



计量经济分析与EViews详解

李 娅 李志鹏 编著



科学出版社



计量经济分析与 EViews 详解

李 娅 李志鹏 编著

本书是云南大学研究生院“精品课程”建设项目成果
并得到其资助出版

科学出版社

北 京

内 容 简 介

本书系统介绍计量经济学的基本理论和常用方法,以经典线性回归模型为主,并引入时间序列和平行数据计量经济学模型,坚持循序渐进,理论联系实际的原则,以各种丰富易懂的例证全面介绍了计量经济学的各种常用回归模型的分析及检验。本书运用 EViews 软件,结合实例分析“无缝式”地展示 EViews 的操作过程,突出计量分析方法应用和 EViews 操作的有机结合,使读者对计量方法的应用与软件的操作有一个全面的了解。

本书可作为本科生及研究生的自学和教学用书,也可作为在经济、统计、金融等领域从事计量分析的工作人员的参考书使用。

图书在版编目(CIP)数据

计量经济分析与 EViews 详解/李娅,李志鹏编著. —北京:科学出版社, 2017

ISBN 977-7-03-050149-3

I. ①计… II. ①李… ②李… III. ①计量经济学-应用软件
IV. ①F224.0-39

中国版本图书馆 CIP 数据核字(2016)第 237284 号

责任编辑:兰 鹏/责任校对:杜子昂
责任印制:张 伟/封面设计:蓝正设计

科 学 出 版 社 出 版

北京东黄城根北街 16 号

邮政编码:100717

<http://www.sciencep.com>

北京康华虎彩印刷有限公司印刷

科学出版社发行 各地新华书店经销

*

2017 年 2 月第 一 版 开本:787×1092 1/16

2017 年 2 月第一次印刷 印张:12 1/4

字数:290 000

定价:42.00 元

(如有印装质量问题,我社负责调换)

前 言

本书是在作者教授的计量经济学教案的基础上编著而形成的。“计量经济学”是一门既难学也难教的课程，如何真正让学生“学懂”并“会用”是作者在教学实践中一直思考的问题。在教学过程中，作者使用过许多国内外经典的教材作为教学用书或参考书，意在取长补短，博采众长，并尝试用最通俗易懂的“讲故事”的方法把计量经济学这门课呈现给学生。多年的教学实践使作者深刻地体会到，要教好计量经济学，最重要的就是要“因人施教”，即明确教学对象、教学目的和教学指导原则。当前，经济学、管理学专业研究生和非计量经济学专业博士研究生学习现代计量经济学的目的是“应用”，而非从事计量经济学理论方法研究。因此，本书在编写过程中始终把应用性和实用性放在首位，着重强调“正确进行经济计量分析”的指导思想。在数学描述方面适当淡化，在详细介绍线性回归模型的数学过程的基础上，各章的重点不是理论方法的数学推导与证明，而是以讲清楚方法、思路为目标，重点放在如何运用各计量经济方法对实际的经济问题进行分析、建模、预测等实际方法的应用和操作上。本书结合 EViews 应用软件，通过系统地讲述应用经济计量分析的相关知识，全面而简洁地介绍了经济计量分析的主要理论和方法，实现经济计量理论与软件的一体化，前后贯通，层次清晰，力求简洁，通俗易懂。

之所以选择 EViews 作为本书的配套教学软件，也是基于对教学对象实施“因人施教”的指导思想。EViews 具有操作简便、界面友好、功能强大等特点，其使用图形交互式用户界面，界面友好且操作简单，可以通过菜单操作和编程两种方式进行分析，使初、中级计量经济学学生较容易地掌握并付诸实践。EViews 提供了与多种应用软件的接口，用户可以方便地把 Excel、SAS、Stata、SPSS 等格式的数据导入 EViews。EViews 拥有统计分析、线性回归分析、非线性单方程模型、联立方程模型、动态回归模型、分布滞后模型、VAR 模型、ARCH/GARCH 模型、离散选择模型、时间序列模型、编程与模拟等分析模块，用户通过 EViews 既可以进行基本的统计和回归分析，也可以完成复杂的计量经济建模。计量经济学是一门实践性要求非常高的课程，对软件的掌握熟练程度的要求非常高，一直强调学生要“干中学”，在实际的数据分析运用中能够切实地解决问题。作者在教学中发现，尽管目前关于计量经济学和 EViews 运用的教材和著作比较多，但是将两者结合起来，尤其是对 EViews 进行“无死角”展示的并不多见，多数教材仅仅就 EViews 的主要步骤进行了展示，对一些中间环节的遗漏造成了学生在实际运用中的知识盲点和运用障碍。本书针对每一个案例对 EViews 的操作进行完整的“无缝”展示，在实际教学中效果极佳，尤其是对于计量经济学零基础的学生和初学者，具有非常好的教学效果。

本书在编写过程中参阅了大量国外有关计量经济学的教材和文献，书中部分案例引自古扎拉蒂、伍德里奇、格林等编著的国外经典教材实例的例题，目的是通过对国外资料的

分析使读者对国外教材有所涉猎，做到与国外教材同步和接轨。同时，本书在编写过程中吸收了一些国内学者的研究成果，在此一并表示感谢。由于作者水平有限，书中难免存在不妥之处，恳请广大读者批评指正。

作者

2016年10月

目 录

第一章 知识准备	1
一、回归分析	1
二、回归模型	1
三、一元线性回归模型	3
四、多元线性回归模型	4
五、随机干扰项	5
第二章 线性回归模型和 OLS	6
一、问题的提出	6
二、解决问题的思路	6
三、解决问题的方法——OLS 估计	7
四、一元线性回归模型的拓展——多元线性回归模型	9
五、高斯-马尔可夫定理	10
六、假设检验	12
七、方差分解	13
八、结构差异检验	15
第三章 异方差与 GLS	34
一、异方差的定义	34
二、异方差的检验	35
三、GLS 法	36
四、异方差的修正	38
第四章 序列相关与 AR	53
一、序列相关的定义	53
二、序列相关的检验	53
三、序列相关的修正	55
第五章 内生解释变量	69
一、引起内生性的原因及其对参数估计的影响	69
二、对内生性的检验	70
三、IV 估计法	71
第六章 多重共线性	87
一、多重共线性的基本概念	87
二、多重共线性产生的原因	87
三、多重共线性的检验方法	88
四、多重共线性的修正	89

第七章 虚拟变量	93
一、虚拟变量定义	93
二、数量因素与变参数模型	94
三、定性因素与变参数模型	95
第八章 离散选择模型	101
一、线性概率模型	101
二、二元离散选择模型	101
三、二元离散选择模型的极大似然估计	104
四、多元离散选择模型	106
第九章 时间序列分析	111
一、平稳时间序列与单位根过程	111
二、协整与误差修正模型	115
第十章 VAR 模型分析	148
一、VAR 模型定义	148
二、VAR 模型的脉冲响应函数和方差分解	149
三、VAR 模型滞后期 k 的选择	152
四、Granger 非因果性检验	153
五、VAR 模型与协整	154
第十一章 面板数据模型分析	169
一、面板数据定义	169
二、面板数据模型分类	169
三、面板数据模型设定的检验方法	172
四、面板数据模型估计方法	173
参考文献	188

第一章 知识准备

本章一句话提示：理解随机干扰项在计量经济学中的重要地位和作用！

一、回归分析

“回归”一词最早来源于“加尔顿普遍回归定律”。加尔顿发现在人口身高的统计中表现出一般规律是：父母高，儿女也高；父母矮，儿女也矮。但是，给定父母的身高，儿女辈的平均身高却趋向于全体人口的平均身高，或者说，是“回归”到全体人口的平均身高。计量经济学的回归分析就是关于研究因变量（被解释变量）与一个或多个自变量（解释变量）之间的因果关系的计算方法和理论。其目的在于通过解释变量（在重复抽样中）的已知或设定值，去估计和预测被解释变量的（总体）均值。

变量间的统计相关关系可以通过相关分析与回归分析来研究。表 1-1 对相关分析和回归分析进行辨析。

表 1-1 相关分析和回归分析

相关分析	回归分析
<p>相关分析主要研究随机变量间的相关形式与相关程度。从变量间相关的形式来看，有线性相关和非线性相关之分。前者往往表现为变量的散点图接近于一条直线。变量间线性相关程度的大小可以通过相关系数来测量，两变量 X 与 Y 之间的总体相关系数为</p> $\rho_{XY} = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)}\sqrt{\text{Var}(Y)}}$ <p>式中，$\text{Cov}(X, Y)$ 是变量 X 和 Y 的协方差，$\text{Var}(X)$、$\text{Var}(Y)$ 分别是 X 和 Y 的方差。</p>	<p>计量经济学的回归分析就是关于变量之间的因果关系的计算方法和理论。其目的在于通过解释变量（在重复抽样中）的已知或设定值，去估计和预测被解释变量的（总体）均值。</p>

回归分析构成计量经济学的方法论基础，其主要内容包括如下三方面。

- (1) 根据样本观察值对计量经济学模型参数进行估计，求得回归方程。
- (2) 对回归方程、参数估计值进行显著性检验。
- (3) 利用回归方程分析、评价及预测。

二、回归模型

(一) 总体回归模型

以不同家庭收入 x_i 和不同消费支出 y_i 为例，两者之间的关系可以表示为

$$y_i = \beta_0 + \beta_1 x_i + u_i$$

式中, y_i 为被解释变量 (因变量); x_i 为解释变量 (自变量); β_0 为常数项 (截距项, 通常未知); β_1 为回归系数 (通常未知); u_i 为随机干扰项, 也称为随机误差项、随机扰动项。

(二) 两分法

为了便于理解, 本书运用两分法的思想, 将总体回归模型分为两个部分。

(1) $E(y_i) = \beta_0 + \beta_1 x_i$ (也称作线性总体回归函数): 对应着计量经济学学习中的一个重要内容“参数估计”。

(2) 随机部分 u_i : 对应着计量经济学学习中的另一个重要内容“假设检验”。与精确的函数关系相比, 回归模型的显著特点是多了随机干扰项 (随机误差项), 计量经济学很多重要而玄妙之处就在于此随机扰动项, 这是计量经济学的学习重点。

(三) 样本回归模型

总体回归模型揭示了所考察总体被解释变量与解释变量间的平均变化规律, 但现实的情况往往是总体回归函数实际上是未知的。一般的做法是从总体样本中取其中的一部分 (抽样), 通过样本的信息来估计总体回归函数。

样本回归模型形式记为

$$y_i = \hat{\beta}_0 + \hat{\beta}_1 x_i + e_i$$

式中, e_i 称为样本残差项, 代表了其他影响 Y 的随机因素的集合, 可以看作 u_i 的估计量 \hat{u}_i 。

$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$, 也称为样本回归函数。

回归分析的主要目的, 就是根据样本回归函数, 估计总体回归函数。

辨析: 总体回归模型、总体回归函数 (方程)、样本回归模型、样本回归函数 (方程), 如下所示。

(1) 总体回归模型: $y_i = \beta_0 + \beta_1 x_i + u_i$ 。

(2) 总体回归方程: $E(y_i) = \beta_0 + \beta_1 x_i$ 。

(3) 样本回归模型: $y_i = \hat{\beta}_0 + \hat{\beta}_1 x_i + e_i$ 。

(4) 样本回归方程: $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$ 。

总体回归方程与样本回归方程、随机干扰项和残差之间的关系如图 1-1 所示。

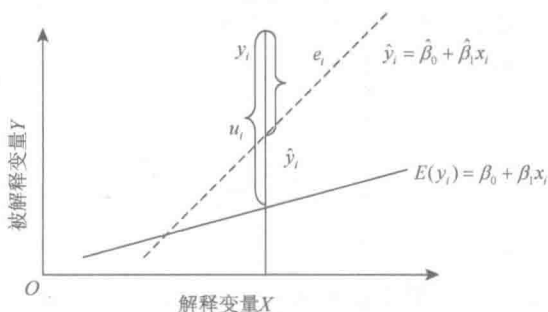


图 1-1 总体回归方程与样本回归方程关系图

三、一元线性回归模型

(一) 一元线性回归模型的模型形式

一元线性回归模型是最简单的计量经济学模型,在模型中只有一个解释变量,其一般形式是

$$Y = \beta_0 + \beta_1 X + u \quad (1-1)$$

式中, Y 为被解释变量; X 为解释变量; β_0 与 β_1 为待估参数; u 为随机干扰项。在有 n 个样本观测点的情况下,式(1-1)也可以写为

$$y_i = \beta_0 + \beta_1 x_i + u_i, \quad i=1,2,\dots,n \quad (1-2)$$

(二) 古典线性回归模型的基本假设

为确保参数估计量具有良好的性质,通常对线性回归模型提出若干基本假设,也称为高斯-马尔可夫假设。按照假设对象的不同,可以分为三个大方面。

(1) 关于回归模型本身的假定:回归模型是正确设定的。

假定:真实模型是

$$Y = \beta_0 + \beta_1 X + u \quad (1-3)$$

有三种情况属于对该假定的违背,如下所示。

- 1) 遗漏了相关的解释变量或者增加了无关的解释变量;
- 2) Y 与 X 间的关系是非线性的;
- 3) β_0 、 β_1 并不是常数。

(2) 关于解释变量 X 的假定。

假定 1: X 是非随机的,即 X 的值是事先固定的。

假定 2: u 和 X 相互独立,即 $\text{Cov}(x_i, u_i) = 0$, 否则分不清 Y 的变化是由 X 引起的,还是由 u 引起的。在重复抽样中, (x_1, x_2, \dots, x_N) 被预先固定下来,即是非随机的,显然,如果解释变量含有随机的测量误差,那么该假定被违背。

假定 3: 对于多元回归模型,假设多个解释变量之间不存在完全共线性。

(3) 关于随机干扰项 u 的假定,也称为“古典模型假设的球形扰动”。

假定 1: 随机误差项 $u_i (i=1, 2, 3, \dots, n)$, 是 n 个随机变量,数学期望为零,即 $E(u_i) = 0$ 。

在模型中,如果能够保证 u_i 中所包含的都是影响 y_i 的微小因素,那么在众多微小因素的作用下,假定就是合理的。

假定 2: 随机误差项的 u_i 方差与 i 无关,为一个常数, $\text{Var}(u_i) = \sigma^2$, 称 u_i 具有同方差性。这一假定的含义是对于任意 u_i , 其分布的方差都是一个常量。反之,称其具有异方差性。

假定 3: 假定不同的随机误差项 u_i 和 u_j 之间相互独立。 $\text{Cov}(u_i, u_j) = 0$, 即所谓的序列不相关假定。

假定 4: u_i 为服从正态分布的随机变量, $u_i \sim N(0, \sigma^2)$ 。

特别强调的是: 对于线性回归模型, 如果扰动项不服从正态分布, 就无法使用小样本普通最小二乘 (ordinary least square, OLS) 法进行统计推断。对于很多估计方法, 如极大似然估计 (maximum likelihood estimate, MLE), 正态分布假定是推导 MLE 的前提。

当以上假定成立时, 在所有线性无偏估计量中, OLS 估计量方差最小。或者说, OLS 估计量是最优线性无偏估计量 (best linear unbiased estimator, BLUE)。这被称为高斯-马尔可夫定理。

说明: 以上线性回归模型的经典假设是一种理想状态, 这些假定同时成立的情况几乎是不可能的。因此: 如何检验这些假定是否成立? 如果一些假定并不成立, 那么 OLS 估计量具有什么性质? 此时应该采取何种估计方法进行修正? 解决以上问题就构成了初级计量经济学的主要内容。

(三) 一元线性回归模型的参数估计

一元线性回归模型的参数估计, 是在一组样本观测值下, 通过一定的参数估计方法, 估计出模型中的未知参数, 即估计样本回归线。常见的估计方法有三种: OLS 估计法及其扩展、MLE 法和广义矩估计法 (generalized method of moments, GMM)。

(四) 一元线性回归模型的统计检验

在利用 OLS 法估计了一元线性回归模型的参数, 并确定了样本回归线后, 要根据经济理论及实际问题中 X 和 Y 的对应关系, 对回归系数的符号、大小及相互关系进行直观判断; 如果上述检验通过, 还需对参数估计值进行统计学检验, 主要包括拟合优度检验、变量的显著性检验以及参数检验的置信区间估计; 另外, 还要进行计量经济学检验, 包括对随机误差项 u_i 的异方差、序列相关检验, 对解释变量的严重多重共线性的检验等, 将在后面章节中一一介绍。

四、多元线性回归模型

(一) 多元线性回归模型的模型形式

多元线性回归模型可以表示为

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_k x_{ik} + u_i, \quad i = 1, 2, \dots, n \quad (1-4)$$

这里 $E(y_i) = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_k x_{ik}$ 为总体多元线性回归方程, 简称总体回归方程。其中, k 表示解释变量个数, β_0 称为截距项, $\beta_1 \beta_2 \cdots \beta_k$ 是总体回归系数。 $\beta_j (j = 1, 2, 3, \dots, k)$ 表示在其他自变量保持不变的情况下, 自变量 x_{ij} 变动一个单位所引起的因变量 y_i 平均变动的数量, 因而也称为偏回归系数。

当给定一个样本 $(y_i, x_{i1}, x_{i2}, \dots, x_{ik}), (i = 1, 2, \dots, n)$ 时, 上述模型可以表示为

$$\begin{cases} y_1 = \beta_0 + \beta_1 x_{11} + \beta_2 x_{12} + \cdots + \beta_k x_{1k} + u_1 \\ y_2 = \beta_0 + \beta_1 x_{21} + \beta_2 x_{22} + \cdots + \beta_k x_{2k} + u_2 \\ y_3 = \beta_0 + \beta_1 x_{31} + \beta_2 x_{32} + \cdots + \beta_k x_{3k} + u_3 \\ \vdots \\ y_n = \beta_0 + \beta_1 x_{n1} + \beta_2 x_{n2} + \cdots + \beta_k x_{nk} + u_n \end{cases} \quad (1-5)$$

y_i 与 x_{ij} 已知, β_i 与 u_i 未知。

其相应的矩阵表达式为

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_n \end{bmatrix}_{(n \times 1)} = \begin{bmatrix} 1 & x_{11} & \cdots & x_{1j} & \cdots & x_{1k} \\ 1 & x_{21} & \cdots & x_{2j} & \cdots & x_{2k} \\ 1 & x_{31} & \cdots & x_{3j} & \cdots & x_{3k} \\ \vdots & \vdots & \cdots & \vdots & \cdots & \vdots \\ 1 & x_{n1} & \cdots & x_{nj} & \cdots & x_{nk} \end{bmatrix}_{(n \times k)} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_k \end{bmatrix}_{(k \times 1)} + \begin{bmatrix} u_1 \\ u_2 \\ u_3 \\ \vdots \\ u_n \end{bmatrix}_{(n \times 1)} \quad (1-6)$$

可以简化为总体回归模型的简化形式:

$$Y = X\beta + u \quad (1-7)$$

(二) 多元线性回归模型中样本容量的问题

1. 最小样本容量

在多元线性回归模型中, 样本容量必须不少于模型中解释变量的数目(包括常数项), 这就是最小样本容量, 即 $n \geq k+1$ 。

2. 满足基本要求的样本容量

一般经验认为, 当 $n \geq 30$ 或者至少 $n \geq 3(k+1)$ 时, 才能说满足模型估计的基本要求。

五、随机干扰项

通过以上一元回归模型和多元回归模型的形式设定可以发现: 与精确的函数关系相比, 回归模型的显著特点是多了随机干扰项(随机误差项)。本书要强调的一个观点就是, 这个被形容为“黑匣子”的“包罗万象”的随机干扰项是现代计量经济学研究的主要内容, 几乎撑起了计量经济学学习的“半壁江山”。无论是经典的计量经济学还是非经典的计量经济学, 很多理论和方法就是从对随机干扰项的设定开始的。理解了这一点, 才能真正理解计量经济学这门学科的内涵。本书将在后续的章节中对随机干扰项进行展开分析。

第二章 线性回归模型和 OLS

本章一句话提示：学会如何在“一堆点里找一条线”！

一、问题的提出

凯恩斯消费理论指出：边际消费倾向（marginal propensity to consume, MPC），即收入每变化一个单位的消费变化率，大于零而小于1， $0 < \text{MPC} < 1$ 。凯恩斯假设了消费与收入之间有正的关系，但没有明确指出两者之间的准确的函数关系。

运用中国 1985~2003 年的经济数据，测算出中国 1985~2003 年的 MPC 约为 0.46，表明在此样本期间，收入每增加一元，平均而言，消费支出将增加 0.46 元。

由此提出的问题是：0.46 是如何测算出来的？它是否合理？

为了弄清楚这个问题，需要从如何收集数据开始，模型的设定、参数的估计和检验等一系列过程的实现，就是计量经济学最基础的任务。

二、解决问题的思路

凯恩斯设定消费与收入之间有正的关系，建设两者之间的函数关系为如下数学模型：

$$Y = \beta_0 + \beta_1 X, \quad 0 < \beta_1 < 1 \quad (2-1)$$

式中， Y 为消费支出（因变量）； X 为收入（自变量）；而被称为模型参数的 β_0 和 β_1 分别是截距和斜率系数。 β_1 是 MPC 的度量。推广到更一般的形式，可描述如下。

假定 Y 与 X 具有近似的线性关系：

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

式中， ε 是随机误差项。对 β_0 、 β_1 这两个参数的值一无所知，任务是利用样本数据去猜测 β_0 、 β_1 的取值。

为了估计得到参数 β_0 和 β_1 的数值，收集了中国经济数据， Y 是消费支出， X 是国内生产总值（gross domestic product, GDP），均以亿元为单位计算。画出这些观察值的散点图（横轴 X ，纵轴 Y ），如图 2-1 所示。

从图 2-1 散点图可以看出，消费支出 Y 与收入 X 具有近似的线性关系，那么就在图中拟合一条直线： $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X$ ，该直线是对 Y 与 X 的真实关系的近似。接下来的问题是，如何确定 $\hat{\beta}_0$ 与 $\hat{\beta}_1$ ，以使图形中的实线最大程度地代表中国 1985~2003 年消费和收入的关系？也就是如何在一堆点里面找到一条能够刻画出其分布规律的线？在一堆样本点中，找

这条“线”的方法就是 OLS 法，其基本思路如图 2-2 所示。

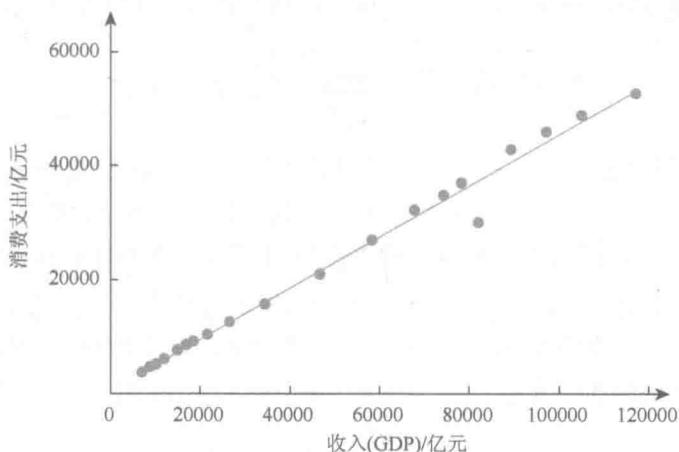


图 2-1 1985~2003 年中国消费支出与收入

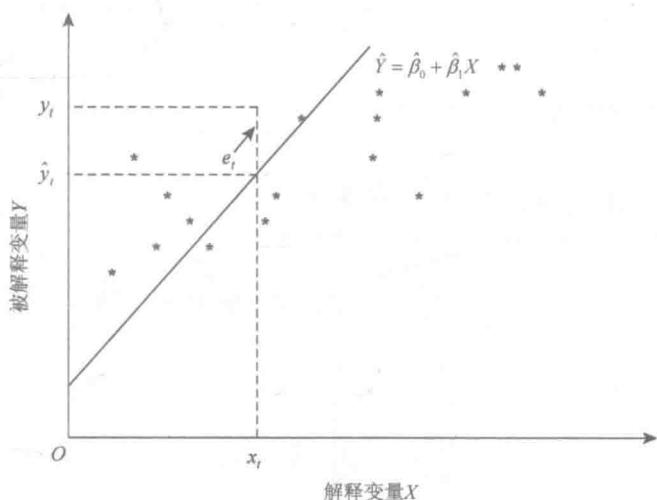


图 2-2 OLS 估计思路图示

三、解决问题的方法——OLS 估计

(一) OLS 估计思路

本书的目标：使拟合出来的直线在某种意义上是最佳的，包括线形、位置。

直观意义：要求估计直线尽可能地靠近各观测点，这意味着应使残差总体上尽可能地小。

方法：要做到这一点，就必须用某种方法将每个点相应的残差加在一起，使其达到最小。

理想的测度：残差平方和，即 $(y_1, y_2, \dots, y_N)'$ 与 $(\hat{y}_1, \hat{y}_2, \dots, \hat{y}_N)'$ 是 N 维空间的两点， $\hat{\beta}_0$ 与 $\hat{\beta}_1$ 的选择应该是这两点的距离最短。这可以归结为求解一个数学问题：

$$\text{Min}_{\hat{\beta}_0, \hat{\beta}_1} \sum_{i=1}^N (y_i - \hat{y}_i)^2 = \text{Min}_{\hat{\beta}_0, \hat{\beta}_1} \sum_{i=1}^N (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 \quad (2-2)$$

由于 $y_i - \hat{y}_i$ 是残差的定义，因此上述获得 $\hat{\beta}_0$ 与 $\hat{\beta}_1$ 的方法即是 $\hat{\beta}_0$ 与 $\hat{\beta}_1$ 的值应该使残差平方和最小。也就是说，给定 x_i ，看起来 y_i 与 \hat{y}_i 越近越好（最近距离是 0）。然而，当选择拟合直线使得 y_i 与 \hat{y}_i 相当近的时候， y_j 与 \hat{y}_j 的距离也许变远了，因此存在一个权衡。一种简单的权衡方式是，给定 x_1, x_2, \dots, x_N ，拟合直线的选择应该使 y_1 与 \hat{y}_1 、 y_2 与 \hat{y}_2 、 \dots 、 y_N 与 \hat{y}_N 的距离的平均值是最小的。距离是一个绝对值，数学处理较为麻烦，因此，把第二种思考方法转化求解数学问题：

$$\text{Min}_{\hat{\beta}_0, \hat{\beta}_1} \sum_{i=1}^N (y_i - \hat{y}_i)^2 / N = \text{Min}_{\hat{\beta}_0, \hat{\beta}_1} \sum_{i=1}^N (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 / N \quad (2-3)$$

由于 N 为常数，式 (2-2) 和式 (2-3) 对于求解 $\hat{\beta}_0$ 与 $\hat{\beta}_1$ 的值是无差异的。

(二) OLS 估计式的推导

定义 $Q = \sum_{i=1}^N (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$ ，利用一阶条件，有

$$\frac{\partial Q}{\partial \hat{\beta}_0} = \sum_{i=1}^N 2(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)(-1) = 0$$

则

$$\sum_{i=1}^N (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0 \quad (2-4)$$

可推出：

$$\sum_{i=1}^N \hat{\varepsilon}_i = 0 \quad (2-5)$$

由式 (2-4) 有

$$\bar{y} = \hat{\beta}_0 + \hat{\beta}_1 \bar{x}$$

式中， $\bar{y} = \frac{1}{N} \sum_{i=1}^N y_i$ ， $\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$ 。方程 (2-4) 与方程 (2-5) 被称为正规方程，把 $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$ 代入方程 (2-5)，有

$$\sum_{i=1}^N [y_i - \bar{y} - \hat{\beta}_1 (x_i - \bar{x})] x_i = 0$$

则

$$\hat{\beta}_1 = \frac{\sum_{i=1}^N (y_i - \bar{y})x_i}{\sum_{i=1}^N (x_i - \bar{x})x_i} \quad (2-6)$$

上述获得 $\hat{\beta}_0$ 、 $\hat{\beta}_1$ 的方法就是 OLS 法。

四、一元线性回归模型的拓展——多元线性回归模型

推广到总体多元回归模型是

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_k x_{ik} + u_i, \quad i=1,2,\dots,n \quad (2-7)$$

即

$$\begin{cases} \sum_{i=1}^n (\hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \hat{\beta}_2 x_{i2} + \cdots + \hat{\beta}_k x_{ik}) = \sum_{i=1}^n y_i \\ \sum_{i=1}^n (\hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \hat{\beta}_2 x_{i2} + \cdots + \hat{\beta}_k x_{ik}) x_{i1} = \sum_{i=1}^n y_i x_{i1} \\ \sum_{i=1}^n (\hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \hat{\beta}_2 x_{i2} + \cdots + \hat{\beta}_k x_{ik}) x_{i2} = \sum_{i=1}^n y_i x_{i2} \\ \vdots \\ \sum_{i=1}^n (\hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \hat{\beta}_2 x_{i2} + \cdots + \hat{\beta}_k x_{ik}) x_{ik} = \sum_{i=1}^n y_i x_{ik} \end{cases}$$

如果用矩阵来描述，首先定义下列向量与矩阵：

$$\begin{bmatrix} n & \sum_{i=1}^n x_{i1} & \cdots & \sum_{i=1}^n x_{ik} \\ \sum_{i=1}^n x_{i1} & \sum_{i=1}^n x_{i1}^2 & \cdots & \sum_{i=1}^n x_{i1} x_{ik} \\ \vdots & \vdots & & \vdots \\ \sum_{i=1}^n x_{ik} & \sum_{i=1}^n x_{ik} x_{i1} & \cdots & \sum_{i=1}^n x_{ik}^2 \end{bmatrix} \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_k \end{bmatrix} = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ x_{11} & x_{12} & \cdots & x_{1n} \\ \vdots & \vdots & & \vdots \\ x_{k1} & x_{k2} & \cdots & x_{kn} \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad (2-8)$$

模型的矩阵表示：

$$(X'X)\hat{\beta} = X'Y \quad (2-9)$$

求解最小化问题： $\text{Min}_{\hat{\beta}} (Y - X\hat{\beta})'(Y - X\hat{\beta})$ ，有

$$\begin{aligned} \frac{\partial [(Y - X\hat{\beta})'(Y - X\hat{\beta})]}{\partial \hat{\beta}} &= \frac{\partial [(Y' - \hat{\beta}'X')(Y - X\hat{\beta})]}{\partial \hat{\beta}} \\ &= \frac{\partial [Y'Y - Y'X\hat{\beta} - \hat{\beta}'X'Y + \hat{\beta}'X'X\hat{\beta}]}{\partial \hat{\beta}} = 0 \end{aligned}$$

而根据矩阵微分的知识, 有

$$\frac{\partial(Y'Y)}{\partial\hat{\beta}}=0, \quad \frac{\partial(Y'X\hat{\beta})}{\partial\hat{\beta}}=(Y'X)'=X'Y$$

$$\frac{\partial(\hat{\beta}'X'Y)}{\partial\hat{\beta}}=X'Y, \quad \frac{\partial(\hat{\beta}'X'X\hat{\beta})}{\partial\hat{\beta}}=X'X\hat{\beta}+(\hat{\beta}'X'X)'=2X'X\hat{\beta}, \quad X'Y=X'X\hat{\beta}$$

则

$$\hat{\beta}=(X'X)^{-1}(X'Y) \quad (2-10)$$

这就是多元回归模型的 OLS 估计。

五、高斯-马尔可夫定理

当高斯-马尔可夫假定成立时, 在所有线性无偏估计量中, OLS 估计量方差最小。或者说, OLS 估计量是最优线性无偏估计量。这被称为高斯-马尔可夫定理。

(一) OLS 估计量是线性估计量

所谓 OLS 估计量是线性估计量, 是指它能够被表示为 y_i 的线性函数。例如:

$$\hat{\beta}_1 = \left[\frac{x_i - \bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2} \right] y_i = \sum_{i=1}^n k_i y_i \quad (2-11)$$

式中, k_i 是非随机的。

(二) OLS 估计量具有无偏性

$$E(\hat{\beta}_1) = \beta_1, \quad E(\hat{\beta}_0) = \beta_0$$

证明 $E(\hat{\beta}_1) = \beta_1$ 。

$$\hat{\beta}_1 = \sum_{i=1}^n k_i y_i = \sum_{i=1}^n k_i (\beta_0 + \beta_1 x_i + u_i) = \beta_0 \sum_{i=1}^n k_i + \beta_1 \sum_{i=1}^n k_i x_i + \sum_{i=1}^n k_i u_i$$

$$E(\hat{\beta}_1) = E\left(\beta_1 + \sum_{i=1}^n k_i u_i\right) = \beta_1 + \sum_{i=1}^n k_i E(u_i) = \beta_1$$

而 $\sum_{i=1}^n k_i = 0$, $\sum_{i=1}^n k_i x_i = 1$ 。因此在重要假定 $E(u_i) = 0$ 下, 有 $E(\hat{\beta}_1) = \beta_1$

$$E(\hat{\beta}_0) = E\left(\beta_0 + \sum_{i=1}^n w_i u_i\right) = E(\beta_0) + \sum_{i=1}^n w_i E(u_i) = \beta_0$$

① 为了保证 $(X'X)^{-1}$ 的存在, OLS 法假设 X 列满秩, 即解释变量不是完全共线的。