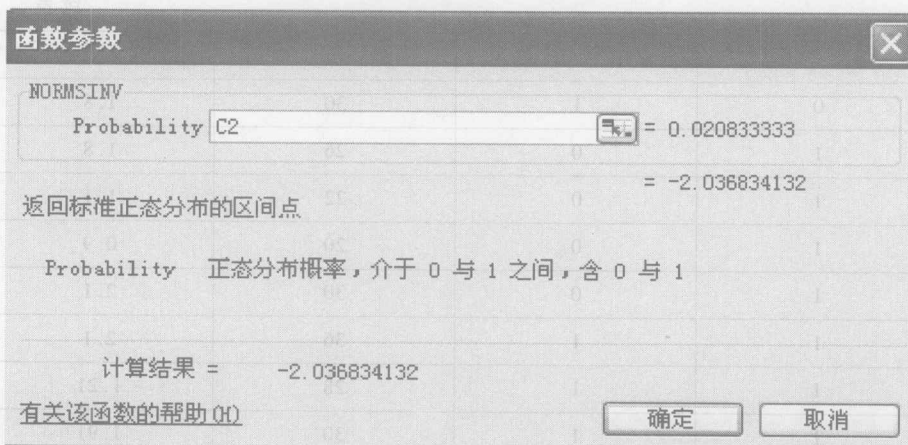


图 1.6

	A	B	C	D
1	j	销售额 $X(j)$	$(j-0.5)/24$	Z_i
2	1	17	0.020833333	-2.036834132
3	2	19	0.0625	-1.534120544
4	3	20	0.104166667	-1.258161561
5	4	22	0.145833333	-1.054472452
6	5	23	0.1875	-0.887146559
7	6	23	0.229166667	-0.741594044
8	7	23	0.270833333	-0.61029461
9	8	24	0.3125	-0.488776411
10	9	25	0.354166667	-0.37409541
11	10	26	0.395833333	-0.264146977
12	11	26	0.4375	-0.157310685
13	12	28	0.479166667	-0.05224518
14	13	28	0.520833333	0.05224518
15	14	30	0.5625	0.157310685

图 1.7

第4步:选择“工具”→“数据分析”,在“分析工具”中选择“回归”,单击“确定”。以 z_i 为纵轴, $X(j)$ 为横轴,绘制标准正态概率图。界面如图1.8所示。

单击“确定”,得到如图1.9所示的标准正态概率图。

其中, $X(j)$ 转化为其对应的百分比排位。从图1.9可看出,由 $(X(j), z_i)$ 形成的点基本上在一条直线附近,可以说该组数据基本上遵从正态分布。

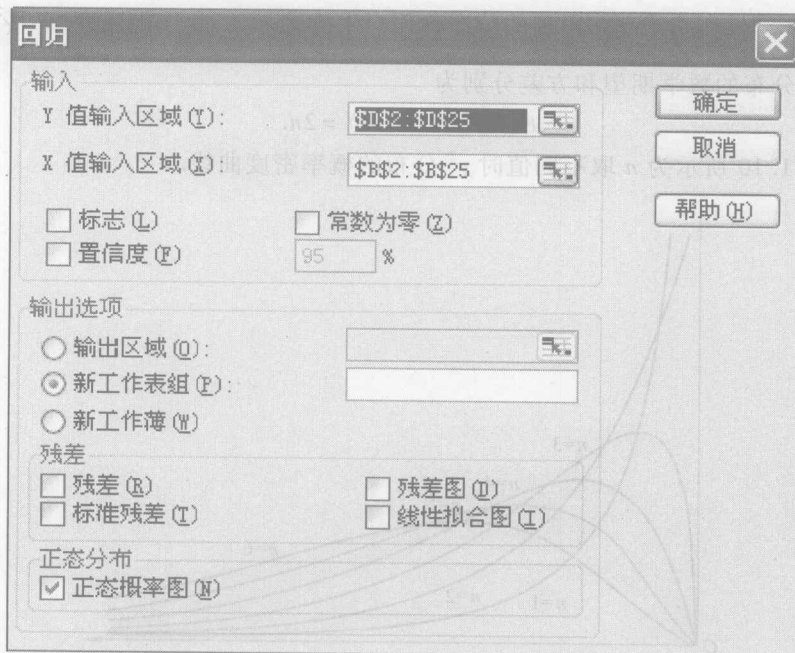


图 1.8

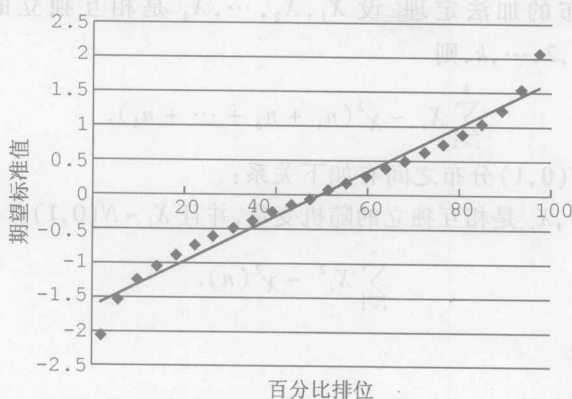


图 1.9 标准正态概率图

三、 χ^2 分布

如果从标准正态分布 $N(0,1)$ 的总体中得到的 n 个随机变量分别为 $X_1,$