

计量经济学

Ender

2023 年 11 月 17 日

EconometricsNote © 2023 by Ender23333 is licensed under CC BY-NC 4.0. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc/4.0/>.

目录

第一章 计量经济学诸论	1
第二章 一元线性回归模型	2
2.1 回归函数	3
2.2 一元线性投影模型的总体参数	10
2.3 一元线性投影模型的样本估计量	17
2.4 一元 OLS 估计量的 Gauss-Markov 定理	23
2.5 一元 OLS 估计量下的拟合值、残差	34
2.6 一元 OLS 估计量下的方差估计量	43
2.7 一元 OLS 估计量的区间估计与假设检验	45
2.8 一元 OLS 估计量的渐进性质	50
2.9 代码	50
第三章 多元线性回归模型	54
3.1 多元线性回归模型的矩阵符号	54
3.2 非简单随机样本的模型假设	56
3.3 多元线性投影模型的总体参数	59
3.4 多元线性投影的样本估计量	62
3.5 多元线性投影模型的分块回归	66
3.6 多元 OLS 估计量下 Gauss-Markov 定理	66
3.7 正交投影与 OLS 估计量	69
3.8 多元 OLS 估计量下拟合值、残差	75
3.9 Leave-One-Out 回归与预测误差	76
3.10 拟合的度量	79
3.11 多元 OLS 估计量下的方差估计量	81
3.12 多元正态回归模型及其区间估计	86
3.13 多元正态回归模型的假设检验	90

3.14 多元线性回归模型的预测	102
3.15 虚拟变量的 OLS 估计量	102
3.16 代码	102
第四章 概率渐进理论	104
4.1 依概率收敛	104
4.2 大数定律与连续映射定理	107
4.3 矩生成函数	113
4.4 渐进分布与中心极限定理	118
4.5 随机过程的渐进性质	119
第五章 OLS 估计量的渐进性质	120
5.1 遍历平稳下 OLS 估计量的渐进性质	121
5.2 代码	121
概念索引	123
模型索引	130

回归模型

$$Y = E(Y|X) + e$$

↑
CEF

↑
回归误差/
CEF误差

$$E(e|X) = 0 \quad Ee = 0$$

$E(Y|X)$ 是未知的

$E(Y|X)$ 为线性函数



线性 CEF 模型

$$Y = X^T \beta + e$$

↑
线性 CEF

$$E(e|X) = 0$$

线性 CEF 依然未知, β 未知

线性投影模型

$$Y = X^T \beta^* + \varepsilon$$

↑
投影误差

性质: $E(X\varepsilon) = 0$ 若有常数项 $E\varepsilon = 0$

Best Line Predictor $P(Y|X) = X^T \beta^*$

依据 MSE 最小得到

$$\beta^* = (E(X^T X))^{-1} E(XY)$$

总体

给定 Random Sample

样本

$$\{(Y_i, \vec{X}_i), \dots, (Y_n, \vec{X}_n)\} \text{ i.i.d.}$$

OLS 估计量 $\hat{\beta}_{OLS} = \underset{\beta \in \mathbb{R}^k}{\operatorname{argmin}} \hat{S}(\beta)$

$\hat{S}(\beta) = \frac{1}{n} \sum_{i=1}^n (Y_i - \vec{X}_i^T \beta)^2$ 为 MSE 的矩估计

第一章 计量经济学诸论

第二章 一元线性回归模型

回归模型源自于数理统计, 其被用于研究具有相关关系的变量, 对于计量经济学而言, 回归分析则是进一步探讨理论上具有因果关系的变量间的关系 (由因果性而具有相关性). 本章将介绍最为简单的一元线性回归模型, 之后的章节则在理解了一元回归模型后, 将进一步拓展至多元的情形.

依据概率论与数理统计的范式, 本章以及之所研究的感兴趣的变量, 均是**随机变量**, 相应地本笔记使用**斜体大写字母** (例如 X, Y, M, D) 或**首字母大小的英文** (例如 $Wage, Post, Treat$) 表示^[1], 而对于多个变量 X_1, \dots, X_n , 本笔记使用**非斜体粗体大写字母**记有随机向量 $\mathbf{X} = (X_1, \dots, X_n)^T$. 特别地, 对于经过函数处理后的变量, 则采用函数名 + 变量的方式表示 (例如 $\ln X$ 或对数工资 $\ln Wage$). 这样表示的变量, 对应了研究问题的总体. 对于研究变量所对应的样本, 本笔记使用形如 \mathbf{X}, \mathbf{Y} 的矩阵符号^[2]记录全部样本的随机变量, 单个样本的随机变量记为 Y_i, X_i 的形式, 而具体的观测值则使用形如 y_1, y_2, \dots, y_n 的**斜体小写字母**来表示.

计量经济学遵循了数理统计的方法, 是利用样本去推断总体, 因而就本笔记所涉及的样本, 不作特别说明时均假定样本是通过**概率抽样**取得的, 进一步地则是假定为简单随机抽样的样本, 故本笔记所探讨的样本个体都是**独立同分布**的. 具体地, 对于样本容量为 K 个的样本, 在仅有一个变量 X 的情况下则记第 k 个样本) 个体 ($k = 1, 2, \dots, K$) 为随机变量 X_k , 其抽样的观测值为 x_k .

明确了符号的使用可以减少后续笔记中出现误解的情况. 由于所研究的变量均是随机变量, 本笔记也采用概率论的符号来表示随机变量的数值特征, 利用对于变量工资 $Wage$ 和受教育年限 $School$, 则

$$\mathbb{E}(Wage|School)$$

意味着给定受教育年限 $School$ 下工资 $Wage$ 的条件期望.

^[1]在 \LaTeX 的数学环境中, 通过斜体命令 `\textit{}` 输入.

^[2]在 \LaTeX 的数学环境中, 命令 `\mathbf{X}` 输出 \mathbf{X} , 而命令 `\boldsymbol{X}` 输出 \mathbf{X} .

2.1 回归函数

倘若在学习概率论前, 读者便已有接触过回归分析, 则很可能错误的将线性回归模型 $Y = \mathbf{X}\beta + \varepsilon$ 等同于“回归函数”. 本节内容即是要纠正这一错误并介绍概率论中回归函数的概念.

随机变量是对不确定性的一种刻画方式^[3], 只有试验发生了, 才能知道结果得到随机变量的观测值. 对于随机变量具体的取值, 可以用概率去进行精细的刻画, 然而在实际应用中即是知道随机变量取值的概率了, 多数情况下便是依据大数定律去理解 (例如重复掷塞子 100 次, 大约有 50 次朝上的点数是奇数), 仍旧不是一个“确切”的结果——只有试验发生了, 我们才能够观测到结果.

幸运的是, 概率论告诉了我们随机变量的数值特征, 例如: 期望、方差、高阶矩等. 只要知道了随机变量的概率分布^[4], 通常而言我们是可以考虑其完全“确定”的数值特征^[5], 而且数值特征亦能够帮助我们了解随机变量所具有的性质. 例如知道了随机变量的期望 μ 和标准差 σ , 依据 Markov 不等式, 便能得到随机变量最可能出现的取值情况 (尽管这个例子还是带有不确定性, 但还是帮助我们理解了数字特征的“确定性”作用).

现考虑两个连续性随机变量 X 和 Y , 要研究二者具有的关系, 用概率论的方法则是考虑联合分布 (X, Y) 的联合密度函数 $f_{XY}(x, y)$, 这样 X 和 Y 的概率密度函数 (亦边缘密度函数) 为

$$f_X(x) = \int_{-\infty}^{+\infty} f_{XY}(x, y) dy, \quad f_Y(y) = \int_{-\infty}^{+\infty} f_{XY}(x, y) dx,$$

但这只包含了单独的 X 或 Y 的信息. 一种能够得知 X 和 Y 间关系的方式, 即是依据条件概率的思想, 计算条件期望 $\mathbb{E}(Y|X)$, 用计量经济学的术语, 则称为回归函数.

定义 2.1.1 (条件期望函数, 总体回归函数). 设随机变量 X 和 Y , 则称 $\mathbb{E}(Y|X) = f(X)$ 为 Y 关于 X 的**条件期望函数** (*conditional expectation function*). 在计量经济学中, 则称 $\mathbb{E}(Y|X)$ 为**总体回归函数** (*population regression function*), 简称**回归函数**.

^[3]在经济学中, 文献 (cite) 将风险 (risk) 定义为有概率分布的随机变量, 而不确定性 (uncertainty) 则是不具有概率分布的随机变量. 但进来的研究也未有遵循这个定义 (cite).

^[4]一般而言, 概率论中具有概率密度函数 (probability density function, PDF) 的随机变量称为连续性随机变量, 然而更多变量不是连续性随机变量, 亦不属于离散型随机变量 (例如随机变量 $X \sim P(\lambda)$, 随机变量 Y 服从标准正态分布, 那么 $X + Y$ 既不是连续性随机变量, 又不是离散型随机变量.) 本笔记为了方便讨论, 所研究的变量均为连续性随机变量, 即其具有 PDF.

^[5]具有 PDF 的连续性随机变量, 其数值特征并非一定存在. 例如服从 Cauchy 分布的随机变量不存在数学期望 (广义积分是发散的). 为了便于讨论, 计量经济学所研究的变量通常都是假定其数字特征是存在的.

对于给定的 $X = x$, 条件期望函数 $\mathbb{E}(Y|X = x)$ 利用了 $X = x$ 提供的信息, 是在此之下随机变量 Y 所具有的期望, 求期望后则消去 Y 的不确定性, 但会保留有 $X = x$, 因而条件期望函数是关于随机变量 X 的观测值 x 的函数. 假设随机变量 X, Y 均为连续性随机变量且存在联合密度函数 $f_{XY}(x, y)$, 那么 Y 关于 X 的条件期望函数为

$$\mathbb{E}(Y|X) = \int_{-\infty}^{+\infty} y f_{Y|X}(y|x) dy = \int_{-\infty}^{+\infty} y \frac{f_{XY}(x, y)}{f_X(x)} dy,$$

其中

$$f_{Y|X}(y|x) = \frac{f_{XY}(x, y)}{f_X(x)}$$

为 Y 关于 X 的条件密度函数.

在没有歧义的情形下, 本笔记将条件期望函数 $\mathbb{E}(Y|X = x)$ 记为 $m(x)$.

为什么计量经济学要以条件期望函数作为回归函数? 在解释这个问题前, 我们首先需要回顾一下条件期望的相关性质.

定理 2.1.1 (简单迭代期望定律, simple law of iterated expectations). 设随机变量 X 和 Y , 若 $\mathbb{E}Y < \infty$, 则有 $\mathbb{E}(\mathbb{E}(Y|X)) = \mathbb{E}Y$.

定理 (2.1.1)^[6]给出了非条件期望一种计算方法, 即是条件期望的期望等于非条件期望. 具体地, 将期望算符 \mathbb{E} 所对应的随机向量标准, 则有 $\boxed{\mathbb{E}_X(\mathbb{E}_Y(Y|X)) = \mathbb{E}Y}$. 可以直观理解, 首先是对 Y 求条件期望, 得到 $\mathbb{E}(Y|X)$ 是关于随机变量 X 的函数, 亦是一个随机变量, 故迭代后即是对 X 求期望. 下面给出定理 (2.1.1) 的证明.

证明. 下面的证明用到了期望的性质

$$\mathbb{E}g(X) = \int_{-\infty}^{+\infty} g(x) f_X(x) dx, \quad \mathbb{E}g(Y) = \sum_{k=1}^{\infty} g(Y) \mathbb{P}(Y = y_k).$$

设随机变量 X 和 Y 均为连续性随机变量, 则有

$$\begin{aligned} & \mathbb{E}(\mathbb{E}(Y|X)) \\ &= \int_{-\infty}^{+\infty} \left(\int_{-\infty}^{+\infty} y \frac{f_{XY}(x, y)}{f_X(x)} dy \right) f_X(x) dx = \int_{-\infty}^{+\infty} \left(\int_{-\infty}^{+\infty} y f_{XY}(x, y) dy \right) \frac{f_X(x)}{f_X(x)} dx \\ &= \int_{-\infty}^{+\infty} \left(\int_{-\infty}^{+\infty} y f_{XY}(x, y) dy \right) dx \xrightarrow{\text{交换积分次序}} \int_{-\infty}^{+\infty} y \left(\int_{-\infty}^{+\infty} f_{XY}(x, y) dx \right) dy \\ &= \int_{-\infty}^{+\infty} y f_Y(y) dy = \mathbb{E}Y. \end{aligned}$$

[6] “简单迭代期望定律”也可称“迭代期望定律”, 强调“简单”这一限定词是为了同定理 (2.1.2) 区别.

再在设随机变量 X 和 Y 为离散型随机变量, 则有

$$\begin{aligned}
 \mathbb{E}(\mathbb{E}(Y|X)) &= \sum_{j=1}^{\infty} \mathbb{E}(Y|X=x_j) \mathbb{P}(X=x_j) \\
 &= \sum_{j=1}^{\infty} \left(\left(\sum_{k=1}^{\infty} y_k \mathbb{P}(Y=y_k|X=x_j) \right) \mathbb{P}(X=x_j) \right) \\
 &= \sum_{j=1}^{\infty} \left(\left(\sum_{k=1}^{\infty} y_k \frac{\mathbb{P}(Y=y_k, X=x_j)}{\mathbb{P}(X=x_j)} \right) \mathbb{P}(X=x_j) \right) \\
 &= \sum_{j=1}^{\infty} \left(\sum_{k=1}^{\infty} y_k \mathbb{P}(Y=y_k, X=x_j) \right) \xrightarrow{\text{交换求和次序}} \sum_{k=1}^{\infty} \left(\sum_{j=1}^{\infty} y_k \mathbb{P}(Y=y_k, X=x_j) \right) \\
 &= \sum_{k=1}^{\infty} y_k \mathbb{P}(Y=y_k) = \mathbb{E}Y.
 \end{aligned}$$

上面的证明中无论积分还是求和, 都涉及到了交换次序, 这实际上需要严格的数学条件, 但本笔记的证明略过这些条件. \square

从定理 (2.1.1) 中, 尤其是离散形式的证明中可以知道, 简单迭代迭代期望定律实际上可以理解为“先求组内均值, 再将各均值加权相加”. 考虑工资和性别, 一个例子则是

$$\begin{aligned}
 \mathbb{E}(\text{工资}) &= \mathbb{E}(\mathbb{E}(\text{工资}|\text{性别})) \\
 &= \mathbb{E}(\text{工资}|\text{性别}=\text{男}) \mathbb{P}(\text{性别}=\text{男}) + \mathbb{E}(\text{工资}|\text{性别}=\text{女}) \mathbb{P}(\text{性别}=\text{女}),
 \end{aligned}$$

即先计算不同性别的平均工资, 然后以性别比例作为权重计算得到平均工资.

定理 2.1.2 (迭代期望定律, law of iterated expectations). 设随机变量 Y, X_1 和 X_2 , 若 $\mathbb{E}Y < \infty$, 则有 $\mathbb{E}(\mathbb{E}(Y|X_1, X_2)|X_1) = \mathbb{E}(Y|X_1)$.

显然, 定理 (2.1.1) 是定理 (2.1.2) 的特例, 这只需有 $X_2 = X$ 及 X_1 为常数. 对于连续性随机变量, 定理 (2.1.2) 中 $\mathbb{E}(Y|X_1, X_2)$ 的一般形式为

$$\mathbb{E}(Y|X_1, X_2) = \int_{-\infty}^{+\infty} y f_{Y|X_1, X_2}(y|x_1, x_2) dy = \int_{-\infty}^{+\infty} y \frac{f_{X_1 X_2 Y}(x_1, x_2, y)}{f_{X_1, X_2}(x_1, x_2)} dy,$$

其中 $f_{X_1 X_2 Y}(x_1, x_2, y)$ 为 (X_1, X_2, Y) 的联合密度函数, $f_{X_1, X_2}(x_1, x_2)$ 为 (X_1, X_2) 的联合密度函数. 一种对定理 (2.1.2) 的精辟概述是: *the smaller information set wins*, 即仅用 X_1 的信息不能得到拥有 (X_1, X_2) 的信息的 $\mathbb{E}(Y|X_1, X_2)$, 只能得到拥有 X_1 的信息的 $\mathbb{E}(Y|X_1)$.

定理 2.1.3 (条件定理, conditioning theorem). 设随机变量 X 和 Y , 若 $\mathbb{E}Y < \infty$, 则有

$\mathbb{E}(g(X)Y|X) = g(X)\mathbb{E}(Y|X)$. 如果 $\mathbb{E}|g(X)| < \infty$, 则有

$$E(g(X)Y) = \mathbb{E}(g(X)\mathbb{E}(Y|X)).$$

定理 (2.1.3) 的前一个命题很好理解, 即是给定随机变量 X 的情况下 $g(X)$ 相当于已知的常数, 依据期望算符 \mathbb{E} 的线性性质便可以提出 $g(X)$. 而后一个命题则是给出了 $E(g(X)Y)$ 的计算方法. $g(X)Y$ 的依然是随机变量, 若求 $g(X)Y$, 先求其概率密度函数并不总是容易, 但利用定理 (2.1.3) 则有

$$\mathbb{E}(g(X)Y) = \mathbb{E}_X(g(X)\mathbb{E}_Y(Y|X)) = \mathbb{E}_X(\mathbb{E}_Y(g(X)Y|X)).$$

有了定理 (2.1.1)(2.1.2)(2.1.3), 我们便可以回答为何用条件期望函数作为回归函数. 计量经济学的核心目的是进行因果推断, 以 $\mathbb{E}(Y|X)$ 作为总体回归函数, 其意味着在已知信息 X 的情况下用条件期望去解释 Y , 或者说用 $\mathbb{E}(Y|X)$ 去预测 (predict) Y . 我们将证明, 以回归方差为标准, 条件期望函数 $\mathbb{E}(Y|X)$ 即是 X 对于 Y 的最优预测 (best predict).

定义 2.1.2 (条件期望函数误差, 回归误差). 设随机变量 X 和 Y , 则称

$$e = Y - \mathbb{E}(Y|X)$$

为条件期望函数误差 (CEF error), 又称为回归误差 (regression error).

依据定义 (2.1.2), 移项则有

$$Y = \mathbb{E}(Y|X) + e, \quad (2.1)$$

式 (2.1) 即是最一般化的计量经济学模型, 我们称之 Y 对 X ^[7] 的为回归方程 (regression equation). 我们常将随机变量 Y 称为被解释变量 (explained variable)、因变量 (dependent variable) 或回归子 (regressand), 将随机变量 X 称为解释变量 (explanatory variable)、自变量 (independent variable) 或回归元 (regressor), 而回归误差 e 又叫随机扰动项 (stochastic disturbance) 或随机误差项 (stochastic error).

模型 2.1.1 (回归模型). 满足以下设定的模型称为回归模型 (regression model).

$$Y = \mathbb{E}(Y|X) + e$$

^[7] 这里是依据英文 “Y on X”, 尽管依据中文含义, 表述 “X 对 Y” 可能更合理. 本笔记后续内容中 “Y 对 X” 均理解为 “Y on X”.

回归误差 e 则是由随机向量 (X, Y) 的联合分布所决定的一个随机变量, 下面我们回归误差 e 的性质.

定理 2.1.4 (回归误差 e 的误差). 回归误差 e 有如下性质:

- a. $\mathbb{E}(e|X) = 0$;
- b. $\mathbb{E}e = 0$;
- c. $\forall r > 1$, 若 $\mathbb{E}|Y|^r < \infty$, 则 $\mathbb{E}|e|^r < \infty$;
- d. $\forall h(X)$, 若 $\mathbb{E}|h(X)e| < \infty$, 则 $\mathbb{E}|h(X)e| = 0$.

证明. 现仅证明性质 a 和 b, 余下证明见 Bruce(cite) 的 Section 2.33.

对于性质 a, 即有

$$\mathbb{E}(e|X) = \mathbb{E}(Y - \mathbb{E}(Y|X)|X) = \mathbb{E}(Y|X) - \mathbb{E}(\mathbb{E}(Y|X)|X) = \mathbb{E}(Y|X) - \mathbb{E}(Y|X) = 0.$$

对于性质 b, 依据定理 (2.1.1) 有 $\mathbb{E}(e|X) = \mathbb{E}(\mathbb{E}(e|X)) = \mathbb{E}0 = 0$. □

在概率论中, 如果随机变量 X 和 Y 满足有 $\mathbb{E}(Y|X = x) = \mathbb{E}Y$, 则称 Y 对于 X **均值独立** (mean independence)^[8], 或 Y 均值独立于 X . 反过来, $\mathbb{E}(X|Y) = \mathbb{E}X$ 则称 X 均值独立于 Y . 容易知道, Y 均值独立于 X 的必要条件是 $\mathbb{E}(Y|X) = 0$, 同理性质 b 的证明, 即得 $\mathbb{E}(Y|X) = 0 = \mathbb{E}$. 因此, 回归误差 e 均值独立于解释变量 X .

回归误差 e 是随机变量, 其以差值的形式刻画了条件期望函数 $\mathbb{E}(Y|X)$ 对被解释变量 Y 的解释或预测的偏离程度, $\mathbb{E}e = 0$ 表明在平均意义上 $\mathbb{E}(Y|X)$ 对 Y 的偏差是不存在的, 仅给定任意的随机变量 $X = x$ 时 $\mathbb{E}(e|X = x) = 0$ 依然说明了条件期望函数的预测在平均意义上是准确的. 回归误差 e 均值独立于解释变量 X , 则进一步表明了回归误差 e 与 X 是线性无关的.

如果仅从 $\mathbb{E}e = 0$ 的角度说明条件期望函数具有最优的预测性质是不过的, 因为在期望的意义上有正有负的回归误差 e 会出现正负相抵的情况 (但 $\mathbb{E}e = 0$ 依然是条件期望函数的重要性质). 为了避免分析中这个问题. 一种简单的方法即是研究回归误差的平方, 由此来定义最优预测.

定义 2.1.3 (最优预测量). 设随机变量 X 和 Y , 记任意的 $g(X)$ 为 Y 的一个预测函数, 我们称使得 $\mathbb{E}(Y - g(X))^2$ 最小的函数为给定 X 下 Y 的**最优预测量** (best predictor).

这样 $\mathbb{E}(Y - g(X))^2$ 是预测误差平方的平均值, 实际上, 这里是借用数理统计中均

^[8]可以证明, 随机变量的独立性是随机变量均值独立的必要条件, 随机变量均值独立是随机变量线性无关的必要条件, 即有 互相独立 \Rightarrow 均值独立 \Rightarrow 不相关.

方误差 (mean square error) 的概念. 对于参数的估计量, 均方误差 MSE 越小则表明其估计的偏差越小^[9], 数理统计中则认为 $\hat{\theta}$ 的估计更加有效. 依据定义 (2.1.3) 我们下面即可证明, 条件期望函数 $\mathbb{E}(Y|X)$ 是最优预测量.

定理 2.1.5 (最优预测). 设随机变量 X 和 Y , 记任意的 $g(X)$ 为 Y 的一个预测函数, 则有不等式

$$\mathbb{E}(Y - g(X))^2 \geq \mathbb{E}(Y - \mathbb{E}(Y|X))^2 = \mathbb{E}e^2$$

成立, 当且仅当 $g(X) = \mathbb{E}(Y|X)$ 时取等.

证明. 注意到

$$\begin{aligned} \mathbb{E}(Y - g(X))^2 &= \mathbb{E}(Y - m(X) + m(X) - g(X))^2 = \mathbb{E}(e + m(X) - g(X))^2 \\ &= \mathbb{E}(e^2 + 2(m(X) - g(X))e + (m(X) - g(X))^2) \\ &= \mathbb{E}e^2 + 2\mathbb{E}((m(X) - g(X))e) + \mathbb{E}(m(X) - g(X))^2, \end{aligned}$$

而依据定理 (2.1.3) 有

$$\mathbb{E}((m(X) - g(X))e) = \mathbb{E}((m(X) - g(X)) \cdot \mathbb{E}(e|X)) = \mathbb{E}((m(X) - g(X)) \cdot 0) = 0,$$

故有

$$\mathbb{E}(Y - g(X))^2 = \mathbb{E}e^2 + \mathbb{E}(m(X) - g(X))^2 \geq \mathbb{E}e^2$$

成立, 当且仅当 $g(X) = m(X) = \mathbb{E}(Y|X)$ 时取等. \square

上述证明, 实际上表明了对于 Y 的任意的关于 X 的预测量 $g(X)$, 其误差平方的期望 $\mathbb{E}(Y - g(X))^2$ 存在一个可达的下界 $\mathbb{E}e^2$, 在非 $g(X) = m(X)$ 的情况下 $\mathbb{E}(Y - g(X))^2$ 总是大于 $\mathbb{E}e^2$ 的. 可以说, 定理 (2.1.5) 即是计量经济学以条件期望函数 $\mathbb{E}(Y|X)$ 作为回归函数的原因.

然而, 仅从均值上考虑随机变量 e 是不够的, Bruce(cite) 在书中讲述了一个经济学家的笑话:

An economist was standing with one foot in a bucket of boiling water and the other foot in a bucket of ice. When asked how he felt, he replied, “*On average I feel just fine.*”

Bruce(cite) 的笑话即是忽略了变量的离散 (dispersion) 情况, 尽管均值为 0, 但观测值的取值则是两个极端. 因此, 我们有必要研究回归误差 e 的方差.

^[9]容易证明, 参数 θ 的估计量 $\hat{\theta}$ 的均方误差 MSE 满足有 $\text{MSE}(\hat{\theta}) = \mathbb{E}(\theta - \hat{\theta})^2 = \text{Var}(\hat{\theta}) + (\mathbb{E}\hat{\theta} - \theta)^2$, 当估计量 $\hat{\theta}$ 为无偏估计时, 其均方误差即为 $\text{Var}(\hat{\theta})$. 特别地, 若估计量 $\hat{\theta}$ 比 θ 的其他任意估计量具有更小的 MSE, 则称 $\hat{\theta}$ 为最小方差无偏估计 (minimum variance unbiased estimate), 简称为 **MVU 估计**.

定义 2.1.4 (误差平方和, 回归方差). 对于回归误差 e , 我们记其方差为

$$\sigma^2 = \text{Var}(e) = \mathbb{E}(e - \mathbb{E}e)^2 = \mathbb{E}e^2,$$

则称 σ^2 为回归方差 (*regression variance*), $\sigma = \sqrt{\text{Var}(e)}$ 为回归标准差. 若 $\mathbb{E}e^2 < \infty$, 我们将回归误差 e 关于 X 的条件方差 (*conditional variance*) 记为

$$\sigma^2(x) = \text{Var}(e|X=x) = \mathbb{E}((e - \mathbb{E}e)^2 | X=x) = \mathbb{E}(e^2 | X=x),$$

我们将 $\sigma^2(X)$ 视为随机变量, 将 $\sigma^2(X=x)$ 视为关于随机变量 $X=x$ 的函数. 同时, 我们将 $\sigma(X) = \sqrt{\text{Var}(e|X)}$ 称为条件标准差 (*conditional standard deviation*).

依据定义 (2.1.4) 以及定理 (2.1.1), 则有

$$\sigma^2 = \mathbb{E}e^2 = \mathbb{E}(\mathbb{E}(e^2|X)) = \mathbb{E}(\sigma^2(X)),$$

即回归方差 σ^2 是条件方差 $\sigma^2(X)$ 的平均值.

利用条件方差函数, 我们可以对回归误差 e 进行标准化, 令

$$u = \frac{e - 0}{\sigma(X)} = \frac{e}{\sigma(X)},$$

则有

$$\begin{aligned} \mathbb{E}u &= \mathbb{E}\left(\frac{e}{\sigma(X)}\right) = \mathbb{E}\left(\mathbb{E}\left(\frac{e}{\sigma(X)} \middle| X\right)\right) = \mathbb{E}(\sigma(X) \mathbb{E}(e|X)) = \mathbb{E}(\sigma(X) \cdot 0) = 0, \\ \text{Var}(u) &= \mathbb{E}(u^2) = \mathbb{E}(\mathbb{E}(u^2|X)) = \mathbb{E}\left(\mathbb{E}\left(\frac{e^2}{\sigma^2(X)} \middle| X\right)\right) = \mathbb{E}\left(\frac{\mathbb{E}(e^2|X)}{\sigma^2(X)}\right) = \mathbb{E}1 = 1, \end{aligned}$$

这样式 (2.1) 便有

$$Y = \mathbb{E}(Y|X) + \sigma(X)u, \quad (2.2)$$

也即条件期望函数有表达式 $m(X) = \mathbb{E}(Y|X) = \mathbb{E}(Y) + \sigma(X)u$.

对于条件方差函数 $\sigma^2(X)$, 我们有如下的重要定义.

定义 2.1.5 (同方差, 异方差). 对于回归误差 e 的条件方差 $\sigma^2(X) = \text{Var}(e|X)$, 若 $\sigma^2(X)$ 不依赖于随机变量 X , 即满足 $\sigma^2(X) = \sigma^2$, 那么称回归误差 e 满足同方差 (*homoskedasticity*); 若 $\sigma^2(X)$ 依赖随机变量 X , 则称回归误差 e 满足异方差 (*heteroskedasticity*).

同方差和异方差在后续讨论 OLS 估计量的性质时将发挥重要作用, 尤其是在同方

差的情形下, OLS 估计量可以成为最优无偏估计量^[10]. 需要指出, 回归误差 e 满足同方差实际上是异方差的一种特例情况, 异方差是一种更加普遍的情况.

由于方差是期望运算, 本节最后补充方差的分解公式以帮助简化求解方差.

定理 2.1.6. 设随机变量 Y 和 X , 若 $\mathbb{E}Y^2 < \infty$ 则有

$$\text{Var}(Y) = \mathbb{E}(\text{Var}(Y|X)) + \text{Var}(\mathbb{E}(Y|X)).$$

证明. 对于条件方差 $\text{Var}(Y|X)$ 有

$$\begin{aligned}\text{Var}(Y|X) &= \mathbb{E}((Y - \mathbb{E}(Y|X))^2 | X) = \mathbb{E}((Y - m(X))^2 | X) \\ &= \mathbb{E}(Y^2 - 2Ym(X) + (m(X))^2 | X) \\ &= \mathbb{E}(Y^2 | X) - 2(m(X))^2 + (m(X))^2 \\ &= \mathbb{E}(Y^2 | X) - (m(X))^2.\end{aligned}$$

将 $m(X)$ 视为随机变量, 利用公式 $\text{Var}(Y|X) = \mathbb{E}(Y^2|X) - (m(X))^2$ 和 $\text{Var}(Y) = \mathbb{E}Y^2 - (\mathbb{E}Y)^2$ 以及定理 (2.1.1), 对 $\text{Var}(Y|X)$ 求期望有

$$\begin{aligned}\mathbb{E}(\text{Var}(Y|X)) &= \mathbb{E}(\mathbb{E}(Y^2|X) - (m(X))^2) = \mathbb{E}(\mathbb{E}(Y^2|X)) - \mathbb{E}[(m(X))^2] \\ &= \mathbb{E}Y^2 - (\text{Var}(\mathbb{E}(Y|X)) + (\mathbb{E}m(X))^2) \\ &= \text{Var}(Y) + (\mathbb{E}Y)^2 - \text{Var}(\mathbb{E}(Y|X)) - (\mathbb{E}m(X))^2 \\ &= \text{Var}(Y) - \text{Var}(\mathbb{E}(Y|X)) + (\mathbb{E}Y)^2 - (\mathbb{E}Y)^2 \\ &= \text{Var}(Y) - \text{Var}(\mathbb{E}(Y|X)),\end{aligned}$$

移项得 $\text{Var}(Y) = \mathbb{E}(\text{Var}(Y|X)) + \text{Var}(\mathbb{E}(Y|X))$. □

定理 (2.1.6) 将无条件方差 $\text{Var}(Y)$ 分解为**组内方差** (within group variance) $\mathbb{E}(\text{Var}(Y|X))$ 和**组间方差** (across group variance) $\text{Var}(\mathbb{E}(Y|X))$. 例如, 如果 X 是教育水平, 那么第一项组内方差就是教育水平条件预期的预期方差, 第二项组间方差是控制教育程度后的方差.

2.2 一元线性投影模型的总体参数

上一节我们介绍 (总体) 回归函数 $\mathbb{E}(Y|X)$, 并由此得到了最为一般的回归模型 $Y = \mathbb{E}(Y|X) + e$. 依据数理统计的方法范式, 统计推断是建立已知总体的情况下利用样本数

^[10]即 Gauss-Markov 定理, 在同方差的假定下, OLS 估计量是最优线性无偏估计量 (BLUE), 而 Bruce(cite) 最新的研究证明了此条件下 OLS 估计量是最优无偏估计量 (BUE), 用数理统计的概念而言即是**最小方差无偏** (minimum variance unbiased, MVU) 估计.

据区推断. 通常而言, 在研究解释变量 X 对被解释变量 Y 的因果关系时, 就模型 (2.1.1) 而言我们只能搜集到样本个体就 X 与 Y 的观测值, 若没有关于总体的理论则我们不能得知条件期望函数 $m(X)$ 的具体形式^[11], 因而不能得到 $\mathbb{E}(Y|X)$ 的观测值, 进而回归误差 e 也是不可观测的. 一种处理方法便是设定 $m(X)$ 的具体形式, 本节我们将研究 $m(X)$ 为线性函数的情况, 并进一步研究最为广泛使用的线性投影模型.

首先我们介绍线性 CEF 模型.

定义 2.2.1 (线性条件期望函数). 对于随机变量 $Y \in \mathbb{R}$ 和随机向量 $\mathbf{X} = (X_1, \dots, X_K)^T \in \mathbb{R}^{k \times 1}$, 若 Y 关于 X 的条件期望函数为

$$m(\mathbf{X}) = \mathbb{E}(Y|\mathbf{X}) = \beta_1 X_1 + \dots + \beta_K X_K, \quad (2.3)$$

则称 $m(X)$ 为**线性条件期望函数** (*linear CEF*).

利用矩阵运算, 记 $\beta = (\beta_1, \dots, \beta_K)^T$ ^[12], 则当模型 (2.1.1) 中的 $\mathbb{E}(Y|X)$ 为式 (2.3) 时, 我们可以得到如下的模型 (2.2.1) 和模型 (2.2.2).

模型 2.2.1 (线性 CEF 模型). 满足如下设定的回归模型称为**线性 CEF 模型** (*linear CEF model*)或**线性回归模型** (*linear regression model*).

$$\begin{aligned} Y &= \mathbb{E}(Y|X) + e \\ \mathbb{E}(Y|X) &= \mathbf{X}^T \beta \\ \mathbb{E}(e|\mathbf{X}) &= 0 \end{aligned}$$

模型 2.2.2 (同方差线性 CEF 模型). 满足如下设定的回归模型称为**同方差线性 CEF 模型** (*homoskedastic linear CEF model*)或**同方差线性回归模型** (*homoskedastic linear regression model*).

$$\begin{aligned} Y &= \mathbb{E}(Y|X) + e \\ \mathbb{E}(Y|X) &= \mathbf{X}^T \beta \\ \mathbb{E}(e|\mathbf{X}) &= 0 \\ \mathbb{E}(e^2|\mathbf{X}) &= \sigma^2 \end{aligned}$$

线性 CEF 模型 (2.2.1) 与回归模型 (2.1.1) 的差别仅就在于, 模型 (2.2.1) 将条件期望函数设定为了关于参数 β 线性函数, 因而线性 CEF 模型 (2.2.1) 只是回归模型 (2.1.1)

^[11]依据样本数据, 可以使用核密度估计等方法去推断总体分布, 但这不是本章要讨论的问题.

^[12]在计量经济学中, 矩阵转置更多的是使用 \mathbf{A}' 表示, 而本笔记遵从线性代数的习惯使用 \mathbf{A}^T 表示矩阵转置.

的特例情形, 故线性 CEF 模型又被称为线性回归模型. 然而, 由于后续的模型 (2.2.3) 和模型 () 也可以被称为线性回归模型, 故为了区别, 本笔记将模型 (2.2.1) 统一称呼为“线性 CEF 模型”.

定理 (2.1.4) 依旧适用于线性 CEF 模型 (2.2.1) 的回归误差.

需要强调, 线性 CEF 模型的“线性”是指条件期望函数是参数 $\beta = (\beta_1, \dots, \beta_K)^T$ 的线性函数, 这样而言,

$$\ln Y = \beta_1 X_1 + \beta_1 X_1^2 + \beta_3 \ln X_2 + \beta_2 X_2 + \beta_3 X_1 X_2 + e$$

依然是线性回归模型. 但形如

$$\ln Y = \beta_1^2 X_1 + \beta_0 + e \quad \text{或} \quad Y = AK^\alpha L^\beta e^\varepsilon$$

的模型不是线性回归模型 (但可以通过变量替换转变为线性回归模型).

本节以及后续将关注最为简单的一元线性 CEF 模型, 其结构为

$$Y = \mathbb{E}(Y|X) = \beta_0 + \beta_1 X + e. \quad (2.4)$$

式 (2.4) 称为一元线性回归方程 (linear regression equation). 这里的“一元”指的是只含有一个解释变量 X , 我们将 $\beta = (\beta_0, \beta_1)^T$ 称为回归系数 (regression coefficient), 具体地, 系数 β_1 被称为被解释变量 Y 对解释变量 X 的回归系数 (regression coefficient of Y on X), 系数 β_0 则称为截距项 (intercept term) 或常数项 (constant term).

对于初学者而言, 式 (2.4) 其实是具有误导性的. 如果依据模型 (2.2.1), 实际上一元线性回归模型的设定应该为

$$Y = \beta_1 X + e,$$

即不含有截距项 β_0 . 更一般的, 式 (2.4) 实际上应该视为特殊的二元线性回归模型, 即

$$Y = \beta_0 X_0 + \beta_1 X + e, \quad X_0 \equiv 1, \quad (2.5)$$

其中解释变量 X_0 为取值为 1 的常数 (也视为随机变量), 这时回归元被称为 **demeaned regressor**. 对于含有常数项的 demeaned regressor, 后续章节中本笔记将会详细探讨其回归模型.

我们采用模型 (2.2.1) 来研究变量 Y 和变量 X 的关系, 我们仍不能依据样本观测值来计算条件期望函数 $\mathbb{E}(Y|X)$, 因为模型 (2.1.1) 中的参数 β 是未知的. 但在此含义上, 我们可以对参数 β 进一步设定以进行统计推断, 从而使得模型更具有可操作性. 定理 (2.1.5) 证明了条件期望函数具有最小的均方误差 MSE, 则以 β 的线性函数作为条件期望函数的 $m(X)$ 也应该满足这项要求 (必要条件), 这样, 我们从线性 CEF 模型衍生出了线性投影模型. 我们记一元线性投影模型为

$$Y = \mathcal{P}(Y|X) + \varepsilon = \beta_0^* + \beta_1^* X + \varepsilon, \quad (2.6)$$

其中 $\mathcal{P}(Y|X) = \beta_0^* + \beta_1^* X$ 称为给定 X 下 Y 的最优线性预测量 (best linear predictor)^[13], 式 (2.6) 称为一元线性投影方程 (linear projection equation). $\mathcal{P}(Y|X)$ 是定义 (2.1.3) 的一种特殊情况, 即是假定了最优预测量为关于参数 $\beta^* = (\beta_0^*, \beta_1^*)^T$ 的线性函数.

需要指出, $\mathcal{P}(Y|X) = \beta_0^* + \beta_1^* X$ 中的 β^* 是为了与式 (2.4) 的系数 β 区别, 前者我们通过后面的推导将给出具体的形式, 而后者属于假定的线性条件期望函数, 是未知的. 我们将 $\varepsilon = Y - \mathcal{P}(Y|X) = Y - \beta_0^* - \beta_1^* X$ 称为投影误差 (projection error), 其不同于定义 (2.1.2) 的回归误差 e , 因为 β^* 并不是模型 (2.1.1) 中未知的 β . 当且仅当 $\mathbb{E}(Y|X)$ 确实为式 (2.4) 中的线性 CEF 时, 投影误差 ε 等于回归误差 e .

同理于式 (2.5), 一元线性投影方程 (2.6) 的真正形式为

$$Y = \beta_0^* X_0 + \beta_1^* X + e, \quad X_0 \equiv 1, \quad (2.7)$$

是具有 demeaned regressor 的二元线性投影模型. 我们之所以在一元回归模型以及一元线性投影模型中引入常数项, 一方面对于具有常数项的线性投影模型, 其具有一些独有的性质, 另一方面在 $(X, Y) \in \mathbb{R}^2$ 中点斜式直线方程的形式为 $y = kx + b$, 一元回归下的 OLS 便是求解了一条距离样本点平均“距离”最小的直线. 不做特别说明, 本章所说的一元线性 CEF 模型或一元线性投影模型均指的是如式 (2.4) 或式 (2.6) 那样的含有截距项的模型.

对于线性投影模型, 该模型的基本假设如下.

假设 2.2.1. 对于随机变量 Y 和随机向量 $\mathbf{X} = (X_1, \dots, X_K)^T$, 假设有

(a) $\mathbb{E}Y^2 < \infty$;

(b) $\mathbb{E}\|\mathbf{X}\|^2 < \infty$;

(c) $\mathbf{Q}_{\mathbf{X}\mathbf{X}} = \mathbb{E}(\mathbf{X}\mathbf{X}^T)$ 是正定矩阵.

我们研究的模型实质上具有的是式 (2.7) 的形式, 对应地假设 (2.2.1) 中后两项是 $\mathbb{E}X^2$ 存在且有限以及对 \mathbf{X} 的 Gramian 矩阵的期望

$$\mathbf{Q}_{\mathbf{X}\mathbf{X}} = \mathbb{E}(\mathbf{x}\mathbf{x}^T) = \mathbb{E} \begin{pmatrix} 1 & X \\ X & X^2 \end{pmatrix}$$

是正定矩阵. 不难发现, 一元线性投影模型的假设只需有随机变量 Y 和随机变量 X 的二阶原点矩存在且有限.

线性回归模型采用了线性函数作为条件期望函数, 但参数 $\beta = (\beta_1, \dots, \beta_K)^T$ 依旧是未知的. 线性投影模型则进一步限定了参数 β 需要满足的条件. 对于形如式 (2.4) 的

^[13]符号 \mathcal{P} 为手写花体字母 P , 在 \LaTeX 中通过命令 `\mathscr{P}` 输入.

一元线性 CEF 模型, 设 MSE 的期望为

$$S(\boldsymbol{\beta}) = S(\beta_0, \beta_1) = \mathbb{E}(Y - \beta_0 - \beta_1 X)^2, \quad (2.8)$$

则一元线性投影模型中 $\mathcal{P}(Y|X)$ 的参数 $\boldsymbol{\beta}^* = (\beta_0^*, \beta_1^*)^T$ 被定义为

$$(\beta_0^*, \beta_1^*) = \underset{(\beta_0, \beta_1) \in \mathbb{R}^2}{\operatorname{argmin}} \mathbb{E}(Y - \beta_0 - \beta_1 X)^2,$$

即能够使得均方误差 MSE 最小. 这样, 一元线性投影模型的参数 $\boldsymbol{\beta}^* = (\beta_0^*, \beta_1^*)^T$ 即是一个最优化问题的解.

现在求解使得式 (2.8) 最小的参数 $\boldsymbol{\beta}^* = (\beta_0^*, \beta_1^*)^T$, 本笔记将介绍两种方法, 一种是基于代数的方法, 另一种则是各教材上常见的微积分方法. 首先介绍代数求解方法. 注意到对一元线性回归方程 (2.4) 求期望有

$$\mathbb{E}Y = \beta_0 + \beta_1 \mathbb{E}X, \quad (2.9)$$

将式 (2.9) 代入式 (2.8) 消去 β_0 , 即有

$$S(\beta_1) = \mathbb{E}(Y - (\mathbb{E}Y - \beta_1 \mathbb{E}X + \beta_1 X))^2 = \mathbb{E}(Y - [\mathbb{E}Y + \beta_1 (X - \mathbb{E}X)])^2,$$

这时 S 变成了关于单参数 β_1 的二次函数, 尝试配方得^[14]

$$\begin{aligned} S(\beta_1) &= \mathbb{E}((X - \mathbb{E}X)\beta_1 - (Y - \mathbb{E}Y))^2 \\ &= \beta_1^2 \mathbb{E}(X - \mathbb{E}X)^2 - 2\beta_1 \mathbb{E}((X - \mathbb{E}X)(Y - \mathbb{E}Y)) + \mathbb{E}(Y - \mathbb{E}Y)^2 \\ &= \beta_1^2 \operatorname{Var}(X) - 2\beta_1 \operatorname{Cov}(X, Y) + \operatorname{Var}(Y) \\ &= \operatorname{Var}(X) \left(\beta_1 - \frac{\operatorname{Cov}(X, Y)}{\operatorname{Var}(X)} \right)^2 - \frac{(\operatorname{Cov}(X, Y))^2}{\operatorname{Var}(X)} + \operatorname{Var}(Y), \end{aligned}$$

则当一元线性 CEF 模型的 MSE 可以取得最小值时, β_1 满足

$$\beta_1 = \frac{\operatorname{Cov}(X, Y)}{\operatorname{Var}(X)}$$

而 MSE 的最小值为

$$\begin{aligned} S_{\min} &= \operatorname{Var}(Y) - \frac{(\operatorname{Cov}(X, Y))^2}{\operatorname{Var}(X)} \\ &= \operatorname{Var}(Y) \left(1 - \frac{(\operatorname{Cov}(X, Y))^2}{\operatorname{Var}(X) \operatorname{Var}(Y)} \right) = \operatorname{Var}(Y) (1 - \rho_{XY}^2). \end{aligned}$$

^[14]假定交换求和及交换积分次序均成立, 这里以及之后涉及到期望算符 \mathbb{E} 时, 需始终牢记算符 \mathbb{E} 是线性的.

利用式 (2.9) 则有

$$\beta_0 = \mathbb{E}Y - \beta_1 \mathbb{E}X = \mathbb{E}Y - \frac{\text{Cov}(X, Y)}{\text{Var}(X)} \mathbb{E}X.$$

以上代数法的推导, 即是一元线性回归方程 (2.4) 满足 MSE 最小这一必要条件时, 未知参数 β 所具有的性质. 则对于一元线性投影方程 (2.6), 在仅保证 MSE 最小的定义下 $\mathcal{P}(Y|X)$ 的参数 $\beta^* = (\beta_0^*, \beta_1^*)^T$ 满足有

$$\beta_1^* = \frac{\text{Cov}(X, Y)}{\text{Var}(X)}, \quad \beta_0^* = \mathbb{E}Y - \beta_1^* \mathbb{E}X = \mathbb{E}Y - \frac{\text{Cov}(X, Y)}{\text{Var}(X)} \mathbb{E}X.$$

下面我们再使用微积分的方法求解参数 $\beta^* = (\beta_0^*, \beta_1^*)^T$. 对于式 (2.8), 求一阶偏导有

$$\frac{\partial}{\partial \beta_0} S(\beta_0, \beta_1) = -2\mathbb{E}(Y - \beta_0 - \beta_1 X), \quad \frac{\partial}{\partial \beta_1} S(\beta_0, \beta_1) = -2\mathbb{E}((Y - \beta_0 - \beta_1 X) X),$$

令偏导数为 0, 得到极值点的一阶条件为

$$\begin{cases} \mathbb{E}(Y - \beta_0 - \beta_1 X) = 0, \\ \mathbb{E}(X(Y - \beta_0 - \beta_1 X)) = 0, \end{cases} \Leftrightarrow \begin{cases} \mathbb{E}e = 0, \\ \mathbb{E}(Xe) = 0, \end{cases} \quad (2.10)$$

方程组 (2.10) 通常被称为一元线性 CEF 模型的正规方程组 (normal equations)^[15]或一元线性回归模型的正规方程组. 当 MSE 取最小值时, 则一元线性投影方程 (2.6) 也满足有方程组 (2.10), 这样对于式投影误差 ε 满足必要条件

$$\begin{cases} \mathbb{E}(Y - \beta_0^* - \beta_1^* X) = 0, \\ \mathbb{E}(X(Y - \beta_0^* - \beta_1^* X)) = 0, \end{cases} \Leftrightarrow \begin{cases} \mathbb{E}\varepsilon = 0, \\ \mathbb{E}(X\varepsilon) = 0, \end{cases} \quad (2.11)$$

成立.

方程组 (2.11) 表明了一元线性投影模型中投影误差 ε 满足有非条件期望为 0 以及与被解释变量 X 正交^[16].

现将正规方程组 (2.11) 展开有

$$\mathbb{E}Y - \beta_0^* - \beta_1^* \mathbb{E}X = 0, \quad (2.12)$$

$$\mathbb{E}(YX) - \beta_0^* \mathbb{E}X - \beta_1^* \mathbb{E}X^2 = 0, \quad (2.13)$$

将式 (2.12) 代入式 (2.13) 中, 消去 β_0 得

$$\mathbb{E}(YX) - (\mathbb{E}Y - \beta_1^* \mathbb{E}X) \mathbb{E}X - \beta_1^* \mathbb{E}X^2 = 0,$$

^[15]也有翻译为“常规方程组”, 见林文夫 (cite).

^[16]这里的正交 (orthogonality) 是概率论中的概念: 若两个随机变量 X, Y , 其积 XY 的数学期望 $\mathbb{E}XY$ 为 0, 那么则称这两个随机变量是正交的.

化简有

$$(\mathbb{E}X^2 - (\mathbb{E}X)^2) \beta_1^* = \mathbb{E}YX - \mathbb{E}Y\mathbb{E}X,$$

解得

$$\beta_1^* = \frac{\mathbb{E}YX - \mathbb{E}Y\mathbb{E}X}{\mathbb{E}X^2 - (\mathbb{E}X)^2} = \frac{\text{Cov}(X, Y)}{\text{Var}(X)},$$

将 β_1 代回式 (2.12) 可得到 β_0^* 的表达式. 可以发现, 代数法和微积分方法所解得的参数 $\beta^* = (\beta_0^*, \beta_1^*)^T$ 是相同的, 由此我们得到了一元线性投影模型.

模型 2.2.3 (一元线性投影模型). 在假设 (2.2.1) 成立时, 满足如下设定的模型称为一元线性投影模型 (*linear projection model*).

$$\begin{aligned} Y &= \mathcal{P}(Y|X) + \varepsilon \\ \mathcal{P}(Y|X) &= \beta_0 + \beta_1 X, \quad \begin{cases} \beta_0 = \mathbb{E}Y - \beta_1 \mathbb{E}X \\ \beta_1 = \frac{\text{Cov}(X, Y)}{\text{Var}(X)} \end{cases} \\ \mathbb{E}(X\varepsilon) &= 0 \\ \mathbb{E}(\varepsilon) &= 0 \quad (\text{需要常数项 } \beta_0) \end{aligned}$$

对比一元线性投影模型 (2.2.3) 与一元线性 CEF 模型 (2.2.1), 两者的最大差别在于对于被解释变量 Y 的解释或预测是不同的. 后者使用条件期望函数 $m(X) = \mathbb{E}(Y|X)$ 作为对被解释变量 Y 的解释或预测, 只不过 $\mathbb{E}(Y|X)$ 使用了线性 CEF(但参数 β 未知, 而前者使用最优线性预测量 $\mathcal{P}(Y|X)$ 来对 Y 进行解释或预测, 而且 $\mathcal{P}(Y|X)$ 的参数是有具体表达式的.

一元线性投影模型 (2.2.3) 又被称为 Y 在 X 上的**线性投影** (*linear projection*), 且和模型 (2.1.1) 一样也被称为 Y 对 X 的**回归**. 这就导致了“回归 (regression)”一词有时是存在歧义的. 本笔记将会严格区分线性 CEF 模型 (2.2.1) 和线性投影模型 (2.2.3), 对于后者仅使用“投影”一词来表述.

另外, 之前推导 $\mathcal{P}(Y|X)$ 时为了与一元线性 CEF 方程 (2.4) 的未知参数 β 区分, 我们将 $\mathcal{P}(Y|X)$ 的参数写为 β^* . 在理解了模型 (2.2.1) 与模型 (2.2.3) 的区别后, 本笔记不再区分式 (2.4) 与式 (2.6) 中参数的符号, 都采用希腊字母 β 表示其参数. 即将一元线性投影方程 (2.6) 记为

$$Y = \mathcal{P}(Y|X) + \varepsilon = \beta_0 + \beta_1 X + \varepsilon. \quad (2.14)$$

之所以称模型 (2.2.3) 为投影模型, 这是因为在给定使得 MSE 最小的线性条件期望函数时, 回归函数 $\mathbb{E}(Y|X)$ 即是将 Y 在随机变量 X 所张成的线性空间中投影, 这一点我们将在多元线性回归模型中详细讨论.

需要指出, 模型 (2.2.3) 中 $\mathbb{E}(e) = 0$ 这一性质, 实际上是因为式 (2.4) 中含有常数项 β_0 才成立, 这是因为对于不含有常数项的一元线性投影模型

$$Y = \mathcal{P}(Y|X) + \varepsilon = \beta_1 X + \varepsilon$$

求解使得 MSE 最小的最优化问题时, 不存在一阶条件

$$\frac{\partial}{\partial \beta_0} S(\beta_0, \beta_1) = -2\mathbb{E}(Y - \beta_0 - \beta_1 X),$$

故在不含有常数项的线性投影模型中并不能够保证 $\mathbb{E}(e) = \mathbb{E}(Y - \beta_1 X) = 0$.

现在我们研究模型 (2.2.3) 的回归方差. 依据定义 (2.1.4), 对式 (2.14) 求方差, 则

$$\text{Var}(Y) = \text{Var}(\beta_0 + \beta_1 X + \varepsilon) = \beta_1^2 \text{Var}(X) + \text{Var}(\varepsilon) + 2\beta_1 \text{Cov}(X, \varepsilon),$$

而模型 (2.2.3) 含有常数项有

$$\text{Cov}(X, \varepsilon) = \mathbb{E}(X\varepsilon) - \mathbb{E}X\mathbb{E}\varepsilon = 0 + \mathbb{E}X \cdot 0 = 0,$$

即 X 与投影误差 ε 不相关, 于是一元线性投影模型的方差 σ^2 满足有

$$\sigma^2 = \text{Var}(\varepsilon) = \text{Var}(Y) - \beta_1^2 \text{Var}(X). \quad (2.15)$$

2.3 一元线性投影模型的样本估计量

在理清了一元线性投影模型的总体参数后, 本节研究使用样本估计对应的估计量. 为了对模型 (2.2.3) 进行参数估计, 设抽取了样本数为 n 的简单随机样本 (simple random sample), 即样本 $\{(X_1, Y_1), \dots, (X_n, Y_n)\}$ 满足独立同分布 (independent and identically distributed), 具体地, 即 (X_i, Y_i) 与 (X_j, Y_j) 是相互独立的 ($i \neq j$), 而样本个体 (X_i, Y_i) 与总体 (X, Y) 具有相同的概率分布 ($i = 1, \dots, n$). 对于总体所具有的概率分布, 在计量经济学中又被称为数据生成过程 (data generating process), 这意味着样本来自于一个无穷大的潜在总体或理论上的总体.

于是, 本节研究的总体与样本分别为

$$\text{总体: } Y = \beta_0 + \beta_1 X + \varepsilon$$

$$\text{样本: } \mathbf{Y} = \beta_0 \mathbf{1} + \beta_1 \mathbf{X} + \boldsymbol{\varepsilon}$$

$$\text{样本个体: } Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i, \quad i = 1, \dots, n$$

其中列向量 $\mathbf{Y}, \mathbf{X}, \boldsymbol{\varepsilon}, \mathbf{1}$ 分别为

$$\mathbf{Y} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_i \\ \vdots \\ Y_n \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} X_1 \\ \vdots \\ X_i \\ \vdots \\ X_n \end{pmatrix}, \quad \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_i \\ \vdots \\ \varepsilon_n \end{pmatrix}, \quad \mathbf{1} = \begin{pmatrix} 1 \\ \vdots \\ 1 \\ \vdots \\ 1 \end{pmatrix}_{n \times 1}.$$

对于模型 (2.2.3), 其参数 $\beta = (\beta_1, \beta_1)^T$ 为

$$\beta_0 = \mathbb{E}Y - \beta_1 \mathbb{E}X, \quad \beta_1 = \frac{\text{Cov}(X, Y)}{\text{Var}(X)},$$

则对于参数 $\beta = (\beta_1, \beta_1)^T$ 的一种估计方法是基于矩估计 (moment estimation) 的方法, 先使用样本矩^[17]去估计 $\text{Var}(X)$ 以及 $\text{Cov}(X, Y)$, 则有参数 β 的估计量 $\hat{\beta}_0$ 和 $\hat{\beta}_1$ (分别读作 β_0 hat 和 β_1 hat, 其余估计量的读法同理) 为

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X},$$

和

$$\hat{\beta}_1 = \frac{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2} = \frac{\frac{1}{n} \sum_{i=1}^n X_i Y_i - \frac{1}{n^2} \left(\sum_{i=1}^n X_i \right) \left(\sum_{i=1}^n Y_i \right)}{\frac{1}{n} \sum_{i=1}^n X_i^2 - \frac{1}{n^2} \left(\sum_{i=1}^n X_i \right)^2}$$

其中 \bar{X} 和 \bar{Y} (分别读作 X bar 和 Y bar, 其余样本均值的读法同理) 为样本均值, 分别满足有

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i, \quad \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i.$$

对于投影误差 ε , 尽管我们知道模型 (2.2.3) 中参数 β 的表达式, 但由于总体的概率分布未知, 我们只能使用样本去估计 β , 因而 ε 是不可观测的. 而对投影误差 ε 的一种合理的估计量为 $\hat{\varepsilon}_i = Y_i - (\beta_0 + \beta_1 X_i)$. 同时, 上一节我们计算了一元线性投影模型的回归方差为式 (2.15), 则对于 $\sigma^2 = \text{Var}(\varepsilon) = \text{Var}(Y) - \beta_1^2 \text{Var}(X)$ 的矩估计为

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - \bar{Y})^2 - \frac{\hat{\beta}_1^2}{n} \sum_{i=1}^n (X_i - \bar{X})^2. \quad (2.16)$$

上面给出了模型 (2.2.3) 的矩估计. 但对于线性投影模型而言, 为了估计参数 β , 更加常用的方法是使用普通最小二乘法.

定义 2.3.1 (最小二乘法估计量). 对于模型 (2.2.3), 设有简单随机样本 $\{(X_1, Y_1), \dots, (X_n, Y_n)\}$, 则均方误差 $\mathbb{E}(Y - \mathbf{X}^T \beta)^2$ 的矩估计量为

$$\hat{S}(\beta) = \frac{1}{n} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i)^2 = \frac{1}{n} \text{SSE}(\beta_0, \beta_1),$$

^[17]对于随机变量 X , 其有简单随机样本 $\{X_1, \dots, X_n\}$, 对于正整数 k . 设

$$a_k = \frac{1}{n} \sum_{i=1}^n X_i^k, \quad m_k = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^k,$$

则称 a_k 为 k 阶样本原点矩 (sample origin moment), m_k 为 k 阶样本中心矩 (sample center moment).

其中

$$\text{SSE}(\beta_0, \beta_1) = \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i)^2$$

被称为误差平方和 (*sum of squared errors*) 函数. 对于使得 $\hat{S}(\beta_0, \beta_1)$ 最小的

$$\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1)^T = \underset{(\beta_0, \beta_1) \in \mathbb{R}^2}{\operatorname{argmin}} \hat{S}(\beta),$$

则称 $(\hat{\beta}_0, \hat{\beta}_1)$ 为最小二乘估计量 (*least squares estimator*) 或普通最小二乘估计量 (*ordinary least squares estimator*), 简称为 **LS** 估计量 (*LS estimator*) 或 **OLS** 估计量 (*OLS estimator*).

定义 (2.3.1) 给出的 $(\hat{\beta}_0, \hat{\beta}_1)^T$ 是对于模型 (2.2.3) 的总体参数 $\beta = (\beta_0, \beta_1)^T$ 的估计量. 需要强调, 参数 $\beta = (\beta_0, \beta_1)^T$ 是由总体确定的非随机的、固定的数值, 而估计量 $(\hat{\beta}_0, \hat{\beta}_1)^T$ 则是随样本变化的随机变量. 另外, 为了表示估计方法, OLS 估计量 $\hat{\beta}$ 时常被加上角标后记为 $\hat{\beta}_{\text{OLS}}$; 为了表示样本容量为 n , 同理 OLS 估计量加上角标后记为 $\hat{\beta}_n$.

现在我们求解 OLS 估计量的具体形式, 与求解模型 (2.2.3) 的总体参数类似, $\hat{\beta}$ 同样有代数和微积分两种方法求解. 同样的, 本笔记将使用两种方法求解. 首先使用代数法, 对于

$$\hat{S}(\beta) = \frac{1}{n} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i)^2,$$

展开即有

$$\begin{aligned} \hat{S}(\beta) &= \frac{1}{n} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i)^2 = \frac{1}{n} \sum_{i=1}^n [Y_i^2 - 2Y_i(\beta_0 + \beta_1 X_i) + (\beta_0 + \beta_1 X_i)^2] \\ &= \frac{1}{n} \sum_{i=1}^n Y_i^2 - \frac{2}{n} \sum_{i=1}^n [Y_i(\beta_0 + \beta_1 X_i)] + \frac{1}{n} \sum_{i=1}^n (\beta_0 + \beta_1 X_i)^2 \\ &= \frac{1}{n} \sum_{i=1}^n Y_i^2 - \frac{2\beta_0}{n} \sum_{i=1}^n Y_i - \frac{2\beta_1}{n} \sum_{i=1}^n Y_i X_i + \left(\beta_0^2 + \frac{2\beta_0\beta_1}{n} \sum_{i=1}^n X_i + \frac{\beta_1^2}{n} \sum_{i=1}^n X_i^2 \right), \end{aligned}$$

分别记随机变量 X, Y 和 XY 的 2 阶样本原点矩为

$$a_2^X = \frac{1}{n} \sum_{i=1}^n X_i^2, \quad a_2^Y = \frac{1}{n} \sum_{i=1}^n Y_i^2, \quad a_2^{XY} = \frac{1}{n} \sum_{i=1}^n Y_i X_i,$$

则 $\hat{S}(\boldsymbol{\beta})$ 有

$$\begin{aligned}
 \hat{S}(\boldsymbol{\beta}) &= a_2^Y - 2\beta_0\bar{Y} - 2\beta_1a_2^{XY} + (\beta_0^2 + 2\beta_0\beta_1\bar{X} + a_2^X\beta_1^2) \\
 &= \beta_0^2 + a_2^X\beta_1^2 + 2\bar{X}\beta_0\beta_1 - 2\bar{Y}\beta_0 - 2a_2^{XY}\beta_1 + a_2^Y \\
 &= \beta_0^2 - 2\beta_0(\bar{Y} - \bar{X}\beta_1) + a_2^X\beta_1^2 - 2a_2^{XY}\beta_1 + a_2^Y \\
 &= [\beta_0 - (\bar{Y} - \bar{X}\beta_1)]^2 - (\bar{Y} - \bar{X}\beta_1)^2 + a_2^X\beta_1^2 - 2a_2^{XY}\beta_1 + a_2^Y \\
 &= [\beta_0 - (\bar{Y} - \bar{X}\beta_1)]^2 - (\bar{Y}^2 - 2\bar{Y}\bar{X}\beta_1 + \bar{X}^2\beta_1^2) + a_2^X\beta_1^2 - 2a_2^{XY}\beta_1 + a_2^Y \\
 &= [\beta_0 - (\bar{Y} - \bar{X}\beta_1)]^2 + (a_2^X - \bar{X}^2)\beta_1^2 - 2(a_2^{XY} - \bar{X}\bar{Y})\beta_1 + a_2^Y - \bar{Y}^2 \\
 &= [\beta_0 - (\bar{Y} - \bar{X}\beta_1)]^2 + (a_2^X - \bar{X}^2)\left(\beta_1 - \frac{a_2^{XY} - \bar{X}\bar{Y}}{a_2^X - \bar{X}^2}\right)^2 - \frac{(a_2^{XY} - \bar{X}\bar{Y})^2}{a_2^X - \bar{X}^2} + a_2^Y - \bar{Y}^2,
 \end{aligned}$$

最终配方得到标准二次型^[18]

$$\begin{aligned}
 \hat{S}(\boldsymbol{\beta}) &= [\beta_0 - (\bar{Y} - \bar{X}\beta_1)]^2 + (a_2^X - \bar{X}^2)\left(\beta_1 - \frac{a_2^{XY} - \bar{X}\bar{Y}}{a_2^X - \bar{X}^2}\right)^2 \\
 &\quad + \left(\frac{a_2^{XY} - \bar{X}\bar{Y}}{a_2^X - \bar{X}^2}\right)^2 + a_2^Y - \bar{Y}^2.
 \end{aligned} \tag{2.17}$$

对于式 (2.17), 前三项分别为平方项, 这三项的取值很不为负, 而第四项

$$a_2^Y - \bar{Y}^2 = \frac{1}{n} \sum_{i=1}^n Y_i^2 - \bar{Y}^2 = \frac{1}{n^2} \left(n \sum_{i=1}^n Y_i^2 - \left(\sum_{i=1}^n Y_i \right)^2 \right),$$

依据 Cauchy-Schwarz 不等式^[19]有

$$\frac{1}{n^2} \left(n \sum_{i=1}^n Y_i^2 - \left(\sum_{i=1}^n Y_i \right)^2 \right) \geq \frac{1}{n^2} \left(n \sum_{i=1}^n Y_i^2 - \left(\sum_{i=1}^n Y_i \right) \left(\sum_{i=1}^n 1^2 \right) \right) = 0,$$

当且仅当 $Y_1 = \cdots = Y_n$ 时取等. 这样, $\hat{S}(\boldsymbol{\beta})$ 的四项均非负数, 且不难得知 $\hat{S}(\beta_0, \beta_1)$ 作为关于 $\boldsymbol{\beta} = (\beta_1, \beta_1)^T$ 存在最小值. 对于式 (2.17), 前两项为关于 $\boldsymbol{\beta} = (\beta_1, \beta_1)^T$ 的二次项, 当这两项取 0 时便有

$$\hat{\beta}_0 = \bar{Y} - \bar{X}\hat{\beta}_1$$

^[18]在线性代数中, 对于二次型的配方问题存在通法, 任意的二次型可以通过矩阵相合 (congruent) 化简为标准二次型.

^[19]设 $a_i, b_i \in \mathbb{R}$ 且 $b_i \neq 0$, 则 Cauchy-Schwarz 不等式为

$$\left(\sum_{i=1}^n a_i b_i \right)^2 \leq \left(\sum_{i=1}^n a_i^2 \right) \left(\sum_{i=1}^n b_i^2 \right),$$

当且仅当 $\frac{a_1}{b_1} = \cdots = \frac{a_n}{b_n}$ 时取等. 在很多问题中, Cauchy-Schwarz 不等式又被称为内积不等式.

以及

$$\hat{\beta}_1 = \frac{a_2^{XY} - \bar{X}\bar{Y}}{a_2^X - \bar{X}^2} = \frac{\frac{1}{n} \sum_{i=1}^n Y_i X_i - \bar{X}\bar{Y}}{\frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}^2} = \frac{\frac{1}{n} \sum_{i=1}^n X_i Y_i - \frac{1}{n^2} \left(\sum_{i=1}^n X_i \right) \left(\sum_{i=1}^n Y_i \right)}{\frac{1}{n} \sum_{i=1}^n X_i^2 - \frac{1}{n^2} \left(\sum_{i=1}^n X_i \right)^2},$$

这与本节开头给出的矩估计是一样的. 同时, 代数法的方便之处在于通过式 (2.17) 我们可以得知 OLS 估计量下均方误差的矩估计量为

$$\hat{S}(\beta_0, \beta_1) = \frac{(a_2^{XY} - \bar{X}\bar{Y})^2}{a_2^X - \bar{X}^2} + a_2^Y - \bar{Y}^2 = \frac{\left(\frac{1}{n} \sum_{i=1}^n Y_i X_i - \bar{X}\bar{Y} \right)^2}{\frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}^2} + \frac{1}{n} \sum_{i=1}^n Y_i^2 - \bar{Y}^2,$$

不难发现, 这正是代数法解出模型 (2.2.3) 的总体参数 $\beta = (\beta_0, \beta_1)^T$ 后所计算的 MSE 的矩估计.

现在我们使用微积分来求解 OLS 估计量. 为了得到

$$\hat{S}(\beta) = \frac{1}{n} \text{SSE}(\beta_0, \beta_1),$$

的最小值, 只需要 $\text{SSE}(\beta_0, \beta_1)$ 取最小值即可. 则对

$$\text{SSE}(\beta_0, \beta_1) = \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i)^2$$

求一阶偏导有

$$\frac{\partial}{\partial \beta_0} \text{SSE} = -2 \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i), \quad \frac{\partial}{\partial \beta_1} \text{SSE} = -2 \sum_{i=1}^n [(Y_i - \beta_0 - \beta_1 X_i) X_i],$$

求二阶偏导数有

$$\frac{\partial^2}{\partial \beta_0^2} \text{SSE} = 2n, \quad \frac{\partial^2}{\partial \beta_1 \partial \beta_0} \text{SSE} = 2 \sum_{i=1}^n X_i,$$

以及

$$\frac{\partial^2}{\partial \beta_0 \partial \beta_1} \text{SSE} = 2 \sum_{i=1}^n X_i, \quad \frac{\partial^2}{\partial \beta_1^2} \text{SSE} = 2 \sum_{i=1}^n X_i^2.$$

SSE 最优化有一阶条件为

$$\begin{cases} \frac{\partial}{\partial \beta_0} \text{SSE} = 0, \\ \frac{\partial}{\partial \beta_1} \text{SSE} = 0, \end{cases}$$

即有

$$\begin{cases} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i) = 0, \\ \sum_{i=1}^n [(Y_i - \beta_0 - \beta_1 X_i) X_i] = 0, \end{cases} \Rightarrow \begin{cases} \sum_{i=1}^n \hat{\varepsilon}_i = 0, \\ \sum_{i=1}^n \hat{\varepsilon}_i X_i = 0, \end{cases} \quad (2.18)$$

其中 $\hat{\varepsilon}_i = Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i$ 称为残差 (residual)^[20], 其为回归误差 e_i 的估计量. 一阶条件得到的式 (2.18) 也称为 OLS 估计量的正规方程组, 其与式 (2.11) 是对应的. 此外, 正规方程组可以用矩阵语言书写为

$$\hat{\varepsilon}^T \mathbf{1} = 0, \quad \hat{\varepsilon}^T \mathbf{X} = 0, \quad (2.19)$$

这里 $\hat{\varepsilon} = (\hat{\varepsilon}_1, \dots, \hat{\varepsilon}_n)^T$ 为残差向量.

将式 (2.18) 两边同时除以 n , 化简得

$$\begin{cases} \bar{Y} - \beta_0 - \bar{X}\beta_1 = 0, \\ \frac{1}{n} \sum_{i=1}^n Y_i X_i - \beta_0 \bar{X} - \frac{\beta_1}{n} \sum_{i=1}^n X_i^2 = 0, \end{cases}$$

这即是关于 β_0 和 β_1 的二元一次方程组^[21], 消去 β_0 得

$$\frac{1}{n} \sum_{i=1}^n Y_i X_i + \bar{X}^2 \beta_1 - \bar{X}\bar{Y} - \frac{\beta_1}{n} \sum_{i=1}^n X_i^2 = 0,$$

解得

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n X_i Y_i - n \bar{X} \bar{Y}}{\sum_{i=1}^n X_i^2 - n \bar{X}^2} = \frac{n \sum_{i=1}^n X_i Y_i - \left(\sum_{i=1}^n X_i \right) \left(\sum_{i=1}^n Y_i \right)}{n \sum_{i=1}^n X_i^2 - \left(\sum_{i=1}^n X_i \right)^2},$$

而 $\hat{\beta}_0 = \bar{Y} - \bar{X} \hat{\beta}_1$.

现在检验极值点的二阶条件^[22], 注意到 Hessian 矩阵为

$$\mathbf{H} = \begin{pmatrix} \frac{\partial^2}{\partial \beta_0^2} \text{SSE} & \frac{\partial^2}{\partial \beta_1 \partial \beta_0} \text{SSE} \\ \frac{\partial^2}{\partial \beta_0 \partial \beta_1} \text{SSE} & \frac{\partial^2}{\partial \beta_1^2} \text{SSE} \end{pmatrix} = \begin{pmatrix} 2n & 2 \sum_{i=1}^n X_i \\ 2 \sum_{i=1}^n X_i & 2 \sum_{i=1}^n X_i^2 \end{pmatrix},$$

则

$$|\mathbf{H}| = 4n \sum_{i=1}^n X_i^2 - 4 \left(\sum_{i=1}^n X_i \right)^2 = 4 \left[n \sum_{i=1}^n X_i^2 - \left(\sum_{i=1}^n X_i \right)^2 \right],$$

^[20] 上述一阶条件中并没有给 β 标上 $\hat{\cdot}$, 这是因为在求解极值过程中视 β 为变量, 而最终求得的极值点 $\hat{\beta}$ 是满足一阶条件的.

^[21] 国内不少计量经济学的教材, 推导到此处后使用 Cramer 法则来求解 $\hat{\beta}_0$ 和 $\hat{\beta}_1$, 这未免有点“高射炮打蚊子”, 令人怀疑教材的编者是否真的动手推导了 OLS 估计量.

^[22] 依据 SSE 的结构, 显然其不存在最大值, 则求解的极值点为极小值点.

同理依据 Cauchy-Schwarz 不等式知 $|\mathbf{H}| \geq 0$, 但且仅当 $X_1 = \cdots = X_n$ 时取等. 通常情况, 简单随机样本不会有所有观察值均相同, 故 $|\mathbf{H}| > 0$, 可以判断 $|\mathbf{H}|$ 为正定矩阵, 则求解的 $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1)^T$ 为极小值点, 也是 SSE 的最小值点.

最终, 无论是代数法还是微积分方法, 我们解得的 OLS 估计量与本节开篇提出的矩估计的形式一样的.

定理 2.3.1. 模型 (2.2.3) 的参数 β_0 和 β_1 的 OLS 估计量分别为

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}, \quad (2.20)$$

$$\begin{aligned} \hat{\beta}_1 &= \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2} = \frac{\sum_{i=1}^n (X_i - \bar{X}) Y_i}{\sum_{i=1}^n (X_i - \bar{X})^2} \\ &= \frac{\sum_{i=1}^n Y_i X_i - n \bar{X} \bar{Y}}{\sum_{i=1}^n X_i^2 - n \bar{X}^2} = \frac{n \sum_{i=1}^n X_i Y_i - \left(\sum_{i=1}^n X_i \right) \left(\sum_{i=1}^n Y_i \right)}{n \sum_{i=1}^n X_i^2 - \left(\sum_{i=1}^n X_i \right)^2}, \end{aligned} \quad (2.21)$$

且 OLS 估计量满足有

$$\begin{cases} \sum_{i=1}^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i) = 0, \\ \sum_{i=1}^n [(Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i) X_i] = 0, \end{cases} \Rightarrow \begin{cases} \sum_{i=1}^n \hat{\varepsilon}_i = 0, \\ \sum_{i=1}^n \hat{\varepsilon}_i X_i = 0. \end{cases}$$

2.4 一元 OLS 估计量的 Gauss-Markov 定理

本节我们将研究 OLS 估计量所具有的有限性质. 在概率论与数理统计中, 对于点估计的估计准则, 可以分为有限样本性质 (或小样本性质) 和大样本性质. 前者是考虑样本容量 n 有限的情况下, 研究估计量的无偏性、有效性, 后者则是在 $n \rightarrow \infty$ 时, 研究估计量相合性、渐进分布等问题. 本节将详细研究模型 (2.2.3) 的 OLS 估计量以及其衍生变量的有限性质, 同时将证明 Gauss-Markov 定理.

在前面的章节中, 模型 (2.2.1) 和模型 (2.2.3) 都可以被称为线性回归模型, 为了区分我们分别称之为一元线性 CEF 模型 (2.2.1) 与一元线性投影模型 (2.2.3). 这样区别的另一个原因, 则是为了将“线性回归模型”的名称留给下面的模型.

模型 2.4.1 (一元线性回归模型). 在假设 (2.2.1) 以及

$$\mathbb{E}(\varepsilon|X) = 0 \quad (2.22)$$

成立时, 满足如下设定的模型称为一元线性回归模型 (*linear regression model*).

$$\begin{aligned}
 Y &= \mathcal{P}(Y|X) + \varepsilon \\
 \mathcal{P}(Y|X) &= \beta_0 + \beta_1 X, \quad \begin{cases} \beta_0 = \mathbb{E}Y - \beta_1 \mathbb{E}X \\ \beta_1 = \frac{\text{Cov}(X, Y)}{\text{Var}(X)} \end{cases} \\
 \mathbb{E}(X\varepsilon) &= 0 \\
 \mathbb{E}(\varepsilon) &= 0 \quad (\text{不需要常数项 } \beta_0)
 \end{aligned}$$

不难发现, 模型 (2.4.1) 即是在一元线性投影模型之上补充了假设式 (2.22), 而 $\mathbb{E}(\varepsilon|X) = 0$ 则是线性 CEF 模型 (2.2.1) 的性质. 实际上, 对于模型 (2.4.1) 的总体回归方程

$$Y = \mathcal{P}(Y|X) + \varepsilon = \beta_0 + \beta_1 X + \varepsilon \quad (2.23)$$

求给定 X 下的条件期望, 则

$$\mathbb{E}(Y|X) = \mathbb{E}(\beta_0 + \beta_1 X + \varepsilon|X) = \beta_0 + \beta_1 X + \mathbb{E}(\varepsilon|X),$$

而模型 (2.4.1) 中补充的假设式 (2.22), 这样模型 (2.4.1) 的条件期望函数即满足有

$$\mathbb{E}(Y|X) = \beta_0 + \beta_1 X,$$

即为符合定义 (2.2.1) 的线性条件期望函数. 这就是说, 若一元线性投影模型 (2.2.3) 满足假设 $\mathbb{E}(\varepsilon|X) = 0$, 那么即为一元线性 CEF 模型 (2.2.1).

模型 (2.4.1) 即是将模型 (2.2.3) 和模型 (2.2.1) 合二为一^[23], 我们也不必再区分 CEF 模型与投影模型的区别, 故将模型 (2.4.1) 称为一元线性回归模型. 对于大多数的计量经济学教材, 其讨论的线性回归模型便是模型 (2.4.1). 由于模型 (2.4.1) 的条件期望函数为线性的, 这样其投影误差 $\varepsilon = Y - \beta_0 - \beta_1 X$ 即是定义 (2.1.2) 中回归误差 (regression error), 我们仍记为 ε 以表明其源自线性投影模型.

在研究 OLS 估计量的有限样本性质前, 我们先研究一下模型 (2.4.1) 简单随机样本 $\{(X_1, Y_1), \dots, (X_n, Y_n)\}$ 下具有的性质. 对于样本个体

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i, \quad i = 1, \dots, n, \quad (2.24)$$

由 (Y_i, X_i) 独立同分布 ($i = 1, \dots, n$) 知投影误差 $\varepsilon_i = Y_i - \beta_0 - \beta_1 X_i$ 也是独立同分布的, 因而 (ε_i, X_i) 也是独立同分布的. 这样, 对于式 (2.24) 而言, 样本个体的被解释变量

^[23]实质上, 从推导的逻辑而言, 线性投影模型 (2.2.3) 的性质也是 (2.2.1) 所具有的. 但我们无法确定研究问题中 Y 与 X 的具体关系, 故分别设定了 CEF 和最优线性预测量.

Y_i 和解释变量 X_i 以及投影残差 ε_i 的非条件期望以及非条件方差满足有

$$\begin{aligned}\mathbb{E}X_i &= \mathbb{E}X, & \text{Var}(X_i) &= \text{Var}(X), \\ \mathbb{E}Y_i &= \mathbb{E}Y, & \text{Var}(Y_i) &= \text{Var}(Y), \\ \mathbb{E}\varepsilon_i &= \mathbb{E}\varepsilon, & \text{Var}(\varepsilon_i) &= \text{Var}(\varepsilon),\end{aligned}$$

这是依据同分布条件得到的, 样本的随机变量的数值特征与总体的数值特征相同. 而对于条件期望 $\mathbb{E}(Y_i|\mathbf{X})$ 则有,

$$\boxed{\mathbb{E}(Y_i|\mathbf{X})} = \mathbb{E}(Y_i|X_i) = \mathbb{E}(Y|X_i) = \boxed{m(x_i)} = \beta_0 + \beta_1 x_i,$$

第一个等号是因为不同 i 的样本是独立的, 而由于样本个体 i 是同分布的, 则 $\mathbb{E}(Y_i|X_i) = \mathbb{E}(Y|X_i)$ 成立, 这意味意味着所以的样本个体 Y_i 具有相同函数形式的条件期望 $m(\cdot)$, 但是其确切的值依赖于样本个体 X_i 的观测值 x_i , 也就是说, 条件期望 $\mathbb{E}(Y_i|\mathbf{X})$ 对于所以的样本个体具有相同的函数形式, 但具体值由索引 i 决定. 用矩阵符号, 则有

$$\mathbb{E}(\mathbf{Y}|\mathbf{X}) = \mathbf{1}_{n \times 1}\beta_0 + \mathbf{X}\beta_1 = (\mathbf{1}_{n \times 1}, \mathbf{X}) \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix} = \begin{pmatrix} \mathbf{1}_{1 \times n} \\ \mathbf{X}^T \end{pmatrix}^T \beta.$$

对于回归误差 ε_i 的条件期望 ($i = 1 \dots, n$), 同理上述推导有

$$\begin{aligned}\mathbb{E}(\varepsilon_i|\mathbf{X}) &= \mathbb{E}(\varepsilon_i|X_i) = \mathbb{E}(\varepsilon|X = x_i) \stackrel{\text{式 (2.22)}}{=} 0, \\ \mathbb{E}(\varepsilon_i\varepsilon_j|\mathbf{X}) &\stackrel{\text{样本 i.i.d.}}{=} \mathbb{E}(\varepsilon_i|\mathbf{X})\mathbb{E}(\varepsilon_j|\mathbf{X}) = 0, \quad i \neq j, \\ \mathbb{E}(\varepsilon_i^2|\mathbf{X}) &= \mathbb{E}(\varepsilon_i^2|X_i) = \mathbb{E}(\varepsilon^2|X = x_i) = \sigma^2(X = x_i),\end{aligned}$$

因为模型 (2.4.1) 中假设了预测误差 ε 均值独立于被解释变量 X , 故 $\mathbb{E}(\varepsilon_i|\mathbf{X})$ 的取值为 0. 依定理 (2.1.2), 这又意味 $\mathbb{E}\varepsilon_i$ 着而 $\mathbb{E}(\varepsilon_i^2|\mathbf{X})$ 的形式为条件方差函数 $\sigma^2(X)$, 但实际的观测值则依赖于索引 i . 若假设

$$\mathbb{E}(\varepsilon_i^2|\mathbf{X}) = \sigma^2(X = x_i) = \sigma^2, \quad (2.25)$$

则模型 (2.4.1) 满足同方差假设.

模型 2.4.2 (同方差一元线性回归模型). 在模型 (2.4.1) 的假设中加入式 (2.25), 则称该模型为同方差线性回归模型 (*homoskedastic linear regression model*).

在同方差假设下, 则被解释变量样本个体 Y_i 的的条件方差为

$$\begin{aligned}\boxed{\text{Var}(Y_i|\mathbf{X})} &= \text{Var}(Y_i|X_i) = \text{Var}(\beta_0 + \beta_1 X_i + \varepsilon_i|X_i) = \boxed{\text{Var}(\varepsilon_i|X_i)} \\ &= \mathbb{E}((\varepsilon_i - \mathbb{E}(\varepsilon_i|X_i))^2|X_i) = \boxed{\mathbb{E}(\varepsilon_i^2|X_i)} = \sigma^2.\end{aligned}$$

对于模型 (2.4.1) 的正交性质 $\mathbb{E}(X\varepsilon) = 0$, 则对于样本有 $\mathbb{E}(X_i\varepsilon_j) = 0$, 这也意味着对于一元线性回归模型, 样本矩阵 (列向量) $\mathbf{X} = (X_1, \dots, X_n)^T$ 与误差向量 $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^T$ 正交, 即 $\mathbb{E}(\mathbf{X}^T \boldsymbol{\varepsilon}) = 0$.

我们现考虑样本均值 \bar{Y} , 在数理统计中我们知道其为总体均值的无偏估计. 现在我们计算 \bar{Y} 的条件均值, 则有

$$\begin{aligned}\mathbb{E}(\bar{Y} | \mathbf{X}) &= \mathbb{E}\left(\frac{1}{n} \sum_{i=1}^n Y_i \middle| \mathbf{X}\right) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}(Y_i | \mathbf{X}) = \boxed{\frac{1}{n} \sum_{i=1}^n \mathbb{E}(Y_i | X_i)} \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}(Y | X_i) = \beta_0 + \frac{\beta_1}{n} \sum_{i=1}^n X_i = \beta_0 + \beta_1 \bar{X},\end{aligned}$$

即样本均值 \bar{Y} 的条件期望是条件期望函数 $m(X)$ 的样本均值, 依据条件期望定理, 不难得知 $\mathbb{E}(\bar{Y} | \mathbf{X})$ 是样本均值 \bar{Y} 的无偏估计. 需要注意的是, 因为样本 $\{(X_1, Y_1), \dots, (X_n, Y_n)\}$ 是独立同分布的, 上面计算中标红的过程才得以成立. 而对于 \bar{Y} 的条件方差, 同样计算有

$$\begin{aligned}\text{Var}(\bar{Y} | \mathbf{X}) &= \text{Var}\left(\frac{1}{n} \sum_{i=1}^n Y_i \middle| \mathbf{X}\right) = \frac{1}{n^2} \text{Var}\left(\sum_{i=1}^n Y_i \middle| \mathbf{X}\right) \\ &\stackrel{\text{样本 i.i.d.}}{=} \frac{1}{n^2} \sum_{i=1}^n \text{Var}(Y_i | \mathbf{X}) = \frac{1}{n^2} \sum_{i=1}^n \mathbb{E}(\varepsilon_i^2 | X_i),\end{aligned}$$

若是满足同方差假设的模型 (2.4.2), 则有 $\text{Var}(\bar{Y} | \mathbf{X}) = \frac{1}{n^2} \sum_{i=1}^n \sigma^2 = n^{-1} \sigma^2$, 样本均值 \bar{Y} 的条件方差与非条件方差相同.

下面我们研究 OLS 估计量的无偏性. 我们首先证明 OLS 估计量是关于 Y_1, \dots, Y_n 的线性组合, 然后证明 $\hat{\boldsymbol{\beta}} = (\hat{\beta}_0, \hat{\beta}_1)^T$ 是 $\boldsymbol{\beta} = (\beta_0, \beta_1)^T$ 的无偏估计, 最后证明 OLS 估计量的 Gauss-Markov 定理.

定理 2.4.1. 模型 (2.2.3) 的 OLS 估计量 $\hat{\beta}_0$ 和 $\hat{\beta}_1$ 是 Y_1, \dots, Y_n 的线性组合.

证明. 我们首先给出使用求和符号的证明, 然后给出基于向量符号的证明.

由定理 (2.3.1) 有

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2},$$

记 $\boxed{A^{-1}(\mathbf{X}) = \sum_{i=1}^n (X_i - \bar{X})^2}$ 以及离差 $\boxed{B_i(\mathbf{X}) = X_i - \bar{X}}$, 其中 $\mathbf{X} = (X_1, \dots, X_n)^T, i =$

$1, \dots, n$, 易知 $\sum_{i=1}^n B_i = 0$, 则

$$\begin{aligned}\hat{\beta}_1 &= A \sum_{i=1}^n B_i (Y_i - \bar{Y}) = A \sum_{i=1}^n B_i Y_i - A \sum_{i=1}^n B_i \bar{Y} \\ &= A \sum_{i=1}^n B_i Y_i - A \bar{Y} \sum_{i=1}^n B_i = A \sum_{i=1}^n B_i Y_i,\end{aligned}$$

由样本均值 \bar{X} 可知 B_i 不全为 0, 故 $\hat{\beta}_0$ 是 Y_1, \dots, Y_n 的线性组合.

而 $\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$, 其中 \bar{Y} 与 $\hat{\beta}_1 \bar{X}$ 均是 Y_1, \dots, Y_n 的线性组合, 故 $\hat{\beta}_0$ 也是 Y_1, \dots, Y_n 的线性组合. 当然, 也计算得

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X} = \frac{1}{n} \sum_{i=1}^n Y_i - A \bar{X} \sum_{i=1}^n B_i Y_i = \sum_{i=1}^n \left(\frac{1}{n} - A \bar{X} B_i \right) Y_i,$$

这在之后的定理的证明中是需要的.

下面我们使用向量符号来书写, 与求和符号相比, 向量更能够体现出 OLS 估计量的几何性质, 这对于之后理解多元线性回归模型的 OLS 估计量中很有帮助.

记离差向量 $\mathbf{b}(\mathbf{X}) = (B_1, \dots, B_n)^T$ 和 $\mathbf{1} = (1, \dots, 1)^T$, 则有 $\mathbf{b}^T \mathbf{1} = 0$, 即向量 \mathbf{b} 同 $\mathbf{1}$ 正交. 记 $\mathbf{Y} = (Y_1, \dots, Y_n)$. 对于 OLS 估计量 $\hat{\beta}_1$, 其用向量乘法可以表示为

$$\hat{\beta}_1 = A \mathbf{b}^T (\mathbf{Y} - \bar{Y} \cdot \mathbf{1}) = A \mathbf{b}^T \mathbf{Y} - A \bar{Y} (\mathbf{b}^T \mathbf{1}) = A \mathbf{b}^T \mathbf{Y},$$

而 $A \mathbf{b}$ 的元素不全为 0, 故 $\hat{\beta}_0$ 是 Y_1, \dots, Y_n 的线性组合. $\hat{\beta}_0$ 也是 Y_1, \dots, Y_n 的线性组合的证明同上, 但我们还是可以计算 $\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X} = \bar{Y} - A \mathbf{b}^T \mathbf{Y} \bar{X} = n^{-1} \mathbf{1}^T \mathbf{Y} - A \bar{X} \mathbf{b}^T \mathbf{Y} = (n^{-1} \mathbf{1} - A \bar{X} \mathbf{b})^T \mathbf{Y}$, 后续的定理的证明需要此. \square

定理 2.4.2. 模型 (2.2.3) 的 OLS 估计量 $\hat{\beta}_0$ 和 $\hat{\beta}_1$ 是参数 β_0 和 β_1 的无偏估计量.

证明. 同样地, 我们先使用求和符号书写证明, 然后再给出向量证明的形式.

首先使用求和符号书写. 为了证明 OLS 估计量是无偏的, 除了直接计算期望 $\mathbb{E} \hat{\beta}$, 另一种方法便是使用迭代期望定律, 通常而言使用后者更简单且能够得到条件期望. 对于 $\hat{\beta}_1$, 计算条件期望有

$$\begin{aligned}\mathbb{E}(\hat{\beta}_1 | \mathbf{X}) &= \mathbb{E} \left(A(\mathbf{X}) \sum_{i=1}^n B_i(\mathbf{X}) Y_i \middle| \mathbf{X} \right) \\ &= A(\mathbf{X}) \mathbb{E} \left(\sum_{i=1}^n B_i(\mathbf{X}) (\beta_0 + \beta_1 X_i + \varepsilon_i) \middle| \mathbf{X} \right) \\ &= A \beta_0 \sum_{i=1}^n B_i + A \beta_1 \sum_{i=1}^n B_i X_i + A \mathbb{E} \left(\sum_{i=1}^n B_i \varepsilon_i \middle| \mathbf{X} \right) \\ &= A \beta_1 \sum_{i=1}^n B_i X_i + A \mathbb{E} \left(\sum_{i=1}^n B_i \varepsilon_i \middle| \mathbf{X} \right),\end{aligned}$$

对于第一项, 回顾数理统计中样本方差的计算可知

$$A \sum_{i=1}^n B_i X_i = A \sum_{i=1}^n (X_i - \bar{X}) X_i = A \left(\sum_{i=1}^n X_i^2 - n\bar{X}^2 \right) = \frac{\sum_{i=1}^n X_i^2 - n\bar{X}^2}{\sum_{i=1}^n (X_i - \bar{X})^2} = 1,$$

而第二项有

$$\begin{aligned} \mathbb{E} \left(\sum_{i=1}^n B_i \varepsilon_i \middle| \mathbf{X} \right) &= \mathbb{E} \left(\sum_{i=1}^n (X_i - \bar{X}) \varepsilon_i \middle| \mathbf{X} \right) = \sum_{i=1}^n \mathbb{E} \left((X_i - \bar{X}) \varepsilon_i \middle| \mathbf{X} \right) \\ &= \sum_{i=1}^n (X_i - \bar{X}) \mathbb{E}(\varepsilon_i | \mathbf{X}) = \sum_{i=1}^n (X_i - \bar{X}) \cdot 0 = 0, \end{aligned}$$

故 $\hat{\beta}_0$ 的条件期望为 $\mathbb{E}(\hat{\beta}_1 | \mathbf{X}) = \beta_1 \cdot 1 + A \cdot 0 = \beta_1$. 依据定理 (2.1.1) 知 $\mathbb{E}(\hat{\beta}_1) = \mathbb{E}(\mathbb{E}(\hat{\beta}_1 | \mathbf{X})) = \beta_1$, 即 OLS 估计量 $\hat{\beta}_1$ 是 β_1 的无偏估计.

对于 $\hat{\beta}_0$, 其条件期望为

$$\mathbb{E}(\hat{\beta}_0 | \mathbf{X}) = \mathbb{E}(\bar{Y} - \hat{\beta}_1 \bar{X} | \mathbf{X}) = \mathbb{E}(\bar{Y} | \mathbf{X}) - \bar{X} \mathbb{E}(\hat{\beta}_1 | \mathbf{X}) = \bar{Y} - \bar{X} \beta_1,$$

同理得

$$\begin{aligned} \mathbb{E} \hat{\beta}_0 &= \mathbb{E}(\mathbb{E}(\hat{\beta}_0 | \mathbf{X})) = \mathbb{E}(\bar{Y} - \bar{X} \beta_1) \\ &= \mathbb{E} \bar{Y} - \beta_1 \mathbb{E} \bar{X} = \mathbb{E} Y - \beta_1 \mathbb{E} X = \beta_0, \end{aligned}$$

即 $\hat{\beta}_0$ 是参数 β_0 的无偏估计.

现在使用向量来证明 OLS 估计量的无偏性. 依据定理 2.4.1 的证明结果, 对于 $\hat{\beta}_1 = A(\mathbf{X}) \mathbf{b}(\mathbf{X})^T \mathbf{Y}$ 求条件期望有

$$\begin{aligned} \mathbb{E}(\hat{\beta}_1 | \mathbf{X}) &= \mathbb{E} \left(A(\mathbf{X}) \mathbf{b}(\mathbf{X})^T \mathbf{Y} \middle| \mathbf{X} \right) \stackrel{\text{向量乘法是线性变换}}{=} A \mathbf{b}^T \mathbb{E}(\mathbf{Y} | \mathbf{X}) \\ &= A \mathbf{b}^T (\mathbf{1} \beta_0 + \mathbf{X} \beta_1) = A \beta_0 \mathbf{b}^T \mathbf{1} + A \beta_1 \mathbf{b}^T \mathbf{X} = A \beta_0 \cdot 0 + A \beta_1 \mathbf{b}^T \mathbf{X} = A \beta_1 \mathbf{b}^T \mathbf{X}, \end{aligned}$$

而 $\mathbf{b}^T \mathbf{X} = \sum_{i=1}^n (X_i - \bar{X}) X_i = A^{-1}$, 故有 $\mathbb{E}(\hat{\beta}_1 | \mathbf{X}) = \beta_1$. 对于同理可证. 故 $\hat{\beta}_0, \hat{\beta}_1$ 分别是 β_0, β_1 的无偏估计. \square

定理 2.4.3 (一元线性回归模型的 Gauss-Markov 定理). 设模型 (2.4.1) 的 OLS 估计量为 $\hat{\beta}_0$ 和 $\hat{\beta}_1$, 且满足有同方差假设 $\mathbb{E}(\varepsilon^2 | \mathbf{X}) = \sigma^2$. 对于模型 (2.4.1) 的任意无偏估计量 $\tilde{\beta}_0$ 和 $\tilde{\beta}_1$, 若 $\tilde{\beta}_0$ 和 $\tilde{\beta}_1$ 是 Y_1, \dots, Y_n 的线性组合, 则有

$$\text{Var}(\tilde{\beta}_0 | \mathbf{X}) \geq \text{Var}(\hat{\beta}_0 | \mathbf{X}), \quad \text{Var}(\tilde{\beta}_1 | \mathbf{X}) \geq \text{Var}(\hat{\beta}_1 | \mathbf{X}).$$

证明. 我们首先需要计算 OLS 估计量的条件方差. 对于 $\hat{\beta}_1$ 有

$$\begin{aligned}\text{Var}(\hat{\beta}_1 | \mathbf{X}) &= \text{Var}\left(A(\mathbf{X}) \sum_{i=1}^n B_i(\mathbf{X}) Y_i \middle| \mathbf{X}\right) = A^2 \text{Var}\left(\sum_{i=1}^n B_i(\mathbf{X}) Y_i \middle| \mathbf{X}\right) \\ &\stackrel{\text{样本 i.i.d.}}{=} A^2 \sum_{i=1}^n \text{Var}(B_i(\mathbf{X}) Y_i | \mathbf{X}) = A^2 \sum_{i=1}^n B_i^2 \text{Var}(Y_i | \mathbf{X}) \\ &= A^2 \sum_{i=1}^n B_i^2 \mathbb{E}(\varepsilon_i^2 | X_i),\end{aligned}$$

而在同方差假设下 $\mathbb{E}(\varepsilon_i^2 | X_i) = \sigma^2$, 则有

$$\begin{aligned}\text{Var}(\hat{\beta}_1 | \mathbf{X}) &= A^2 \sum_{i=1}^n B_i^2 \sigma^2 = \sigma^2 A^2 \sum_{i=1}^n B_i^2 = \sigma^2 A^2 \sum_{i=1}^n (X_i - \bar{X})^2 \\ &= \sigma^2 A^2 \cdot A^{-1} = \boxed{A\sigma^2} = \boxed{\frac{\sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2}}.\end{aligned}$$

对于 $\hat{\beta}_0$, 由定理 (2.4.1) 的证明有

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X} = \frac{1}{n} \sum_{i=1}^n Y_i - A\bar{X} \sum_{i=1}^n B_i Y_i = \sum_{i=1}^n \left(\frac{1}{n} - A\bar{X} B_i\right) Y_i,$$

则 $\hat{\beta}_0$ 的条件方差为

$$\begin{aligned}\text{Var}(\hat{\beta}_0 | \mathbf{X}) &= \text{Var}\left(\sum_{i=1}^n \left(\frac{1}{n} - A\bar{X} B_i\right) Y_i \middle| \mathbf{X}\right) \stackrel{\text{样本 i.i.d.}}{=} \sum_{i=1}^n \text{Var}\left(\left(\frac{1}{n} - A\bar{X} B_i\right) Y_i \middle| \mathbf{X}\right) \\ &= \sum_{i=1}^n \left(\frac{1}{n} - A\bar{X} B_i\right)^2 \text{Var}(Y_i | \mathbf{X}) = \sum_{i=1}^n \left(\frac{1}{n} - A\bar{X} B_i\right)^2 \mathbb{E}(\varepsilon_i^2 | X_i),\end{aligned}$$

而同方差假设 $\mathbb{E}(\varepsilon_i^2 | X_i) = \sigma^2$ 下, 则有

$$\begin{aligned}\text{Var}(\hat{\beta}_0 | \mathbf{X}) &= \sum_{i=1}^n \left(\frac{1}{n} - A\bar{X} B_i\right)^2 \sigma^2 = \sigma^2 \sum_{i=1}^n \left(\frac{1}{n} - A\bar{X} B_i\right)^2 \\ &= \sigma^2 \sum_{i=1}^n \left(\frac{1}{n^2} - \frac{2}{n} A\bar{X} B_i + A^2 \bar{X}^2 B_i^2\right) = \sigma^2 \left(\frac{1}{n} - \frac{2}{n} A\bar{X} \sum_{i=1}^n B_i + A^2 \bar{X}^2 \sum_{i=1}^n B_i^2\right) \\ &= \sigma^2 \left(\frac{1}{n} - \frac{2}{n} A\bar{X} \cdot 0 + A^2 \bar{X}^2 A^{-1}\right) = \boxed{\sigma^2 \left(\frac{1}{n} + A\bar{X}^2\right)}.\end{aligned}$$

这也可以化简得

$$\begin{aligned}
 \text{Var}(\hat{\beta}_0 | \mathbf{X}) &= \sigma^2 \left(\frac{1}{n} + A\bar{X}^2 \right) \\
 &= \sigma^2 \left(\frac{1}{n} + \frac{\left(\sum_{i=1}^n X_i \right)^2}{n^2 \sum_{i=1}^n (X_i - \bar{X})^2} \right) = \sigma^2 \cdot \frac{n \sum_{i=1}^n (X_i - \bar{X})^2 + \left(\sum_{i=1}^n X_i \right)^2}{n^2 \sum_{i=1}^n (X_i - \bar{X})^2} \\
 &= \sigma^2 \cdot \frac{n \left(\sum_{i=1}^n X_i^2 - \frac{1}{n} \left(\sum_{i=1}^n X_i \right)^2 \right) + \left(\sum_{i=1}^n X_i \right)^2}{n^2 \sum_{i=1}^n (X_i - \bar{X})^2} = \sigma^2 \cdot \frac{\sum_{i=1}^n X_i^2}{n \sum_{i=1}^n (X_i - \bar{X})^2}.
 \end{aligned}$$

这样我们便求出了 OLS 估计量 $\hat{\beta}_0$ 以及 $\hat{\beta}_1$ 的条件方差分别为

$$\text{Var}(\hat{\beta}_1 | \mathbf{X}) = A\sigma^2 = \frac{\sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2}, \quad (2.26)$$

$$\text{Var}(\hat{\beta}_0 | \mathbf{X}) = \sigma^2 \left(\frac{1}{n} + A\bar{X}^2 \right) = \frac{\sigma^2 \sum_{i=1}^n X_i^2}{n \sum_{i=1}^n (X_i - \bar{X})^2}. \quad (2.27)$$

现考虑模型 (2.4.1) 的任意无偏估计量 $\tilde{\beta}_0$ 和 $\tilde{\beta}_1$, 其满足有

$$\tilde{\beta}_0 = \bar{Y} - \tilde{\beta}_1 \bar{X}, \quad \tilde{\beta}_1 = \sum_{i=1}^n C_i(\mathbf{X}) Y_i = \mathbf{c}^T \mathbf{Y},$$

其中列向量 $\mathbf{c}(\mathbf{X}) = (C_1(\mathbf{X}), \dots, C_n(\mathbf{X}))^T$ 的元素不全为 0. 无偏估计量 $\tilde{\beta}_0$ 和 $\tilde{\beta}_1$ 的条件期望有

$$\begin{aligned}
 \mathbb{E}(\tilde{\beta}_1 | \mathbf{X}) &= \mathbb{E} \left(\sum_{i=1}^n C_i(\mathbf{X}) Y_i \middle| \mathbf{X} \right) = \sum_{i=1}^n \mathbb{E}(C_i(\mathbf{X}) Y_i | \mathbf{X}) \\
 &= \sum_{i=1}^n C_i \mathbb{E}(Y_i | \mathbf{X}) = \mathbf{c}^T \mathbb{E}(\mathbf{Y} | \mathbf{X}) = \mathbf{c}^T \mathbb{E}(\mathbf{Y} | \mathbf{X}) \\
 &= \mathbf{c}^T (\mathbf{1}\beta_0 + \mathbf{X}\beta_1) = \beta_0 \mathbf{c}^T \mathbf{1} + \beta_1 \mathbf{c}^T \mathbf{X}
 \end{aligned}$$

及

$$\mathbb{E}(\tilde{\beta}_0 | \mathbf{X}) = \mathbb{E}(\bar{Y} - \tilde{\beta}_1 \bar{X} | \mathbf{X}) = \mathbb{E}(\bar{Y} | \mathbf{X}) - \bar{X} \mathbb{E}(\tilde{\beta}_1 | \mathbf{X}) = \bar{Y} - \bar{X} \tilde{\beta}_1,$$

则只需有 $\mathbb{E}(\mathbb{E}(\tilde{\beta}_1 | \mathbf{X})) = \beta_1$ 便可以保证 $\tilde{\beta}_0$ 和 $\tilde{\beta}_1$ 为参数 β_0 和 β_1 的无偏估计, 因而依据

$$\mathbb{E}\tilde{\beta}_1 = \mathbb{E}(\mathbb{E}(\tilde{\beta}_1 | \mathbf{X})) = \mathbb{E}(\beta_0 \mathbf{c}^T \mathbf{1} + \beta_1 \mathbf{c}^T \mathbf{X}) = \beta_0 \mathbb{E}(\mathbf{c}^T \mathbf{1}) + \beta_1 \mathbb{E}(\mathbf{c}^T \mathbf{X}) = \beta_1$$

成立待定系数得到必要条件^[24]

$$\begin{cases} \mathbb{E}(\mathbf{c}^T \mathbf{1}) = 0, \\ \mathbb{E}(\mathbf{c}^T \mathbf{X}) = 1, \end{cases} \Leftrightarrow \boxed{\begin{cases} \mathbf{c}^T \mathbf{1} = 0, \\ \mathbf{c}^T \mathbf{X} = 1. \end{cases}} \quad (2.28)$$

现计算任意线性无偏估计量 $\tilde{\beta}_1$ 与 $\tilde{\beta}_0$ 的条件方差为

$$\begin{aligned} \text{Var}(\tilde{\beta}_1 | \mathbf{X}) &= \text{Var}\left(\sum_{i=1}^n C_i(\mathbf{X}) Y_i \middle| \mathbf{X}\right) \stackrel{\text{样本 i.i.d.}}{=} \sum_{i=1}^n \text{Var}(C_i(\mathbf{X}) Y_i | \mathbf{X}) \\ &= \sum_{i=1}^n C_i^2 \text{Var}(Y_i | \mathbf{X}) = \sum_{i=1}^n C_i^2 \mathbb{E}(\varepsilon_i^2 | X_i) \end{aligned}$$

和

$$\begin{aligned} \text{Var}(\tilde{\beta}_0 | \mathbf{X}) &= \text{Var}\left(\frac{1}{n} \sum_{i=1}^n Y_i - \bar{X} \sum_{i=1}^n C_i(\mathbf{X}) Y_i \middle| \mathbf{X}\right) \\ &= \text{Var}\left(\sum_{i=1}^n \left(\frac{1}{n} - \bar{X} C_i(\mathbf{X})\right) Y_i \middle| \mathbf{X}\right) \stackrel{\text{样本 i.i.d.}}{=} \sum_{i=1}^n \text{Var}\left(\left(\frac{1}{n} - \bar{X} C_i(\mathbf{X})\right) Y_i \middle| \mathbf{X}\right) \\ &= \sum_{i=1}^n \left(\frac{1}{n} - \bar{X} C_i\right)^2 \text{Var}(Y_i | \mathbf{X}) = \sum_{i=1}^n \left(\frac{1}{n} - \bar{X} C_i\right)^2 \text{Var}(Y_i | \mathbf{X}) \\ &= \sum_{i=1}^n \left(\frac{1}{n} - \bar{X} C_i\right)^2 \mathbb{E}(\varepsilon_i^2 | X_i), \end{aligned}$$

同方差假设有 $\mathbb{E}(\varepsilon_i^2 | X_i) = \sigma^2$, 而

$$\sum_{i=1}^n \left(\frac{1}{n} - \bar{X} C_i\right)^2 = \frac{1}{n} - \frac{2\bar{X}}{n} \sum_{i=1}^n C_i + \bar{X}^2 \sum_{i=1}^n C_i^2 = \frac{1}{n} + \bar{X}^2 \sum_{i=1}^n C_i^2,$$

于是 $\tilde{\beta}_1$ 与 $\tilde{\beta}_0$ 的条件方差分别为

$$\text{Var}(\tilde{\beta}_1 | \mathbf{X}) = \sigma^2 \sum_{i=1}^n C_i^2, \quad (2.29)$$

$$\text{Var}(\tilde{\beta}_0 | \mathbf{X}) = \sigma^2 \left(\frac{1}{n} + \bar{X}^2 \sum_{i=1}^n C_i^2 \right). \quad (2.30)$$

现在我们来比较 OLS 估计量与任意线性无偏估计量的非条件方差, 将式 (2.29) 减去式 (2.26)、式 (2.30) 减去式 (2.27) 得

$$\text{Var}(\tilde{\beta}_1 | \mathbf{X}) - \text{Var}(\hat{\beta}_1 | \mathbf{X}) = \sigma^2 \sum_{i=1}^n C_i^2 - \sigma^2 A = \sigma^2 \left(\sum_{i=1}^n C_i^2 - A \right)$$

^[24]严格的来说, 式 (2.28) 的必要条件只有 $\begin{cases} \mathbb{E}(\mathbf{C}^T \mathbf{1}) = 0, \\ \mathbb{E}(\mathbf{C}^T \mathbf{X}) = 1, \end{cases}$, 这时利用定理 (2.1.6) 我们可以证明非条

件方差形式的 Gauss-Markov 定理.

与

$$\begin{aligned}\text{Var}\left(\tilde{\beta}_0 \mid \mathbf{X}\right) - \text{Var}\left(\hat{\beta}_0 \mid \mathbf{X}\right) &= \sigma^2 \left(\frac{1}{n} + \bar{X}^2 \sum_{i=1}^n C_i^2 \right) - \sigma^2 \left(\frac{1}{n} + A \bar{X}^2 \right) \\ &= \sigma^2 \bar{X}^2 \left(\sum_{i=1}^n C_i^2 - A \right),\end{aligned}$$

这样, 我们只需要研究 $\sum_{i=1}^n C_i^2 - A$. 由 $A^{-1} = \sum_{i=1}^n (X_i - \bar{X})^2 = \mathbf{b}^T \mathbf{b}$, $B_i = X_i - \bar{X}$, 不难注意到式 (2.28) 有

$$\mathbf{c}^T \mathbf{b} = \mathbf{c}^T (\mathbf{X} - \mathbf{1} \bar{X}) = \mathbf{C}^T \mathbf{X} - \mathbf{C}^T \mathbf{1} \cdot \bar{X} = 1,$$

那么

$$\begin{aligned}\sum_{i=1}^n C_i^2 - A &= \mathbf{c}^T \mathbf{c} - A = (\mathbf{c} - \mathbf{A} \mathbf{b} + \mathbf{A} \mathbf{b})^T (\mathbf{c} - \mathbf{A} \mathbf{b} + \mathbf{A} \mathbf{b}) - A \\ &= (\mathbf{c} - \mathbf{A} \mathbf{b})^T (\mathbf{c} - \mathbf{A} \mathbf{b}) + 2(\mathbf{c} - \mathbf{A} \mathbf{b})^T \mathbf{A} \mathbf{b} + (\mathbf{A} \mathbf{b})^T (\mathbf{A} \mathbf{b}) - A \\ &= (\mathbf{c} - \mathbf{A} \mathbf{b})^T (\mathbf{c} - \mathbf{A} \mathbf{b}) + 2A \mathbf{c}^T \mathbf{b} - 2(\mathbf{A} \mathbf{b})^T \mathbf{A} \mathbf{b} + (\mathbf{A} \mathbf{b})^T (\mathbf{A} \mathbf{b}) - A \\ &= (\mathbf{c} - \mathbf{A} \mathbf{b})^T (\mathbf{c} - \mathbf{A} \mathbf{b}) + 2A - A^2 \mathbf{b}^T \mathbf{b} - A \\ &= (\mathbf{c} - \mathbf{A} \mathbf{b})^T (\mathbf{c} - \mathbf{A} \mathbf{b}) + 2A - A^2 A^{-1} - A \\ &= (\mathbf{c} - \mathbf{A} \mathbf{b})^T (\mathbf{c} - \mathbf{A} \mathbf{b}) \geq 0,\end{aligned}$$

当且仅当 $\mathbf{c} = \mathbf{A} \mathbf{b}$ 即 $\tilde{\beta}$ 为 OLS 估计量时取等, 故有

$$\text{Var}\left(\tilde{\beta}_0 \mid \mathbf{X}\right) \geq \text{Var}\left(\hat{\beta}_0 \mid \mathbf{X}\right), \quad \text{Var}\left(\tilde{\beta}_1 \mid \mathbf{X}\right) \geq \text{Var}\left(\hat{\beta}_1 \mid \mathbf{X}\right).$$

□

依据定理 (2.4.3), OLS 估计量又被称为**最优线性无偏估计量** (best linear unbiased estimator), 简称为 **BLUE**, 该性质可以说是 OLS 估计量最为重要的性质之一, 其意味着对于一元线性回归模型 (2.4.1), 从 MSE 的角度出发, 其他任何**线性**估计量均没有 OLS 估计量有效. 一个自然的问题便是是否存在模型 (2.4.1) 的非线性估计量, 其比 OLS 估计量更为有效? 长久以来, 这是一个开放问题 (open problem), 但 Bruce 在 2021 年证明了现代的 Gauss-Markov 定理: OLS 估计量是线性回归模型的 MVU 估计. 也因此, OLS 估计量又可以被称为**最优无偏估计量** (best unbiased estimator), 简称为 **BUE**.

下面我们利用 Monte Carlo 模拟, 进一步说明 OLS 估计量的有限性质, 代码见本章最后一节. 假设被解释变量 Y 与随机变量 X 具有线性条件期望函数 $E(Y|X) = 10x + 25$, 数据生成过程为 $Y = E(Y|X) + e$, 其中 e 服从正态分布 $N(0, 5^2)$. 现在我们用 OLS 去回归方程 $Y_1 = \beta_0 + \beta_1 X_i + e_i$ 的参数, 我们进行重复抽样 $k = 1000$ 次, 每次抽样的样本

容量为 $n = 50$, 抽样方法则是有放回简单随机抽样. 这样, 我们每抽样一次便都能得到由该样本所决定的 $\hat{\beta}_0$ 和 $\hat{\beta}_1$.

定理 (2.4.2) 说明了 OLS 估计量是参数 β 的无偏估计, 而大数定律则指出对于任意随机变量 X 的样本均值 \bar{X} , 其依概率收敛于非条件期望 $\mathbb{E}X$, 这意味着为了估计 $\mathbb{E}X$ 我们可以通过反复抽样求平均值来消除概率抽样所导致的误差, 尽管不能完全消除. 而 OLS 估计量 $\hat{\beta}_{OLS}$ 也是由样本所决定的随机变量, 其无偏性便表明为了得到 β 的准确估计, 重复抽样计算均值便可以提升估计的准确性.

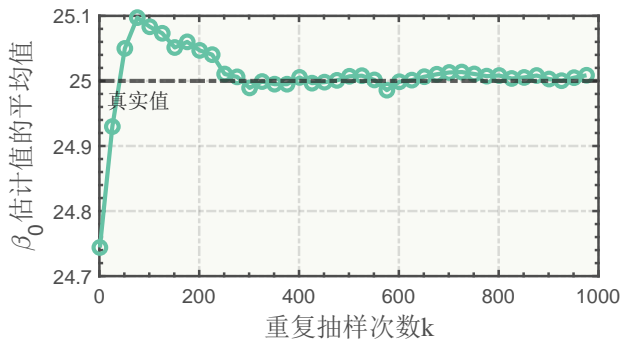
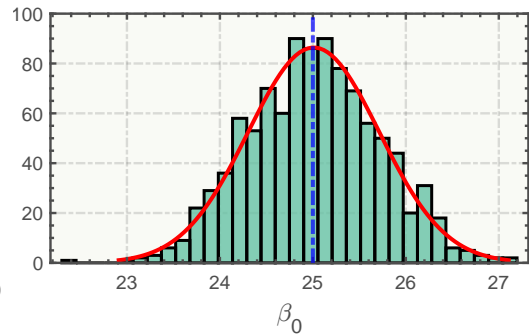
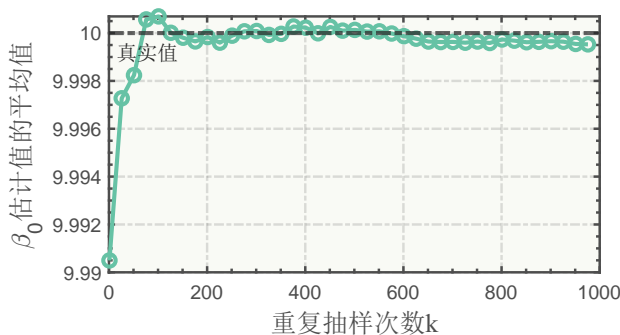
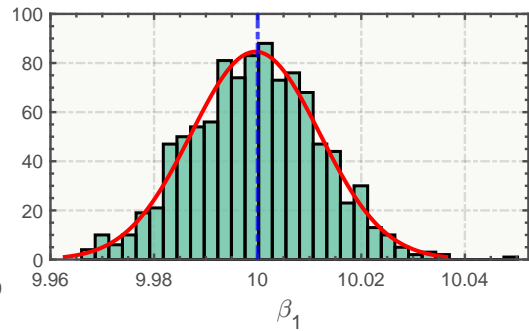
图 2.1: $\hat{\beta}_0$ 的样本均值图 2.2: $\hat{\beta}_0$ 的直方图图 2.3: $\hat{\beta}_1$ 的样本均值图 2.4: $\hat{\beta}_1$ 的直方图

图 (2.1) 和图 (2.3) 绘制了第 k 次抽样后计算的 $\hat{\beta}_0$ 以及 $\hat{\beta}_1$ 的样本均值. 可以发现, 与设定的 $E(Y|X) = 10x + 25$ 的真实值相比较, 首次抽样计算的 OLS 估计量与真实值还是有较大的差距, 但随着抽样次数 k 的增大, 估计量的样本均值便趋于真实值. 图 (2.1) 中 $\hat{\beta}_0$ 的样本均值大约在 300 次抽样后便稳定在真实值 $\beta_0 = 25$ 附近, 而图 (2.3) 中 $\hat{\beta}_1$ 的样本均值大约在 150 次抽样后对于真实值 $\beta_1 = 10$ 的偏离情况非常小了.

图 (2.2) 和图 (2.4) 则是绘制了 $k = 1000$ 次重复抽样后 OLS 估计量的直方图以及拟合的概率密度函数, 图中的蓝色竖线为对应的真实值. 不难发现, 蓝线与拟合的概率密度函数的期望值是非常接近的. 为了进一步地感受这点, 图 (2.5) 和图 (2.6) 这是依次绘制了 $k = 25$ 次、 $k = 100$ 次以及 $k = 1000$ 次下拟合的 OLS 估计量的概率密度函数,

可以发现, 随着 K 的增大, 拟合的 OLS 估计量的概率密度函数最终其期望趋向于真实值^[25]. 这便通过 Monte Carlo 模拟验证了 OLS 估计量的无偏性.

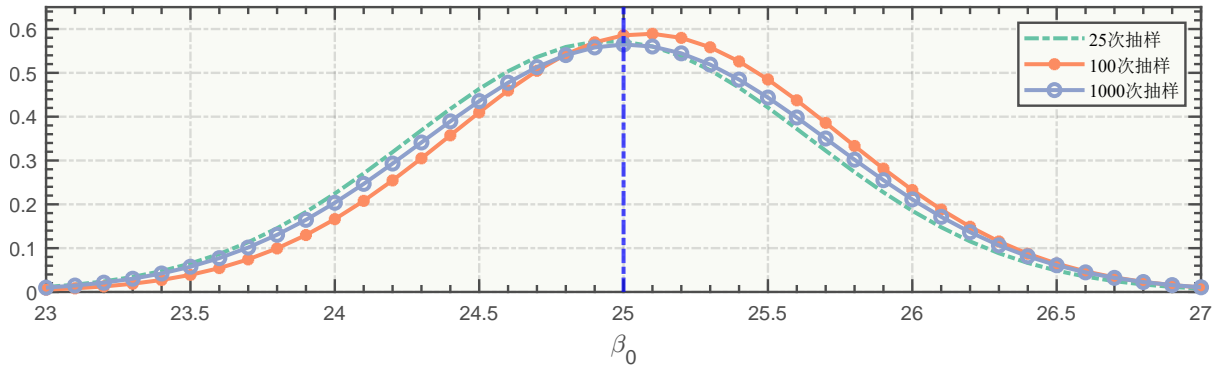


图 2.5: 重复抽样下 $\hat{\beta}_0$ 的概率密度函数

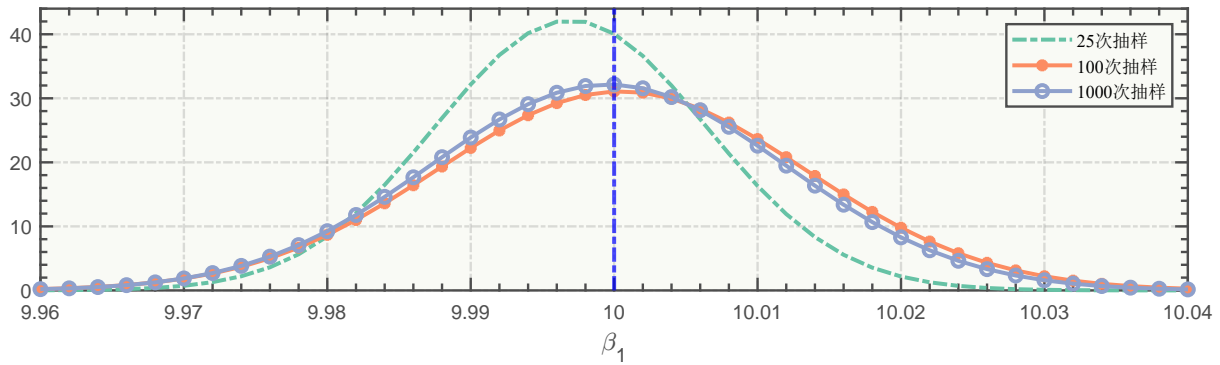


图 2.6: 重复抽样下 $\hat{\beta}_1$ 的概率密度函数

2.5 一元 OLS 估计量下的拟合值、残差

本节我们介绍由 OLS 估计量所导出的拟合值以及残差, 并介绍一元线性回归模型下的投影矩阵 (projection matrix) 与归零矩阵 (annihilator matrix), 并借此研究拟合值和残差具有有限样本性质.

依据普通最小二乘法, 我们同样可以得到模型 (2.4.1) 的 OLS 估计量, 这样我们使用 $\hat{\beta}_0$ 和 $\hat{\beta}_1$ 可以得到模型 (2.2.3) 的其他部分的估计量.

定义 2.5.1 (拟合值, 残差). 对于模型 (2.4.1), 设其参数估计为 $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1)^T$, 则记 $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$ 为 Y_i 的拟合值 (fitted value), 拟合值 \hat{Y}_i 为 Y_i 的估计量. 记 $\hat{\varepsilon}_i = \hat{Y}_i - Y_i$

^[25]这里绘制的其实是样本均值的分布. 对于分布, 中心极限定理说明了样本均值的分布依概率收敛于正态分布, 同时大数定律说明了样本均值依概率收敛于期望.

为残差 (*residual*), 残差 $\hat{\varepsilon}_i$ 为回归误差 $\varepsilon_i = Y_i - \mathbb{E}(Y|X_i)$ 的估计量.

我们使用向量符号, 记全部样本个体的拟合值、残差分别为拟合值向量 $\hat{\mathbf{Y}} = (\hat{Y}_1, \dots, \hat{Y}_n)^T$ 和残差向量 $\hat{\boldsymbol{\varepsilon}} = (\hat{\varepsilon}_1, \dots, \hat{\varepsilon}_n)^T$.

由 $\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$ 及 $\hat{\beta}_1 = A \sum_{i=1}^n B_i Y_i = \mathbf{A} \mathbf{b}^T \mathbf{Y}$, 我们现研究 OLS 估计量下拟合值的表达式. 对于单独的样本个体, 其拟合值 \hat{Y}_i 有

$$\begin{aligned}\hat{Y}_i &= \hat{\beta}_0 + \hat{\beta}_1 X_i = \bar{Y} + \hat{\beta}_1 (X_i - \bar{X}) = \bar{Y} + \hat{\beta}_1 B_i \\ &= \frac{1}{n} \sum_{j=1}^n Y_j - A \mathbf{B}_i \sum_{j=1}^n B_j Y_j = \sum_{j=1}^n \left(\frac{1}{n} - A \mathbf{B}_i B_j \right) Y_j,\end{aligned}$$

其中的 B_i 是上一节用以表示离差的符号, 可见 \hat{Y}_i 的形式谈不上直观. 但我们考虑拟合值向量 $\hat{\mathbf{Y}}$, 同样变形有

$$\begin{aligned}\hat{\mathbf{Y}} &= \mathbf{1} \hat{\beta}_0 + \mathbf{X} \hat{\beta}_1 = \mathbf{1} \bar{Y} + (\mathbf{X} - \mathbf{1} \bar{X}) \hat{\beta}_1 = \mathbf{1} \bar{Y} + \mathbf{b} \hat{\beta}_1 \\ &= \frac{1}{n} \mathbf{1} (\mathbf{1}^T \mathbf{Y}) + \mathbf{b} A \mathbf{b}^T \mathbf{Y} = \left(\frac{1}{n} (\mathbf{1} \mathbf{1}^T) + A \mathbf{b} \mathbf{b}^T \right) \mathbf{Y},\end{aligned}$$

注意到 $n = \mathbf{1}^T \mathbf{1}$, $A^{-1} = \mathbf{b}^T \mathbf{b}$, 这里的 \mathbf{b} 为上一节的离差向量, 于是利用矩阵的分块运算有

$$\begin{aligned}\hat{\mathbf{Y}} &= \left((\mathbf{1}^T \mathbf{1})^{-1} (\mathbf{1} \mathbf{1}^T) + (\mathbf{b}^T \mathbf{b})^{-1} \mathbf{b} \mathbf{b}^T \right) \mathbf{Y} \xrightarrow{\text{内积视为 } 1 \times 1 \text{ 矩阵}} \left(\mathbf{1} (\mathbf{1}^T \mathbf{1})^{-1} \mathbf{1}^T + \mathbf{b} (\mathbf{b}^T \mathbf{b})^{-1} \mathbf{b}^T \right) \mathbf{Y} \\ &= \left(\mathbf{1} (\mathbf{1}^T \mathbf{1})^{-1}, \mathbf{b} (\mathbf{b}^T \mathbf{b})^{-1} \right) \begin{pmatrix} \mathbf{1}^T \\ \mathbf{b}^T \end{pmatrix} \mathbf{Y} = (\mathbf{1}, \mathbf{b}) \begin{pmatrix} (\mathbf{1}^T \mathbf{1})^{-1} & 0 \\ 0 & (\mathbf{b}^T \mathbf{b})^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{1}^T \\ \mathbf{b}^T \end{pmatrix} \mathbf{Y},\end{aligned}$$

不难发现, 中间的矩阵 $\begin{pmatrix} (\mathbf{1}^T \mathbf{1})^{-1} & 0 \\ 0 & (\mathbf{b}^T \mathbf{b})^{-1} \end{pmatrix}$ 即为矩阵 $\begin{pmatrix} \mathbf{1}^T \\ \mathbf{b}^T \end{pmatrix} (\mathbf{1}, \mathbf{b})$ 的逆, 这是因为

$$\begin{aligned}& \left(\begin{pmatrix} \mathbf{1}^T \\ \mathbf{b}^T \end{pmatrix} (\mathbf{1}, \mathbf{b}) \right)^{-1} = \begin{pmatrix} \mathbf{1}^T \mathbf{1} & \mathbf{1}^T \mathbf{b} \\ \mathbf{b}^T \mathbf{1} & \mathbf{b}^T \mathbf{b} \end{pmatrix}^{-1} \\ &= \begin{pmatrix} \mathbf{1}^T \mathbf{1} & 0 \\ 0 & \mathbf{b}^T \mathbf{b} \end{pmatrix}^{-1} \xrightarrow{\text{对角矩阵的逆}} \begin{pmatrix} (\mathbf{1}^T \mathbf{1})^{-1} & 0 \\ 0 & (\mathbf{b}^T \mathbf{b})^{-1} \end{pmatrix}.\end{aligned}$$

这样, 拟合值向量 $\hat{\mathbf{Y}}$ 便具有了对称性的形式, 记

$$\mathbf{X}^* = (\mathbf{1}, \mathbf{b}) = (\mathbf{1}, \mathbf{X}) - (\mathbf{0}, \mathbf{1}) \bar{X} = \begin{pmatrix} 1 & X_1 - \bar{X} \\ \vdots & \vdots \\ 1 & X_n - \bar{X} \end{pmatrix},$$

则矩阵 \mathbf{X}^* 包含了符合式 (2.7) 形式的样本个体 (也即 demeaned regressor 的样本个体), 我们记 n 阶方阵

$$\mathbf{P} = (\mathbf{1}, \mathbf{b})_{n \times 2} \begin{pmatrix} (\mathbf{1}^T \mathbf{1})^{-1} & 0 \\ 0 & (\mathbf{b}^T \mathbf{b})^{-1} \end{pmatrix}_{2 \times 2} \begin{pmatrix} \mathbf{1}^T \\ \mathbf{b}^T \end{pmatrix}_{2 \times n} = \mathbf{X}^* \left((\mathbf{X}^*)^T \mathbf{X}^* \right)^{-1} (\mathbf{X}^*)^T, \quad (2.31)$$

则称矩阵 \mathbf{P} 为投影矩阵 (projection matrix).

投影矩阵 \mathbf{P} 的含义十分清晰, $\hat{\mathbf{Y}} = \mathbf{P}\mathbf{Y}$ 即将被解释变量 $\hat{\mathbf{Y}}$ 映射为 $\hat{\mathbf{Y}}$, 但其“投影”几何含义我们将在下一章介绍.

对于投影矩阵 \mathbf{P} , 其满足有如下性质.

定理 2.5.1 (一元 OLS 估计量的投影矩阵 \mathbf{P} 的性质). 对于模型 (2.4.1), 则 OLS 估计量下投影矩阵 $\mathbf{P} = \mathbf{X}^* \left((\mathbf{X}^*)^T \mathbf{X}^* \right)^{-1} (\mathbf{X}^*)^T$ 具有如下性质:

- (a) $\mathbf{P}\mathbf{X}^* = \mathbf{X}^*, \mathbf{P}\mathbf{Y} = \hat{\mathbf{Y}};$
- (b) 投影矩阵 \mathbf{P} 是对称矩阵, 即 $\mathbf{P}^T = \mathbf{P};$
- (c) 投影矩阵 \mathbf{P} 为幂等矩阵, 即 $\mathbf{P}^2 = \mathbf{P};$
- (d) 投影矩阵的 \mathbf{P} 的迹为 2;
- (e) 投影矩阵 \mathbf{P} 有 2 个特征值为 1, 余下的 $n - 2$ 个特征值为 0;
- (f) 投影矩阵 \mathbf{P} 的秩为 2.

证明. 我们仅证明性质 (a)(b)(c)(d), 余下的性质见下一章定理 (3.7.1) 的证明. 依据式 (2.31), 对于性质 (a) 直接计算有 $\mathbf{P}\mathbf{X}^* = \mathbf{X}^* \left((\mathbf{X}^*)^T \mathbf{X}^* \right)^{-1} (\mathbf{X}^*)^T \mathbf{X}^* = \mathbf{X}^*.$

性质 (b) 计算得

$$\begin{aligned} \mathbf{P}^T &= \left(\mathbf{X}^* \left((\mathbf{X}^*)^T \mathbf{X}^* \right)^{-1} (\mathbf{X}^*)^T \right)^T = \left((\mathbf{X}^*)^T \right)^T \left(\left((\mathbf{X}^*)^T \mathbf{X}^* \right)^T \right)^{-1} (\mathbf{X}^*)^T \\ &= \mathbf{X}^* \left((\mathbf{X}^*)^T \mathbf{X}^* \right)^{-1} (\mathbf{X}^*)^T = \mathbf{P}. \end{aligned}$$

性质 (c) 计算得

$$\begin{aligned} \mathbf{P}^2 &= \mathbf{X}^* \left((\mathbf{X}^*)^T \mathbf{X}^* \right)^{-1} (\mathbf{X}^*)^T \mathbf{X}^* \left((\mathbf{X}^*)^T \mathbf{X}^* \right)^{-1} (\mathbf{X}^*)^T \\ &= \mathbf{X}^* \left((\mathbf{X}^*)^T \mathbf{X}^* \right)^{-1} (\mathbf{X}^*)^T = \mathbf{P}. \end{aligned}$$

性质 (d) 计算得

$$\begin{aligned}\operatorname{tr} \mathbf{P} &= \operatorname{tr} \left(\mathbf{X}^* \left((\mathbf{X}^*)^T \mathbf{X}^* \right)^{-1} (\mathbf{X}^*)^T \right) \stackrel{\text{迹的性质}}{=} \operatorname{tr} \left(\left((\mathbf{X}^*)^T \mathbf{X}^* \right)^{-1} \mathbf{X}^* (\mathbf{X}^*)^T \right) \\ &= \operatorname{tr} \left(\left[(\mathbf{1}, \mathbf{b})^T (\mathbf{1}, \mathbf{b}) \right]^{-1} (\mathbf{1}, \mathbf{b})^T (\mathbf{1}, \mathbf{b}) \right) = \operatorname{tr} (\mathbf{I}_2) = 2.\end{aligned}$$

□

我们现在研究残差. 对于样本个体而言, 其残差 $\hat{\varepsilon}_i$ 即为

$$\hat{\varepsilon}_i = \hat{Y}_i - Y_i = \sum_{j=1}^n \left(\frac{1}{n} - A \mathbf{B}_i \mathbf{B}_j \right) Y_j - Y_i,$$

单独来看依旧缺少对称性. 考虑残差向量 $\hat{\varepsilon}$, 则

$$\hat{\varepsilon} = \mathbf{Y} - \hat{\mathbf{Y}} = \mathbf{Y} - \mathbf{P}\mathbf{Y} = (\mathbf{I}_n - \mathbf{P})\mathbf{Y},$$

我们记 n 阶方阵 $\mathbf{M} = \mathbf{I}_n - \mathbf{P}$, 则称矩阵 \mathbf{M} 为归零矩阵 (annihilator matrix), 这样残差向量也具有对称性的形式 $\hat{\varepsilon} = \mathbf{M}\mathbf{Y}$. 对于归零矩阵 \mathbf{M} , 其具有投影矩阵 \mathbf{P} 类似的代数性质.

定理 2.5.2 (一元 OLS 估计量的阵归零矩阵 \mathbf{M} 的性质). 对于模型 (2.4.1), 则 OLS 估计量下归零矩阵 $\mathbf{M} = \mathbf{I}_n - \mathbf{P}$ 具有如下性质:

- (a) $\mathbf{M}\mathbf{X}^* = \mathbf{O}, \mathbf{M}\mathbf{Y} = \mathbf{M}\varepsilon = \hat{\varepsilon}$;
- (b) 归零矩阵 \mathbf{M} 是对称矩阵, 即 $\mathbf{M}^T = \mathbf{M}$;
- (c) 归零矩阵 \mathbf{P} 为幂等矩阵, 即 $\mathbf{M}^2 = \mathbf{M}$;
- (d) 归零矩阵的 \mathbf{M} 的迹为 $n - 2$.

证明. 这些性质都是简单的代数计算便可得的. 对于性质 (a), 即有 $\mathbf{M}\mathbf{X}^* = (\mathbf{I}_n - \mathbf{P})\mathbf{X}^* = \mathbf{X}^* - \mathbf{P} \cdot \mathbf{X}^* = \mathbf{X}^* - \mathbf{X}^* = \mathbf{O}_{n \times 2}$ 及 $\mathbf{M}\mathbf{Y} = \mathbf{M}(\mathbf{X}^*\beta + \varepsilon) = \mathbf{M}\mathbf{X}^*\beta + \mathbf{M}\varepsilon = \mathbf{M}\varepsilon = \hat{\mathbf{Y}}$.

对于性质 (b), 则 $\mathbf{M}^T = (\mathbf{I}_n - \mathbf{P})^T = (\mathbf{I}_n - \mathbf{P}) = \mathbf{M}$.

对于性质 (c), 即 $\mathbf{M}^2 = (\mathbf{I}_n - \mathbf{P})(\mathbf{I}_n - \mathbf{P}) = \mathbf{I}_n - 2\mathbf{P} + \mathbf{P} = \mathbf{I}_n - \mathbf{P} = \mathbf{M}$.

对于性质 (d), 这是因为 $\operatorname{tr} \mathbf{M} = \operatorname{tr} (\mathbf{I}_n - \mathbf{P}) = \operatorname{tr} \mathbf{I}_n - \operatorname{tr} \mathbf{P} = n - 2$. □

有了投影矩阵 \mathbf{P} 及归零矩阵 \mathbf{M} , 我们对拟合值 $\hat{\mathbf{Y}}$ 和残差 $\hat{\varepsilon}$ 相关性质的研究将容易很多. 现在, 我们对回归方差进行分解. 由于

$$\mathbf{Y} = \hat{\mathbf{Y}} + \hat{\varepsilon},$$

则考虑被解释变量的样本平方和 $\mathbf{Y}^T \mathbf{Y}$ 有

$$\mathbf{Y}^T \mathbf{Y} = (\hat{\mathbf{Y}} + \hat{\boldsymbol{\varepsilon}})^T (\hat{\mathbf{Y}} + \hat{\boldsymbol{\varepsilon}}) = \hat{\mathbf{Y}}^T \hat{\mathbf{Y}} + 2\hat{\mathbf{Y}}^T \hat{\boldsymbol{\varepsilon}} + \hat{\boldsymbol{\varepsilon}}^T \hat{\boldsymbol{\varepsilon}} = \hat{\mathbf{Y}}^T \hat{\mathbf{Y}} + \hat{\boldsymbol{\varepsilon}}^T \hat{\boldsymbol{\varepsilon}},$$

注意到

$$\hat{\mathbf{Y}}^T \hat{\boldsymbol{\varepsilon}} = (\mathbf{P}\mathbf{Y})^T (\mathbf{M}\mathbf{Y}) = \mathbf{Y}^T \mathbf{P}\mathbf{M}\mathbf{Y} = \mathbf{Y}^T [\mathbf{P}(\mathbf{I} - \mathbf{P})] \mathbf{Y} = 0,$$

这也即意味着拟合值 $\hat{\mathbf{Y}}$ 与残差 $\hat{\boldsymbol{\varepsilon}}$ 正交, 故样本平方和 $\mathbf{Y}^T \mathbf{Y}$ 可以被分解为

$$\mathbf{Y}^T \mathbf{Y} = \hat{\mathbf{Y}}^T \hat{\mathbf{Y}} + \hat{\boldsymbol{\varepsilon}}^T \hat{\boldsymbol{\varepsilon}}, \Leftrightarrow \sum_{i=1}^n Y_i^2 = \sum_{i=1}^n \hat{Y}_i^2 + \sum_{i=1}^n \hat{\varepsilon}_i^2, \quad (2.32)$$

式 (2.32) 称为平方和分解公式, 这个等式其仅仅利用了投影矩阵 \mathbf{P} 和归零矩阵 \mathbf{M} 的性质. 由于模型 (2.4.1) 含义常数项 β_0 , 另一种更常用的分解公式是考虑

$$\mathbf{Y} - \mathbf{1}\bar{Y} = (\hat{\mathbf{Y}} - \mathbf{1}\bar{Y}) + \hat{\boldsymbol{\varepsilon}},$$

则离差平方和 $(\mathbf{Y} - \mathbf{1}\bar{Y})^T (\mathbf{Y} - \mathbf{1}\bar{Y})$ 满足

$$\begin{aligned} (\mathbf{Y} - \mathbf{1}\bar{Y})^T (\mathbf{Y} - \mathbf{1}\bar{Y}) &= [(\hat{\mathbf{Y}} - \mathbf{1}\bar{Y}) + \hat{\boldsymbol{\varepsilon}}]^T [(\hat{\mathbf{Y}} - \mathbf{1}\bar{Y}) + \hat{\boldsymbol{\varepsilon}}] \\ &= (\hat{\mathbf{Y}} - \mathbf{1}\bar{Y})^T (\hat{\mathbf{Y}} - \mathbf{1}\bar{Y}) + 2(\hat{\mathbf{Y}} - \mathbf{1}\bar{Y})^T \hat{\boldsymbol{\varepsilon}} + \hat{\boldsymbol{\varepsilon}}^T \hat{\boldsymbol{\varepsilon}} \\ &= (\hat{\mathbf{Y}} - \mathbf{1}\bar{Y})^T (\hat{\mathbf{Y}} - \mathbf{1}\bar{Y}) + \hat{\boldsymbol{\varepsilon}}^T \hat{\boldsymbol{\varepsilon}}, \end{aligned}$$

即得到了

$$(\mathbf{Y} - \mathbf{1}\bar{Y})^T (\mathbf{Y} - \mathbf{1}\bar{Y}) = (\hat{\mathbf{Y}} - \mathbf{1}\bar{Y})^T (\hat{\mathbf{Y}} - \mathbf{1}\bar{Y}) + \hat{\boldsymbol{\varepsilon}}^T \hat{\boldsymbol{\varepsilon}}, \quad (2.33)$$

或

$$\sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 + \sum_{i=1}^n \hat{\varepsilon}_i^2, \quad (2.34)$$

式 (2.33) 或 (2.34) 被称为普通最小二乘法的方差分析 (analysis-of-variance). 需要指出, 为了得到 (2.33), 我们使用了含有截距项时残差和为 0 的性质, 则是正规方程组 (2.18) 所得到的. OLS 的方差分析, 最为广泛的引用便是定义了统计量拟合优度 R^2 .

定义 2.5.2 (一元线性回归模型的拟合优度). 对于模型 (2.4.1), 依据式 (2.34), 我们称 $\sum_{i=1}^n (Y_i - \bar{Y})^2$ 为总体平方和 (total sum of square) 或离差平方和 (sum of squared deviations), $\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2$ 称为解释平方和 (explained sum of square) 或回归平方和 (re-

gression sum of square), $\sum_{i=1}^n \hat{\varepsilon}_i^2$ 为残差平方和 (residual sum of square), 记

$$\text{TSS} = \sum_{i=1}^n (Y_i - \bar{Y})^2, \quad \text{ESS} = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2, \quad \text{RSS} = \sum_{i=1}^n \hat{\varepsilon}_i^2 = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2,$$

则定义统计量 R^2 为

$$R^2 = \frac{\text{ESS}}{\text{TSS}} = 1 - \frac{\text{RSS}}{\text{TSS}} = \frac{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} = 1 - \frac{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2},$$

我们称 R^2 为拟合优度 (goodness of fit) 或可决系数 (coefficient of determination).

对于定义 (2.5.2) 的拟合优度, 依据式 (2.34) 不难得知 $R^2 \in [0, 1]$. 拟合优度 R^2 通常被描述为“由最小二乘法拟合值 \hat{Y} 所解释的 Y 的样本方差的百分比”^[26]. 需要指出, 拟合优度仅度量了样本回归函数 $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X$ 对于样本的拟合程度, R^2 数值越接近于 1, 说明 RSS 越小, 即拟合值 \hat{Y} 与样本观察值 Y 的偏差越小——这并不表明计量模型是否更具有合理性.

拟合优度的定义依赖于回归模型有常数项, 没有常数项时则式 (2.33) 或 (2.34) 不会成立. 此时可以使用式 (2.32) 来类似地定义非中心化的 R_{uc}^2 , 除此之外, 更多的有关度量拟合的内容见第3.10节.

需要强调, 定义 (2.5.2) 中 RSS 是残差平方和, 在一些计量经济学教材中有记为误差平方和 SSE, 但本笔记不使用这个说法, 因为我们已经在定义中 (2.3.1) 将 $\sum_{i=1}^n e_i^2 = \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i)^2$ 称为 SSE, 其是回归误差 e 的平方和, 是不可观察的.

本节的最后我们研究一下拟合值 \hat{Y}_i 与残差 $\hat{\varepsilon}_i$ 的有限样本性质. 倘若使用投影矩阵 \mathbf{P} 以及 \mathbf{M} , 则对于拟合值向量 $\hat{\mathbf{Y}}$ 及残差向量 $\hat{\boldsymbol{\varepsilon}}$ 其条件期望、条件方差将有非常简洁的推导和最终结果, 但笔者将这些内容放在下一章多元线性回归分析里. 在下一章, 我们会发现拟合值和残差值的条件期望以及条件方差都依赖于投影矩阵 \mathbf{P} 的对角线元素, 这些元素在统计学中被称为杠杆值.

本笔记现使用求和符号来推导拟合值 \hat{Y}_i 与残差 $\hat{\varepsilon}_i$ 的条件期望、条件方差.

对于样本 (Y_i, X_i) , 其拟合值 \hat{Y}_i 就给定 $\mathbf{X} = (X_1, \dots, X_n)^T$ 的条件期望为

$$\begin{aligned} \mathbb{E}(\hat{Y}_i | \mathbf{X}) &= \mathbb{E}(\hat{\beta}_0 + \hat{\beta}_1 X_i | \mathbf{X}) = \mathbb{E}(\hat{\beta}_0 | \mathbf{X}) + X_i \mathbb{E}(\hat{\beta}_1 | \mathbf{X}) \\ &= \boxed{\beta_0 + \beta_1 X_i} = \mathbb{E}(Y | X_i), \end{aligned}$$

^[26]Bruce(cite) 上原文:...the fraction of the sample variance of Y which is explained by the least squares-fit. R^2 is a crude measure of regression fit.

故 $\mathbb{E}\hat{Y}_i = \mathbb{E}(\mathbb{E}(\hat{Y}_i | \mathbf{X})) = \mathbb{E}(\beta_0 + \beta_1 X_i) = \boxed{\beta_0 + \beta_1 \mathbb{E}X}$, 也就是说, 模型 (2.4.1) 的拟合值 \hat{Y}_i 并不是样本个体 $Y_i = \beta_0 + \beta_1 X_i$ 的无偏估计. 除非解释变量 X 为常数, 即满足有 $X = X_i = \mathbb{E}X$ 时, 这时一元线性回归模型的拟合值便是对于被解释变量 Y_i 的无偏估计. 我们再考虑拟合值 \hat{Y}_i 的条件方差, 则依据公式

$$\text{Var}(aX + bY | Z) = a^2 \text{Var}(X | Z) + b^2 \text{Var}(Y | Z) + 2ab \text{Cov}(X, Y | Z),$$

则有

$$\begin{aligned} \text{Var}(\hat{Y}_i | \mathbf{X}) &= \text{Var}(\hat{\beta}_0 + \hat{\beta}_1 X_i | \mathbf{X}) \\ &= \text{Var}(\hat{\beta}_0 | \mathbf{X}) + X_i^2 \text{Var}(\hat{\beta}_1 | \mathbf{X}) + 2X_i \text{Cov}(\hat{\beta}_0, \hat{\beta}_1 | \mathbf{X}), \end{aligned}$$

在定理 (2.4.3) 的证明中我们已经计算了 OLS 估计量的条件方差, 现在计算条件协方差 $\text{Cov}(\hat{\beta}_0, \hat{\beta}_1 | \mathbf{X})$. 由方差的分配律

$$\text{Cov}(aX_1 + bX_2, Y | Z) = a \text{Cov}(X_1, Y | Z) + b \text{Cov}(X_2, Y | Z)$$

得

$$\begin{aligned} \text{Cov}(\hat{\beta}_0, \hat{\beta}_1 | \mathbf{X}) &= \text{Cov}(\bar{Y} - \hat{\beta}_1 \bar{X}, \hat{\beta}_1 | \mathbf{X}) = \text{Cov}(\bar{Y}, \hat{\beta}_1 | \mathbf{X}) - \bar{X} \text{Cov}(\hat{\beta}_1, \hat{\beta}_1 | \mathbf{X}) \\ &= \text{Cov}(\bar{Y}, \hat{\beta}_1 | \mathbf{X}) - \bar{X} \text{Var}(\hat{\beta}_1 | \mathbf{X}), \end{aligned}$$

而依据定理 (2.4.3) 的证明中计算的 $\text{Var}(\hat{\beta}_0 | \mathbf{X}) = \left(\frac{1}{n} + A\bar{X}^2\right) \sigma^2$, 这样我们由

$$\begin{aligned} \text{Var}(\hat{\beta}_0 | \mathbf{X}) &= \text{Var}(\bar{Y} - \hat{\beta}_1 \bar{X} | \mathbf{X}) \\ &= \text{Var}(\bar{Y} | \mathbf{X}) + \bar{X}^2 \text{Var}(\hat{\beta}_1 | \mathbf{X}) - 2\bar{X} \text{Cov}(\bar{Y}, \hat{\beta}_1 | \mathbf{X}) \end{aligned}$$

可以推导出

$$\begin{aligned} -2\bar{X} \text{Cov}(\bar{Y}, \hat{\beta}_1 | \mathbf{X}) &= \left(\frac{1}{n} + A\bar{X}^2\right) \sigma^2 - \text{Var}(\bar{Y} | \mathbf{X}) + \bar{X}^2 \text{Var}(\hat{\beta}_1 | \mathbf{X}) \\ &= \left(\frac{1}{n} + A\bar{X}^2\right) \sigma^2 - \frac{1}{n} \sigma^2 - \bar{X}^2 A \sigma^2 = 0, \end{aligned}$$

即样本均值 \bar{Y} 与 OLS 估计量 $\hat{\beta}_1$ 完全不相干, 故

$$\boxed{\text{Cov}(\hat{\beta}_0, \hat{\beta}_1 | \mathbf{X}) = -\bar{X} \text{Var}(\hat{\beta}_1 | \mathbf{X}) = -A\bar{X} \sigma^2}.$$

回到拟合值 \hat{Y}_i 的条件方差的计算, 于是

$$\begin{aligned} \text{Var}(\hat{Y}_i | \mathbf{X}) &= \text{Var}(\hat{\beta}_0 | \mathbf{X}) + X_i^2 \text{Var}(\hat{\beta}_1 | \mathbf{X}) + 2X_i \text{Cov}(\hat{\beta}_0, \hat{\beta}_1 | \mathbf{X}) \\ &= \left(\frac{1}{n} + A\bar{X}^2\right) \sigma^2 + X_i^2 A \sigma^2 - 2X_i A \bar{X} \sigma^2 \\ &= \left(\frac{1}{n} + A\bar{X}^2 + AX_i^2 - 2AX_i \bar{X}\right) \sigma^2 = \boxed{\left[\frac{1}{n} + A(X_i - \bar{X})^2\right] \sigma^2}. \end{aligned}$$

注意到

$$\mathbf{P} = \frac{1}{n} (\mathbf{1}\mathbf{1}^T) + A\mathbf{b}\mathbf{b}^T = \frac{1}{n} \begin{pmatrix} 1 & 1 & \cdots & 1 \\ 1 & 1 & \cdots & 1 \\ \vdots & \vdots & & \vdots \\ 1 & 1 & \cdots & 1 \end{pmatrix} + A \begin{pmatrix} B_1^2 & B_1B_2 & \cdots & B_1B_n \\ B_2B_1 & B_2^2 & \cdots & B_2B_n \\ \vdots & \vdots & & \vdots \\ B_nB_1 & B_nB_2 & \cdots & B_n^2 \end{pmatrix},$$

则投影矩阵对角线上第 i 个元素 h_{ii} 为

$$h_{ii} = \frac{1}{n} + A\bar{X}^2 + AX_i^2 - 2AX_i\bar{X} = \frac{1}{n} + A(X_i - \bar{X})^2, \quad i = 1, \dots, n,$$

在统计学中, 我们称 h_{ii} 为杠杆值 (leverage value), 因此同方差假设下拟合值 \hat{Y}_i 的条件方差为 $\boxed{\text{Var}(\hat{Y}_i | \mathbf{X}) = h_{ii}\sigma^2}$.

现在我们来计算残差 $\hat{\varepsilon}_i = Y_i - \hat{Y}_i$ 的条件期望和条件方差, 容易知道

$$\boxed{\mathbb{E}(\hat{\varepsilon}_i | \mathbf{X})} = \mathbb{E}(Y_i - \hat{Y}_i | \mathbf{X}) = \mathbb{E}(Y_i | \mathbf{X}) - \mathbb{E}(\hat{Y}_i | \mathbf{X}) = \mathbb{E}(Y | X_i) - \mathbb{E}(Y | X_i) = 0,$$

这样可以推知, 残差 $\hat{\varepsilon}_i$ 即是回归误差 ε_i 的无偏估计. 同时 $\boxed{\text{Var}(\hat{\varepsilon}_i | \mathbf{X})} = \mathbb{E}(\hat{\varepsilon}_i^2 | \mathbf{X}) = \mathbb{E}\left[(Y_i - \hat{Y}_i)^2 | \mathbf{X}\right]$. 由于

$$\begin{aligned} (Y_i - \hat{Y}_i)^2 &= [\beta_0 + \beta_1 X_i + \varepsilon_i - (\hat{\beta}_0 + \hat{\beta}_1 X_i)]^2 = [(\beta_0 - \hat{\beta}_0) + (\beta_1 - \hat{\beta}_1) X_i + \varepsilon_i]^2 \\ &= (\beta_0 - \hat{\beta}_0)^2 + X_i^2 (\beta_1 - \hat{\beta}_1)^2 + \varepsilon_i^2 + 2X_i (\beta_0 - \hat{\beta}_0) (\beta_1 - \hat{\beta}_1) \\ &\quad + 2\varepsilon_i (\beta_0 - \hat{\beta}_0) + 2X_i \varepsilon_i (\beta_1 - \hat{\beta}_1), \end{aligned}$$

则

$$\begin{aligned} &\mathbb{E}\left[(Y_i - \hat{Y}_i)^2 | \mathbf{X}\right] \\ &= \text{Var}(\hat{\beta}_0 | \mathbf{X}) + X_i^2 \text{Var}(\hat{\beta}_1 | \mathbf{X}) + \mathbb{E}(\varepsilon_i^2 | \mathbf{X}) + 2X_i \text{Cov}(\hat{\beta}_0, \hat{\beta}_1 | \mathbf{X}) \\ &\quad + 2\mathbb{E}(\varepsilon_i (\beta_0 - \hat{\beta}_0) | \mathbf{X}) + 2X_i \mathbb{E}(\varepsilon_i (\beta_1 - \hat{\beta}_1) | \mathbf{X}) \\ &= \text{Var}(\hat{\beta}_0 | \mathbf{X}) + X_i^2 \text{Var}(\hat{\beta}_1 | \mathbf{X}) + \mathbb{E}(\varepsilon_i^2 | \mathbf{X}) + 2X_i \text{Cov}(\hat{\beta}_0, \hat{\beta}_1 | \mathbf{X}) \\ &\quad - 2\mathbb{E}(\varepsilon_i \hat{\beta}_0 | \mathbf{X}) - 2X_i \mathbb{E}(\varepsilon_i \hat{\beta}_1 | \mathbf{X}) \end{aligned}$$

前四项我们已经计算过了, 现计算后两项. 首先计算

$$\mathbb{E}(\varepsilon_i \hat{\beta}_1 | \mathbf{X}) = \mathbb{E}\left(\varepsilon_i A \sum_{j=1}^n B_j Y_j \middle| \mathbf{X}\right) = A \sum_{j=1}^n B_j \mathbb{E}(Y_j \varepsilon_i | \mathbf{X}),$$

而

$$\begin{aligned}
 \mathbb{E}(Y_j \varepsilon_i | \mathbf{X}) &\stackrel{\text{样本 i.i.d.}}{=} \mathbb{E}(Y_j | \mathbf{X}) \mathbb{E}(\varepsilon_i | \mathbf{X}) = \mathbb{E}(Y_j | \mathbf{X}) \cdot 0 = 0, \quad i \neq j, \\
 \boxed{\mathbb{E}(Y_i \varepsilon_i | \mathbf{X})} &= \text{Cov}(Y_i, \varepsilon_i | \mathbf{X}) + \mathbb{E}(Y_i | \mathbf{X}) \mathbb{E}(\varepsilon_i | \mathbf{X}) = \boxed{\text{Cov}(Y_i, \varepsilon_i | \mathbf{X})} \\
 &= \text{Cov}(\beta_0 + \beta_1 X_i + \varepsilon_i, \varepsilon_i | \mathbf{X}) \\
 &= \underbrace{\text{Cov}(\beta_0, \varepsilon_i | \mathbf{X})}_{\beta_0 \text{ 为常数}} + \beta_1 \underbrace{\text{Cov}(X_i, \varepsilon_i | \mathbf{X})}_{\text{不相关}} + \text{Cov}(\varepsilon_i, \varepsilon_i | \mathbf{X}) \\
 &= \text{Var}(\varepsilon_i | \mathbf{X}) \stackrel{\text{同方差假设}}{=} \sigma^2,
 \end{aligned}$$

故 $\boxed{\mathbb{E}(\varepsilon_i \hat{\beta}_1 | \mathbf{X}) = AB_i \sigma^2}$. 再者,

$$\begin{aligned}
 \mathbb{E}(\varepsilon_i \hat{\beta}_0 | \mathbf{X}) &= \mathbb{E}(\varepsilon_i (\bar{Y} - \hat{\beta}_1 \bar{X}) | \mathbf{X}) = \mathbb{E}(\varepsilon_i \bar{Y} | \mathbf{X}) - \bar{X} \mathbb{E}(\varepsilon_i \hat{\beta}_1 | \mathbf{X}) \\
 &= \text{Cov}(\bar{Y}, \varepsilon_i | \mathbf{X}) + \mathbb{E}(\bar{Y} | \mathbf{X}) \mathbb{E}(\varepsilon_i | \mathbf{X}) - \bar{X} \mathbb{E}(\varepsilon_i \hat{\beta}_1 | \mathbf{X}) \\
 &= \text{Cov}\left(\frac{1}{n} \sum_{j=1}^n Y_j, \varepsilon_i | \mathbf{X}\right) - \bar{X} \mathbb{E}(\varepsilon_i \hat{\beta}_1 | \mathbf{X}) = \frac{1}{n} \sum_{j=1}^n \underbrace{\text{Cov}(Y_j, \varepsilon_i | \mathbf{X})}_{\text{样本 i.i.d.}} - \bar{X} \mathbb{E}(\varepsilon_i \hat{\beta}_1 | \mathbf{X}) \\
 &= \frac{1}{n} \text{Cov}(Y_i, \varepsilon_i | \mathbf{X}) - \bar{X} \mathbb{E}(\varepsilon_i \hat{\beta}_1 | \mathbf{X}) = \frac{1}{n} \sigma^2 - \bar{X} AB_i \sigma^2,
 \end{aligned}$$

即有 $\boxed{\mathbb{E}(\varepsilon_i \hat{\beta}_0 | \mathbf{X}) = \left(\frac{1}{n} - \bar{X} AB_i\right) \sigma^2}$. 这样, 残差 $\hat{\varepsilon}_i$ 的条件方差为

$$\begin{aligned}
 &\text{Var}(\hat{\varepsilon}_i | \mathbf{X}) \\
 &= \text{Var}(\hat{\beta}_0 | \mathbf{X}) + X_i^2 \text{Var}(\hat{\beta}_1 | \mathbf{X}) + \mathbb{E}(\varepsilon_i^2 | \mathbf{X}) + 2X_i \text{Cov}(\hat{\beta}_0, \hat{\beta}_1 | \mathbf{X}) \\
 &\quad - 2\mathbb{E}(\varepsilon_i \hat{\beta}_0 | \mathbf{X}) - 2X_i \mathbb{E}(\varepsilon_i \hat{\beta}_1 | \mathbf{X}) \\
 &= \left(\frac{1}{n} + A\bar{X}^2\right) \sigma^2 + X_i^2 A \sigma^2 + \sigma^2 + 2X_i (-A\bar{X} \sigma^2) - 2\left(\frac{1}{n} - \bar{X} AB_i\right) \sigma^2 - 2AB_i \sigma^2 \\
 &= \left[1 - \frac{1}{n} + A(\bar{X}^2 - 2X_i \bar{X} + X_i^2 + 2\bar{X} B_i - 2X_i B_i)\right] \sigma^2 \\
 &= \left[1 - \frac{1}{n} + A((X_i - \bar{X})^2 + 2B_i(\bar{X} - X_i))\right] \sigma^2 \\
 &= \left[1 - \frac{1}{n} + A((X_i - \bar{X})^2 - 2(X_i - \bar{X})^2)\right] \sigma^2 \\
 &= \boxed{\left[1 - \frac{1}{n} - A(X_i - \bar{X})^2\right] \sigma^2} = \boxed{(1 - h_{ii}) \sigma^2}.
 \end{aligned}$$

也就是说, 同方差假设下, $\text{Var}(\hat{\varepsilon}_i | \mathbf{X}) = (1 - h_{ii}) \sigma^2 < \sigma^2 = \text{Var}(\varepsilon_i | \mathbf{X})$, 即残差的条件方差要小于回归误差的条件方差.

2.6 一元 OLS 估计量下的方差估计量

本节我们将介绍同方差假设下, 由 OLS 估计量所导出的回归方差 σ^2 的估计量. 对于一元线性回归模型 (2.4.1), 因为

$$\sigma^2 = \text{Var}(\varepsilon) = \mathbb{E}\varepsilon^2,$$

则理想的情况下, 回归误差 ε_i 可以被观察, 则回归方差 σ^2 的矩估计 σ_{ME}^2 为

$$\sigma_{\text{ME}}^2 = \frac{1}{n} \sum_{i=1}^n \varepsilon_i^2 = \frac{1}{n} \boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon},$$

且

$$\mathbb{E}\sigma_{\text{ME}}^2 = \mathbb{E}\left(\frac{1}{n} \sum_{i=1}^n \varepsilon_i^2\right) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}\varepsilon_i^2 = \frac{1}{n} \sum_{i=1}^n \sigma^2 = \sigma^2,$$

即 σ_{ME}^2 是回归方差的无偏估计. 但可惜的是, 回归误差 ε_i 不可被观察, 必须选择其他的估计方法.

上一节中, 我们证明了残差 $\hat{\varepsilon}_i$ 是 ε_i 的无偏估计, 则用 $\hat{\varepsilon}_i$ 替换 ε_i , 得到估计量

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \hat{\varepsilon}_i^2 = \frac{1}{n} \hat{\boldsymbol{\varepsilon}}^T \hat{\boldsymbol{\varepsilon}} = \frac{1}{n} (\mathbf{M}\boldsymbol{\varepsilon})^T \mathbf{M}\boldsymbol{\varepsilon} = \frac{1}{n} \boldsymbol{\varepsilon}^T \mathbf{M}^T \mathbf{M} \boldsymbol{\varepsilon} = \frac{1}{n} \boldsymbol{\varepsilon}^T \mathbf{M} \boldsymbol{\varepsilon},$$

但由于

$$\sigma_{\text{ME}}^2 - \hat{\sigma}^2 = \frac{1}{n} \boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon} - \frac{1}{n} \boldsymbol{\varepsilon}^T \mathbf{M} \boldsymbol{\varepsilon} = \frac{1}{n} \boldsymbol{\varepsilon}^T (\mathbf{I} - \mathbf{M}) \boldsymbol{\varepsilon} = \frac{1}{n} \boldsymbol{\varepsilon}^T \mathbf{P} \boldsymbol{\varepsilon} \geq 0,$$

这里的投影矩阵 \mathbf{P} 是半正定, 取期望则有 $\sigma^2 = \mathbb{E}\sigma_{\text{ME}}^2 \geq \mathbb{E}\hat{\sigma}^2$, 即直接使用 $\hat{\varepsilon}_i$ 替换 ε_i 得到的方差估计量是有偏的. 这样, 为了得到回归方差的无偏估计, 我们需要计算残差平方和的期望. 我们将采用两种方法来计算 RSS 的条件期望.

第一种方法则是利用上一节计算的 $\text{Var}(\hat{\varepsilon}_i | \mathbf{X}) = \mathbb{E}(\hat{\varepsilon}_i^2 | \mathbf{X}) = (1 - h_{ii}) \sigma^2$, 则

$$\begin{aligned} \mathbb{E}\left(\sum_{i=1}^n \hat{\varepsilon}_i^2 \middle| \mathbf{X}\right) &= \sum_{i=1}^n \mathbb{E}(\hat{\varepsilon}_i^2 | \mathbf{X}) = \sum_{i=1}^n (1 - h_{ii}) \sigma^2 \\ &= \sigma^2 \text{tr}(\mathbf{I}_n - \mathbf{P}) = \sigma^2 \text{tr} \mathbf{M} = (n - 2) \sigma^2, \end{aligned}$$

这里利用了杠杆值 h_{ii} 的是投影矩阵 \mathbf{P} 的对角线元素以及定理 (2.5.2). 由此定义

$$s^2 = \frac{1}{n - 2} \sum_{i=1}^n \hat{\varepsilon}_i^2 = \frac{\text{RSS}}{n - 2},$$

则有 $\mathbb{E}s^2 = \mathbb{E}(s^2 | \mathbf{X}) = \sigma^2$, 即 s^2 是回归方差 σ^2 的无偏估计, 我们称 s^2 为 σ^2 的偏差校正估计量 (bias-corrected estimator).

我们再使用另一个方法计算 RSS 的条件期望, 利用式 (2.32) 则有

$$\begin{aligned}
 \mathbb{E} \left(\sum_{i=1}^n \hat{\varepsilon}_i^2 \middle| \mathbf{X} \right) &= \mathbb{E} \left(\sum_{i=1}^n Y_i^2 \middle| \mathbf{X} \right) - \mathbb{E} \left(\sum_{i=1}^n \hat{Y}_i^2 \middle| \mathbf{X} \right) \\
 &= \sum_{i=1}^n \mathbb{E} (Y_i^2 | \mathbf{X}) - \sum_{i=1}^n \mathbb{E} (\hat{Y}_i^2 | \mathbf{X}) \\
 &= \sum_{i=1}^n [\text{Var} (Y_i | \mathbf{X}) + (\mathbb{E} (Y_i | \mathbf{X}))^2] - \sum_{i=1}^n [\text{Var} (\hat{Y}_i | \mathbf{X}) + (\mathbb{E} (\hat{Y}_i | \mathbf{X}))^2] \\
 &= \sum_{i=1}^n [\text{Var} (Y_i | \mathbf{X}) + (\mathbb{E} (Y | X_i))^2] - \sum_{i=1}^n [\text{Var} (\hat{Y}_i | \mathbf{X}) + (\mathbb{E} (Y | X_i))^2] \\
 &= \sum_{i=1}^n \text{Var} (Y_i | \mathbf{X}) - \sum_{i=1}^n \text{Var} (\hat{Y}_i | \mathbf{X}) = n\sigma^2 - \sum_{i=1}^n h_{ii}\sigma^2 = \left(n - \sum_{i=1}^n h_{ii} \right) \sigma^2 \\
 &= (n - \text{tr} \mathbf{P}) \sigma^2 = (n - 2) \sigma^2.
 \end{aligned}$$

这样, 我们同样可以推导出校正偏差的 s^2 作为 σ^2 .

s^2 是最广泛使用的回归方差的无偏估计量, 除此之外, 还有一种被广泛使用的 σ^2 的无偏估计量则是使用标准化残差还代替矩估计 σ_{ME}^2 中不可观察的 ε_i . 由 $\mathbb{E}(\hat{\varepsilon}_i | \mathbf{X}) = 0$ 以及 $\text{Var}(\hat{\varepsilon}_i | \mathbf{X}) = (1 - h_{ii}) \sigma^2$, 记

$$\bar{\varepsilon}_i = \frac{\hat{\varepsilon}_i - 0}{\sqrt{1 - h_{ii}}},$$

则称 $\bar{\varepsilon}_i$ 为标准化残差 (standardized residual). 显然 $\bar{\varepsilon}_i$ 满足有 $\mathbb{E}(\bar{\varepsilon}_i | \mathbf{X}) = 0$ 及

$$\text{Var}(\bar{\varepsilon}_i | \mathbf{X}) = \mathbb{E}(\bar{\varepsilon}_i^2 | \mathbf{X}) = \mathbb{E} \left(\frac{\hat{\varepsilon}_i^2}{1 - h_{ii}} \middle| \mathbf{X} \right) = \frac{\mathbb{E}(\hat{\varepsilon}_i^2 | \mathbf{X})}{1 - h_{ii}} = \frac{(1 - h_{ii}) \sigma^2}{1 - h_{ii}} = \sigma^2,$$

这样对于 $\sum_{i=1}^n \bar{\varepsilon}_i^2 = \sum_{i=1}^n \frac{\hat{\varepsilon}_i^2}{1 - h_{ii}}$ 便有

$$\mathbb{E} \left(\sum_{i=1}^n \bar{\varepsilon}_i^2 \middle| \mathbf{X} \right) = \sum_{i=1}^n \mathbb{E}(\bar{\varepsilon}_i^2 | \mathbf{X}) = \sum_{i=1}^n \sigma^2 = n\sigma^2,$$

容易知道, 使用标准化残差代替回归误差的矩估计

$$\bar{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \bar{\varepsilon}_i^2 = \frac{1}{n} \sum_{i=1}^n \frac{\hat{\varepsilon}_i^2}{1 - h_{ii}}$$

是回归方差 σ^2 的无偏估计.

2.7 一元 OLS 估计量的区间估计与假设检验

本节将介绍一元正态回归模型, 其为一元线性回归模型的一种特殊情形. 在之前的模型 (2.4.1) 中, 我们尚无设定被解释变量 Y 以及回归误差 ε 的概率分布, 而在本节我们将研究具有正态分布下的线性回归模型——有了概率分布, 我们便可以引入区间估计以及假设检验.

模型 2.7.1 (一元正态回归模型). 在一元线性回归模型 (2.4.1) 的假设之上, 补充假设

$$\varepsilon|X \sim N(0, \sigma^2) \quad (2.35)$$

成立时, 满足如下设定的模型称为一元正态回归模型 (*normal regression model*).

$$\begin{aligned} Y &= \mathcal{P}(Y|X) + \varepsilon \\ \mathcal{P}(Y|X) &= \beta_0 + \beta_1 X, \quad \begin{cases} \beta_0 = \mathbb{E}Y - \beta_1 \mathbb{E}X \\ \beta_1 = \frac{\text{Cov}(X, Y)}{\text{Var}(X)} \end{cases} \\ \mathbb{E}(X\varepsilon) &= 0 \\ \mathbb{E}(\varepsilon) &= 0 \quad (\text{不需要常数项 } \beta_0) \\ \text{Var}(\varepsilon) &= \mathbb{E}\varepsilon^2 = \mathbb{E}(\varepsilon|X) = \sigma^2 \end{aligned}$$

一元正态回归模型较模型 (2.4.1) 的差别, 仅在于假设了回归误差 ε 在给定解释变量 X 下的条件分布服从正态分布, 而式 (2.35) 通常被称为正态假设 (normality assumption). 总体的回归误差 $\varepsilon|X \sim N(0, \sigma^2)$, 这里的正态分布 $N(0, \sigma^2)$ 是不依赖于解释变量 X , 也就是说条件随机变量 $\varepsilon|X$ 的各种性质均与 X 无关. 这样, 简单随机样本下样本个体的回归误差 ε_i 独立同分布, 且也满足 $\varepsilon_i|X_i \sim N(0, \sigma^2)$, 因而对于误差向量便有 $\boldsymbol{\varepsilon}|\mathbf{X} \sim N(0, \sigma^2 \mathbf{I}_n)$, $\boldsymbol{\varepsilon}|\mathbf{X}$ 的联合分布不受到 \mathbf{X} 的影响.

于是对于一元正态回归模型 (2.7.1), 同方差假设

$$\mathbb{E}(\varepsilon_i^2|\mathbf{X}) = \mathbb{E}\varepsilon_i^2 = \sigma^2$$

直接成立.

因此, 在引入条件分布服从正态分布的回归误差后, 模型 (2.7.1) 的总体和样本便是

$$\text{总体: } Y = \beta_0 + \beta_1 X + \varepsilon, \quad \varepsilon|X \sim N(0, \sigma^2),$$

$$\text{样本: } \mathbf{Y} = \beta_0 \mathbf{1} + \beta_1 \mathbf{X} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon}|\mathbf{X} \sim N(0, \sigma^2 \mathbf{I}_n),$$

$$\text{样本个体: } Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i, \quad \varepsilon_i|X_i \text{ i.i.d. } \sim N(0, \sigma^2), i = 1, \dots, n,$$

一元正态回归模型至少假定了回归误差的条件分布, 而对于研究的变量 (Y, X) 并没有限定其联合分布. 如果在模型 (2.4.1) 上假定了 (Y, X) 服从联合正态分布, 那么得到的模型则称为线性回归模型的经典模型.

模型 2.7.2 (一元经典线性回归模型). 在模型 (2.4.1) 的假设中加入 (Y, X) 服从联合正态分布, 则称该模型为一元经典线性回归模型 (*classic linear regression model*).

由于回归误差 ε 是 (Y, X) 的线性变换, 模型 (2.7.2) 中 (Y, X) 服从联合正态分布, 则可以得到 ε 服从正态分布 $N(0, \sigma^2)$. 同时在线性回归模型中, X 与 ε 不相关, 对于两个服从正态分布的随机变量而言不相关等价于独立, 因此则模型 (2.7.2) 也满足有 $\varepsilon|X \sim N(0, \sigma^2)$. 故一元经典线性回归模型 (2.7.2) 只是一元正态回归模型 (2.7.1) 的特例情形.

一元正态回归模型给定了回归误差的条件分布, 现在我们可以推导其他随机变量的条件分布. 显然, 被解释变量 Y 的总体具有条件分布 $Y|X \sim N(\beta_0 + \beta_1 X, \sigma^2)$, 相应地 Y 的样本个体具有条件分布 $Y_i|X \sim N(\beta_0 + \beta_1 X_i, \sigma^2)$.

对于参数的 OLS 估计量

$$\hat{\beta}_1 = A(\mathbf{X}) \sum_{i=1}^n B_i(\mathbf{X}) Y_i, \quad \hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X},$$

给定样本 $\mathbf{X} = (X_1, \dots, X_n)^T$ 的情形下, 则 $\hat{\beta}_0$ 与 $\hat{\beta}_1$ 均是 Y_1, \dots, Y_n 的线性组合, 而正态分布的线性组合仍是线性组合, 依据此前章节计算的条件期望和条件方差, 这样 OLS 估计量便有

$$\hat{\beta}_1 | \mathbf{X} \sim N(\beta_1, A\sigma^2), \quad \hat{\beta}_0 | \mathbf{X} \sim N\left(\beta_0, \left(\frac{1}{n} + A\bar{X}^2\right)\sigma^2\right). \quad (2.36)$$

同理可以推知, 拟合值 \hat{Y}_i 与残差 $\hat{\varepsilon}_i$ 具有条件正态分布, 即

$$\hat{Y}_i | \mathbf{X} \sim N(\beta_0 + \beta_1 X_i, h_{ii}\sigma^2), \quad \hat{\varepsilon}_i | \mathbf{X} \sim N(0, (1 - h_{ii})\sigma^2). \quad (2.37)$$

我们现在介绍回归方差估计量 $s^2 = \frac{1}{n-2} \sum_{i=1}^n \hat{\varepsilon}_i^2$ 的条件分布, 利用线性代数的知识可以证明

$$\frac{(n-2)s^2}{\sigma^2} | \mathbf{X} \sim \chi^2(n), \quad (2.38)$$

但这个证明并不算轻松, 我们将在下一章的多元线性回归模型里去详细介绍. 在此, 我们给出另一种 RSS 的分解以说明式 (2.38) 成立.

首先, 对于样本个体可以推知

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i \Rightarrow \bar{Y} = \beta_0 + \beta_1 \bar{X} + \bar{\varepsilon} \Rightarrow Y_i - \bar{Y} = \beta_1 (X_i - \bar{X}) + (\varepsilon_i - \bar{\varepsilon}),$$

其中

$$\bar{\varepsilon} = \frac{1}{n} \sum_{i=1}^n \varepsilon_i$$

为回归误差 ε_i 的样本均值 (不可观察). 同样使用前面的符号记离差 $B_i = X_i - \bar{X}, i = 1, \dots, n$, 则有模型 (2.4.1) 即有消去截距 β_0 的**中心化**形式

$$Y_i - \bar{Y} = \beta_1 B_i + (\varepsilon_i - \bar{\varepsilon}). \quad (2.39)$$

我们再注意到 OLS 估计量 $\hat{\beta}_1$ 其较参数 β_1 的差值为

$$\begin{aligned} \hat{\beta}_1 - \beta_1 &= \beta_1 - A \sum_{i=1}^n B_i Y_i = A \sum_{i=1}^n B_i (\beta_1 B_i + (\varepsilon_i - \bar{\varepsilon}) + \bar{Y}) - \beta_1 \\ &= A \beta_1 \sum_{i=1}^n B_i^2 + A \sum_{i=1}^n B_i (\varepsilon_i - \bar{\varepsilon}) - \beta_1 = A \beta_1 A^{-1} + A \sum_{i=1}^n B_i (\varepsilon_i - \bar{\varepsilon}) - \beta_1 \\ &= A \sum_{i=1}^n B_i (\varepsilon_i - \bar{\varepsilon}), \end{aligned}$$

即 $\hat{\beta}_1 - \beta_1$ 是回归误差 $\varepsilon_1, \dots, \varepsilon_n$ 的线性组合, 而 $\varepsilon_i | X_i \text{ i.i.d. } \sim N(0, \sigma^2), i = 1, \dots, n$, 于是

$$\hat{\beta}_1 - \beta_1 | \mathbf{X} \sim N(0, A\sigma^2),$$

当然, $\hat{\beta}_1 - \beta_1$ 的条件分布也可以用式 (2.36) 得到, 我们这里是为了得到 $\hat{\beta}_1 - \beta_1$ 的表达式. 一般地, 我们称参数的估计量与参数的差值为**抽样误差** (sampling error).

现在我们利用式 (2.39) 对 $\text{RSS} = \sum_{i=1}^n \hat{\varepsilon}_i^2$ 进行分解, 不难算得

$$\begin{aligned} \sum_{i=1}^n \hat{\varepsilon}_i^2 &= \sum_{i=1}^n \left[(\beta_1 - \hat{\beta}_1) B_i + (\varepsilon_i - \bar{\varepsilon}) \right]^2 \\ &= \sum_{i=1}^n \left[(\beta_1 - \hat{\beta}_1)^2 B_i^2 - 2(\hat{\beta}_1 - \beta_1) B_i (\varepsilon_i - \bar{\varepsilon}) + (\varepsilon_i - \bar{\varepsilon})^2 \right] \\ &= A^{-1} (\beta_1 - \hat{\beta}_1)^2 + 2(\beta_1 - \hat{\beta}_1) \sum_{i=1}^n B_i (\varepsilon_i - \bar{\varepsilon}) + \sum_{i=1}^n (\varepsilon_i - \bar{\varepsilon})^2 \\ &= A^{-1} (\beta_1 - \hat{\beta}_1)^2 - 2A^{-1} (\beta_1 - \hat{\beta}_1)^2 + \sum_{i=1}^n (\varepsilon_i - \bar{\varepsilon})^2 \\ &= \sum_{i=1}^n (\varepsilon_i - \bar{\varepsilon})^2 - A^{-1} (\beta_1 - \hat{\beta}_1)^2, \end{aligned}$$

这里用到了 $A^{-1} = \sum_{i=1}^n B_i^2$, 故有

$$\sum_{i=1}^n (\varepsilon_i - \bar{\varepsilon})^2 = \sum_{i=1}^n \hat{\varepsilon}_i^2 + A^{-1} (\beta_1 - \hat{\beta}_1)^2. \quad (2.40)$$

而对于样本 2 阶中心距, 我们熟知 $\sum_{i=1}^n (\varepsilon_i - \bar{\varepsilon})^2 = \sum_{i=1}^n \varepsilon_i^2 - n\bar{\varepsilon}^2$, 代入式 (2.40) 中整理得

$$\sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n \hat{\varepsilon}_i^2 + A^{-1} (\beta_1 - \hat{\beta}_1)^2 + n\bar{\varepsilon}^2, \quad (2.41)$$

这即是对于误差平方和 SSE 的分解. 我们对于式 (2.41) 两边同除以 σ^2 , 整理即得到

$$\sum_{i=1}^n \left(\frac{\varepsilon_i}{\sigma} \right)^2 = \frac{1}{\sigma^2} \sum_{i=1}^n \hat{\varepsilon}_i^2 + \left(\frac{\beta_1 - \hat{\beta}_1}{\sqrt{A\sigma^2}} \right)^2 + \left(\frac{\bar{\varepsilon}}{\sqrt{\sigma^2/n}} \right)^2. \quad (2.42)$$

不难发现, 式 (2.42) 中回归误差 ε_i 、抽样误差 $\hat{\beta}_1 - \beta_1$ 和样本均值 $\bar{\varepsilon}$ 均进行了标准化处理. 在没有正态假设下, 式 (2.40)、(2.41) 及 (2.42) 也是成立的, 且可以证明式 (2.42) 右边三项是互相独立的. 这样, 对式 (2.42) 求条件期望, 也可以由此推导偏差校正估计量 s^2 . 而加入了正态假设, 回归误差 ε_i 、抽样误差 $\hat{\beta}_1 - \beta_1$ 和样本均值 $\bar{\varepsilon}$ 这三者的条件分布均服从正态分布, 标准化再平方后则有

$$\sum_{i=1}^n \left(\frac{\varepsilon_i}{\sigma} \right)^2 \Big| \mathbf{X} \sim \chi^2(n), \quad \left(\frac{\beta_1 - \hat{\beta}_1}{\sqrt{A\sigma^2}} \right)^2 \Big| \mathbf{X} \sim \chi^2(1), \quad \left(\frac{\bar{\varepsilon}}{\sqrt{\sigma^2/n}} \right)^2 \Big| \mathbf{X} \sim \chi^2(1)$$

成立. 这样式 (2.42) 的左边是自由度为 n 的卡方分布, 右边后两项都是自由度为 1 的卡方分布成立, 利用概率论中矩生成函数 (moment generating function, 见第 4.3 节) 可以证明

$$\frac{1}{\sigma^2} \sum_{i=1}^n \hat{\varepsilon}_i^2 \Big| \mathbf{X} \sim \chi^2(n-2),$$

也即证明了式 (2.38).

现在我们便得到了正态假设下 OLS 估计量下统计量的概率分布, 由此我们即可进行区间估计和假设检验.

定义 2.7.1 (一元正态回归模型的 t 统计量). 对于模型 (2.7.1) 的 OLS 估计量 $\hat{\beta}_0$ 和 $\hat{\beta}_1$, 其标准化后有

$$\frac{\beta_1 - \hat{\beta}_1}{\sqrt{A\sigma^2}} \Big| \mathbf{X} \sim N(0, 1), \quad \frac{\beta_0 - \hat{\beta}_0}{\sqrt{\left(\frac{1}{n} + A\bar{X}^2\right)\sigma^2}} \Big| \mathbf{X} \sim N(0, 1),$$

由于 σ^2 为未知的厌恶参数 (nuisance parameter), 使用 σ^2 的偏差校正估计量 s^2 代替, 记

$$\text{SE}(\hat{\beta}_1) = \sqrt{As^2}, \quad \text{SE}(\hat{\beta}_0) = \sqrt{\left(\frac{1}{n} + A\bar{X}^2\right)s^2}, \quad (2.43)$$

我们称式 (2.43) 为 OLS 估计量的标准误差 (standard error) 或标准误, 则我们称

$$T = \frac{\beta_i - \hat{\beta}_i}{\text{SE}(\hat{\beta}_i)} \quad (2.44)$$

为 t 统计量 (t -statistic).

对于 $\hat{\beta}_1$ 的 t 统计量, 由于

$$T = \frac{\beta_1 - \hat{\beta}_1}{\text{SE}(\hat{\beta}_1)} = \frac{\beta_1 - \hat{\beta}_1}{\sqrt{A\sigma^2}} \frac{\sqrt{A\sigma^2}}{\text{SE}(\hat{\beta}_1)} = \frac{\beta_1 - \hat{\beta}_1}{\sqrt{A\sigma^2}} \bigg/ \sqrt{\frac{(\text{SE}(\hat{\beta}_1))^2}{A\sigma^2}},$$

而

$$\frac{\beta_1 - \hat{\beta}_1}{\sqrt{A\sigma^2}} \bigg| \mathbf{X} \sim N(0, 1), \quad \frac{(\text{SE}(\hat{\beta}_1))^2}{A\sigma^2} \bigg| \mathbf{X} = \frac{(n-2)s^2/\sigma^2}{n-2} \bigg| \mathbf{X} \sim \frac{\chi^2(n-2)}{n-2},$$

故依据数理统计的知识, 在给定 \mathbf{X} 的条件下 $\hat{\beta}_1$ 的 t 统计量服从自由度为 $n-2$ 的 t 分布. 同理可以证明 $\hat{\beta}_0$ 的 t 统计量 T 也满足有 $T | \mathbf{X} \sim t(n-2)$.

利用 t 统计量作为枢轴变量 (pivotal variable), 我们便可以对总体参数 β 进行区间估计. 对于参数 β , 若存在区间 \hat{I} 使得

$$\mathbb{P}(\beta \in \hat{I}) = 1 - \alpha$$

成立, 则称区间估计量 \hat{I} 是参数 β 的置信水平 (confidence level) 为 $1 - \alpha$ 的置信区间 (confidence interval). 对于概率 $\mathbb{P}(\beta \in \hat{I})$ 的理解, 应该是将置信区间 \hat{I} 视为随机的, 而总体参数 β 应该视为“已知”或固定的, 这样, 置信水平为 $1 - \alpha$ 的置信区间 \hat{I} 意味着重复抽样下区间 \hat{I} 含义参数 β 的真实值的概率为 $1 - \alpha$.

在计量经济学中, 对于回归系数 β 而言我们通常构造如下的置信区间:

$$\hat{I} = [\hat{\beta}_i - c \cdot \text{SE}(\hat{\beta}_i), \hat{\beta}_i + c \cdot \text{SE}(\hat{\beta}_i)],$$

其中 c 是某个大于 0 的常数. 如果 \hat{I} 的置信水平为 $1 - \alpha$, 则有

$$\begin{aligned} \mathbb{P}(\hat{\beta}_i - c \cdot \text{SE}(\hat{\beta}_i) \leq \beta \leq \hat{\beta}_i + c \cdot \text{SE}(\hat{\beta}_i)) &= 1 - \alpha \\ \Leftrightarrow \mathbb{P}\left(-c \leq \frac{\beta - \hat{\beta}_i}{\text{SE}(\hat{\beta}_i)} \leq c\right) &= 1 - \alpha, \end{aligned}$$

也即是 t 统计量满足

$$\mathbb{P}(-c \leq t \leq c) = F(c) - F(-c) \stackrel{\text{对称性}}{=} F(c) - (1 - F(c)) = 2F(c) - 1 = 1 - \alpha,$$

即有 $F(c) = 1 - \frac{\alpha}{2}$, 这里 F 是 t 分布的累计密度函数. 对于连续的累计密度函数函数 F , 其反函数 F^{-1} 一定是存在的, 故 $c = F^{-1}\left(1 - \frac{\alpha}{2}\right)$. 在数理统计中我们已经学过, 这里的 c 即是上 $\alpha/2$ 分位数 $t_{\alpha/2}$.

一条有用的经验法则是, 当 $n - k > 61$ (一元 OLS 即是 $k = 2$) 时, 参数 β 的 95% 水平的置信区间近似为 $\left[\hat{\beta}_i - 2\text{SE}(\hat{\beta}_i), \hat{\beta}_i + 2\text{SE}(\hat{\beta}_i)\right]$.

同样, 以卡方分布作为枢轴变量, 考虑回归方差 σ^2 的置信水平为 $1 - \alpha$ 的置信区间 \hat{I} 有

$$\begin{aligned}\mathbb{P}\left(c_1 \leq \frac{(n-2)s^2}{\sigma^2} \leq c_2\right) &= F(c_2) - F(c_1) = 1 - \alpha \\ \Leftrightarrow \mathbb{P}\left(\frac{(n-2)s^2}{c_2} \leq \sigma^2 \leq \frac{(n-2)s^2}{c_1}\right) &= 1 - \alpha,\end{aligned}$$

即 σ^2 的置信区间为

$$\hat{I} = \left[\frac{(n-2)s^2}{c_2}, \frac{(n-2)s^2}{c_1}\right],$$

其中 $c_1 = F^{-1}(\alpha/2)$, $c_2 = F^{-1}(1 - \alpha/2)$, 也即是卡方分布的上 $\alpha/2$ 分位数和上 $1 - \alpha/2$ 分位数.

关于假设检验问题, 请见第节.

2.8 一元 OLS 估计量的渐进性质

简单介绍什么是“大样本”性质, 具体参加后面章节.

2.9 代码

```
1 %% OLS估计量无偏性
2 %设定E(Y|X)=10x+25,e~N(0,5^2)
3 %考虑重复抽样k次,单次样本容量为n
4 clear
5 clc
6 k=1000;%抽样次数
7 n=50;%单次抽样容量
8 s=RandStream('mcg16807','Seed',114514);%设定复现的随机流以及种子
9 Beta=zeros(2,k);%存储每次估计的系数beta
10 for i=1:k
11 X=[ones(n,1),randsample(s,-100:100,n)'];%生成样本X
12 e=5*randn(s,n,1);%生成扰动项
13 Y=X*[25;10]+e;%数据生成过程
```

```

14 beta=X\Y;%OLS估计量
15 Beta(:,i)=beta;%保存OLS估计量
16 end
17 beta_hat=cumsum(Beta,2)./(1:k);%OLS估计的平均值
18 selfGrootDefault(2)%绘图美化函数,没有请注释
19 %绘制估计值均值变化
20 figure(1)
21 set(gcf,'unit','centimeters',...
22 'Position',[0 0 21*1.25 29.7*0.35*1.25]);%设置figure为A4纸宽度
23 tiledlayout(2,2);%创建2*2画布
24 nexttile
25 histfit(Beta(1,:))%beta0的直方图和拟合的pdf
26 xline(25,'-.b','LineWidth',2)
27 xlabel('\beta_0','FontSize',12.5)
28
29 nexttile
30 plot(1:25:k,beta_hat(1,1:25:k),'-o')
31 xlabel('\fontname{宋体}重复抽样次数\fontname{Times New ...
    Roman}k','FontSize',12.5)
32 ylabel('\beta_0\fontname{宋体}估计值的平均值','FontSize',12.5);
33 yline(25,'-.','\fontname{宋体}真实值','LineWidth',2,'LabelHorizontalAlignment' ...
    ...
34 , 'left','LabelVerticalAlignment','bottom')
35
36 nexttile
37 histfit(Beta(2,:))%beta1的直方图和拟合的pdf
38 xline(10,'-.b','LineWidth',2)
39 xlabel('\beta_1','FontSize',12.5)
40
41 nexttile
42 plot(1:25:k,beta_hat(2,1:25:k),'-o')
43 xlabel('\fontname{宋体}重复抽样次数\fontname{Times New ...
    Roman}k','FontSize',12.5)
44 ylabel('\beta_0\fontname{宋体}估计值的平均值','FontSize',12.5);
45 yline(10,'-.','\fontname{宋体}真实值','LineWidth',2,'LabelHorizontalAlignment' ...
    ...
46 , 'left','LabelVerticalAlignment','bottom')
47 %绘制重复抽样pdf
48 figure(2)
49 set(gca,'FontName','Times New Roman');
50 set(gcf,'unit','centimeters',...
51 'Position',[0 0 21*1.25 29.7*0.4*1.25]);%设置figure为A4纸宽度

```



```

52 tiledlayout(2,1);%创建2*2画布
53 nexttile
54 pd1=fitdist(Beta(1,1:25)', 'Normal');%beta0拟合的pdf_25次抽样
55 pd2=fitdist(Beta(1,1:100)', 'Normal');%beta0拟合的pdf_100次抽样
56 pd3=fitdist(Beta(1,1:k)', 'Normal');%beta0拟合的pdf_k次抽样
57 x_values=20:0.1:30;
58 y1=pdf(pd1, x_values);%计算pdf_25次抽样
59 y2=pdf(pd2, x_values);%计算pdf_100次抽样
60 y3=pdf(pd3, x_values);%计算pdf_k次抽样
61 plot(x_values, y1, '-. ')
62 hold on
63 plot(x_values, y2, '-* ')
64 plot(x_values, y3, '-o ')
65 xline(25, '-.b', 'LineWidth', 2, 'DisplayName', '\beta_0 ...
    \fontname{宋体}真实值')
66 xlabel('\beta_0', 'FontSize', 12.5)
67 legend('\fontname{Times New ...
    Roman}25\fontname{宋体}次抽样', '\fontname{Times New ...
    Roman}100\fontname{宋体}次抽样', '\fontname{Times New ...
    Roman}1000\fontname{宋体}次抽样')
68 xlim([23, 27])
69 ylim([0, 0.65])
70
71 nexttile
72 pd1=fitdist(Beta(2,1:25)', 'Normal');%beta1拟合的pdf_25次抽样
73 pd2=fitdist(Beta(2,1:100)', 'Normal');%beta1拟合的pdf_100次抽样
74 pd3=fitdist(Beta(2,1:k)', 'Normal');%beta1拟合的pdf_k次抽样
75 x_values=9.5:0.002:10.5;
76 y1=pdf(pd1, x_values);%计算pdf_25次抽样
77 y2=pdf(pd2, x_values);%计算pdf_100次抽样
78 y3=pdf(pd3, x_values);%计算pdf_k次抽样
79 plot(x_values, y1, '-. ')
80 hold on
81 plot(x_values, y2, '-* ')
82 plot(x_values, y3, '-o ')
83 xline(10, '-.b', 'LineWidth', 2, 'DisplayName', '\beta_1 ...
    \fontname{宋体}真实值')
84 xlabel('\beta_1', 'FontSize', 12.5)
85 legend('\fontname{Times New ...
    Roman}25\fontname{宋体}次抽样', '\fontname{Times New ...
    Roman}100\fontname{宋体}次抽样', '\fontname{Times New ...
    Roman}1000\fontname{宋体}次抽样')

```

```
86 xlim([9.96,10.04])  
87 ylim([0,44])
```

第三章 多元线性回归模型

本章我们将一元线性回归模型拓展至多元的情况, 即从单个解释变量 X 拓展至 K 个解释变量 $\mathbf{X} = (X_1, \dots, X_K)^T$. 本章内容结构上与上一节并无多大差异, 利用矩阵符号, 一元 OLS 的各类问题都可以自然对应至多元情形, 且有着更为更一般的视角.

3.1 多元线性回归模型的矩阵符号

本节介绍多元线性回归模型将使用的矩阵符号. 我们所研究的被解释变量依旧记为 Y , 而感兴趣的解释变量则共有 K 个, 分别是随机变量 X_1, \dots, X_K , 我们用随机向量 $\mathbf{X} = (X_1, \dots, X_K)^T$ 一概表示.

本章考虑的依旧是定义 (2.2.1) 的线性 CEF, 我们使用与模型 (2.2.1) 相同的矩阵符号, 即记回归系数为 $K \times 1$ 的列向量 $\boldsymbol{\beta} = (\beta_1, \dots, \beta_K)^T$, 那么线性 CEF 可以用向量内积表示为

$$m(\mathbf{X}) = \mathbb{E}(Y|\mathbf{X}) = \beta_1 X_1 + \dots + \beta_K X_K = \mathbf{X}^T \boldsymbol{\beta}.$$

在上一章中, 我们区别了一元投影模型 (2.2.3) 以及一元线性回归模型 (2.4.1) 的区别. 相应地, 我们将这两个模型拓展为多元情形. 对于投影模型, 给定 \mathbf{X} 下 Y 的最优线性预测量记为

$$\mathcal{P}(Y|\mathbf{X}) = \beta_1 X_1 + \dots + \beta_n X_K = \mathbf{X}^T \boldsymbol{\beta},$$

则投影模型方程为 $Y = \mathcal{P}(Y|\mathbf{X}) + \varepsilon$, 其中 ε 为投影误差. 当回归模型 $Y = \mathbb{E}(Y|\mathbf{X}) + e$ 中条件期望函数确实为线性 CEF 时, 投影模型即为线性回归模型, 则是投影误差 ε 等价于回归误差 e . 这样, 线性回归模型的总体方程为

$$Y = \beta_1 X_1 + \dots + \beta_n X_K + \varepsilon = \mathbf{X}^T \boldsymbol{\beta} + \varepsilon. \quad (3.1)$$

现在考虑样本容量为 n 的抽样样本, 对于样本个体 $i = 1, \dots, n$, 则其线性回归方程为

$$Y_i = \beta_1 X_{i1} + \dots + \beta_n X_{iK} + \varepsilon_i = \mathbf{X}_i^T \boldsymbol{\beta} + \varepsilon_i, \quad (3.2)$$

其中 $\mathbf{X}_i = (X_{i1}, \dots, X_{iK})^T$ 为个体 i 对应的被解释变量 \mathbf{X} . 我们将 n 个样本回归方程以线性方程组的形式写出, 即有

$$\begin{cases} Y_1 = \beta_1 X_{11} + \dots + \beta_n X_{K1} + \varepsilon_1, \\ \dots \\ Y_i = \beta_1 X_{1i} + \dots + \beta_n X_{Ki} + \varepsilon_i, \\ \dots \\ Y_n = \beta_1 X_{1n} + \dots + \beta_n X_{Kn} + \varepsilon_n, \end{cases}$$

依据线性代数的知识, 我们将线性方程组以使用矩阵符号书写为

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad (3.3)$$

其中

$$\mathbf{Y}_{n \times 1} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_i \\ \vdots \\ Y_n \end{pmatrix}, \quad \mathbf{X}_{n \times K} = \begin{pmatrix} X_{11} & X_{12} & \dots & X_{1K} \\ \vdots & \vdots & & \vdots \\ X_{i1} & X_{i2} & \dots & X_{iK} \\ \vdots & \vdots & & \vdots \\ X_{n1} & X_{n2} & \dots & X_{nK} \end{pmatrix}, \quad \boldsymbol{\varepsilon}_{n \times 1} = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_i \\ \vdots \\ \varepsilon_n \end{pmatrix}.$$

通常而言, n 维列向量 \mathbf{Y} 被称为**数据向量** (data vector), 而 $n \times k$ 维矩阵 \mathbf{X} 被称为**数据矩阵** (data matrix). 数据矩阵 \mathbf{X} 满足有

$$\mathbf{X}_{n \times K} = \begin{pmatrix} \mathbf{X}_1^T \\ \vdots \\ \mathbf{X}_i^T \\ \vdots \\ \mathbf{X}_n^T \end{pmatrix} = (\mathbf{X}_1, \dots, \mathbf{X}_i, \dots, \mathbf{X}_n)^T,$$

这用矩阵的分块运算

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \Leftrightarrow \begin{pmatrix} Y_1 \\ \vdots \\ Y_i \\ \vdots \\ Y_n \end{pmatrix} = \begin{pmatrix} \mathbf{X}_1^T \boldsymbol{\beta} \\ \vdots \\ \mathbf{X}_i^T \boldsymbol{\beta} \\ \vdots \\ \mathbf{X}_n^T \boldsymbol{\beta} \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_i \\ \vdots \\ \varepsilon_n \end{pmatrix} \Leftrightarrow \begin{pmatrix} Y_1 \\ \vdots \\ Y_i \\ \vdots \\ Y_n \end{pmatrix} = \begin{pmatrix} \mathbf{X}_1^T \\ \vdots \\ \mathbf{X}_i^T \\ \vdots \\ \mathbf{X}_n^T \end{pmatrix} \boldsymbol{\beta} + \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_i \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

是显而易见的.

多元线性回归模型的总体方程 (3.1) 不含有截距, 但我们只需令 $X_1 \equiv 1$, 相应地样本个体满足 $\mathbf{X}_i = (1, X_{i2}, \dots, X_{iK})^T$, 即得到了含有截距的线性回归模型为

$$\text{总体: } Y = \beta_1 + \beta_2 X_2 + \dots + \beta_n X_K + \varepsilon,$$

$$\text{样本个体: } Y_i = \beta_1 + \beta_2 X_{i2} + \dots + \beta_n X_{iK} + \varepsilon_i,$$

此时数据矩阵 \mathbf{X} 满足有

$$\mathbf{X}_{n \times K} = \begin{pmatrix} 1 & X_{12} & \cdots & X_{1K} \\ 1 & X_{22} & \cdots & X_{2K} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & X_{n2} & \cdots & X_{nK} \end{pmatrix},$$

因此, 式 (3.3) 便已经包含了含有截距的线性回归模型.

3.2 非简单随机样本的模型假设

在一元线性回归模型中, 我们选择了简单随机样本以推导 OLS 估计量, 但简单随机样本并非是一种普遍的情形. 就线性回归模型而言, 计量经济学对于样本的假设更为一般化.

假设 3.2.1 (线性回归模型的一般假设). 对于线性回归模型, 我们假设样本服从如下假设:

(a) **线性假设** (*linearity assumption*): 对于样本个体 i , 给定样数据矩阵 \mathbf{X} 下 Y_i 条件期望函数为定义 (2.2.1) 的线性函数, 即有

$$\mathbb{E}(Y_i | \mathbf{X}) = \mathbf{X}_i^T \boldsymbol{\beta}, \quad i = 1, \dots, n; \quad (3.4)$$

(b) **严格外生性** (*strict exogeneity*): 样本个体的回归误差 ε_i 均值独立于数据矩阵 \mathbf{X} , 即有

$$\mathbb{E}(\varepsilon_i | \mathbf{X}) = 0, \quad i = 1, \dots, n; \quad (3.5)$$

(c) **无多重共线性** (*no multicollinearity*): 数据矩阵 \mathbf{X} 以概率为 1 列满秩, 即

$$\mathbb{P}(\text{rank}(\mathbf{X}_{n \times K}) = K) = 1; \quad (3.6)$$

(d) **球形误差方差** (*spherical error variance*): 样本个体的回归误差 ε_i 满足条件同方差 (*conditional homoskedasticity*) 假设

$$\mathbb{E}(\varepsilon_i^2 | \mathbf{X}) = \sigma^2 > 0, \quad i = 1, \dots, n \quad (3.7)$$

和不存在自相关 (autocorrelation) 或序列相关 (serial correlation), 即

$$\mathbb{E}(\varepsilon_i \varepsilon_j | \mathbf{X}) = 0, \quad i, j = 1, \dots, n \text{ 且 } i \neq j. \quad (3.8)$$

我们现在将对假设 (3.2.1) 的四条假设逐一分析.

对于线性假设式 (3.4), 此假设使得线性投影模型等价于线性回归模型. 可以说假设 (3.2.1)(a) 是线性回归模型相关推导的基础, 否则这些推导都只能归类至线性投影模型 (尽管对于线性回归模型是成立的).

对于严格外生性式 (3.5), 利用简单迭代期望定律 (2.1.1) 便有

$$\mathbb{E}\varepsilon_i = \mathbb{E}(\mathbb{E}(\varepsilon_i | \mathbf{X})) = 0, \quad i = 1, \dots, n$$

成立. 这条假设则是对应了总体的回归误差 ε 的性质 (见定理 (2.1.4)). 进一步, 我们由严格外生性可以推出回归误差 ε_i 与任意的样本观测值 X_{ij} 是完全不相关的, 这是因为对于索引 $i, j = 1, \dots, n, k = 1, \dots, K$, 依据迭代期望定律 (2.1.2) 有

$$\mathbb{E}(\varepsilon_i | X_{jk}) = \mathbb{E}(\mathbb{E}(\varepsilon_i | \mathbf{X}) | X_{jk}) = \mathbb{E}(0 | X_{jk}) = 0,$$

再依据定理 (2.1.3) 有

$$\mathbb{E}(X_{jk}\varepsilon_i) = \mathbb{E}(X_{jk}\mathbb{E}(\varepsilon_i | X_{jk})) = \mathbb{E}(X_{jk} \cdot 0) = 0,$$

即回归误差 ε_i 与任意的样本观测值 X_{jk} 是正交的, 故有

$$\text{Cov}(X_{jk}, \varepsilon_i) = \mathbb{E}(X_{jk}\varepsilon_i) - \mathbb{E}(X_{jk})\mathbb{E}(\varepsilon_i) = 0 - \mathbb{E}(X_{jk}) \cdot 0 = 0.$$

特别地, 当 $i = j$ 时 $\text{Cov}(X_{ik}, \varepsilon_i)$ 意味着样本个体 i 自身的各解释变量与回归误差完全无关.

式 (3.4) 和式 (3.5) 用矩阵符号则有 $\mathbb{E}(\mathbf{Y} | \mathbf{X}) = \mathbf{X}\boldsymbol{\beta}$ 和 $\mathbb{E}(\boldsymbol{\varepsilon} | \mathbf{X}) = \mathbf{0}_{n \times 1}$.

对于无多重共线性式 (3.6), \mathbf{X} 列满秩也即以为着其列向量组是线性无关的, 而 \mathbf{X} 的 K 个列向量代表了总体被解释变量 $\mathbf{X} = (X_1, \dots, X_K)^T$ 的各个分量, 也就是说式 (3.6) 假设了是各个解释变量 X_i 互相线性无关以概率为 1 成立. 之所以用概率考虑则是考虑到了解释变量为随机变量. 为了确保式 (3.6) 成立, 则有必要条件 $n \geq \text{rank}(\mathbf{X}_{n \times K}) = K$ 成立.

对于球形误差方差式 (3.7) 和式 (3.8), 这是从协方差矩阵的角度考虑. 对于任意的随机向量 $\mathbf{X} = (X_1, \dots, X_n)^T$, 给定条件 Z 下其条件协方差矩阵 (conditional covariance

matrix) 定义为

$$\begin{aligned}\text{Var}(\mathbf{X}|Z) &= \mathbb{E} \left[(\mathbf{X} - \mathbb{E}(\mathbf{X}|Z)) (\mathbf{X} - \mathbb{E}(\mathbf{X}|Z))^T | Z \right] \\ &= \begin{pmatrix} \text{Var}(X_1|Z) & \text{Cov}(X_1, X_2|Z) & \cdots & \text{Cov}(X_1, X_n|Z) \\ \text{Cov}(X_2, X_1|Z) & \text{Var}(X_2|Z) & \cdots & \text{Cov}(X_2, X_n|Z) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(X_n, X_1|Z) & \text{Cov}(X_n, X_2|Z) & \cdots & \text{Var}(X_n|Z) \end{pmatrix}.\end{aligned}$$

由此误差向量 $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^T$ 的条件协方差为

$$\boxed{\text{Var}(\boldsymbol{\varepsilon}|\mathbf{X})} = \mathbb{E} \left[(\boldsymbol{\varepsilon} - \mathbb{E}(\boldsymbol{\varepsilon}|\mathbf{X})) (\boldsymbol{\varepsilon} - \mathbb{E}(\boldsymbol{\varepsilon}|\mathbf{X}))^T | \mathbf{X} \right] \xrightarrow{\text{式 (3.5)}} \boxed{\mathbb{E}(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^T | \mathbf{X})},$$

这里已经使用了严格外生性假设, 进一步的在式 (3.7) 和式 (3.8) 下误差向量的协方差矩阵为

$$\mathbb{E}(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^T | \mathbf{X}) = \begin{pmatrix} \mathbb{E}(\varepsilon_1^2 | \mathbf{X}) & \mathbb{E}(\varepsilon_1 \varepsilon_2 | \mathbf{X}) & \cdots & \mathbb{E}(\varepsilon_1 \varepsilon_n | \mathbf{X}) \\ \mathbb{E}(\varepsilon_2 \varepsilon_1 | \mathbf{X}) & \mathbb{E}(\varepsilon_2^2 | \mathbf{X}) & \cdots & \mathbb{E}(\varepsilon_2 \varepsilon_n | \mathbf{X}) \\ \vdots & \vdots & \ddots & \vdots \\ \mathbb{E}(\varepsilon_n \varepsilon_1 | \mathbf{X}) & \mathbb{E}(\varepsilon_n \varepsilon_2 | \mathbf{X}) & \cdots & \mathbb{E}(\varepsilon_n^2 | \mathbf{X}) \end{pmatrix} = \begin{pmatrix} \sigma^2 & 0 & \cdots & 0 \\ 0 & \sigma^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma^2 \end{pmatrix},$$

即有 $\text{Var}(\boldsymbol{\varepsilon}|\mathbf{X}) = \sigma^2 \mathbf{I}$, “球形” 便是对此的具象化描述.

现在我们对假设 (3.2.1) 和简单随机样本的下的假设, 由于简单随机样本满足有独立同分布, 正如在前一章的计算一样, 对于样本个体 $i = 1, \dots, n$, 回归误差 ε_i 满足有

$$\begin{aligned}\mathbb{E}(\varepsilon_i | \mathbf{X}) &= \mathbb{E}(\varepsilon_i | \mathbf{X}_i), \\ \mathbb{E}(\varepsilon_i^2 | \mathbf{X}) &= \mathbb{E}(\varepsilon_i^2 | \mathbf{X}_i), \\ \mathbb{E}(\varepsilon_i \varepsilon_j | \mathbf{X}) &= \mathbb{E}(\varepsilon_i | \mathbf{X}_i) \mathbb{E}(\varepsilon_j | \mathbf{X}_j), \quad i \neq j,\end{aligned}$$

同分布保证了上述的期望具有相同的形式, 但具体的取值由 \mathbf{X}_i 的具体观测值 \mathbf{x}_i 所决定. 因此, 严格外生性以及球形误差方差假设仍然需要, 但可以简化. 严格外生性式 (3.5) 退化为 $\mathbb{E}(\varepsilon_i | \mathbf{X}_i) = 0$, 这样无自相关式 (3.8) 便直接成立, 条件同方差式 (3.7) 也可以退化为 $\mathbb{E}(\varepsilon_i^2 | \mathbf{X}_i) = \sigma^2$. 假设 (3.2.1) 中的其余假设仍是必须的.

对于线性回归模型, 假设 (3.2.1) 是 OLS 估计量具有良好性质的必要条件. 这在之后的章节中, 我们也会逐步放宽上述假设要求. 如果严格外生性式 (3.5) 不成立, 则计量经济学模型具有内生性 (endogeneity) 问题——这可以说是实证研究中最常见的问题了——计量经济学家为此建立了反事实的现代因果推断方法, 包括工具变量法、双重差分法、断点回归、合成控制法等等方法, 我们将在后续的章节中介绍这些方法. 球形误

差方差具有式 (3.7) 和 (3.8) 两条假设, 倘若仅有无自相关 (3.8) 成立, 则回归误差 ε_i 存在异方差 (heteroskedasticity) 问题, 这可以写成

$$\mathbb{E}(\varepsilon_i^2 | \mathbf{X}) = \sigma_i^2 > 0, \quad i = 1, \dots, n, \quad (3.9)$$

这时误差向量 $\boldsymbol{\varepsilon}$ 的条件协方差矩阵满足有

$$\mathbb{E}(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^T | \mathbf{X}) = \begin{pmatrix} \sigma_1^2 & 0 & \cdots & 0 \\ 0 & \sigma_2^2 & \cdots & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & \sigma_n^2 \end{pmatrix} = \text{diag}(\sigma_1^2, \dots, \sigma_n^2). \quad (3.10)$$

我们将式 (3.9) 称为条件异方差 (conditional heteroskedasticity), 同时记对角阵 $\mathbf{D} = \text{diag}(\sigma_1^2, \dots, \sigma_n^2)$. 对角阵 \mathbf{D} 被称为条件异方差矩阵. 容易知道, 条件同方差式 (3.7) 实际上是条件异方差 (3.9) 的特殊情形, 即 $\mathbf{D} = \sigma^2 \mathbf{I}$.

本节的最后, 我们介绍条件协方差矩阵的有用性质.

定理 3.2.1 (条件协方差矩阵的性质). 对于任意的随机向量 $\mathbf{X} = (X_1, \dots, X_n)^T$, 给定条件 Z 下, 则其协方差矩阵 $\text{Var}(\mathbf{X} | Z)$ 满足如下性质:

- (a) 协方差矩阵计算公式为 $\text{Var}(\mathbf{X} | Z) = \mathbb{E}(\mathbf{X}\mathbf{X}^T | Z) - \mathbb{E}(\mathbf{X} | Z)\mathbb{E}(\mathbf{X} | Z)^T$,
- (b) 条件协方差矩阵 $\text{Var}(\mathbf{X} | Z)$ 为对称方阵;
- (c) 条件协方差矩阵 $\text{Var}(\mathbf{X} | Z)$ 是半正定的;
- (d) 对于任意的非随机矩阵 $\mathbf{A} \in \mathbb{R}^{m \times n}$ 和 m 维常数列向量 \mathbf{a} , 则有 $\text{Var}(\mathbf{A}\mathbf{X} + \mathbf{a} | Z) = \mathbf{A}\text{Var}(\mathbf{X} | Z)\mathbf{A}^T$.

这些性质的证明仅涉及一些基本的矩阵运算, 留给读者自证不难.

3.3 多元线性投影模型的总体参数

本节我们介绍线性投影模型的总体参数, 尽管在线性假设式 (3.4) 成立的情况下, 线性投影模型即是线性回归模型, 但这只是一个假设 (assumption, 而非 hypothesis).

多元线性投影模型的总体方程为

$$Y = \mathcal{P}(Y | \mathbf{X}) + \varepsilon, \quad (3.11)$$

其中 $\mathcal{P}(Y | \mathbf{X}) = \beta_1 X_1 + \cdots + \beta_n X_K = \mathbf{X}^T \boldsymbol{\beta}$ 为给定

$$\mathbf{X} = (X_1, \dots, X_K)^T$$

下 Y 的最优线性预测量 (best linear predictor).

由于 $\mathcal{P}(Y|\mathbf{X})$ 使得投影误差 ε 具有最小的 MSE, 据此我们定义了线性投影模型的回归系数 β 为

$$\beta = \underset{\beta \in \mathbb{R}^{K \times 1}}{\operatorname{argmin}} \operatorname{MSE}(\beta) = \underset{\beta \in \mathbb{R}^{K \times 1}}{\operatorname{argmin}} \mathbb{E} (Y - \mathbf{X}^T \beta)^2. \quad (3.12)$$

我们来求解这个最优化问题, 将 $\operatorname{MSE}(\beta)$ 展开有

$$\begin{aligned} \operatorname{MSE}(\beta) &= \mathbb{E} (Y - \mathbf{X}^T \beta)^2 = \mathbb{E} (Y^2 - 2Y\beta^T \mathbf{X} + (\beta^T \mathbf{X}) (\mathbf{X}^T \beta)) \\ &= \mathbb{E} Y^2 - 2\beta^T \mathbb{E} (\mathbf{X}Y) + \beta^T \mathbb{E} (\mathbf{X}\mathbf{X}^T) \beta, \end{aligned}$$

利用矩阵微积分的知识^[1], 对参数 β 求导得到一阶条件

$$\frac{\partial}{\partial \beta} \operatorname{MSE}(\beta) = -2\mathbb{E} (\mathbf{X}Y) + 2\mathbb{E} (\mathbf{X}\mathbf{X}^T) \beta = 0,$$

即有 $\mathbb{E} (\mathbf{X}\mathbf{X}^T) \beta = \mathbb{E} (\mathbf{X}Y)$ 成立. 在假设 (2.2.1) 下便可以解出线性投影模型的总体参数为

$$\boxed{\beta = (\mathbb{E} (\mathbf{X}\mathbf{X}^T))^{-1} \mathbb{E} (\mathbf{X}Y)}. \quad (3.13)$$

在式 (3.13) 中, $(\mathbb{E} (\mathbf{X}\mathbf{X}^T))^{-1}$ 为 K 阶方阵, $\mathbb{E} (\mathbf{X}Y)$ 为 $K \times 1$ 的列向量, 因此矩阵乘法的结果即是 $K \times 1$ 的列向量.

由于 $\operatorname{MSE}(\beta)$ 是非负的, 当其取极值时只能为极小值, 由此保障了式 (3.13) 满足式 (3.12) 的最优化要求.

在一元线性投影模型中, 我们知道投影误差 ε 与解释变量 X 正交, 这对于多元的情形同样是成立的. 我们将一阶条件 $\mathbb{E} (\mathbf{X}\mathbf{X}^T) \beta = \mathbb{E} (\mathbf{X}Y)$ 中的 Y 替换有

$$\mathbb{E} (\mathbf{X}\mathbf{X}^T) \beta = \mathbb{E} (\mathbf{X} (\mathbf{X}^T \beta + \varepsilon)) \Leftrightarrow \mathbb{E} (\mathbf{X}\mathbf{X}^T) \beta = \mathbb{E} (\mathbf{X}\mathbf{X}^T) \beta + \mathbb{E} (\mathbf{X}\varepsilon),$$

即有 $\boxed{\mathbb{E} (\mathbf{X}\varepsilon) = \mathbf{0}_{K \times 1}}$ 成立, 这等价于 $\boxed{\mathbb{E} (X_i \varepsilon) = 0}, i = 1, \dots, n$. 也就是说, 对于多元线性投影模型, K 个解释变量均与投影误差 ε 正交.

特别地, 当最优线性预测了含有常数时, 即 $X_1 \equiv 1$, 则 $\mathbb{E} (X_i \varepsilon) = \mathbb{E} (1 \times \varepsilon) = 0$, 即包含常数项时线性投影模型满足有 $\mathbb{E} \varepsilon = 0$. 注意, 线性投影模型不需要假设 (3.2.1), 而只需要假设 (2.2.1), 因此截距项对于线性投影模型具有非凡的意义.

在求解一元线性投影模型时, 我们还使用配方法来解决. 实际上, 对于任意的关于列向量 \mathbf{a} 的函数 $f(\mathbf{a}) = c - 2\mathbf{b}^T \mathbf{a} + \mathbf{a}^T \mathbf{S} \mathbf{a}$, 其中 c 为常数, \mathbf{b} 为常数列向量, \mathbf{S} 为正定

^[1]详见参考资料 (cite)

的实对称方阵, 则 $f(\mathbf{a})$ 可以配方^[2]为

$$\begin{aligned} f(\mathbf{a}) &= c - 2\mathbf{b}^T \mathbf{a} + \mathbf{a}^T \mathbf{S} \mathbf{a} - \mathbf{b}^T \mathbf{S}^{-1} \mathbf{b} + \mathbf{b}^T \mathbf{S}^{-1} \mathbf{b} \\ &= (c - \mathbf{b}^T \mathbf{S}^{-1} \mathbf{b}) + \mathbf{a}^T \mathbf{S} \mathbf{a} - 2\mathbf{b}^T \mathbf{a} + \mathbf{b}^T \mathbf{S}^{-1} \mathbf{b} \\ &= (c - \mathbf{b}^T \mathbf{S}^{-1} \mathbf{b}) + \mathbf{a}^T \mathbf{S} \mathbf{a} - 2\mathbf{b}^T (\mathbf{S}^{-1} \mathbf{S}) \mathbf{a} + \mathbf{b}^T \mathbf{S}^{-1} \mathbf{b} \\ &= (c - \mathbf{b}^T \mathbf{S}^{-1} \mathbf{b}) + (\mathbf{a} - \mathbf{S}^{-1} \mathbf{b})^T \mathbf{S} (\mathbf{a} - \mathbf{S}^{-1} \mathbf{b}), \end{aligned}$$

这样, $\text{MSE}(\boldsymbol{\beta}) = \mathbb{E}Y^2 - 2\boldsymbol{\beta}^T \mathbb{E}(\mathbf{X}Y) + \boldsymbol{\beta}^T \mathbb{E}(\mathbf{X}\mathbf{X}^T) \boldsymbol{\beta}$ 亦可以配方为如此的形式, 同样 MSE 取最小值时即可得到式 (3.13), 我们利用此还可以得到 MSE 取最小值为

$$\mathbb{E}Y^2 - (\mathbb{E}(\mathbf{X}Y))^T (\mathbb{E}(\mathbf{X}\mathbf{X}^T))^{-1} \mathbb{E}(\mathbf{X}Y).$$

上述推导的线性投影模型并没有直接包含有常数项, 我们可以令 $X_1 \equiv 1$ 来得到含义常数项的线性投影模型, 另一种计算含有常数项的投影模型的方法则是进行中心化处理. 设含义常数项的线性投影模型为

$$Y = \alpha + \beta_1 X_1 + \cdots + \beta_n X_n + \varepsilon = \alpha + \mathbf{X}^T \boldsymbol{\beta} + \varepsilon, \quad (3.14)$$

对式 (3.14) 求期望则有

$$\mathbb{E}Y = \alpha + \mathbb{E}(\mathbf{X}^T) \boldsymbol{\beta} + \mathbb{E}\varepsilon,$$

当含有常数项 α 时 $\mathbb{E}\varepsilon = 0$, 故有

$$\mathbb{E}Y = \alpha + \mathbb{E}(\mathbf{X}^T) \boldsymbol{\beta}, \quad (3.15)$$

将式 (3.11) 减去式 (3.15), 即可消去常数项 α 并得到新的线性投影模型

$$Y - \mathbb{E}Y = (\mathbf{X} - \mathbb{E}(\mathbf{X}))^T \boldsymbol{\beta} + \varepsilon. \quad (3.16)$$

式 (3.16) 的形式与式 (3.11) 是一致的, 于是依据 MSE 最小的要求我们便得到式 (3.16) 中参数 $\boldsymbol{\beta}$ 的表达式为

$$\begin{aligned} \boldsymbol{\beta} &= \left[\mathbb{E} \left((\mathbf{X} - \mathbb{E}(\mathbf{X})) (\mathbf{X} - \mathbb{E}(\mathbf{X}))^T \right) \right]^{-1} \mathbb{E}(\mathbf{X} - \mathbb{E}(\mathbf{X})) (Y - \mathbb{E}Y) \\ &= (\text{Var}(\mathbf{X}))^{-1} \text{Cov}(\mathbf{X}, Y), \end{aligned}$$

这与我们求解的一元线性投影模型的参数 $\beta_1 = \frac{\text{Cov}(X, Y)}{\text{Var}(X)}$ 具有相同的形式, 只不过换成了协方差矩阵以及矩阵的逆. 利用式 (3.15) 还可以得到常数项 α 的表达式为 $\alpha = \mathbb{E}Y - \mathbb{E}(\mathbf{X}^T) \boldsymbol{\beta}$.

我们将本节的多元线性投影模型总结如下.

^[2]这个配方过程依然能够通过矩阵相合实现.

模型 3.3.1 (多元线性投影模型). 在假设 (2.2.1) 成立时, 满足如下设定的模型称为多元线性投影模型 (*linear projection model*).

$$\begin{aligned} Y &= \mathcal{P}(Y|\mathbf{X}) + \varepsilon \\ \mathcal{P}(Y|\mathbf{X}) &= \mathbf{X}^T \boldsymbol{\beta}, \quad \boldsymbol{\beta} = (\mathbb{E}(\mathbf{X}\mathbf{X}^T))^{-1} \mathbb{E}(\mathbf{X}Y) \\ \mathbb{E}(\mathbf{X}\varepsilon) &= \mathbf{0} \\ \mathbb{E}(\varepsilon) &= 0 \quad (\text{需要常数项}) \end{aligned}$$

模型 3.3.2 (多元线性投影模型 (含有截距项)). 在假设 (2.2.1) 成立时, 满足如下设定的模型称为含有截距项的多元线性投影模型 (*linear projection model*).

$$\begin{aligned} Y &= \mathcal{P}(Y|\mathbf{X}) + \varepsilon \\ \mathcal{P}(Y|\mathbf{X}) &= \alpha + \mathbf{X}^T \boldsymbol{\beta}, \quad \begin{cases} \boldsymbol{\beta} = (\text{Var}(\mathbf{X}))^{-1} \text{Cov}(\mathbf{X}, Y) \\ \alpha = \mathbb{E}Y - \mathbb{E}(\mathbf{X}^T) \end{cases} \\ \mathbb{E}(\mathbf{X}\varepsilon) &= \mathbf{0} \\ \mathbb{E}(\varepsilon) &= 0 \quad (\text{需要常数项}) \end{aligned}$$

3.4 多元线性投影的样本估计量

这一节我们介绍多元线性投影模型的 OLS 方法, 并由此得到 OLS 估计量.

依据本章第一节介绍的符号, 对于线性投影方程 (3.11), 设有样本容量为 n 的样本 $\{(Y_1, \mathbf{X}_1), \dots, (Y_n, \mathbf{X}_n)\}$, 则样本个体的线性投影方程为

$$Y_i = \mathcal{P}(Y_i|\mathbf{X}_i) + \varepsilon = \mathbf{X}_i^T \boldsymbol{\beta} + \varepsilon_i,$$

同样可以矩阵形式表示 $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$.

在第一章中我们介绍了普通最小二乘法, 对于多元的情形我们仍旧依据定义 (2.3.1), 即依据样本对于总体的 MSE 的矩估计为

$$\hat{S}(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n (Y_i - \mathbf{X}_i^T \boldsymbol{\beta})^2 = \frac{1}{n} \text{SSE}(\boldsymbol{\beta}),$$

而误差平方和 $\text{SSE}(\boldsymbol{\beta}) = \sum_{i=1}^n (Y_i - \mathbf{X}_i^T \boldsymbol{\beta})^2$. 这样, OLS 估计量 (OLS estimator) 即被定义为

$$\hat{\boldsymbol{\beta}}_{\text{OLS}} = \underset{\boldsymbol{\beta} \in \mathbb{R}^{K \times 1}}{\text{argmin}} \hat{S}(\boldsymbol{\beta}). \quad (3.17)$$

为了求解式 (3.17) 的最优化问题, 同样只需要研究 $\text{SSE}(\boldsymbol{\beta})$ 的极值. 由于

$$\text{SSE}(\boldsymbol{\beta}) = \sum_{i=1}^n (Y_i - \mathbf{X}_i^T \boldsymbol{\beta})^2 = \sum_{i=1}^n [Y_i - (\beta_1 X_{i1} + \cdots + \beta_n X_{iK})]^2,$$

则对参数 $\boldsymbol{\beta}$ 的分量依次求导, 得到包含有 K 个方程的一阶条件

$$\begin{cases} \frac{\partial}{\partial \beta_1} \text{SSE}(\boldsymbol{\beta}) = -2 \sum_{i=1}^n X_{i1} [Y_i - (\beta_1 X_{i1} + \cdots + \beta_n X_{iK})] = 0, \\ \cdots \\ \frac{\partial}{\partial \beta_k} \text{SSE}(\boldsymbol{\beta}) = -2 \sum_{i=1}^n X_{ik} [Y_i - (\beta_1 X_{i1} + \cdots + \beta_n X_{iK})] = 0, \\ \cdots \\ \frac{\partial}{\partial \beta_K} \text{SSE}(\boldsymbol{\beta}) = -2 \sum_{i=1}^n X_{iK} [Y_i - (\beta_1 X_{i1} + \cdots + \beta_n X_{iK})] = 0, \end{cases}$$

也即有

$$\begin{cases} \sum_{i=1}^n X_{i1} \hat{\varepsilon}_i = 0, \\ \cdots \\ \sum_{i=1}^n X_{ik} \hat{\varepsilon}_i = 0, \\ \cdots \\ \sum_{i=1}^n X_{iK} \hat{\varepsilon}_i = 0, \end{cases} \Leftrightarrow \begin{cases} (X_{11}, \cdots, X_{n1}) \cdot \hat{\boldsymbol{\varepsilon}} = 0, \\ \cdots \\ (X_{1k}, \cdots, X_{nk}) \cdot \hat{\boldsymbol{\varepsilon}} = 0, \\ \cdots \\ (X_{1K}, \cdots, X_{nK}) \cdot \hat{\boldsymbol{\varepsilon}} = 0, \end{cases}$$

这里的 $\varepsilon_i = Y_i - \mathbf{X}_i^T \hat{\boldsymbol{\beta}}$ 为样本残差, 其依赖于 OLS 估计量, 而 OLS 估计量满足一阶条件. 注意到 $\{(X_{11}, \cdots, X_{n1})^T, \cdots, (X_{1k}, \cdots, X_{nk})^T, \cdots, (X_{1K}, \cdots, X_{nK})^T\}$ 为数据矩阵 \mathbf{X} 的列向量组, 于是一阶条件等价于

$$((X_{11}, \cdots, X_{n1}), \cdots, (X_{1k}, \cdots, X_{nk}), \cdots, (X_{1K}, \cdots, X_{nK}))^T \cdot \hat{\boldsymbol{\varepsilon}} = \mathbf{0}_{K \times 1}$$

即有

$$\mathbf{X}^T \hat{\boldsymbol{\varepsilon}} = \mathbf{0}_{K \times 1} \quad (3.18)$$

成立. 式 (3.18) 即是多元线性投影模型的 OLS 估计量的正规方程组. 我们将式 (3.18) 中的残差向量 $\hat{\boldsymbol{\varepsilon}}$ 替换掉有

$$\mathbf{X}^T (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}}) = \mathbf{0} \Leftrightarrow \mathbf{X}^T \mathbf{Y} = \mathbf{X}^T \mathbf{X} \hat{\boldsymbol{\beta}}$$

于是解得多元线性投影模型的 OLS 估计量为

$$\hat{\boldsymbol{\beta}}_{\text{OLS}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}. \quad (3.19)$$

同样 $\text{SSE}(\beta)$ 是非负的, 一阶条件求得的极值点 $\hat{\beta}$ 便是极小值点, 因而式 (3.19) 符合式 (3.17) 的要求.

注意到

$$\mathbf{X}_{n \times K} = \begin{pmatrix} \mathbf{X}_1^T \\ \vdots \\ \mathbf{X}_i^T \\ \vdots \\ \mathbf{X}_n^T \end{pmatrix}, \quad \mathbf{X}^T = (\mathbf{X}_1, \cdots, \mathbf{X}_i, \cdots, \mathbf{X}_n),$$

则

$$\underbrace{\mathbf{X}^T \mathbf{X}}_{K \text{ 阶方阵}} = (\mathbf{X}_1, \cdots, \mathbf{X}_i, \cdots, \mathbf{X}_n) \begin{pmatrix} \mathbf{X}_1^T \\ \vdots \\ \mathbf{X}_i^T \\ \vdots \\ \mathbf{X}_n^T \end{pmatrix} = \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^T,$$

$$\underbrace{\mathbf{X}^T \mathbf{Y}}_{K \times 1 \text{ 列向量}} = (\mathbf{X}_1, \cdots, \mathbf{X}_i, \cdots, \mathbf{X}_n) \begin{pmatrix} Y_1 \\ \vdots \\ Y_i \\ \vdots \\ Y_n \end{pmatrix} = \sum_{i=1}^n \mathbf{X}_i Y_i,$$

其中 $\mathbf{X}_i \mathbf{X}_i^T$ 为 $K \times 1$ 列向量与 $1 \times K$ 行向量的积, $\mathbf{X}_i Y_i$ 是 $K \times 1$ 列向量与 1×1 标量的积. 这样 OLS 估计量便可以求和符号表达式

$$\hat{\beta}_{\text{OLS}} = \left(\sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^T \right)^{-1} \sum_{i=1}^n \mathbf{X}_i Y_i. \quad (3.20)$$

由于 $\text{SSE}(\beta)$ 可以用向量内积表示为

$$\text{SSE}(\beta) = \sum_{i=1}^n (Y_i - \mathbf{X}_i^T \beta)^2 = (\mathbf{Y} - \mathbf{X}\beta)^T (\mathbf{Y} - \mathbf{X}\beta),$$

熟悉矩阵微积分的读者则可以直接对参数 β 求梯度矩阵, 得到一阶条件

$$\frac{\partial}{\partial \beta} \text{SSE}(\beta) = -2\mathbf{X}^T (\mathbf{Y} - \mathbf{X}\beta) = \mathbf{0}_{K \times 1},$$

这样亦可推导出式 (3.19) 的 OLS 估计量.

对于正规方程组 (3.18), 其等价于

$$\mathbf{X}^T \hat{\epsilon} = \mathbf{0} \Leftrightarrow (\mathbf{X}_1, \cdots, \mathbf{X}_i, \cdots, \mathbf{X}_n) \hat{\epsilon} = \mathbf{0} \Leftrightarrow \boxed{\sum_{i=1}^n \mathbf{X}_i \hat{\epsilon}_i = \mathbf{0}} \Leftrightarrow \sum_{i=1}^n X_{ik} \hat{\epsilon}_i = 0,$$

依据严格外生性假设式 (3.5) 我们推知回归误差 ε_i 同任意的解释变量样本 X_{jk} 是正交的, 因而可以推知 $\mathbb{E}(\mathbf{X}\varepsilon_i) = 0$, 其在 OLS 估计量的对应形式便是 $\mathbf{X}^T \hat{\varepsilon} = \mathbf{0}$. 特别地, 当线性投影模型含有常数项 $X_{i1} \equiv 1$ 时, 即有残差和 $\sum_{i=1}^n X_{i1} \hat{\varepsilon}_i = \sum_{i=1}^n \hat{\varepsilon}_i = 0$.

我们现在推导含有常数项时的 OLS 估计量, 我们依旧是利用中心化的方法. 设含有常数项 α 的样本个体有投影方程

$$Y_i = \alpha + \mathbf{X}_i^T \boldsymbol{\beta} + \varepsilon_i,$$

则求样本均值得 $\bar{Y} = \alpha + \bar{\mathbf{X}}^T \boldsymbol{\beta} + \bar{\varepsilon}$, 其中 $\bar{\mathbf{X}} = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i$, $\bar{\varepsilon} = \frac{1}{n} \sum_{i=1}^n \varepsilon_i$, 于是有样本个体有中心化的投影方程

$$Y_i - \bar{Y} = (\mathbf{X}_i - \bar{\mathbf{X}})^T \boldsymbol{\beta} + (\varepsilon_i - \bar{\varepsilon}) = (\mathbf{X}_i - \bar{\mathbf{X}})^T \boldsymbol{\beta} + \mu_i,$$

这里我们使用 $\mu_i = \varepsilon_i - \bar{\varepsilon}$ 定义了新的投影误差, 于是依据式 (3.20) 得到含有常数项 α 时的 OLS 估计量为

$$\hat{\boldsymbol{\beta}}_{\text{OLS}} = \left(\sum_{i=1}^n (\mathbf{X}_i - \bar{\mathbf{X}}) (\mathbf{X}_i - \bar{\mathbf{X}})^T \right)^{-1} \sum_{i=1}^n (\mathbf{X}_i - \bar{\mathbf{X}}) (Y_i - \bar{Y}), \quad (3.21)$$

由于含有常数项时残差和为 0, 故常数项 α 的 OLS 估计量为 $\hat{\alpha} = \bar{Y} - \bar{\mathbf{X}}^T \hat{\boldsymbol{\beta}}$. 这与定理 (2.3.1) 是完全一致的.

定理 3.4.1. 模型 (3.3.1) 的参数 $\boldsymbol{\beta}$ 的 OLS 估计量为

$$\hat{\boldsymbol{\beta}}_{\text{OLS}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \quad \text{或} \quad \hat{\boldsymbol{\beta}}_{\text{OLS}} = \left(\sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^T \right)^{-1} \sum_{i=1}^n \mathbf{X}_i Y_i$$

且 OLS 估计量满足有

$$\hat{\varepsilon} = \mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}}, \quad \mathbf{X}^T \hat{\varepsilon} = \mathbf{0}.$$

若是含有常数项 α 的模型 (3.3.2), 则 OLS 估计量为

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= \left(\sum_{i=1}^n (\mathbf{X}_i - \bar{\mathbf{X}}) (\mathbf{X}_i - \bar{\mathbf{X}})^T \right)^{-1} \sum_{i=1}^n (\mathbf{X}_i - \bar{\mathbf{X}}) (Y_i - \bar{Y}) \\ \hat{\alpha} &= \bar{Y} - \bar{\mathbf{X}}^T \hat{\boldsymbol{\beta}}, \end{aligned}$$

且 OLS 估计量满足有

$$\mathbf{X}^T \hat{\varepsilon} = \mathbf{0}, \quad \mathbf{1}^T \hat{\varepsilon} = 0.$$

3.5 多元线性投影模型的分块回归

3.6 多元 OLS 估计量下 Guass-Markov 定理

本节我们将研究多元 OLS 估计量的有限样本性质, 本节以及之后, 我们假设线性投影模型 (3.3.1) 满足有假设 (3.2.1), 这时线性投影模型便等价于线性回归模型.

模型 3.6.1 (多元线性回归模型). 在假设 (2.2.1) 以及

$$\mathbb{E}(Y|\mathbf{X}) = \mathcal{P}(Y|\mathbf{X}) \quad \text{或} \quad \mathbb{E}(\varepsilon|\mathbf{X}) = 0$$

成立时, 满足如下设定的模型称为多元线性回归模型 (*linear regression model*).

$$\begin{aligned} Y &= \mathcal{P}(Y|\mathbf{X}) + \varepsilon \\ \mathcal{P}(Y|\mathbf{X}) &= \mathbf{X}^T \boldsymbol{\beta}, \quad \boldsymbol{\beta} = (\mathbb{E}(\mathbf{X}\mathbf{X}^T))^{-1} \mathbb{E}(\mathbf{X}Y) \\ \mathbb{E}(\mathbf{X}\varepsilon) &= 0 \\ \mathbb{E}(\varepsilon) &= 0 \quad (\text{不需要常数项}) \end{aligned}$$

模型 (3.6.1) 是多元线性回归模型的统计总体的设定, 而样本的假设则遵从假设 (3.2.1). 本节以及后续, 我们并不设定样本为独立同分布的简单随机样本, 而是依据假设 (3.2.1) 这样更一般的设定, 两者的区别我们已经在 3.2 节辨析过了.

我们首先研究 OLS 估计量的代数性质. 由于 OLS 估计量为式 (3.19), 于是多元线性回归模型 (3.6.1) 的抽样误差 (sampling error) 即为

$$\begin{aligned} \hat{\boldsymbol{\beta}} - \boldsymbol{\beta} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} - \boldsymbol{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}) - \boldsymbol{\beta} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X}\boldsymbol{\beta} - \boldsymbol{\beta} + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\varepsilon} = \boxed{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\varepsilon}}, \end{aligned}$$

这便表明抽样误差 $\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}$ 是误差向量 $\boldsymbol{\varepsilon}$ 的线性组合.

我们再来计算 OLS 估计量的协方差矩阵, 由于 $\hat{\boldsymbol{\beta}}_{\text{OLS}} = (\beta_1, \dots, \beta_K)^T$ 是一个 $K \times 1$ 维的随机向量, 则依据定理 (3.2.1) 得到

$$\begin{aligned} \text{Var}(\hat{\boldsymbol{\beta}}|\mathbf{X}) &= \text{Var}\left((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \middle| \mathbf{X}\right) \\ &= \left((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\right) \text{Var}(\mathbf{Y}|\mathbf{X}) \left((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\right)^T \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \text{Var}(\mathbf{Y}|\mathbf{X}) \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}, \end{aligned}$$

而条件协方差矩阵 $\text{Var}(\mathbf{Y}|\mathbf{X})$ 同样计算有

$$\boxed{\text{Var}(\mathbf{Y}|\mathbf{X})} = \text{Var}(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}|\mathbf{X}) = \boxed{\text{Var}(\boldsymbol{\varepsilon}|\mathbf{X})} = \boxed{\mathbb{E}(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^T|\mathbf{X})}.$$

这样, 对于模型 (3.6.1) 而言, OLS 估计量 $\hat{\beta}$ 的条件协方差矩阵为

$$\text{Var}(\hat{\beta} | \mathbf{X}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\sigma^2 \mathbf{I}) \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} = \boxed{(\mathbf{X}^T \mathbf{X})^{-1} \sigma^2},$$

这一结果依赖于条件同方差式 (3.7). 如果模型 (3.6.1) 具有异方差问题, 则依据式 (3.10) 得到

$$\text{Var}(\hat{\beta} | \mathbf{X}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{D} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1},$$

由于矩阵乘法中乘以对角矩阵仅作用于对角线元素, 因此对于 $\mathbf{X}^T = (\mathbf{X}_1, \dots, \mathbf{X}_n)$ 不难得到

$$\mathbf{X}^T \mathbf{D} \mathbf{X} = \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^T \sigma_i^2,$$

故异方差下 $\hat{\beta}$ 的条件协方差矩阵也可以写作

$$\text{Var}(\hat{\beta} | \mathbf{X}) = \underbrace{\left(\sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^T \right)^{-1}}_{K\text{阶方阵}} \underbrace{\left(\sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^T \sigma_i^2 \right)}_{K\text{阶方阵}} \underbrace{\left(\sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^T \right)^{-1}}_{K\text{阶方阵}}. \quad (3.22)$$

定理 3.6.1 (OLS 估计量的协方差矩阵). 对于模型 (3.6.1) 的 OLS 估计量 $\hat{\beta}$, 我们记 $\hat{\beta}$ 的条件协方差为 $\mathbf{V}_{\hat{\beta}}$, 则条件异方差式 (3.9) 下有

$$\mathbf{V}_{\hat{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{D} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1},$$

当条件同方差式 (3.7) 成立时则有 $\mathbf{V}_{\hat{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \sigma^2$.

现在我们研究 OLS 估计量的有限性质. 对于一元的情形我们首先证明了 $\hat{\beta}$ 是 Y_1, \dots, Y_n 的线性组合, 在多元情形下 $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$ 这是平凡的. 现在我们研究 OLS 估计量的有效性.

定理 3.6.2 (OLS 估计量的无偏性). 模型 (3.6.1) 的 OLS 估计量 $\hat{\beta}$ 是参数 β 的无偏估计量.

证明. 对抽样误差 $\hat{\beta} - \beta = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\varepsilon}$ 求条件期望有

$$\begin{aligned} \mathbb{E}(\hat{\beta} - \beta | \mathbf{X}) &= \mathbb{E}\left((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\varepsilon} | \mathbf{X}\right) \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{E}(\boldsymbol{\varepsilon} | \mathbf{X}) \stackrel{\text{严格外生性式 (3.5)}}{=} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \cdot \mathbf{0}_{n \times 1} = \mathbf{0}_{K \times 1}, \end{aligned}$$

则有 $\mathbb{E}(\hat{\beta} | \mathbf{X}) = \mathbb{E}(\beta | \mathbf{X}) = \beta$. 当然亦可直接计算

$$\begin{aligned} \mathbb{E}(\hat{\beta} | \mathbf{X}) &= \mathbb{E}\left((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} | \mathbf{X}\right) \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{E}(\mathbf{Y} | \mathbf{X}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \beta = \beta. \end{aligned}$$

于是, 依据定理 (2.1.1) 得 $\mathbb{E}\hat{\beta} = \mathbb{E}\left(\mathbb{E}\left(\hat{\beta} \middle| \mathbf{X}\right)\right) = \beta$. \square

我们再证明 OLS 估计量是模型 (3.6.1) 的最优线性无偏估计量 (best linear unbiased estimator), 即 Gauss-Markov 定理.

定理 3.6.3 (Gauss-Markov 定理). 设模型 (3.6.1) 的 OLS 估计量为 $\hat{\beta}$, 且满足有条件同方差式 (3.7). 对于模型 (3.6.1) 的任意线性无偏估计量 $\tilde{\beta} = \mathbf{A}(\mathbf{X})\mathbf{Y}$, 则有

$$\text{Var}\left(\tilde{\beta} \middle| \mathbf{X}\right) \geq \mathbf{V}_{\hat{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \sigma^2,$$

这里 \geq 表示 $\text{Var}\left(\tilde{\beta} \middle| \mathbf{X}\right) - \mathbf{V}_{\hat{\beta}}$ 为半正定矩阵. 这也意味着 $\text{Var}\left(\tilde{\beta}_k \middle| \mathbf{X}\right) \geq \text{Var}\left(\hat{\beta}_k \middle| \mathbf{X}\right), k = 1, \dots, K$.

证明. 估计量 $\tilde{\beta} = \mathbf{A}(\mathbf{X})\mathbf{Y}$ 是参数 β 的无偏估计量, 而

$$\mathbb{E}\left(\tilde{\beta} \middle| \mathbf{X}\right) = \mathbb{E}\left(\mathbf{A}(\mathbf{X})\mathbf{Y} \middle| \mathbf{X}\right) = \mathbf{A}\mathbb{E}(\mathbf{Y} \middle| \mathbf{X}) = \mathbf{A}\mathbf{X}\beta,$$

故 $\tilde{\beta}$ 满足有必要条件^[3] $\mathbf{A}\mathbf{X} = \mathbf{I}_K$.

这样, 线性无偏估计量 $\tilde{\beta}$ 的条件协方差为

$$\begin{aligned} \text{Var}\left(\tilde{\beta}_k \middle| \mathbf{X}\right) &= \text{Var}(\mathbf{A}\mathbf{Y} \middle| \mathbf{X}) = \mathbf{A}\text{Var}(\mathbf{Y} \middle| \mathbf{X})\mathbf{A}^T \\ &= \mathbf{A}\mathbf{D}\mathbf{A}^T \stackrel{\text{条件同方差式 (3.7)}}{=} \mathbf{A}\mathbf{A}^T \sigma^2. \end{aligned}$$

于是, $\text{Var}\left(\tilde{\beta} \middle| \mathbf{X}\right) - \mathbf{V}_{\hat{\beta}} = \mathbf{A}\mathbf{A}^T \sigma^2 - (\mathbf{X}^T \mathbf{X})^{-1} \sigma^2 = \left(\mathbf{A}\mathbf{A}^T - (\mathbf{X}^T \mathbf{X})^{-1}\right) \sigma^2$.

Guass 在 1823 年时便注意到可以将 \mathbf{A} 分解为 $\mathbf{A} - (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$, 令 $\mathbf{C} = \mathbf{A} - (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$, 则有

$$\begin{aligned} \mathbf{A}\mathbf{A}^T - (\mathbf{X}^T \mathbf{X})^{-1} &= \left(\mathbf{C} + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\right) \left(\mathbf{C} + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\right)^T - (\mathbf{X}^T \mathbf{X})^{-1} \\ &= \mathbf{C}\mathbf{C}^T + \mathbf{C}\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{C}^T + (\mathbf{X}^T \mathbf{X})^{-1} - (\mathbf{X}^T \mathbf{X})^{-1} \\ &= \mathbf{C}\mathbf{C}^T + \mathbf{A}\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} - (\mathbf{X}^T \mathbf{X})^{-1} + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{A}^T - (\mathbf{X}^T \mathbf{X})^{-1} \\ &= \mathbf{C}\mathbf{C}^T + (\mathbf{X}^T \mathbf{X})^{-1} - (\mathbf{X}^T \mathbf{X})^{-1} + (\mathbf{X}^T \mathbf{X})^{-1} - (\mathbf{X}^T \mathbf{X})^{-1} \\ &= \mathbf{C}\mathbf{C}^T, \end{aligned}$$

这样便将 $\mathbf{A}\mathbf{A}^T - (\mathbf{X}^T \mathbf{X})^{-1}$ 转化为了 $K \times n$ 矩阵 $\mathbf{C} = \mathbf{A} - (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ 的 Gram 矩阵. 对于任意矩阵 \mathbf{A} , 其 Gram 矩阵 $\mathbf{G} = \mathbf{A}\mathbf{A}^T$ 是半正定的, 当且仅当 \mathbf{A} 的列向量组线性无关时 Gram 矩阵 \mathbf{G} 是正定的. 于是一定有

$$\text{Var}\left(\tilde{\beta} \middle| \mathbf{X}\right) - \mathbf{V}_{\hat{\beta}} = \left(\mathbf{A}\mathbf{A}^T - (\mathbf{X}^T \mathbf{X})^{-1}\right) \sigma^2 \geq 0$$

^[3]严格的说, 依据 $\mathbb{E}\tilde{\beta} = \mathbb{E}\left(\mathbb{E}\left(\tilde{\beta} \middle| \mathbf{X}\right) \middle| \mathbf{X}\right) = \mathbb{E}(\mathbf{A}\mathbf{X})\beta = \beta$ 得到的必要条件为 $\mathbb{E}(\mathbf{A}\mathbf{X}) = \mathbf{I}_K$.

成立. □

定理 (3.6.3) 表明了 $\hat{\beta}_{OLS}$ 是最优线性无偏估计量, 而 Bruce(cite) 给出了现代形式的 Gauss-Markov 定理, 在条件同方差式 (3.7) 的假设下, 他证明了 OLS 估计量即是最优无偏估计量 (best unbiased estimator), 是模型 (3.6.1) 参数 β 的 MVU 估计.

对于 OLS 有限样本的性质, 我们已经在一元形式讨论过了其意义, 并进行了 Monte Carlo 模拟以验证.

3.7 正交投影与 OLS 估计量

到目前为止, 本笔记所涉及的不少模型和概念, 都与“投影 (Projection)”这个词语相关——实质上, “投影”这个术语源自于线性代数. 本节将介绍线性代数中的正交投影的概念, 并介绍其与 OLS 估计量之间的关联.

在线性代数中, 对于向量空间 (或线性空间) V , 若线性变换 $f: V \rightarrow V$ 满足有 $f \circ f = f$ ($f(f(\cdot)) = f$), 则称线性变换 f 为投影 (projection). 如果 V 定义了内积 $\langle \cdot, \cdot \rangle$, 对于任意的 $x, y \in V$, 投影 f 满足有 $\langle f(x), y \rangle = \langle x, f(y) \rangle$, 则称投影 f 为正交投影 (orthogonal projection). 显然, 这里的投影是个一般化的抽象概念, 它反映了投影在代数上的本质, 但对于我们理解线性模型中出现了投影基本上没有帮助.

为了理解上述的抽象定义, 我们回顾中学数学中的向量投影的概念, 但加以线性代数的语言来叙述. 设想一个如图 (3.1) 的情景, 已知二维平面 \mathbb{R}^2 中的一条过原点的直线 l 和一点 x , 由于他们都有显式的二维坐标, 故我们可以直接使用列向量 l 和 x 来表示二者. 在中学数学中, 我们知道点 x 在直线 l 上的投影, 便是 l 上的与点 x 距离最小的一点. 在二维平面中, 我们用内积定义距离, 也即是 Euclid 距离.

依据勾股定理, 不难得知距离最小意味着点 x 与其投影点的连线垂直于直线 l , 否则, 总可以找到使得距离更小的投影点. 利用平面几何的知识, 我们可以证明将点 x 投影至直线 l 是一个线性变换, 而线性变换可以用矩阵乘法表示, 这样我们不妨记投影矩阵 $P_l \in \mathbb{R}^{2 \times 2}$, 则点 x 在直线 l 上的投影点为 $P_l x$.

于是我们描绘了如图 (3.1) 所示的 (正交) 投影, 这一情形是可以自然推广至更高维的情形. 现思考一个 n 维的向量空间 \mathbb{R}^n , 则其中的点 $x \in \mathbb{R}^n$ 包含有 n 个分量 (视为列向量). 在高维时与二维平面中过原点直线 l 所对应的存在, 数学中称之为超平面 (hyperplane), 其是 \mathbb{R}^n 的一个子空间^[4], 我们仍旧可以记为列向量 $l \in \mathbb{R}^n$.

为了确定投影矩阵 P_l 的具体形式, 我们定义

$$P_l = \operatorname{argmin}_{P_l \in \mathbb{R}^{n \times n}} \|P_l x - x\| = \operatorname{argmin}_{P_l \in \mathbb{R}^{n \times n}} \sqrt{(P_l x - x)^T (P_l x - x)}, \quad (3.23)$$

$$s.t. \quad P_l x = l t, \quad t \in \mathbb{R},$$

^[4] l 过原点确保了子空间是封闭的.

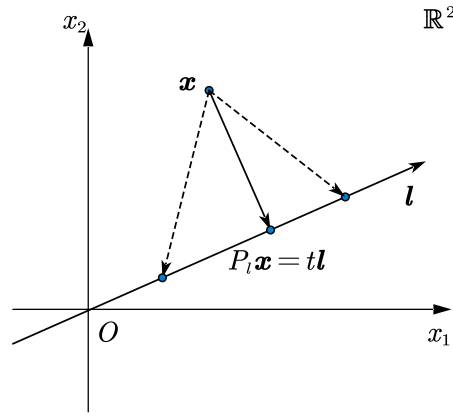


图 3.1: 二维平面中的正交投影

其中 $P_l x = tl$ 为约束条件, 使得投影点在 l 之上. 不难发现, 式 (3.23) 与 OLS 估计量式 (3.17) 在形式上是相似的. 为了求解式 (3.23), 一种思路即是 Lagrange 乘数法, 但式 (3.23) 的约束条件非常简单, 直接代入优化函数得

$$s(t) = (tl - x)^T (tl - x) = t^2 l^T l - 2tl^T x + x^T x,$$

对 $s(t)$ 求导得到一阶条件

$$\frac{d}{dt} s(t) = 2tl^T l - 2l^T x = 0 \quad \Rightarrow, \quad tl^T l = l^T x,$$

解得 $t = (l^T l)^{-1} l^T x$. 故有 $P_l x = tl = l(l^T l)^{-1} l^T x$, 则将点 x 投影至 l 上的线性变换可以用投影矩阵

$$P_l = l(l^T l)^{-1} l^T \quad (3.24)$$

表示. 不难发现, 投影点 $P_l x$ 至点 x 的向量 $x - P_l x$ 与 l 是正交的, 即

$$l^T (x - P_l x) = l^T x - l^T P_l x = l^T x - l^T l (l^T l)^{-1} l^T x = l^T x - l^T x = 0,$$

因此, $P_l x$ 称为点 x 在 l 上的正交投影.

对于超平面 l , 我们进一步考虑其上的 K 个向量 l_1, \dots, l_K , 由于 l 是向量空间 \mathbb{R}^n 的子空间, 则向量组 $\{l_1, \dots, l_k\}$ 的线性组合仍属于超平面 l . 这一过程也是线性变换, 记矩阵 $L = (l_1, \dots, l_k) \in \mathbb{R}^{n \times K}$, 则 L 与任意的 k 维列向量的积仍属于超平面内, 但乘积蕴含了向量组 $\{l_1, \dots, l_k\}$ 所张成 (span) 的向量空间. 我们将式 (3.24) 中的 l 替换为 L , 则得到 n 阶方阵

$$P_L = L(L^T L)^{-1} L^T, \quad (3.25)$$

则式 (3.25) 是点 x 投影至 $\{l_1, \dots, l_k\}$ 所张成向量空间的投影矩阵 (projection matrix).

这样, 我们便可以理解为什么 OLS 估计量与投影存在关联了. 模型 (3.6.1) 的 OLS 估计量为 $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$, 对应的拟合值为

$$\hat{\mathbf{Y}} = \mathbf{X} \hat{\beta} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} = \mathbf{P}_X \mathbf{Y},$$

依据上文分析即可得知, 拟合值 $\hat{\mathbf{Y}}$ 实际上是被解释变量的样本 \mathbf{Y} 在数据矩阵 \mathbf{X} 的列向量组所张成的线性空间上的正交投影. 实际上, 投影是一种简化 n 维信息的方法, 其将信息压缩到一个 (超) 平面上. 这在社会科学领域尤其有用, 因为我们所研究的现象非常复杂, 不可能做出准确的预测. 相反, 我们往往希望构建一些模型, 由繁杂多变的数据得到更简单、更合理的解释, 投影便是实现这一目标的统计方法.

我们通过实例来理解 OLS 估计量的投影性质. 考虑一元线性回归方程 (2.23), 且有如表 (3.1) 所示的样本观测值, 则有

$$\mathbf{Y} = \begin{pmatrix} -0.5 \\ 1.25 \\ 2 \end{pmatrix}, \quad \mathbf{X}^* = \begin{pmatrix} 1 & -1 \\ 1 & 1 \\ 1 & 0 \end{pmatrix} = (\mathbf{X}, \mathbf{1}).$$

表 3.1: 正交投影示例观测值

y	x	常数项
-0.5	-1	1
1.25	1	1
2	0	1

我们将数据向量 \mathbf{Y} 和数据矩阵 \mathbf{X}^* 的列向量组都视为 \mathbb{R}^3 中的列向量, 则数据矩阵 \mathbf{X}^* 可以视为一个线性变换, 其即是对应了列向量 \mathbf{X} 和 $\mathbf{1}$ 所张成的过 $O(0, 0, 0)$ 的平面 L . 由于 OLS 估计量和拟合值为

$$\begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_0 \end{pmatrix} = ((\mathbf{X}^*)^T \mathbf{X}^*)^{-1} (\mathbf{X}^*)^T \mathbf{Y} = \begin{pmatrix} 0.8750 \\ 0.9167 \end{pmatrix}, \quad \hat{\mathbf{Y}} = \mathbf{X}^* \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_0 \end{pmatrix} = \begin{pmatrix} 0.0417 \\ 1.7917 \\ 0.9167 \end{pmatrix}$$

我们在图 (3.2) 中绘制了样本点以及样本回归直线. 依据投影的视角, 拟合值向量 $\hat{\mathbf{Y}}$ 是平面 L 中的向量, 残差向量 $\hat{\epsilon}$ 即是连接了 $\hat{\mathbf{Y}}$ 和 \mathbf{Y} 两者终点的向量, 且残差向量 $\hat{\epsilon}$ 同平面 L 是垂直的, 这样结果如图 (3.3) 所示.

下面我们研究式 (3.25) 的投影矩阵的性质.

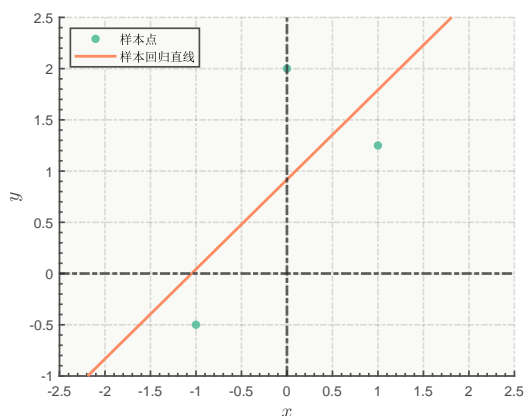


图 3.2: 一元线性回归模型的拟合直线

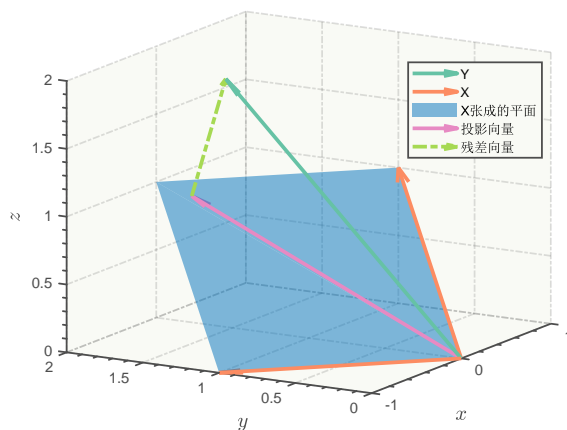


图 3.3: 一元线性回归模型的正交投影

定理 3.7.1 (投影矩阵的性质). 对于式 (3.25) 的投影矩阵 $P_L = L(L^T L)^{-1} L^T, L \in \mathbb{R}^{n \times K}, n \geq K$, 其满足有如下性质:

- (a) 投影矩阵 P_L 为幂等矩阵, 即 $P^2 = P$, 且 $P_L L = L$;
- (b) 投影矩阵 P_L 是对称矩阵, 即 $P_L^T = P_L$;
- (c) 投影矩阵的 P_L 的迹为 K ;
- (d) 投影矩阵 P_L 有 K 个特征值为 1, 余下的 $n - K$ 个特征值为 0;
- (e) 投影矩阵 P_L 的秩为 K .

证明. (a) 计算有 $P_L^2 = L(L^T L)^{-1} L^T L(L^T L)^{-1} L^T = P_L$, 这即是符合本节开头所讲的 $f \circ f = f$, 复合映射两次后的结果相同. 而 $P_L L = L$ 也是平凡的, 其含义便是将 L 投影至 L 上, 在投影后距离最小的约束条件下则 $P_L L$ 自然是 L 本身.

(b) 这也是平凡的, 这条性质保证了投影是正交的, 注意到

$$L^T(x - P_L x) = (P_L L)^T(x - P_L x) = L^T P_L^T x - L^T P_L^T P_L x = L^T P_L^T (I_n - P_L) x,$$

则 $P_L^T = P_L$ 是 $L^T(x - P_L x) = 0$ 的充分条件.

性质 (c) 至 (e) 则是线性代数中多个命题的综合, 本笔记下面证明中所使用的结论, 相应地都可以在线性代数的基础教材中找到.

首先, 对于投影矩阵 P_L , 其特征值与特征向量即是关于实数 λ 和非零列向量 x 的矩阵方程 $P_L x = \lambda x$ 的解, 由于 P_L 是幂等矩阵, 则有

$$P_L x = x\lambda, \Rightarrow P_L P_L x = P_L x\lambda, \Rightarrow P_L x = x\lambda \cdot \lambda, \Rightarrow x\lambda = x\lambda^2, \Rightarrow x\lambda(1 - \lambda) = 0,$$

由于 \mathbf{x} 是非零向量, 故投影矩阵 \mathbf{P}_L 的特征值只能是 0 或 1.

这样, 我们计算投影矩阵 \mathbf{P}_L 的迹为

$$\text{tr}(\mathbf{P}_L) = \text{tr}\left(\mathbf{L}(\mathbf{L}^T \mathbf{L})^{-1} \mathbf{L}^T\right) = \text{tr}\left((\mathbf{L}^T \mathbf{L})^{-1} \mathbf{L}^T \mathbf{L}\right) = \text{tr}(\mathbf{I}_K) = K,$$

即证性质 (c). 而方阵的特征值的和又满足有

$$\sum_{i=1}^n \lambda_i = \text{tr}(\mathbf{P}_L) = K,$$

则投影矩阵 \mathbf{P}_L 的特征值必有 K 个等于 1, 余下的 $n - K$ 个等于 0, 即证性质 (d).

最后, 由于投影矩阵 \mathbf{P}_L 为实对称方阵, 且其特征值非负, 这等价于 \mathbf{P}_L 为半正定矩阵, 而半正定矩阵的秩等于其非负特征值的数量, 即得 $\text{rank} \mathbf{P}_L = K$, 性质 (e) 证毕. \square

在一元线性回归模型 (2.4.1), 我们粗略地得到投影矩阵后便定义了回归矩阵, 在此我们同样定义

$$\mathbf{M}_L = \mathbf{I}_n - \mathbf{P}_L \quad (3.26)$$

为归零矩阵 (annihilator matrix), 其具有如下的性质.

定理 3.7.2 (归零矩阵的性质). 对于形如式 (3.26) 的归零矩阵 $\mathbf{M}_L = \mathbf{I}_n - \mathbf{P}_L$, 其具有如下性质:

- (a) $\mathbf{M}_L \mathbf{L} = \mathbf{O}_{n \times K}, \mathbf{M} \mathbf{P}_L = \mathbf{O}_n$;
- (b) 归零矩阵 \mathbf{M}_L 为幂等矩阵, 即 $\mathbf{M}^2 = \mathbf{M}$;
- (c) 归零矩阵 \mathbf{M}_L 是对称矩阵, 即 $\mathbf{M}^T = \mathbf{M}$;
- (d) 归零矩阵的 \mathbf{M}_L 的迹为 $n - K$;
- (e) 归零矩阵 \mathbf{M}_L 有 $n - K$ 个特征值为 1, 余下的 n 个特征值为 0;
- (f) 归零矩阵 \mathbf{M}_L 的秩为 $n - K$.

证明. 性质 (a) 至 (c) 的是平凡的, 余下的性质与定理 (3.7.1) 后三条性质的证明同理. 相信读者自证不难. \square

最后, 本节介绍线性代数中实对称方阵的谱分解以及应用, 这在后续章节中会使用.

定理 3.7.3 (实对称方阵的谱分解). 设任意的实对称方阵 $\mathbf{A} \in \mathbb{R}^{n \times n}$, 则其一定可以通过正交矩阵 (orthogonal matrix) 相似至对角阵 $\mathbf{\Lambda} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$, 即有

$$\mathbf{A} = \mathbf{H}^{-1} \mathbf{\Lambda} \mathbf{H},$$

其中 Λ 的元素为 A 的特征值, 正交矩阵 H 满足有 $H^{-1} = H^T$, 我们称此为实对称方阵 A 的或谱分解 (spectral decomposition).

定理 (3.7.3) 可以在基础的线性代数教材上找到 (cite), 其一个有用的应用便是考虑幂等的 n 阶实对称方阵 A , 设 $\text{rank} A = r$, 则有

$$A = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n) = \underbrace{\text{diag}(1, \dots, 1, 0, \dots, 0)}_{r \uparrow 1, n-r \uparrow 0},$$

于是

$$A = A^2 = H^{-1} \Lambda H H^{-1} \Lambda H = H^{-1} \Lambda^2 H = H^{-1} \Lambda H,$$

亦即有

$$A = H^{-1} \begin{pmatrix} I_r & O_{r \times (n-r)} \\ O_{(n-r) \times r} & O_{n-r} \end{pmatrix} H. \quad (3.27)$$

定理 (3.7.3) 的另一个应用便是正定矩阵的平方根. 对于任意的方阵 A , 若存在方阵 B 使得 $A = B^2 = B \cdot B$, 则称方阵 B 是矩阵 A 的平方根 (square root of the matrix A), 特别地, 如果方阵 A 是正定的, 则有如下命题成立.

定理 3.7.4 (正定矩阵的平方根). 对于任意的正定矩阵 $A \in \mathbb{R}^{n \times n}$, 则存在唯一的平方根 B , 使得 $A = B^2$ 成立. 记有 $B = A^{1/2}$. 且有

$$A^{1/2} = H^T \Lambda^{1/2} H,$$

这里的 H 为使得 A 相似至对角阵 Λ 的正交矩阵, $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$ 为 A 的特征值构成的对角阵, $\Lambda^{1/2} = \text{diag}(\sqrt{\lambda_1}, \sqrt{\lambda_2}, \dots, \sqrt{\lambda_n})$.

证明. 对于任意的 n 阶正定的实对称方阵 A , 即有 $\text{rank} A = n$, 那么其所有的特征值都是正数, 依据定理 (3.7.3) 即有 $A = H^{-1} \Lambda H$, $H^{-1} = H^T$. 不难注意到

$$\begin{aligned} A &= H^{-1} \Lambda^{1/2} \Lambda^{1/2} H = H^T \Lambda^{1/2} I_n \Lambda^{1/2} H = H^T \Lambda^{1/2} H H^T \Lambda^{1/2} H \\ &= (H^T \Lambda^{1/2} H) (H^T \Lambda^{1/2} H) = (H^T \Lambda^{1/2} H)^2, \end{aligned}$$

令 $B = H^T \Lambda^{1/2} H$, 则有 $A = B^2$ 成立, 亦即有 $A^{1/2} = H^T \Lambda^{1/2} H$. 显然, 正定矩阵 A 的平方根是实对称方阵. 我们还需要证明 $A^{1/2}$ 的唯一性, 这个问题留给读者思考. \square

平凡地, 对于任意的正定的对角矩阵 $\Lambda = \text{diag}(a_1, \dots, a_n)$, 其所有对角元都是正数, 则 Λ 的平方根为 $\Lambda^{1/2} = \text{diag}(\sqrt{a_1}, \dots, \sqrt{a_n})$.

3.8 多元 OLS 估计量下拟合值、残差

本节将利用上一节介绍的投影矩阵, 研究 OLS 估计量的拟合值 (fitted value) 向量 $\hat{\mathbf{Y}}$ 和残差 (residual) 向量 $\hat{\boldsymbol{\varepsilon}}$ 的性质.

正如我们在上一章研究过 $\hat{\mathbf{Y}}$ 和残差向量 $\hat{\boldsymbol{\varepsilon}}$, 利用投影矩阵 $\mathbf{P} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ 和归零矩阵 $\mathbf{M} = \mathbf{I} - \mathbf{P}$ 得

$$\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} = \boxed{\mathbf{PY}},$$

$$\hat{\boldsymbol{\varepsilon}} = \hat{\mathbf{Y}} - \mathbf{Y} = \mathbf{PY} - \mathbf{Y} = \boxed{\mathbf{MY}} = \mathbf{M}(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}) = \mathbf{0} \cdot \boldsymbol{\beta} + \mathbf{M}\boldsymbol{\varepsilon} = \boxed{\mathbf{M}\boldsymbol{\varepsilon}},$$

且拟合值向量 $\hat{\mathbf{Y}}$ 和残差向量 $\hat{\boldsymbol{\varepsilon}}$ 是正交的, 这是因为

$$\hat{\mathbf{Y}}^T \hat{\boldsymbol{\varepsilon}} = (\mathbf{PY})^T (\mathbf{MY}) = \mathbf{Y}^T (\mathbf{PM}) \mathbf{Y} = \mathbf{Y}^T \mathbf{O}_n \mathbf{Y} = 0.$$

对于投影矩阵 \mathbf{P} , 其第 i 个对角线元素我们记为 h_{ii} , 则由

$$\mathbf{P} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = \begin{pmatrix} \mathbf{X}_1^T \\ \vdots \\ \mathbf{X}_n^T \end{pmatrix} \underbrace{(\mathbf{X}^T \mathbf{X})^{-1}}_{n\text{阶方阵}} (\mathbf{X}_1, \dots, \mathbf{X}_n)$$

不难推知

$$h_{ii} = \mathbf{X}_i^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_i,$$

我们将 h_{ii} 称为杠杆值 (leverage value).

定理 3.8.1 (杠杆值的性质). 对于投影矩阵 \mathbf{P} 得对角线元素 h_{ii} , 杠杆值 h_{ii} 满足如下性质.

(a) $0 \leq h_{ii} \leq 1$;

(b) 当模型 (3.6.1) 含有常数项时, $h_{ii} \geq n^{-1}$;

(c) $\sum_{i=1}^n h_{ii} = \text{tr} \mathbf{P} = K$.

现在我们来研究拟合值向量 $\hat{\mathbf{Y}}$ 和残差向量 $\hat{\boldsymbol{\varepsilon}}$ 的有限样本性质. 首先, 对于 $\hat{\mathbf{Y}}$ 其条件期望、方差分别为

$$\mathbb{E}(\hat{\mathbf{Y}} | \mathbf{X}) = \mathbb{E}(\mathbf{PY} | \mathbf{X}) = \mathbf{P} \mathbb{E}(\mathbf{Y} | \mathbf{X}),$$

$$\text{Var}(\hat{\mathbf{Y}} | \mathbf{X}) = \mathbf{P} \text{Var}(\mathbf{Y} | \mathbf{X}) \mathbf{P}^T \stackrel{\text{异方差假设式 (3.9)}}{=} \mathbf{P} \mathbf{D} \mathbf{P},$$

在条件同方差式 (3.7) 成立时, 便可以得到

$$\text{Var}(\hat{\mathbf{Y}} | \mathbf{X}) = \mathbf{P} (\mathbf{I} \sigma^2) \mathbf{P} = \mathbf{P} \sigma^2,$$

则有 $\boxed{\text{Var}(\hat{Y}_i | \mathbf{X}) = h_{ii} \sigma^2}$.

对于残差向量 $\hat{\boldsymbol{\varepsilon}}$, 同样计算条件期望、方差得

$$\mathbb{E}(\hat{\boldsymbol{\varepsilon}} | \mathbf{X}) = \mathbb{E}(\mathbf{M}\boldsymbol{\varepsilon} | \mathbf{X}) = \mathbf{M}\mathbb{E}(\boldsymbol{\varepsilon} | \mathbf{X}) \xrightarrow{\text{严格外生性式 (3.5)}} \mathbf{0},$$

$$\text{Var}(\hat{\boldsymbol{\varepsilon}} | \mathbf{X}) = \text{Var}(\mathbf{M}\boldsymbol{\varepsilon} | \mathbf{X}) = \mathbf{M}\text{Var}(\boldsymbol{\varepsilon} | \mathbf{X})\mathbf{M}^T \xrightarrow{\text{异方差假设式 (3.9)}} \mathbf{M}\mathbf{D}\mathbf{M},$$

当条件同方差式 (3.7) 成立时, 则有

$$\text{Var}(\hat{\boldsymbol{\varepsilon}} | \mathbf{X}) = \mathbf{M}(\mathbf{I}\sigma^2)\mathbf{M} = \mathbf{M}\sigma^2,$$

此时残差 $\hat{\varepsilon}_i$ 的条件方差为 $\boxed{\text{Var}(\hat{\varepsilon}_i | \mathbf{X}) = (1 - h_{ii})\sigma^2}$.

与第一章中计算模型 (2.4.2) 的拟合值、残差相比, 我们刚才利用矩阵工具来运算, 要便捷的多. 我们在第一章中提出了**标准化残差** (standardized residual) 为 $\bar{\varepsilon}_i = (1 - h_{ii})^{-\frac{1}{2}} \hat{\varepsilon}_i$, 则标准化残差的向量形式可以表示成

$$\bar{\boldsymbol{\varepsilon}} = \text{diag}\left(\frac{1}{\sqrt{1 - h_{11}}}, \dots, \frac{1}{\sqrt{1 - h_{11}}}\right) \mathbf{M}\boldsymbol{\varepsilon},$$

记 $\mathbf{M}^* = \text{diag}\left(\frac{1}{\sqrt{1 - h_{11}}}, \dots, \frac{1}{\sqrt{1 - h_{11}}}\right)$, 则标准化残差为 $\bar{\boldsymbol{\varepsilon}} = \mathbf{M}^*\mathbf{M}\boldsymbol{\varepsilon}$. 对于标准化残差, 计算条件期望、方差得

$$\mathbb{E}(\bar{\boldsymbol{\varepsilon}} | \mathbf{X}) = \mathbb{E}(\mathbf{M}^*\mathbf{M}\boldsymbol{\varepsilon} | \mathbf{X}) = \mathbf{M}^*\mathbf{M}\mathbb{E}(\boldsymbol{\varepsilon} | \mathbf{X}) = \mathbf{0}$$

与

$$\begin{aligned} \text{Var}(\bar{\boldsymbol{\varepsilon}} | \mathbf{X}) &= \text{Var}(\mathbf{M}^*\mathbf{M}\boldsymbol{\varepsilon} | \mathbf{X}) = \mathbf{M}^*\mathbf{M}\text{Var}(\boldsymbol{\varepsilon} | \mathbf{X})\mathbf{M}\mathbf{M}^* \\ &\xrightarrow{\text{异方差假设式 (3.9)}} \mathbf{M}^*\mathbf{M}\mathbf{D}\mathbf{M}\mathbf{M}^* \xrightarrow{\text{同方差假设式 (3.7)}} \mathbf{M}^*\mathbf{M}\mathbf{M}^*\sigma^2, \end{aligned}$$

则

$$\text{Var}(\bar{\varepsilon}_i | \mathbf{X}) = \frac{1}{\sqrt{1 - h_{ii}}} \cdot (1 - h_{ii})\sigma^2 \cdot \frac{1}{\sqrt{1 - h_{ii}}} = \sigma^2,$$

可见标准化残差 $\bar{\varepsilon}_i$ 与回归误差 ε_i 具有相同的期望和方差.

3.9 Leave-One-Out 回归与预测误差

如果读者在此前有学习过机器学习的相关内容, 则应该了解机器学习中不少方法都是基于**子样本** (sub-sample) 来实现的, 例如自助法 Bootstrap、k-折交叉验证等方法. 本节则介绍基于**差一法** (leave one out)^[5]的线性回归模型, 也即 LOO 回归.

^[5]国内就 leave one out 并没有统一的译法, 最常见的翻译即是“留一法”, 但显然这个翻译是 leave one out 的相反含义. 本笔记后文将 leave one out 简称为 LLO, 而尽量不使用中文译名.

差一法, 即是剔除某个单独的样本后得一个子样本, 然后再将某种统计方法运用于子样本之上, 这即契合了其英文名称 “leave one out”. 据此定义, 显然对于样本容量为 n 的情形 LLO 需要重复统计方法 n 次, 这是十分繁琐的. 但就线性回归模型 (3.6.1), LLO 的结果有着非常简洁的形式, 并且可以利用此去对回归结果做一些细致分析.

模型 3.9.1 (leave-one-out 回归). 在线性回归模型 (3.6.1) 的基础上, 采用 LLO 的回归模型称为 **LLO 回归** (LLO regression). 具体地, 设第 i 个样本为 Y_i 及 $\mathbf{X}_i = (X_{i1}, \dots, X_{ik}, \dots, X_{iK})^T$, 则将利用剔除了第 i 个样本得到的 OLS 估计量记为 $\hat{\beta}_{(-i)}$, 并称第 i 个样本的 **LLO 预测值** (LLO prediction value) 为 $\tilde{Y}_i = \mathbf{X}_i^T \hat{\beta}_{(-i)}$, 同时我们记

$$\tilde{\varepsilon}_i = Y_i - \tilde{Y}_i = Y_i - \mathbf{X}_i^T \hat{\beta}_{(-i)}, \quad (3.28)$$

则称 $\tilde{\varepsilon}_i$ 为第 i 个样本的**预测误差** (prediction error)、**预测残差** (prediction residual) 或 **LLO 残差** (LLO residual).

依据式 (3.20), 则 $\hat{\beta}_{(-i)}$ 可以表示为

$$\hat{\beta}_{(-i)} = \left(\sum_{j \neq i} \mathbf{X}_j \mathbf{X}_j^T \right)^{-1} \sum_{j \neq i} \mathbf{X}_j Y_j,$$

而式 (3.20) 又有依据矩阵分块下的 “内积” 运算得出的, 这样我们便可将 $\hat{\beta}_{(-i)}$ 写作

$$\hat{\beta}_{(-i)} = (\mathbf{X}^T \mathbf{X} - \mathbf{X}_i \mathbf{X}_i^T)^{-1} (\mathbf{X}^T \mathbf{Y} - \mathbf{X}_i Y_i). \quad (3.29)$$

式 (3.29) 可以化简为更简洁的形式, 并与 OLS 估计量联系起来. 将式 (3.29) 左乘以 $\mathbf{X}^T \mathbf{X} - \mathbf{X}_i \mathbf{X}_i^T$, 得到

$$(\mathbf{X}^T \mathbf{X} - \mathbf{X}_i \mathbf{X}_i^T) \hat{\beta}_{(-i)} = \mathbf{X}^T \mathbf{Y} - \mathbf{X}_i Y_i,$$

化简有

$$\mathbf{X}^T \mathbf{X} \hat{\beta}_{(-i)} - \underbrace{\mathbf{X}_i \mathbf{X}_i^T \hat{\beta}_{(-i)}}_{\tilde{Y}_i} = \mathbf{X}^T \mathbf{Y} - \mathbf{X}_i Y_i \quad \Rightarrow \quad \mathbf{X}^T \mathbf{X} \hat{\beta}_{(-i)} = \mathbf{X}^T \mathbf{Y} - \underbrace{\mathbf{X}_i (Y_i - \tilde{Y}_i)}_{\tilde{\varepsilon}_i},$$

故 $\hat{\beta}_{(-i)} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} - (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_i \tilde{\varepsilon}_i$, 注意到 OLS 估计量为 $\hat{\beta}_{\text{OLS}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$, 于是有剔除了第 i 个样本的 LLO 回归结果为

$$\hat{\beta}_{(-i)} = \hat{\beta}_{\text{OLS}} - (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_i \tilde{\varepsilon}_i. \quad (3.30)$$

进一步地, 依据式 (3.30) 则预测误差 $\tilde{\varepsilon}_i$ 为

$$\begin{aligned} \tilde{Y}_i &= \mathbf{X}_i^T \hat{\beta}_{(-i)} = \mathbf{X}_i^T \left(\hat{\beta}_{\text{OLS}} - (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_i \tilde{\varepsilon}_i \right) \\ &= \mathbf{X}_i^T \hat{\beta}_{\text{OLS}} - \mathbf{X}_i^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_i \tilde{\varepsilon}_i \\ &= \hat{Y}_i - h_{ii} \tilde{\varepsilon}_i, \end{aligned}$$

亦即有

$$(Y_i - \hat{Y}_i) + h_{ii}\tilde{\varepsilon}_i = (Y_i - \tilde{Y}_i) \Leftrightarrow \hat{\varepsilon}_i + h_{ii}\tilde{\varepsilon}_i = \tilde{\varepsilon}_i,$$

最终化简得到 OLL 回归的预测误差为

$$\tilde{\varepsilon}_i = \frac{1}{1 - h_{ii}} \hat{\varepsilon}_i. \quad (3.31)$$

需要指出,OLL 方法估计出的回归系数 $\hat{\beta}_{(-i)}$ 是不含有第 i 个样本的信息, 而依据式 (3.30) 和式 (3.28), 则我们也可以从全样本的 OLS 估计量中“剔除”第 i 个样本的信息, 得到 $\hat{\beta}_{(-i)}$. 这样计算 $\hat{\beta}_{(-i)}$ 实则又使用了第 i 个样本的信息.

依据式 (3.31), 我们将预测误差用矩阵符号表示为

$$\tilde{\varepsilon} = \text{diag} \left(\frac{1}{1 - h_{11}}, \dots, \frac{1}{1 - h_{nn}} \right) \hat{\varepsilon} = \text{diag} \left(\frac{1}{1 - h_{11}}, \dots, \frac{1}{1 - h_{nn}} \right) \mathbf{M} \varepsilon,$$

记

$$\mathbf{M}^* = \text{diag} \left(\frac{1}{\sqrt{1 - h_{11}}}, \dots, \frac{1}{\sqrt{1 - h_{nn}}} \right),$$

和

$$\mathbf{M}^{*2} = (\mathbf{M}^*)^2 = \text{diag} \left(\frac{1}{1 - h_{11}}, \dots, \frac{1}{1 - h_{nn}} \right),$$

则预测误差向量满足有 $\tilde{\varepsilon} = \mathbf{M}^{*2} \hat{\varepsilon} = \mathbf{M}^{*2} \mathbf{M} \varepsilon$.

我们现在计算预测误差向量 $\tilde{\varepsilon}$ 的条件期望、方差, 即有

$$\mathbb{E}(\tilde{\varepsilon} | \mathbf{X}) = \mathbb{E}(\mathbf{M}^{*2} \mathbf{M} \varepsilon | \mathbf{X}) = \mathbf{M}^{*2} \mathbf{M} \mathbb{E}(\varepsilon | \mathbf{X}) = 0$$

和

$$\begin{aligned} \text{Var}(\tilde{\varepsilon} | \mathbf{X}) &= \text{Var}(\mathbf{M}^{*2} \hat{\varepsilon} | \mathbf{X}) = \mathbf{M}^{*2} \text{Var}(\hat{\varepsilon} | \mathbf{X}) \mathbf{M}^{*2} \\ &\stackrel{\text{异方差式 (3.9)}}{=} \mathbf{M}^{*2} \mathbf{M} \mathbf{D} \mathbf{M} \mathbf{M}^{*2} \stackrel{\text{同方差式 (3.7)}}{=} \mathbf{M}^{*2} \mathbf{M} \mathbf{M}^{*2} \sigma^2, \end{aligned}$$

也不难得到

$$\text{Var}(\tilde{\varepsilon}_i | \mathbf{X}) = \frac{\text{Var}(\hat{\varepsilon}_i | \mathbf{X})}{(1 - h_{ii})^2} = \frac{(1 - h_{ii}) \sigma^2}{(1 - h_{ii})^2} = \frac{\sigma^2}{1 - h_{ii}}.$$

待补充的内容:

拟合 \neq 插值 \neq 预测

LLO 回归的一个重要用途是用于寻找具有有影响的观测值.

3.10 拟合的度量

我们在上一章介绍了一元线性回归模型 (2.4.1) 的拟合优度 (见定义 (2.5.2)), 本节我们将系统的解释有关拟合的度量的更多概念.

我们在解释拟合优度时, 是基于线性回归模型的方差分析公式对平方和进行了分解, 这在多元的情形自然是成立的. 利用投影矩阵 \mathbf{P} 和归零矩阵 \mathbf{M} , 则由

$$\hat{\mathbf{Y}} = \mathbf{P}\mathbf{Y}, \quad \hat{\boldsymbol{\varepsilon}} = \mathbf{M}\boldsymbol{\varepsilon}, \quad \hat{\mathbf{Y}}^T \hat{\boldsymbol{\varepsilon}} = 0, \quad \hat{\boldsymbol{\varepsilon}}^T \mathbf{1} = 0 \text{ (需要常数项)},$$

可得

$$\mathbf{Y}^T \mathbf{Y} = (\hat{\mathbf{Y}} + \hat{\boldsymbol{\varepsilon}})^T (\hat{\mathbf{Y}} + \hat{\boldsymbol{\varepsilon}}) = \hat{\mathbf{Y}}^T \hat{\mathbf{Y}} + \hat{\boldsymbol{\varepsilon}}^T \hat{\boldsymbol{\varepsilon}} \Leftrightarrow \sum_{i=1}^n Y_i^2 = \sum_{i=1}^n \hat{Y}_i^2 + \sum_{i=1}^n \hat{\varepsilon}_i^2 \quad (3.32)$$

和

$$\begin{aligned} & (\mathbf{Y} - \bar{Y}\mathbf{1})^T (\mathbf{Y} - \bar{Y}\mathbf{1}) \\ &= (\hat{\mathbf{Y}} - \bar{Y}\mathbf{1} + \hat{\boldsymbol{\varepsilon}})^T (\hat{\mathbf{Y}} - \bar{Y}\mathbf{1} + \hat{\boldsymbol{\varepsilon}}) \Leftrightarrow \sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 + \sum_{i=1}^n \hat{\varepsilon}_i^2 \\ &= (\hat{\mathbf{Y}} - \bar{Y}\mathbf{1})^T (\hat{\mathbf{Y}} - \bar{Y}\mathbf{1}) + \hat{\boldsymbol{\varepsilon}}^T \hat{\boldsymbol{\varepsilon}} \end{aligned} \quad (3.33)$$

成立. 式 (3.32) 的分解用到了拟合值向量 $\hat{\mathbf{Y}}$ 与残差向量 $\hat{\boldsymbol{\varepsilon}}$ 正交的性质, 式 (3.32) 在此之上还依赖于模型 (3.6.1) 含义常数项时残差和为 0 的性质. 利用式 (3.32) 和 (3.32), 我们依次定义非中心化 R^2 (uncentered R^2) 为

$$R_{uc}^2 = 1 - \frac{\hat{\boldsymbol{\varepsilon}}^T \hat{\boldsymbol{\varepsilon}}}{\mathbf{Y}^T \mathbf{Y}} = 1 - \frac{\sum_{i=1}^n \hat{Y}_i^2}{\sum_{i=1}^n \hat{\varepsilon}_i^2}, \quad (3.34)$$

R^2 , 亦即拟合优度 (goodness of fit) 或可决系数 (coefficient of determination) 为

$$R^2 = 1 - \frac{\hat{\boldsymbol{\varepsilon}}^T \hat{\boldsymbol{\varepsilon}}}{(\mathbf{Y} - \bar{Y}\mathbf{1})^T (\mathbf{Y} - \bar{Y}\mathbf{1})} = 1 - \frac{\sum_{i=1}^n \hat{\varepsilon}_i^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2}. \quad (3.35)$$

对于拟合优度 R^2 或是非中心化的 R_{uc}^2 , 其都是利用残差平方和 RSS 度量了模型 (3.6.1) 的可解释能力——RSS 越小, 则意味着拟合值 $\hat{\mathbf{Y}}$ 与真实值 \mathbf{Y} 越接近. 依据式 (3.32) 和 (3.32) 可知, R^2 或 R_{uc}^2 的取值范围都是 $[0, 1]$, 取值为 0 时表明解释变量 $\mathbf{X} = (X_1, \dots, X_K)^T$ 与被解释变量 Y 线性无关 (这几乎不会出现), 或出现在被解释变量只有常数项的情形; 取值为 1 时, 则有 $\text{RSS} = 0$, 这意味着样本数据即是纯粹的线性关系, 不存在有误差 $\boldsymbol{\varepsilon}$.

拟合优度 R^2 或是非中心化的 R_{uc}^2 存在的一个较为严重的问题, 则是当解释变量的总数 K 增大时, R^2 和 R_{uc}^2 也会响应的增加. 这是因为在样本容量 n 不变时残差平方和 RSS 是关于 K 的非增函数, 因而 R^2 和 R_{uc}^2 便是 K 的增函数. 这一特性意味解释变量越多, 拟合的效果自然越好. 对此 Theil(1961) 提出了使用**调整的** R^2 (adjusted R^2) 或 \bar{R}^2 来代替 R^2 , 经调整的 \bar{R}^2 被定义为

$$\bar{R}^2 = 1 - \frac{(n-K)^{-1} \hat{\boldsymbol{\varepsilon}}^T \hat{\boldsymbol{\varepsilon}}}{(n-1)^{-1} (\mathbf{Y} - \bar{Y}\mathbf{1})^T (\mathbf{Y} - \bar{Y}\mathbf{1})} = 1 - \frac{\frac{1}{n-K} \sum_{i=1}^n \hat{\varepsilon}_i^2}{\frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2}, \quad (3.36)$$

其即是对残差平方和 RSS 和总体平方和 TSS 予以其自由度的约束 (也即是替换为了 s^2 和 Y 的样本方差), 经过了除以自由度的惩罚后, 在 n 不变的情况下, \bar{R}^2 中的分子不再是关于 K 的非增函数, 这在一定程度上校正了 R^2 和 R_{uc}^2 的问题. 需要指出, 式 (3.36) 可能出现负值, 这是其缺点.

较上述的三种 R^2 而言, Bruce(cite) 给出了一种更为推荐的 \tilde{R}^2 , 其使用预测误差 $\tilde{\varepsilon}_i$ 定义了

$$\tilde{R}^2 = 1 - \frac{\tilde{\boldsymbol{\varepsilon}}^T \tilde{\boldsymbol{\varepsilon}}}{(\mathbf{Y} - \bar{Y}\mathbf{1})^T (\mathbf{Y} - \bar{Y}\mathbf{1})} = 1 - \frac{\sum_{i=1}^n (1 - h_{ii})^{-2} \hat{\varepsilon}_i^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2}, \quad (3.37)$$

依据第节的内容, 可以推知 \tilde{R}^2 可以作为预测方差中被能回归模型所够预测的百分比的估计量. \tilde{R}^2 能够完全避免由于解释变量导致的数值提升问题, 其取值也可能出现负值的情形, 但这仅在模型的预测能力比仅有常数项的模型还差时出现.

需要指出, 本节介绍的四种 R^2 仅仅度量了模型对于样本数据的拟合情况, 其不能说明回归模型的经济含义. 对于同一个经济学模型, 将之做适当变换后得到的计量模型的形式往往不止一个, 这些形式都可以被用于正确地估计出经济学模型中的参数, 但相应的回归结果往往具有不同的 R^2 . 此外, 本节介绍的四种 R^2 在机器学习中也广泛使用, 且在处理高维数据和非参数问题时 \tilde{R}^2 被广泛采用.

最后, 我们给出 R^2 与 \bar{R}^2 和 R_{uc}^2 直接的换算关系. 对于调整的 R^2 . 由于

$$\begin{aligned} \bar{R}^2 &= 1 - \frac{\boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon} / (n-K)}{(\mathbf{Y} - \mathbf{1}\bar{Y})^T (\mathbf{Y} - \mathbf{1}\bar{Y}) / (n-1)} = 1 - \frac{\boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon}}{(\mathbf{Y} - \mathbf{1}\bar{Y})^T (\mathbf{Y} - \mathbf{1}\bar{Y})} \cdot \frac{n-1}{n-K} \\ &= \frac{n-1}{n-K} \left(\frac{n-K}{n-1} - \frac{\boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon}}{(\mathbf{Y} - \mathbf{1}\bar{Y})^T (\mathbf{Y} - \mathbf{1}\bar{Y})} \right) \\ &= \frac{n-1}{n-K} \left(\frac{n-K}{n-1} - 1 + 1 - \frac{\boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon}}{(\mathbf{Y} - \mathbf{1}\bar{Y})^T (\mathbf{Y} - \mathbf{1}\bar{Y})} \right) \\ &= \frac{n-1}{n-K} \left(\frac{n-K}{n-1} - 1 + R^2 \right) = \frac{n-1}{n-K} \left(\frac{1-K}{n-1} + R^2 \right), \end{aligned}$$

故化简有

$$1 - R^2 = \frac{n - K}{n - 1} (1 - \bar{R}^2). \quad (3.38)$$

对于 R^2 与非中心化 R^2 , 则有

$$\begin{aligned} 1 - R^2 &= \frac{\hat{\epsilon}^T \hat{\epsilon}}{(\mathbf{Y} - \mathbf{1}\bar{Y})^T (\mathbf{Y} - \mathbf{1}\bar{Y})} = \frac{\mathbf{Y}^T \mathbf{Y}}{(\mathbf{Y} - \mathbf{1}\bar{Y})^T (\mathbf{Y} - \mathbf{1}\bar{Y})} \cdot \frac{\hat{\epsilon}^T \hat{\epsilon}}{\mathbf{Y}^T \mathbf{Y}} \\ &= \frac{\mathbf{Y}^T \mathbf{Y}}{(\mathbf{Y} - \mathbf{1}\bar{Y})^T (\mathbf{Y} - \mathbf{1}\bar{Y})} \cdot \left[1 - \left(1 - \frac{\hat{\epsilon}^T \hat{\epsilon}}{\mathbf{Y}^T \mathbf{Y}} \right) \right] = \frac{\mathbf{Y}^T \mathbf{Y}}{(\mathbf{Y} - \mathbf{1}\bar{Y})^T (\mathbf{Y} - \mathbf{1}\bar{Y})} \cdot (1 - R_{uc}^2) \\ &= \frac{(\mathbf{Y} - \bar{Y}\mathbf{1})^T (\mathbf{Y} - \bar{Y}\mathbf{1}) + n\bar{Y}^2}{(\mathbf{Y} - \mathbf{1}\bar{Y})^T (\mathbf{Y} - \mathbf{1}\bar{Y})} \cdot (1 - R_{uc}^2) = \left[1 + \frac{n\bar{Y}^2}{(\mathbf{Y} - \mathbf{1}\bar{Y})^T (\mathbf{Y} - \mathbf{1}\bar{Y})} \right] (1 - R_{uc}^2) \end{aligned}$$

即有

$$1 - R^2 = \left[1 + \frac{n\bar{Y}^2}{(\mathbf{Y} - \mathbf{1}\bar{Y})^T (\mathbf{Y} - \mathbf{1}\bar{Y})} \right] (1 - R_{uc}^2). \quad (3.39)$$

3.11 多元 OLS 估计量下的方差估计量

本节我们将研究模型 (3.6.1) 的各种方差估计量, 除了在上一章已经有涉及的回归方差 σ^2 的估计量外, 本节我们还将介绍 OLS 估计量 $\hat{\beta}$ 的方差估计量.

现在我们研究总体参数回归方差 $\mathbb{E}\epsilon^2 = \sigma^2$ 的估计量, 同样地, 结合条件同方差 $\mathbb{E}(\epsilon_i^2 | \mathbf{X}) = \sigma^2$ 依据矩估计的思想有

$$\hat{\sigma}_{\text{ideal}}^2 = \frac{1}{n} \sum_{i=1}^n \epsilon_i^2 = \frac{1}{n} \boldsymbol{\epsilon}^T \boldsymbol{\epsilon},$$

且

$$\hat{\sigma}_{\text{ideal}}^2 = \frac{1}{n} \mathbb{E}(\boldsymbol{\epsilon}^T \boldsymbol{\epsilon} | \mathbf{X}) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}(\epsilon_i^2 | \mathbf{X}) \xrightarrow{\text{条件同方差式 (3.7)}} \frac{1}{n} \sum_{i=1}^n \sigma^2 = \sigma^2,$$

即可以推知 $\hat{\sigma}_{\text{ideal}}^2$ 是 σ^2 的无偏估计. 但回归误差 ϵ_i 是不可观察的, 使用残差 $\hat{\epsilon}_i$ 替换有

$$\frac{1}{n} \hat{\epsilon}^T \hat{\epsilon} = \frac{1}{n} \boldsymbol{\epsilon}^T \mathbf{M} \mathbf{M} \boldsymbol{\epsilon} = \frac{1}{n} \boldsymbol{\epsilon}^T \mathbf{M} \boldsymbol{\epsilon} = \frac{1}{n} \boldsymbol{\epsilon}^T (\mathbf{I} - \mathbf{P}) \boldsymbol{\epsilon} \leq \frac{1}{n} \boldsymbol{\epsilon}^T \boldsymbol{\epsilon} = \hat{\sigma}_{\text{ideal}}^2,$$

可以知悉直接替换的得到的 $\frac{1}{n} \hat{\epsilon}^T \hat{\epsilon}$ 是有偏的. 为此, 计算条件期望

$$\begin{aligned} \mathbb{E}(\hat{\epsilon}^T \hat{\epsilon} | \mathbf{X}) &= \mathbb{E}(\boldsymbol{\epsilon}^T \mathbf{M} \boldsymbol{\epsilon} | \mathbf{X}) = \mathbb{E}(\text{tr}(\boldsymbol{\epsilon}^T \mathbf{M} \boldsymbol{\epsilon}) | \mathbf{X}) = \mathbb{E}(\text{tr}(\mathbf{M} \boldsymbol{\epsilon} \boldsymbol{\epsilon}^T) | \mathbf{X}) \\ &\xrightarrow{\text{迹是线性变换}} \text{tr}(\mathbb{E}(\mathbf{M} \boldsymbol{\epsilon} \boldsymbol{\epsilon}^T | \mathbf{X})) = \text{tr}(\mathbf{M} \mathbb{E}(\boldsymbol{\epsilon} \boldsymbol{\epsilon}^T | \mathbf{X})) \xrightarrow{\text{异方差式 (3.9)}} \text{tr}(\mathbf{M} \mathbf{D}), \end{aligned}$$

则在同方差假设式 (3.7) 有

$$\mathbb{E}(\hat{\epsilon}^T \hat{\epsilon} | \mathbf{X}) = \text{tr}(\mathbf{M} \mathbf{I} \sigma^2) = (n - K) \sigma^2,$$

同样地, 我们得到 σ^2 的偏差校正估计量 (bias-corrected estimator) 为

$$s^2 = \frac{1}{n-K} \hat{\boldsymbol{\varepsilon}}^T \hat{\boldsymbol{\varepsilon}} = \frac{1}{n-K} \sum_{i=1}^n \varepsilon_i^2. \quad (3.40)$$

另一种回归方差方式则是使用具体相同方差的标准化残差 $\bar{\varepsilon}_i$ 替换回归误差 ε_i , 即有方差估计量

$$\bar{\sigma}^2 = \frac{1}{n} \bar{\boldsymbol{\varepsilon}}^T \bar{\boldsymbol{\varepsilon}} = \frac{1}{n} \sum_{i=1}^n \bar{\varepsilon}_i^2 = \frac{1}{n} \sum_{i=1}^n \frac{\hat{\varepsilon}_i^2}{\sqrt{1-h_{ii}}}, \quad (3.41)$$

不难计算

$$\begin{aligned} \mathbb{E}(\bar{\boldsymbol{\varepsilon}}^T \bar{\boldsymbol{\varepsilon}} | \mathbf{X}) &= \mathbb{E}\left((\mathbf{M}^* \mathbf{M} \boldsymbol{\varepsilon})^T (\mathbf{M}^* \mathbf{M} \boldsymbol{\varepsilon}) \middle| \mathbf{X}\right) = \mathbb{E}(\boldsymbol{\varepsilon}^T (\mathbf{M} \mathbf{M}^*) (\mathbf{M}^* \mathbf{M}) \boldsymbol{\varepsilon} | \mathbf{X}) \\ &= \mathbb{E}(\text{tr}(\boldsymbol{\varepsilon}^T \mathbf{M} (\mathbf{M}^*)^2 \mathbf{M} \boldsymbol{\varepsilon}) | \mathbf{X}) = \mathbb{E}(\text{tr}(\boldsymbol{\varepsilon}^T (\mathbf{M}^*)^2 \mathbf{M} \boldsymbol{\varepsilon}) | \mathbf{X}) \\ &= \mathbb{E}(\text{tr}((\mathbf{M}^*)^2 \mathbf{M} \boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^T) | \mathbf{X}) = \text{tr}(\mathbb{E}((\mathbf{M}^*)^2 \mathbf{M} \boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^T | \mathbf{X})) \\ &= \text{tr}((\mathbf{M}^*)^2 \mathbf{M} \mathbb{E}(\boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^T | \mathbf{X})) \stackrel{\text{异方差式 (3.9)}}{=} \text{tr}((\mathbf{M}^*)^2 \mathbf{M} \mathbf{D}) \\ &= \text{tr}(\text{diag}(1-h_{11}, \dots, 1-h_{nn}) \mathbf{M} \mathbf{D}) \\ &\stackrel{\text{同方差式 (3.7)}}{=} \sigma^2 \text{tr}(\text{diag}(1-h_{11}, \dots, 1-h_{nn}) \mathbf{M}) \\ &= \sigma^2 \sum_{i=1}^n \frac{1-h_{ii}}{1-h_{ii}} = n\sigma^2, \end{aligned}$$

由此可以推知, 式 (3.41) 的 $\bar{\sigma}^2$ 是 σ^2 的无偏估计, 因而也是 σ^2 的偏差校正估计量.

线性回归模型 (3.3.1) 的方差中除了 σ^2 之外, 我们还关注 OLS 估计量的协方差矩阵, 因为由其可以衍生出的各个 $\hat{\beta}_i$ 的方差. 依据定理 (3.6.1), 我们已经计算出 OLS 估计量的协方差矩阵为

$$\mathbf{V}_{\hat{\boldsymbol{\beta}}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{D} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1},$$

在同方差式 (3.7) 的情况下则协方差矩阵化简为了 $\mathbf{V}_{\hat{\boldsymbol{\beta}}} = (\mathbf{X}^T \mathbf{X})^{-1} \sigma^2$. 由于 σ^2 未知, 则需要使用偏差校正估计量去代替. 对于同方差的情形, 最为简单的协方差矩阵估计量即为

$$\hat{\mathbf{V}}_{\hat{\boldsymbol{\beta}}}^0 = (\mathbf{X}^T \mathbf{X})^{-1} s^2, \quad (3.42)$$

对协方差矩阵 $\hat{\mathbf{V}}_{\hat{\boldsymbol{\beta}}}^0$ 计算条件期望得

$$\mathbb{E}(\hat{\mathbf{V}}_{\hat{\boldsymbol{\beta}}}^0 | \mathbf{X}) = \mathbb{E}((\mathbf{X}^T \mathbf{X})^{-1} s^2 | \mathbf{X}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbb{E}(s^2 | \mathbf{X}) = (\mathbf{X}^T \mathbf{X})^{-1} \sigma^2 = \mathbf{V}_{\hat{\boldsymbol{\beta}}},$$

则可以进一步推知在同方差的条件下 $\hat{\mathbf{V}}_{\hat{\boldsymbol{\beta}}}^0$ 是 $\mathbf{V}_{\hat{\boldsymbol{\beta}}}$ 的无偏估计. 实际上, 通过这里的计算也不难知悉, 用任意的偏差校正估计量替换 σ^2 , 得到的协方差矩阵估计量都是无偏的. $\hat{\mathbf{V}}_{\hat{\boldsymbol{\beta}}}^0$ 是多数计量经济学软件 (包括 Stata、eViews) 的默认的协方差矩阵估计量.

倘若我们的模型确实满足假设 (3.2.1) 而不存在异方差问题, 利用式 (3.42) 去计算 OLS 估计量的方差, 通常而言不会有较大的偏差. 但一旦模型 (3.6.1) 是存在异方差式 (3.9), 则式 (3.2.1) 将与 $V_{\hat{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{D} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}$ 存在较大的偏离, 在 Bruce(cite) 中这种偏离可以高达 30 倍! 研究异方差下如何计算 OLS 估计量的协方差矩阵的估计值, 这是一个非常重要的问题.

记对角矩阵 $\mathbf{D}_{\varepsilon} = \text{diag}(\varepsilon_1^2, \dots, \varepsilon_n^2)$, 则有

$$\mathbb{E}(\mathbf{D}_{\varepsilon} | \mathbf{X}) = \mathbb{E}(\boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^T | \mathbf{X}) = \text{diag}(\sigma_1^2, \dots, \sigma_n^2) = \mathbf{D},$$

故 \mathbf{D}_{ε} 是条件异方差矩阵 \mathbf{D} 的条件无偏估计.

异方差的情况下, OLS 估计量的协方差矩阵 $V_{\hat{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{D} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}$, 由于

$$\mathbf{X}^T \mathbf{D} \mathbf{X} = \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \sigma_i^2,$$

则我们需要对 $\sigma_1^2, \dots, \sigma_n^2$ 进行估计. 显然, 理想的情况即是使用回归误差的平方 $\varepsilon_1^2, \dots, \varepsilon_n^2$ 去估计, 也即是使用 \mathbf{D}_{ε} 去估计 \mathbf{D} . 记

$$\hat{V}_{\hat{\beta}}^{\text{ideal}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{D}_{\varepsilon} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} = (\mathbf{X}^T \mathbf{X})^{-1} \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \varepsilon_i^2 \right) (\mathbf{X}^T \mathbf{X})^{-1}, \quad (3.43)$$

则

$$\begin{aligned} \mathbb{E}(\hat{V}_{\hat{\beta}}^{\text{ideal}} | \mathbf{X}) &= \mathbb{E}((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{D}_{\varepsilon} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} | \mathbf{X}) \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{E}(\mathbf{D}_{\varepsilon} | \mathbf{X}) \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{D} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} = V_{\hat{\beta}}, \end{aligned}$$

即 $\hat{V}_{\hat{\beta}}^{\text{ideal}}$ 是协方差矩阵 $V_{\hat{\beta}}$ 的无偏估计量. 可惜, 回归误差 $\boldsymbol{\varepsilon}$ 不可观察, 因此我们需要其他的估计方法.

最基本的方法即是使用残差 $\hat{\varepsilon}_i$ 替代回归误差 ε_i , 即有

$$\hat{V}_{\hat{\beta}}^{\text{HC0}} = (\mathbf{X}^T \mathbf{X})^{-1} \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \hat{\varepsilon}_i^2 \right) (\mathbf{X}^T \mathbf{X})^{-1}, \quad (3.44)$$

这里的上标 HC 是 heteroskedasticity-consistent 的缩写, 而 HC0 则是代表 $\hat{V}_{\hat{\beta}}^{\text{HC0}}$ 是异方差下协方差矩阵估计量的基准. $\hat{V}_{\hat{\beta}}^{\text{HC0}}$ 是由 Eicker(1963,cite) 提出的, 最早由 White(1980,cite) 引入计量经济学中, 因此 $\hat{V}_{\hat{\beta}}^{\text{HC1}}$ 又被称为 **Eicker-White** 协方差矩阵估计量 (Eicker-White covariance matrix estimator) 或 **White** 协方差矩阵估计量 (White covariance matrix estimator).

依据我们多次计算的经验和, 显然, 由于 $\hat{\varepsilon}_i$ 的条件方差较 ε_i 偏小, 则式 (3.44) 会是有偏的, 仿照校正直接替换 $\frac{1}{n} \boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon}$ 的方法, 则 $\hat{V}_{\hat{\beta}}^{\text{HC0}}$ 的校正估计量为

$$\hat{V}_{\hat{\beta}}^{\text{HC1}} = \frac{n}{n-K} \hat{V}_{\hat{\beta}}^{\text{HC0}} = \frac{n}{n-K} (\mathbf{X}^T \mathbf{X})^{-1} \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \hat{\varepsilon}_i^2 \right) (\mathbf{X}^T \mathbf{X})^{-1}, \quad (3.45)$$

式 (3.45) 比式 (3.44) 更为推荐, 其予以了 $\hat{\mathbf{V}}_{\beta}^{\text{HC0}}$ 的 $n - K$ 的自由度的惩罚, 这种形式的协方差矩阵估计量由 Hinkley(1977,cite) 提出, 在 Stata 软件中输入可选命令 “,r” 时, Stata 便使用 $\hat{\mathbf{V}}_{\beta}^{\text{HC1}}$ 作为协方差矩阵的估计量.

另几种常见的方法便是使用标准化残差 $\bar{\varepsilon}_i$ 或预测误差 $\tilde{\varepsilon}_i$ 去替换残差, 即有

$$\begin{aligned}\hat{\mathbf{V}}_{\beta}^{\text{HC2}} &= (\mathbf{X}^T \mathbf{X})^{-1} \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \bar{\varepsilon}_i^2 \right) (\mathbf{X}^T \mathbf{X})^{-1} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \left(\sum_{i=1}^n (1 - h_{ii})^{-1} \mathbf{x}_i \mathbf{x}_i^T \hat{\varepsilon}_i^2 \right) (\mathbf{X}^T \mathbf{X})^{-1},\end{aligned}\quad (3.46)$$

$$\begin{aligned}\hat{\mathbf{V}}_{\beta}^{\text{HC3}} &= (\mathbf{X}^T \mathbf{X})^{-1} \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \tilde{\varepsilon}_i^2 \right) (\mathbf{X}^T \mathbf{X})^{-1} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \left(\sum_{i=1}^n (1 - h_{ii})^{-2} \mathbf{x}_i \mathbf{x}_i^T \hat{\varepsilon}_i^2 \right) (\mathbf{X}^T \mathbf{X})^{-1},\end{aligned}\quad (3.47)$$

我们将式 (3.44)、(3.45)、(3.46) 和 (3.47) 统称为稳健的协方差矩阵估计量 (robust covariance matrix estimator)、异方差一致的协方差矩阵估计量 (heteroskedasticity-consistent covariance matrix estimator) 或异方差稳健的协方差矩阵估计量 (heteroskedasticity-robust covariance matrix estimator). 使用使用标准化残差 $\bar{\varepsilon}_i$ 替换得到的 $\hat{\mathbf{V}}_{\beta}^{\text{HC2}}$ 是 Horn 和 Duncan(1975,cite) 提出的, 在 Stata 中可以使用可选项 “,vce(2)” 来汇报依据 $\hat{\mathbf{V}}_{\beta}^{\text{HC2}}$ 计算的稳健的回归结果. $\hat{\mathbf{V}}_{\beta}^{\text{HC3}}$ 则是最晚提出的协方差矩阵估计量, 其依据 jackknife 原则由 MacKinnon and (1985,cite) 提出, 在 Stata 中通过可选命令 “,vce(3)” 来汇报相关结果.

由于 $1 < \frac{1}{1 - h_{ii}} < \frac{1}{(1 - h_{ii})^2}$, $0 < h_{ii} < 1$, 则可以推知

$$\hat{\mathbf{V}}_{\beta}^{\text{HC0}} < \hat{\mathbf{V}}_{\beta}^{\text{HC2}} < \hat{\mathbf{V}}_{\beta}^{\text{HC3}}, \quad (3.48)$$

这里 $<$ 是针对二次型而言的, 即 $\mathbf{A} < \mathbf{B}$ 等价于 $\mathbf{B} - \mathbf{A}$ 是正定的二次型.

对于 $\hat{\mathbf{V}}_{\beta}^{\text{HC0}}$ 、 $\hat{\mathbf{V}}_{\beta}^{\text{HC2}}$ 和 $\hat{\mathbf{V}}_{\beta}^{\text{HC3}}$, 我们可以尝试计算条件期望以判断其无偏性. 以 $\hat{\mathbf{V}}_{\beta}^{\text{HC0}}$ 为例, 即有

$$\begin{aligned}\mathbb{E} \left(\hat{\mathbf{V}}_{\beta}^{\text{HC0}} \middle| \mathbf{X} \right) &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbb{E} \left(\left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \hat{\varepsilon}_i^2 \right) \middle| \mathbf{X} \right) (\mathbf{X}^T \mathbf{X})^{-1} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \mathbb{E}(\hat{\varepsilon}_i^2 | \mathbf{X}) \right) (\mathbf{X}^T \mathbf{X})^{-1},\end{aligned}$$

对于 $\hat{\mathbf{V}}_{\beta}^{\text{HC2}}$ 和 $\hat{\mathbf{V}}_{\beta}^{\text{HC3}}$ 的计算是同理的, 也就是说, $\hat{\mathbf{V}}_{\beta}^{\text{HC0}}$ 、 $\hat{\mathbf{V}}_{\beta}^{\text{HC2}}$ 和 $\hat{\mathbf{V}}_{\beta}^{\text{HC3}}$ 的条件期望依次取决于条件二阶原点矩 $\mathbb{E}(\hat{\varepsilon}_i^2 | \mathbf{X})$ 、 $\mathbb{E}(\bar{\varepsilon}_i^2 | \mathbf{X})$ 和 $\mathbb{E}(\tilde{\varepsilon}_i^2 | \mathbf{X})$. 由于残差、标准化残差和预测误差的条件期望为 0, 则其条件二阶矩即等于其期望, 依据此前的计算有

$$\text{Var}(\hat{\varepsilon} | \mathbf{X}) \xrightarrow{\text{异方差}} \mathbf{MDM} \Rightarrow \mathbb{E}(\hat{\varepsilon}_i^2 | \mathbf{X}) = \text{Var}(\hat{\varepsilon}_i | \mathbf{X}) = (\mathbf{MDM})_{ii} \xrightarrow{\text{同方差}} (1 - h_{ii}) \sigma^2,$$

以及

$$\begin{aligned}\mathbb{E}(\tilde{\varepsilon}_i^2 | \mathbf{X}) &= \text{Var}(\tilde{\varepsilon}_i | \mathbf{X}) = \text{Var}\left(\frac{\hat{\varepsilon}_i}{\sqrt{1-h_{ii}}} \middle| \mathbf{X}\right) = \frac{\text{Var}(\hat{\varepsilon}_i | \mathbf{X})}{1-h_{ii}} \frac{\text{同方差}}{\sigma^2}, \\ \mathbb{E}(\tilde{\varepsilon}_i^2 | \mathbf{X}) &= \text{Var}(\tilde{\varepsilon}_i | \mathbf{X}) = \text{Var}\left(\frac{\hat{\varepsilon}_i}{1-h_{ii}} \middle| \mathbf{X}\right) = \frac{\text{Var}(\hat{\varepsilon}_i | \mathbf{X})}{(1-h_{ii})^2} \frac{\text{同方差}}{\sigma^2} \frac{\sigma^2}{1-h_{ii}},\end{aligned}$$

异方差的情形下需要计算 \mathbf{MDM} , 其形式较为复杂, 在此暂不展开讨论. 但在同方差式 (3.7) 的条件下, 对于 $\hat{\mathbf{V}}_{\beta}^{\text{HC0}}$ 、 $\hat{\mathbf{V}}_{\beta}^{\text{HC2}}$ 和 $\hat{\mathbf{V}}_{\beta}^{\text{HC3}}$ 这三个异方差稳健的估计量, 则容易得到

$$\mathbb{E}(\hat{\mathbf{V}}_{\beta}^{\text{HC0}} | \mathbf{X}) < \mathbf{V}_{\beta} < \mathbb{E}(\hat{\mathbf{V}}_{\beta}^{\text{HC3}} | \mathbf{X}), \quad \mathbf{V}_{\beta} = \mathbb{E}(\hat{\mathbf{V}}_{\beta}^{\text{HC2}} | \mathbf{X}), \quad (3.49)$$

这里的 $<$ 依然是针对二次型的正定性而言的. 式 (3.49) 与式 (3.48) 的结果是相同的. 尽管这三个估计量是针对异方差情形的, 但在同方差的情况下我们还是能够借此对三者进行直观的比较.

在具体的回归分析中, 无论我们选择了哪一个估计量去估计 \mathbf{V}_{β} , 我们总是关系协方差矩阵的对角线元素, 其是 OLS 估计量的方差估计值.

定义 3.11.1 (标准误差, 标准差). 对于 OLS 估计量协方差矩阵 \mathbf{V}_{β} 的估计量 $\hat{\mathbf{V}}_{\beta}$, 我们称 $\hat{\mathbf{V}}_{\beta}$ 的对角线元素的平方根

$$\text{SE}(\hat{\beta}_k) = \sqrt{(\hat{\mathbf{V}}_{\beta})_{kk}} \quad (3.50)$$

为 OLS 估计量 $\hat{\beta}_k$ 的标准误差 (standard error) 或标准误.

依据定义 (3.11.1), 标准误 $\text{SE}(\hat{\beta}_k)$ 实际上是 OLS 估计量 $\hat{\beta}_k$ 的标准差的估计量. 当模型 (3.6.1) 满足有同方差假设式 (3.7) 时, 则依据式 (3.42) 则式 (3.50) 等于

$$\text{SE}(\hat{\beta}_k) = s \sqrt{((\mathbf{X}^T \mathbf{X})^{-1})_{kk}}. \quad (3.51)$$

下面我们通过数据实例来说明 $\hat{\mathbf{V}}_{\beta}^0$ 、 $\hat{\mathbf{V}}_{\beta}^{\text{HC0}}$ 、 $\hat{\mathbf{V}}_{\beta}^{\text{HC1}}$ 、 $\hat{\mathbf{V}}_{\beta}^{\text{HC2}}$ 和 $\hat{\mathbf{V}}_{\beta}^{\text{HC3}}$ 这五个协方差矩阵估计量的区别. 考虑一个如下的存在异方差的一元线性回归模型: 假设解释变量 X 仅有取值 $X = 0, 5, 10, 15$, 数据生成过程为 $Y_i = 10 + 25X_i + \varepsilon_i$ 且有条件异方差

$$\mathbb{E}(\varepsilon_i | X = 0) = 25^2, \quad \mathbb{E}(\varepsilon_i | X = 5) = 30^2, \quad \mathbb{E}(\varepsilon_i | X = 15) = 15^2, \quad \mathbb{E}(\varepsilon_i | X = 20) = 5^2,$$

成立. 设 X 取值 0, 5, 10, 15 的样本数分别为 $n_1 = 25, n_2 = 30, n_3 = 15, n_4 = 20$, 则样本容量 n 为 $n = n_1 + n_2 + n_3 + n_4 = 90$, 这样, 我们便可以得到该模型的协方差矩阵为

$$\mathbf{D} = \begin{pmatrix} 25^2 \mathbf{I}_{25} & & & \\ & 30^2 \mathbf{I}_{30} & & \\ & & 15^2 \mathbf{I}_{15} & \\ & & & 5^2 \mathbf{I}_{20} \end{pmatrix}.$$

我们对这个异方差的回归模型使用 Monte Carlo 模拟以求解 OLS 估计量. 依据式 (3.19) 则有

$$\hat{\beta}_{\text{OLS}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} = \begin{pmatrix} 21.3049 \\ 10.2696 \end{pmatrix},$$

即有样本回归直线为 $Y = 10.2696 + 21.3049X$, 可见异方差下 $\hat{\beta}_1$ 与真实值 $\beta_1 = 25$ 有较大的偏差. 现我们依据本节的推导, 计算 OLS 估计量的标准误, 结果如表 (3.2) 所示. 不难发现, 与真实的 $\mathbf{V}_{\hat{\beta}}$ 所计算的 OLS 估计量的标准差相比, 同方差假设下 $\hat{\mathbf{V}}_{\hat{\beta}}^0$ 计算的标准误出现了非常大的偏差, 对于 $\hat{\beta}_0$ 的标准差估计的误差为 11.53%, 而 $\hat{\beta}_1$ 的标准差估计的误差更是高达 -29.54%! 而余下的四种考虑了异方差的协方差矩阵估计量, 其计算的标准误较真实的标准差而言, 则误差控制在了 5% 以下, 且不难验证式 (3.48) 的二次型关系 (对角元的大小).

表 3.2: 协方差矩阵估计量

	$\hat{\beta}_0$ 的标准误	SE($\hat{\beta}_0$) 与 真实值的误差	$\hat{\beta}_1$ 的标准误	SE($\hat{\beta}_1$) 与 真实值的误差
真实值	4.2300	0.00%	0.3336	0.00%
$\hat{\mathbf{V}}_{\hat{\beta}}^0$	3.7425	11.53%	0.4321	-29.54%
$\hat{\mathbf{V}}_{\hat{\beta}}^{\text{HC0}}$	4.3036	-1.74%	0.3394	-1.74%
$\hat{\mathbf{V}}_{\hat{\beta}}^{\text{HC1}}$	4.3522	-2.89%	0.3433	-2.89%
$\hat{\mathbf{V}}_{\hat{\beta}}^{\text{HC2}}$	4.3526	-2.90%	0.3439	-3.10%
$\hat{\mathbf{V}}_{\hat{\beta}}^{\text{HC3}}$	4.4024	-4.08%	0.3485	-4.47%

3.12 多元正态回归模型及其区间估计

本节我们将介绍多元正态线性回归模型, 其是引入了正态假设的模型 (3.6.1), 是一种特例的线性回归模型. 在此之上, 我们便可以给出 OLS 估计量的区间估计与假设检验.

模型 3.12.1 (多元正态回归模型). 在多元线性回归模型 (3.6.1) 的假设之上, 补充假设

$$\text{总体: } \varepsilon | X \sim N(0, \sigma^2), \quad (3.52)$$

$$\text{样本: } \varepsilon | \mathbf{X} \sim N(\mathbf{0}, \mathbf{I}_n \sigma^2), \quad (3.53)$$

成立时, 满足如下设定的模型称为多元正态回归模型 (*normal regression model*).

$$\begin{aligned} Y &= \mathcal{P}(Y|\mathbf{X}) + \varepsilon \\ \mathcal{P}(Y|\mathbf{X}) &= \mathbf{X}^T \boldsymbol{\beta}, \quad \boldsymbol{\beta} = (\mathbb{E}(\mathbf{X}\mathbf{X}^T))^{-1} \mathbb{E}(\mathbf{X}Y) \\ \mathbb{E}(\mathbf{X}\varepsilon) &= 0 \\ \mathbb{E}(\varepsilon) &= 0 \quad (\text{不需要常数项}) \\ \text{Var}(\varepsilon) &= \mathbb{E}\varepsilon^2 = \mathbb{E}(\varepsilon|\mathbf{X}) = \sigma^2 \end{aligned}$$

与在上一章的分析一样, 当总体满足正态假设 (3.52) 时, 则 $Y = \mathbf{X}^T \boldsymbol{\beta} + \varepsilon$ 满足有条件分布 $Y|\mathbf{X} \sim N(\mathbf{X}^T \boldsymbol{\beta}, \sigma^2)$. 相应地, 正态假设 (3.53) 使得样本 $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ 有条件分布 $\mathbf{Y}|\mathbf{X} \sim N(\mathbf{X}\boldsymbol{\beta}, \mathbf{I}_n \sigma^2)$ 成立, 亦即有 $Y_i|\mathbf{X}_i \sim N(\mathbf{X}_i^T \boldsymbol{\beta}, \sigma^2)$. 注意, 正态假设的式 (3.52) 和 (3.53) 是给出了回归误差的条件分布为正态分布——正态分布仅仅依赖于期望和方差——此时的回归误差便不受到解释变量的影响, 即模型 (3.12.1) 自然有同方差假设成立.

在此之上, 我们现在研究 OLS 估计量所具有的概率分布. 由于抽样误差为 $\hat{\boldsymbol{\beta}} - \boldsymbol{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\varepsilon}$, 则其是服从多元正态分布的回归误差 $\boldsymbol{\varepsilon}$ 向量的线性变换, 故抽样误差依然服从正态分布, 且满足有

$$\hat{\boldsymbol{\beta}} - \boldsymbol{\beta} | \mathbf{X} \sim N(\mathbf{0}, (\mathbf{X}^T \mathbf{X})^{-1} \sigma^2),$$

由此推出 OLS 估计量的条件分布为 $\boxed{\hat{\boldsymbol{\beta}} | \mathbf{X} \sim N(\boldsymbol{\beta}, (\mathbf{X}^T \mathbf{X})^{-1} \sigma^2)}$, 据此则可以得到

$$\hat{\beta}_k | \mathbf{X} \sim N\left(\beta_k, \left((\mathbf{X}^T \mathbf{X})^{-1} \sigma^2\right)_{kk}\right). \quad (3.54)$$

对于拟合值向量 $\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$ 和残差向量 $\hat{\boldsymbol{\varepsilon}} = \mathbf{M}\boldsymbol{\varepsilon}$, 同理可以推知

$$\hat{\mathbf{Y}} | \mathbf{X} \sim N(\mathbf{X}\boldsymbol{\beta}, \mathbf{P}\sigma^2), \quad \hat{\boldsymbol{\varepsilon}} | \mathbf{X} \sim N(\mathbf{0}, \mathbf{M}\sigma^2). \quad (3.55)$$

特别地, 由于抽样误差 $\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}$ 和残差向量 $\hat{\boldsymbol{\varepsilon}}$ 均是误差向量 $\boldsymbol{\varepsilon}$ 的线性变换, 故我们考虑矩阵分块有

$$\begin{pmatrix} \hat{\boldsymbol{\beta}} - \boldsymbol{\beta} \\ \hat{\boldsymbol{\varepsilon}} \end{pmatrix} = \begin{pmatrix} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\varepsilon} \\ \mathbf{M}\boldsymbol{\varepsilon} \end{pmatrix} = \begin{pmatrix} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \\ \mathbf{M} \end{pmatrix} \boldsymbol{\varepsilon},$$

这样我们计算 $\left((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T, \mathbf{M}\right)^T \boldsymbol{\varepsilon}$ 的条件期望、方差为

$$\begin{aligned} \mathbb{E}\left(\begin{pmatrix} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \\ \mathbf{M} \end{pmatrix} \boldsymbol{\varepsilon} \middle| \mathbf{X}\right) &= \begin{pmatrix} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \\ \mathbf{M} \end{pmatrix} \mathbb{E}(\boldsymbol{\varepsilon} | \mathbf{X}) \\ &= \begin{pmatrix} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \\ \mathbf{M} \end{pmatrix} \cdot \mathbf{0} = \mathbf{0}_{(K+n) \times 1} \end{aligned}$$

及

$$\begin{aligned}
 & \text{Var} \left(\begin{pmatrix} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \\ \mathbf{M} \end{pmatrix} \varepsilon \middle| \mathbf{X} \right) \\
 &= \begin{pmatrix} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \\ \mathbf{M} \end{pmatrix} \text{Var}(\varepsilon | \mathbf{X}) \begin{pmatrix} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}, \mathbf{M} \end{pmatrix} \\
 &= \begin{pmatrix} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \\ \mathbf{M} \end{pmatrix} I_n \sigma^2 \begin{pmatrix} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}, \mathbf{M} \end{pmatrix} = \begin{pmatrix} (\mathbf{X}^T \mathbf{X})^{-1} \sigma^2 & \mathbf{O}_{K \times n} \\ \mathbf{O}_{n \times K} & \mathbf{M} \sigma^2 \end{pmatrix},
 \end{aligned}$$

故有

$$\begin{pmatrix} \hat{\beta} - \beta \\ \hat{\varepsilon} \end{pmatrix} \middle| \mathbf{X} \sim N \left(\mathbf{0}_{(K+n) \times 1}, \begin{pmatrix} (\mathbf{X}^T \mathbf{X})^{-1} \sigma^2 & \mathbf{O}_{K \times n} \\ \mathbf{O}_{n \times K} & \mathbf{M} \sigma^2 \end{pmatrix} \right). \quad (3.56)$$

式 (3.56) 表明, 抽样误差 $\hat{\beta} - \beta$ 和残差向量 $\hat{\varepsilon}$ 的联合分布是一个多元正态分布, 且依据协方差矩阵可知 $\hat{\beta} - \beta$ 与 $\hat{\varepsilon}$ 是完全不相关的, 而对于服从正态分布的随机向量而言, 完全不相关与相互独立是等价的, 即式 (3.56) 表明了多元正态回归模型 (3.12.1) 下, 抽样误差与残差向量独立. 这还可以进一步的推出, OLS 估计量 $\hat{\beta}$ 与残差向量 $\hat{\varepsilon}$ 是相互独立的.

现在我们来研究方差估计量 s^2 的条件分布, 在一元正态回归模型 (2.7.1) 中我们说明了其经过变形后服从于卡方分布. 现在我们利用线性代数的知识来证明这一点. 对于 RRS 有 $\hat{\varepsilon}^T \hat{\varepsilon} = (\mathbf{M} \varepsilon)^T \mathbf{M} \varepsilon = \varepsilon^T \mathbf{M} \varepsilon$, 而由于投影矩阵 \mathbf{M} 是幂等的实对称方阵, 则依据式 (3.27) 得到

$$\mathbf{M} = \mathbf{H}^{-1} \begin{pmatrix} \mathbf{I}_{n-K} & \mathbf{O}_{(n-K) \times n} \\ \mathbf{O}_{K \times (n-K)} & \mathbf{I}_K \end{pmatrix} \mathbf{H},$$

其中 \mathbf{H} 是依赖于 \mathbf{X} 的正交方阵, 满足有 $\mathbf{H}^{-1} = \mathbf{H}^T$. 故 RSS 有

$$\hat{\varepsilon}^T \hat{\varepsilon} = \varepsilon^T \mathbf{H}^{-1} \begin{pmatrix} \mathbf{I}_{n-K} & \mathbf{O}_{(n-K) \times n} \\ \mathbf{O}_{K \times (n-K)} & \mathbf{I}_K \end{pmatrix} \mathbf{H} \varepsilon = (\mathbf{H} \varepsilon)^T \begin{pmatrix} \mathbf{I}_{n-K} & \mathbf{O}_{(n-K) \times n} \\ \mathbf{O}_{K \times (n-K)} & \mathbf{O}_K \end{pmatrix} \mathbf{H} \varepsilon$$

成立. 令 $\boldsymbol{\mu} = \mathbf{H} \varepsilon = (\mu_1, \dots, \mu_n)^T$, 即是对回归误差向量 ε 做正交变换, 这自然也是线性变换, 则 $\boldsymbol{\mu}$ 服从多元正态分布. 不难计算有

$$\mathbb{E}(\boldsymbol{\mu} | \mathbf{X}) = \mathbb{E}(\mathbf{H} \varepsilon | \mathbf{X}) = \mathbf{H} \mathbb{E}(\varepsilon | \mathbf{X}) = \mathbf{H} \mathbf{0} = \mathbf{0},$$

$$\text{Var}(\boldsymbol{\mu} | \mathbf{X}) = \text{Var}(\mathbf{H} \varepsilon | \mathbf{X}) = \mathbf{H} \text{Var}(\varepsilon | \mathbf{X}) \mathbf{H}^T = \mathbf{H} (I_n \sigma^2) \mathbf{H}^T = I_n \sigma^2,$$

故 $\boldsymbol{\mu} | \mathbf{X} \sim N(\mathbf{0}, I_n \sigma^2)$, 且 $\boldsymbol{\mu}$ 的分量满足有 $\mu_i | \mathbf{X} \sim N(0, \sigma^2)$.

我们将 $\boldsymbol{\mu}$ 进行分块处理, 设 $\boldsymbol{\mu} = (\boldsymbol{\mu}_1^T, \boldsymbol{\mu}_2^T)^T$, 则不难得知 $\boldsymbol{\mu}_1 | \mathbf{X} \sim N(\mathbf{0}, I_{n-K} \sigma^2)$. 再

计算 RSS 得

$$\begin{aligned}\hat{\varepsilon}^T \hat{\varepsilon} &= \boldsymbol{\mu}^T \begin{pmatrix} \mathbf{I}_{n-K} & \mathbf{O}_{(n-K) \times n} \\ \mathbf{O}_{K \times (n-K)} & \mathbf{O}_K \end{pmatrix} \boldsymbol{\mu} \\ &= (\boldsymbol{\mu}_1^T, \boldsymbol{\mu}_2^T) \begin{pmatrix} \mathbf{I}_{n-K} & \mathbf{O}_{(n-K) \times n} \\ \mathbf{O}_{K \times (n-K)} & \mathbf{O}_K \end{pmatrix} \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix} \\ &= (\boldsymbol{\mu}_1^T, \mathbf{0}_{1 \times K}) \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix} = \boldsymbol{\mu}_1^T \boldsymbol{\mu}_1,\end{aligned}$$

故可以推知方差估计量 s^2 满足有

$$\frac{(n-K)s^2}{\sigma^2} \Big| \mathbf{X} = \frac{\boldsymbol{\mu}_1^T \boldsymbol{\mu}_1}{\sigma^2} \Big| \mathbf{X} = \sum_{i=1}^{n-K} \left(\frac{\mu_i - 0}{\sigma} \right)^2 \Big| \mathbf{X} \sim \chi^2(n-K). \quad (3.57)$$

由于正态假设 (3.52) 和 (3.53) 保证了模型 (3.12.1) 是满足条件同方差式 (3.7), 我们在本章无需考虑异方差问题, 对于需要 σ^2 的情形便使用偏差校正估计量 s^2 去代替.

现在我们便可以对 OLS 估计量进行区间估计. 与定义 (2.7.1) 完全一致, 我们使用 t 统计量进行区间估计.

定义 3.12.1 (多元正态回归模型的 t 统计量). 对于模型 (3.12.1) 的 OLS 估计量 $\hat{\boldsymbol{\beta}}$, 其标准化后有

$$\frac{\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}}{\sqrt{(\mathbf{X}^T \mathbf{X})^{-1} \sigma^2}} \Big| \mathbf{X} \sim \mathbf{N}(\mathbf{0}, \mathbf{I}),$$

由于 σ^2 为未知的厌恶参数 (nuisance parameter), 使用 s^2 的偏差校正估计量 s^2 代替, 则对于 OLS 估计量的第 k 个参数估计量, 我们称

$$T = \frac{\hat{\beta}_k - \beta_k}{\text{SE}(\hat{\beta}_k)} = \frac{\hat{\beta}_k - \beta_k}{\sqrt{s^2 ((\mathbf{X}^T \mathbf{X})^{-1})_{kk}}} \quad (3.58)$$

为 t 统计量 (t -statistic).

显然, $\hat{\beta}_k$ 的 t 统计量 T 有

$$T = \frac{\beta_k - \hat{\beta}_k}{\sqrt{\sigma^2 ((\mathbf{X}^T \mathbf{X})^{-1})_{kk}}} / \frac{\sqrt{s^2 ((\mathbf{X}^T \mathbf{X})^{-1})_{kk}}}{\sqrt{\sigma^2 ((\mathbf{X}^T \mathbf{X})^{-1})_{kk}}} = \frac{\beta_k - \hat{\beta}_k}{\sqrt{\sigma^2 ((\mathbf{X}^T \mathbf{X})^{-1})_{kk}}} / \sqrt{\frac{(n-K)s^2/\sigma^2}{n-K}},$$

而

$$\frac{\beta_k - \hat{\beta}_k}{\sqrt{\sigma^2 ((\mathbf{X}^T \mathbf{X})^{-1})_{kk}}} \Big| \mathbf{X} \sim \mathbf{N}(0, 1), \quad \frac{(n-K)s^2}{\sigma^2} \Big| \mathbf{X} \sim \chi^2(n-K)$$

故 t 统计量 T 的条件分布服从自由度为 $n - K$ 的 t 分布.

在上一章第2.7节中, 我们已经详细讲解了以 t 统计量作为枢轴变量去构造 OLS 估计量的置信区间, 本笔记不再重复叙述. 在此, 我们列出如下定理.

定理 3.12.1 (OLS 估计量的置信区间). 在正态假设 (3.52) 和 (3.53) 下, 模型 (3.12.1) 总体参数 β 的第 k 个回归系数 β_k 的置信水平为 $1 - \alpha$ 的置信区间为

$$\hat{I} = \left[\hat{\beta}_k - c \cdot \text{SE}(\hat{\beta}_k), \hat{\beta}_k + c \cdot \text{SE}(\hat{\beta}_k) \right], \quad (3.59)$$

其中 $c = F^{-1}\left(1 - \frac{\alpha}{2}\right)$, F 是自由度为 $n - K$ 的 t 分布的累计密度函数, F^{-1} 为 F 的反函数. 常数 c 亦即是自由度为 $n - K$ 的 t 分布的上 $\alpha/2$ 分位数 $t_{\alpha/2}$.

对于回归方差 σ^2 , 其置信水平为 $1 - \alpha$ 的置信区间为

$$\hat{I} = \left[\frac{(n-2)s^2}{c_2}, \frac{(n-2)s^2}{c_1} \right], \quad (3.60)$$

其中 $c_1 = F^{-1}(\alpha/2)$, $c_2 = F^{-1}(1 - \alpha/2)$, F 是自由度为 $n - K$ 的卡方分布的累计密度函数, F^{-1} 为 F 的反函数. 常数 c_1 和 c_2 亦即是自由度为 $n - K$ 的卡方分布的上 $\alpha/2$ 分位数和上 $1 - \alpha/2$ 分位数.

特别地, 对于置信区间式 (3.59), 一条有用的经验公式为当 $n - K \geq 61$ 时, β_k 的置信水平为 95% 的置信区间近似为

$$\hat{I} = \left[\hat{\beta}_k - 2 \cdot \text{SE}(\hat{\beta}_k), \hat{\beta}_k + 2 \cdot \text{SE}(\hat{\beta}_k) \right]. \quad (3.61)$$

3.13 多元正态回归模型的假设检验

作为区间估计的置信区间同点估计一样, 其都是被视为由抽取的样本所决定的随机变量, 因此其具有的概率性质, 必须从重复抽样的前提下去理解——这即是大数定律赋予随机变量的现实意义. 然而, 在多数情况下, 计量经济学所要研究的问题是不能重复调查的, 相应地问题所需的样本仅有唯一一次的抽样结果. 就单独的一次样本而言, OLS 估计量的无偏性、置信区间的 $1 - \alpha$ 的置信水平似乎都没有意义——就我们为了确保估计结果的准确性而言.

但假设检验可以就单次样本的估计结果予以可信度的说明. 就模型 (3.12.1) 的回归结果而言, 其结果是否正确我们可以给出自己的判断, 这种判断可以源自于经济学的理论, 由此我们便对于回归结果提出了假设 (hypothesis). 通常而言, 对于某个总体参数 θ , 我们认为其具有一定的限制, 我们用集合 B_0 表示该限制, 则我们称

$$H_0 : \theta \in B_0 \quad (3.62)$$

为原假设 (null hypothesis), 相应地, 我们称原假设不成立的对应情况

$$H_1 : \theta \notin B_0 \quad (3.63)$$

为备择假设 (alternative hypothesis). 对于原假设 H_0 为真, 我们记为 $\theta \in H_0$, 反之则记为 $\theta \in H_1$.

注意, 原假设 H_0 与备择假设 H_1 的选取取决于我们对于所谓“限制”的判断. 例如, 我们要研究一个命题 $\theta \neq \theta_0$, 则原假设 H_0 与备择假设 H_1 的选取通常有如下两种抉择:

- (i) 如果我们要证明命题 $\theta \neq \theta_0$ 为真, 依据反证法的思想, 令其逆命题为原假设 $H_0 : \theta_0 = 0$, 而备择假设即是 $H_1 : \theta \neq \theta_0$;
- (ii) 依据实际问题考虑, 选取犯第一类错误 (见后文) 后危害更小的情形, 将其作为原假设 H_0 .

通常而言, 原假设 H_0 与备择假设 H_1 的选择一般都是遵从 (i) 的方法, 这种选择为假设检验的讨论提供了一个公认的范式.

假设检验 (hypothesis test) 即是通过样本数据构造 **检验统计量** (test statistic), 据此判断原假设 H_0 是否正确. 对于依据检验统计量的实现值来论断, 原假设 H_0 是否成立的方法, 我们称为 **检验** (test), 且常记为 ϕ . 例如, 我们针对回归系数 β_k 原假设 $H_0 : \beta_k > c$, 则以 OLS 估计量为检验统计量的

$$\phi : \hat{\beta}_k > C_0$$

便是原假设 H_0 的一个检验, 其意味当 $\hat{\beta}_k > C_0$ 时我们接受 (不拒绝) 原假设 $H_0 : \beta_k > c$. 我们将检验 ϕ 中使得接受 (不拒绝) 原假设 H_0 的检验统计量的集合, 称为检验 ϕ 的 **接受域** (acceptance region), 反之使得拒绝原假设 H_0 的集合称为检验 ϕ 的 **拒绝域** (rejection region). 接受域与拒绝域的边界通常称为 **临界值** (critical value). 例如上例检验 ϕ 中的接受域是集合 $\{\hat{\beta}_k \mid \hat{\beta}_k > C_0\}$, 拒绝域是 $\{\hat{\beta}_k \mid \hat{\beta}_k \leq C_0\}$, 而 C_0 即是临界值.

对于一个原假设 H_0 , 其检验 ϕ 是多种的, 而 ϕ 所使用的检验统计量是依赖于抽取的样本, 因而 检验 ϕ 实际上是具有随机性的. 为了避免源自抽样的误差的影响, 我们需要对于检验 ϕ 提出一个优劣判断准则.

定义 3.13.1 (功效函数). 设关于总体参数 θ 的原假设 H_0 , 其有基于样本构造检验统计量的检验 ϕ , 我们称

$$\beta_\phi(\theta) = \mathbb{P}(\text{检验 } \phi \text{ 下 } H_0 \text{ 被拒绝}) \quad (3.64)$$

为检验 ϕ 的 **功效函数** (power function).

不难知晓, 当 $\theta \in H_0$ 时功效函数 $\beta_\phi(\theta)$ 应该尽可能小; 当 $\theta \in H_1$ 时功效函数应该尽可能大 (体现了“功效 (power)”一词). 故我们可以以 $\beta_\phi(\theta)$ 作为检验 ϕ 的判断准则.

在假设检验中, 由于检验 ϕ 具有随机性, 对于真实情形 $\theta \in H_0$ 或 $\theta \in H_1$ 便可能出现判断错误. 假设检验存在有两类错误:

- 第一类错误 (type I error): H_0 正确但被拒绝 (弃真);
- 第二类错误 (type II error): H_0 错误但被接受 (取伪).

这两类错误可以用表 (3.3) 直观展示.

表 3.3: 假设检验的判断结果

	$\theta \in H_0$	$\theta \in H_1$
接受 (不拒绝)	正确	第二类错误
拒绝	第一类错误	正确

由于 $\theta \in H_0$ 和 $\theta \in H_1$ 两者只能出现一种情形, 不难计算, 对于检验 ϕ , 记其出现第一、二类错误的概率分布为 $\alpha_1(\theta|\phi)$, $\alpha_2(\theta|\phi)$, 则有

$$\alpha_1(\theta|\phi) = \begin{cases} \beta_\phi(\theta), & \theta \in H_0, \\ 0, & \theta \in H_1, \end{cases} \quad \alpha_2(\theta|\phi) = \begin{cases} 0, & \theta \in H_0, \\ 1 - \beta_\phi(\theta), & \theta \in H_1 \end{cases}$$

成立. 为了保障假设检验的准确性, 被广泛接受的一种方法即是控制犯第一类错误的概率 $\alpha_1(\theta|\phi)$ 在某一较小的水平下. 在前文的按照 (i) 去选择原假设 H_0 时, 即依据反证法的思想将要证明的命题的逆命题设定为原假设, 控制了 $\alpha_1(\theta|\phi)$ 较小, 则意味着 $1 - \beta_\phi(\theta)$ 有较大的值, 这就是说检验 ϕ 下接受 (不拒绝) 原假设的概率较大. 这也意味着, 如果检验 ϕ 能够拒绝原假设 H_0 , 其有检验统计量确定的概率是较小的, 故需要有充分的证据才能拒绝原假设, 也由此避免了检验 ϕ 的随机误差.

定义 3.13.2 (显著性水平). 设 ϕ 是原假设 H_0 的一个检验, 对于常数 $0 \leq \alpha \leq 1$, 若 $\forall \theta \in H_0$ 总有 $\beta_\phi(\theta) \leq \alpha$, 我们称检验 ϕ 是原假设 H_0 的一个显著性水平 (significance level) 为 α 的检验.

显然, 一个显著性水平为 α 的检验 ϕ 表保障了犯第一类错误的概率小于等于常数 α , 通常而言, 尝试 α 的取值为 10%、5% 和 1%.

经过了上述的理论分析, 现在我们就 OLS 估计量进行假设检验的实际操作. 对于单个回归系数 β_k , 最为常见的假设检验问题便是

$$H_0 : \beta_k = c \quad v.s. \quad H_1 : \beta_k \neq c,$$

由于定理 (3.6.3) 指出 OLS 估计量 $\hat{\beta}_k$ 是参数 β_k 的最优线性无偏估计量, 我们便可使用 $\hat{\beta}_k$ 作为 β_k 的近似值. 一种检验 ϕ 即是认为, 如果 OLS 估计量 $\hat{\beta}_k$ 与原假设的 c 偏

离较大, 则应该拒绝 H_0 , 这种检验方法被称为 **Wald 检验** (Wald test). 由于我们还没有进行检验, 则不能拒绝原假设 $H_0: \beta_k = c$ 为假, 则我们便可以依据抽样误差的绝对值 $|\hat{\beta}_k - \beta_k| = |\beta_k - c|$ 来设立检验 ϕ , 但为了具有可比性, 还需要考虑系数自身大小影响, 故我们考虑对抽样误差做标准化处理——也就是说, 我们使用 $\hat{\beta}_k$ 的 t 统计量

$$T = \frac{\hat{\beta}_k - \beta_k}{\text{SE}(\hat{\beta}_k)} \stackrel{H_0}{=} \frac{\beta_k - c}{\text{SE}(\hat{\beta}_k)}$$

作为检验 ϕ 的检验统计量, 且检验 ϕ 为

$$\phi: \text{若 } |T| > C \text{ 则拒绝 } H_0, \text{ 其中 } C > 0.$$

现在我们来研究检验 ϕ 可以具有的置信水平. 考虑其功效函数 $\beta_\phi(\beta_k)$ 为

$$\begin{aligned} \beta_\phi(\beta_k) &= \mathbb{P}(\text{检验 } \phi \text{ 下 } H_0 \text{ 被拒绝}) = \mathbb{P}(|T| > C) \\ &= \mathbb{P}(T < -C) + \mathbb{P}(T > C) = F(-C) + 1 - F(C) \\ &= 1 - F(C) + 1 - F(C) = 2(1 - F(C)), \end{aligned}$$

故检验 ϕ 犯第一类错误的概率为 $\alpha_1(\theta|\phi) = \begin{cases} 2(1 - F(C)), & \beta_k = c, \\ 0, & \beta_k \neq c. \end{cases}$ 这样, 为了保障检验 ϕ 具有 α 的显著性水平, 只需有

$$\alpha_1(\theta|\phi) \leq \alpha \Rightarrow 2(1 - F(C)) \leq \alpha \Rightarrow F(C) \geq 1 - \frac{1}{2}\alpha \Rightarrow C = F^{-1}\left(1 - \frac{1}{2}\alpha\right),$$

这里的 F 是自由度为 $n - K$ 的 t 分布的累计密度函数, F^{-1} 为 F 的函数. 上面的推导过程用到了 t 分布为对称分布以及 F 为增函数这两条性质. $C = F^{-1}(1 - \frac{1}{2}\alpha)$ 便是 t 分布的上 $\alpha/2$ 分位数.

定理 3.13.1 (回归模型单个系数的 t 检验). 对于模型 (3.12.1) 的回归系数 β_k , 就假设

$$H_0: \beta_k = c \quad v.s. \quad H_1: \beta_k \neq c,$$

而言, 以 $\hat{\beta}_k$ 的 t 统计量

$$T = \frac{\hat{\beta}_k - \beta_k}{\text{SE}(\hat{\beta}_k)} \stackrel{H_0}{=} \frac{\beta_k - c}{\text{SE}(\hat{\beta}_k)}$$

作为检验统计量, 则有显著性水平为 α 的检验

$$\phi: \text{若 } |T| > F^{-1}\left(1 - \frac{1}{2}\alpha\right) \text{ 则拒绝 } H_0,$$

其中 F 是自由度为 $n - K$ 的 t 分布的累计密度函数, F^{-1} 为 F 的函数. 我们称检验 ϕ 为回归系数 β_k 的 **t 检验** (t test).

特别地, 当定理 (3.13.1) 中假设常数 c 取 0 时, 则此时的 t 检验称为回归系数 β_k 的显著性检验 (significance test). 如果检验不能通过, 则意味着我们接受 (不能拒绝) $H_0: \beta_k = 0$, 这意味着在 α 的显著性水平下, β_k 对应的解释变量 X_k 对于被解释变量 Y 没有解释力 (因为系数为 0).

除了计算检验统计量 T 外, 我们在概率论与数理统计中学过另一种判断假设是否成立的方法即是计算 p 值.

定义 3.13.3 (p 值). 对于一个假设检验, 设检验统计量为 S , 若检验 ϕ 为

$$\phi: \text{若 } S > C \text{ 则拒绝 } H_0, C \in \mathbb{R},$$

则检验 ϕ 的 p 值 (p -value) 为 $p = 1 - G(S)$, 其中 G 是 H_0 下检验统计量的累计密度函数.

p 值是原假设 H_0 为真时依据检验统计量的累计密度函数所计算的, 进一步计算有

$$p = 1 - G(S) = 1 - \mathbb{P}(s \leq S_{\text{实现值}} | H_0) = \mathbb{P}(s \geq S_{\text{实现值}} | H_0)$$

即 p 值是原假设成立的情形下, 检验统计量 S 的取值不小于由样本计算得到的实现值的概率. 显然, 对于检验 ϕ 是 S 的实现值越大则越有可能拒绝 H_0 , 则计算的 p 值越小, 则出现大于 S 的实现值的情况越少. 若检验 ϕ 具有 α 的显著性水平, 则有

$$\beta_\phi(\beta_k) = \mathbb{P}(\text{检验 } \phi \text{ 下 } H_0 \text{ 被拒绝}) = \mathbb{P}(S > C) \leq \alpha$$

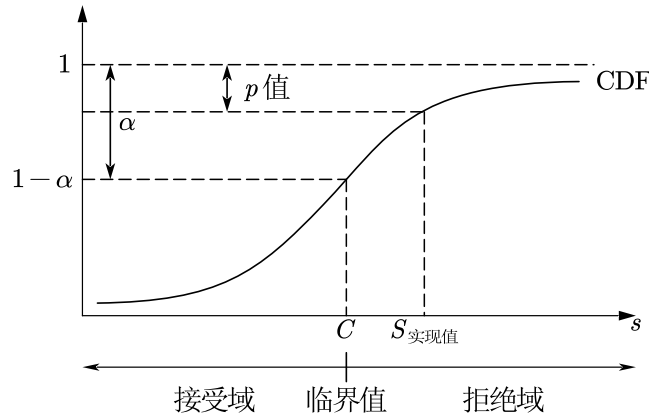
成立, 这实际上也就是说检验统计量 S 的取值出现大于 C 的概率要小于 α . 若我们计算出 p 值满足有 $p \leq \alpha$, 即检验统计量 S 出现了由样本决定的实现值的概率小于等于 α , 那么必然有 S 的现实值大于 C , 这论断可以通过如图 (3.4) 所示的累计密度函数函数形象说明.

一种常见的误解是认为 p 值是原假设成立的概率, 这是错误的. p 值仅仅是一个由样本计算的检验统计量的现实值所决定的随机变量, p 值越小, 则表明样本观测值越不支持原假设 H_0 , 且 p 值与显著性水平 α 是相关联的.

对于定理 (3.13.1), 则其 p 值为

$$\begin{aligned} p &= \mathbb{P}(s \geq S_{\text{实现值}} | H_0) = \mathbb{P}(|t| \geq |T|) = \mathbb{P}(t \geq |T|) + \mathbb{P}(t \leq -|T|) \\ &= 1 - F(|T|) + F(-|T|) = 1 - F(|T|) + 1 - F(|T|) \\ &= 2(1 - F(|T|)), \end{aligned}$$

其中 F 是自由度服从 $n - K$ 的 t 分布的累计密度函数.

图 3.4: 检验 ϕ 的 p 值

既然我们能够对某个单个系数 β_k 进行假设检验, 自然对于模型 (3.12.1) 我们需要进行 K 次 t 检验才能得到

$$\begin{cases} H_{01} : \beta_1 = c_1 & v.s. & H_{11} : \beta_1 \neq c_1, \\ \dots & & \\ H_{0k} : \beta_k = c_k & v.s. & H_{1k} : \beta_k \neq c_k, \\ \dots & & \\ H_{0K} : \beta_K = c_K & v.s. & H_{1K} : \beta_K \neq c_K \end{cases}$$

这 K 个原假设 H_{01}, \dots, H_{0K} 的检验结果. 我们是否能够据此得到合并这 K 个原假设, 依据 K 次 t 检验来给出

$$H_0 : \boldsymbol{\beta} = \mathbf{c} = (c_1, \dots, c_K)^T$$

的结果? 答案是否定的, 因为单个系数的 t 检验都是独立进行的, 没有考虑参数 $\boldsymbol{\beta}$ 的联合分布. 如果把 K 次 t 检验的结果作为 $H_0 : \boldsymbol{\beta} = \mathbf{c}$ 的结果, 还会积累这 K 次 t 检验的第一类错误. 倘若每次 t 检验都在 α 的显著性水平上拒绝了 t 检验的原假设, 以此结果认为能够拒绝 $H_0 : \boldsymbol{\beta} = \mathbf{c}$, 那么累计的犯第一类错误概率便有上界 α^K . 故此, 我们需要对于 $H_0 : \boldsymbol{\beta} = \mathbf{c}$ 需要新的检验方法—— F 检验.

假设模型 (3.12.1) 的回归线性 $\boldsymbol{\beta}$ 有假设

$$H_0 : \mathbf{R}\boldsymbol{\beta} = \mathbf{r} \quad v.s. \quad H_1 : \mathbf{R}\boldsymbol{\beta} \neq \mathbf{r}$$

需要检验, 其中矩阵 $\mathbf{R} \in \mathbb{R}^{r \times K}$ 满足有行满秩 $\text{rank} \mathbf{R} = r$, 而列向量 $\mathbf{r} \in \mathbb{R}^{r \times 1}$, 则原假设

即为

$$\begin{aligned} \mathbf{R}\boldsymbol{\beta} = \mathbf{r} &\Leftrightarrow \begin{pmatrix} R_{11} & R_{12} & \cdots & R_{1K} \\ R_{21} & R_{22} & \cdots & R_{2K} \\ \vdots & \vdots & & \vdots \\ R_{r1} & R_{r2} & \cdots & R_{rK} \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_K \end{pmatrix} = \begin{pmatrix} r_1 \\ r_2 \\ \vdots \\ r_r \end{pmatrix} \\ &\Leftrightarrow \begin{cases} R_{11}\beta_1 + R_{12}\beta_2 + \cdots + R_{1K}\beta_K = r_1, \\ R_{21}\beta_1 + R_{22}\beta_2 + \cdots + R_{2K}\beta_K = r_2, \\ \cdots \\ R_{r1}\beta_1 + R_{r2}\beta_2 + \cdots + R_{rK}\beta_K = r_r, \end{cases} \end{aligned}$$

这便是 r 个关于参数 β_1, \dots, β_K 的线性方程, 其构成对于参数 $\boldsymbol{\beta}$ 的一个线性约束条件. 由于要检验的假设等价于

$$H_0: \mathbf{R}\boldsymbol{\beta} - \mathbf{r} = \mathbf{0} \quad v.s. \quad H_0: \mathbf{R}\boldsymbol{\beta} - \mathbf{r} \neq \mathbf{0},$$

我们同样使用 OLS 估计量 $\hat{\boldsymbol{\beta}}$ 去代替总体参数 $\boldsymbol{\beta}$, 依据 Wald 检验的思想, 如果原假设成立则向量模长 $\|\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r}\|$ 应该有较小的取值, 反之若 $\|\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r}\|$ 的取值较大, 则应该拒绝原假设 H_0 .

令 $\mathbf{u} = \mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r} \in \mathbb{R}^{r \times 1}$, 则 \mathbf{u} 为 OLS 估计量的一个线性变换, 其条件分布依然是多元正态分布. 在 $H_0: \mathbf{R}\boldsymbol{\beta} - \mathbf{r} = \mathbf{0}$ 成立的情况下, 计算 \mathbf{u} 的条件期望、方差为

$$\begin{aligned} \mathbb{E}(\mathbf{u} | \mathbf{X}) &= \mathbb{E}(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r} | \mathbf{X}) = \mathbf{R}\mathbb{E}(\hat{\boldsymbol{\beta}} | \mathbf{X}) - \mathbf{r} = \mathbf{R}\boldsymbol{\beta} - \mathbf{r} \stackrel{H_0}{=} \mathbf{0}, \\ \text{Var}(\mathbf{u} | \mathbf{X}) &= \text{Var}(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r} | \mathbf{X}) = \mathbf{R}\text{Var}(\hat{\boldsymbol{\beta}} | \mathbf{X})\mathbf{R}^T = \mathbf{R}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{R}^T\sigma^2, \end{aligned}$$

故 r 维列向量 $\mathbf{u} | \mathbf{X} \sim N(\mathbf{0}, \mathbf{R}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{R}^T\sigma^2)$.

为了度量 $\|\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r}\| = \|\mathbf{u}\|$ 偏离原假设的程度, 我们还需要对 \mathbf{u} 进行标准化, 将其转化为标准的多元正态分布, 显然这是可以通过线性变换来实现的. 由于假设式 (3.6) 保证了 $(\mathbf{X}^T\mathbf{X})^{-1}$ 的存在性, 则可以推知协方差矩阵 $\text{Var}(\mathbf{u} | \mathbf{X}) = \mathbf{R}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{R}^T\sigma^2$ 是正定的. 依据定理 (3.7.4), 则可将协方差矩阵写作 $\text{Var}(\mathbf{u} | \mathbf{X}) = \mathbf{B}^2$, 其中 r 阶方阵 \mathbf{B} 为 $\text{Var}(\mathbf{u} | \mathbf{X})$ 的平方根. 不难注意到

$$\begin{aligned} \mathbf{u}^T (\text{Var}(\mathbf{u} | \mathbf{X}))^{-1} \mathbf{u} &= \mathbf{u}^T (\mathbf{B}^2)^{-1} \mathbf{u} = \mathbf{u}^T (\mathbf{B}\mathbf{B})^{-1} \mathbf{u} \\ &= \mathbf{u}^T \mathbf{B}^{-T} \mathbf{B}^{-1} \mathbf{u} = (\mathbf{B}^{-1}\mathbf{u})^T (\mathbf{B}^{-1}\mathbf{u}), \end{aligned}$$

在原假设下, 对于 $\mathbf{B}^{-1}\mathbf{u}$ 有

$$\begin{aligned} \mathbb{E}(\mathbf{B}^{-1}\mathbf{u} | \mathbf{X}) &= \mathbf{B}^{-1}\mathbb{E}(\mathbf{u} | \mathbf{X}) = \mathbf{0}, \\ \text{Var}(\mathbf{B}^{-1}\mathbf{u} | \mathbf{X}) &= \text{Var}(\mathbf{B}^{-1}\mathbf{u} | \mathbf{X}) = \mathbf{B}^{-1}\text{Var}(\mathbf{u} | \mathbf{X})\mathbf{B}^{-1} = \mathbf{B}^{-1}(\mathbf{B}^2)\mathbf{B}^{-1} = \mathbf{I}_r, \end{aligned}$$

故有 $B^{-1}\mathbf{u}|\mathbf{X} \sim N(\mathbf{0}, \mathbf{I}_r)$ 成立. 我们以 r 维列向量 $B^{-1}\mathbf{u}$ 的模长 $\|B^{-1}\mathbf{u}\|$ 来度量原假设的偏离情况, 也通过列向量 $B^{-1}\mathbf{u}$ 的内积来衡量. 注意到 $B^{-1}\mathbf{u} = (u_1, \dots, u_r)^T$ 的每个分量都是独立同分布的服从标准正态分布, 那么其内积便满足有

$$\mathbf{u}^T (\text{Var}(\mathbf{u}|\mathbf{X}))^{-1} \mathbf{u} = (B^{-1}\mathbf{u})^T (B^{-1}\mathbf{u}) = \sum_{i=1}^r u_i^2,$$

故有 $\mathbf{u}^T (\text{Var}(\mathbf{u}|\mathbf{X}))^{-1} \mathbf{u}|\mathbf{X} \sim \chi^2(r)$, 也即有

$$\left(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r} \right)^T \left(\mathbf{R}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{R}^T \sigma^2 \right)^{-1} \left(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r} \right) \Big| \mathbf{X} \sim \chi^2(r). \quad (3.65)$$

对于形如式 (3.65) 的具有二次型的统计量

$$W = (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)^T \hat{\mathbf{V}}_{\hat{\boldsymbol{\theta}}} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0), \quad (3.66)$$

我们将式 (3.66) 称为 **Wald 统计量** (Wald statistic). 式 (3.65) 即是回归系数在线性约束条件下的 Wald 统计量. 对于 Wald 统计量, 在学习了渐进理论后我们将对其进行详细讨论.

我们现在继续研究 $H_0: \mathbf{R}\boldsymbol{\beta} - \mathbf{r} = \mathbf{0}$ 的检验问题, 需要指出式 (3.65) 较式 (3.66) 仍是有区别的, 因为式 (3.65) 中有未知的冗余参数 σ^2 , 但真正的 Wald 统计量是使用了协方差矩阵的估计量 $\hat{\mathbf{V}}_{\hat{\boldsymbol{\theta}}}$. 由于正态回归模型 (3.12.1) 是自然满足同方差式 (3.7), 于是我们使用方差估计量 s^2 来替代 σ^2 , 则有

$$\begin{aligned} W &= (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})^T \left(\mathbf{R}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{R}^T s^2 \right)^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r}) \\ &= \frac{(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})^T \left(\mathbf{R}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{R}^T \right)^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})}{s^2} \\ &= \frac{(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})^T \left(\mathbf{R}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{R}^T \sigma^2 \right)^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})}{s^2 / \sigma^2} \\ &= \frac{(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})^T \left(\mathbf{R}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{R}^T \sigma^2 \right)^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})}{\frac{(n-K)s^2}{\sigma^2} / (n-K)}, \end{aligned}$$

记

$$W^0 = \frac{(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})^T \left(\mathbf{R}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{R}^T \sigma^2 \right)^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})}{\frac{(n-K)s^2}{\sigma^2} / (n-K)}, \quad (3.67)$$

我们称式 (3.67) 为同方差的 **Wald 统计量** (homoskedastic Wald statistic). 注意到 s^2 是残差向量 $\hat{\boldsymbol{\varepsilon}}$ 的函数, 式 (3.65) 是 OLS 估计量 $\hat{\boldsymbol{\beta}}$ 的函数, 而 $\hat{\boldsymbol{\varepsilon}}$ 和 $\hat{\boldsymbol{\beta}}$ 是独立的, 故其两者

函数也是独立的. 依据式 (3.57), 不难得知, 令 $F = W^0/r$, 则有

$$F = \frac{W^0}{r} = \frac{\chi^2(r)/r}{\chi^2(r)/(n-K)} \sim F(r, n-K), \quad (3.68)$$

于是我们便得到了一个不依赖冗余参数 σ^2 的检验估计量 F , 其条件分布服从自由度为 $(r, n-K)$ 的 F 分布.

由于 $F = W^0/r, r$ 为线性约束条件的方程个数, 其是固定的, 则我们依旧遵循有 F 越大, 相应地 W^0 越大, 则应该拒绝原假设 $H_0: \mathbf{R}\boldsymbol{\beta} - \mathbf{r} = \mathbf{0}$ 的检验方法. 这从 F 分布的角度而言也是如此. 设检验

$$\phi: \text{若 } F > C \text{ 则拒绝 } H_0,$$

则在 α 的显著性水平下功效函数满足有

$$\beta_\phi(\boldsymbol{\beta}) = \mathbb{P}(\text{检验 } \phi \text{ 下 } H_0 \text{ 被拒绝}) = \mathbb{P}(F > C) = 1 - F(C) \leq \alpha,$$

即可推知

$$F(C) \geq 1 - \alpha \Rightarrow C \geq F^{-1}(1 - \alpha),$$

故有

$$\phi: \text{若 } F > F^{-1}(1 - \alpha) \text{ 则拒绝 } H_0,$$

其中第一个 F 是检验统计量式 (3.68) 的实现值, 第二个 F 为自由度服从 $(r, n-K)$ 的 F 分布的累计密度函数, F^{-1} 为 F 的反函数, $F^{-1}(1 - \alpha)$ 也便是 F 分布的上 α 分位数.

定理 3.13.2 (回归系数线性约束的 F 检验). 对于模型 (3.12.1) 的回归系数 $\boldsymbol{\beta}$, 就假设

$$H_0: \mathbf{R}\boldsymbol{\beta} = \mathbf{r} \quad v.s. \quad H_1: \mathbf{R}\boldsymbol{\beta} \neq \mathbf{r}$$

而言, 以同方差的 Wald 统计量式 (3.67) 得到

$$F = \frac{W^0}{r}$$

作为检验统计量, 则有显著性水平为 α 的检验

$$\phi: \text{若 } F > F^{-1}(1 - \alpha) \text{ 则拒绝 } H_0,$$

其中第一个 F 是检验统计量式 (3.68) 的实现值, 第二个 $F(\cdot)$ 为自由度服从 $(r, n-K)$ 的 F 分布的累计密度函数. 我们称检验 ϕ 为回归系数 $\boldsymbol{\beta}$ 线性约束的 **F 检验** (F test).

定理 (3.13.2) 的 p 值为

$$p = \mathbb{P}(f \geq F_{\text{实现值}} | H_0) = \mathbb{P}(f \geq F_{\text{实现值}}) = 1 - F(F_{\text{实现值}}),$$

其中 $F(\cdot)$ 为自由度服从 $(r, n-K)$ 的 F 分布的累计密度函数.

对于原假设 $H_0: \mathbf{R}\boldsymbol{\beta} = \mathbf{r}$, 定理 (3.13.2) 中统计检验量 F 是依赖于样本所确定的 OLS 估计量 $\hat{\boldsymbol{\beta}}$, 而 $\hat{\boldsymbol{\beta}}$ 则是定义于式 (3.17) 的最优化问题——这是一个无约束条件的优化问题. 我们在推导统计检验量 F 中是用到了原假设 $H_0: \mathbf{R}\boldsymbol{\beta} = \mathbf{r}$ 成立的条件, 但是 OLS 估计量的推导没有用到, 故我们还可以从受约束的最优化问题

$$\begin{aligned} \hat{\boldsymbol{\beta}}_{\text{CLS}} &= \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^{K \times 1}} \hat{S}(\boldsymbol{\beta}), \\ \text{s.t. } \mathbf{R}\boldsymbol{\beta} &= \mathbf{r}, \end{aligned}$$

来推导总体参数 $\boldsymbol{\beta}$ 的估计量 $\hat{\boldsymbol{\beta}}_{\text{CLS}}$, 此时我们得到的估计量称为**受约束的最小二乘估计量** (constrained least squares estimator). 我们可以以 $\hat{\boldsymbol{\beta}}_{\text{CLS}}$ 来研究 $H_0: \mathbf{R}\boldsymbol{\beta} = \mathbf{r}$ 的检验问题, 但这我们留在本笔记的第 () 节. 我们仅在此介绍依据 $\hat{\boldsymbol{\beta}}_{\text{CLS}}$ 得到的 F 统计量的似然比检验, 且可以证明, 似然比检验与定理 (3.13.2) 是等价的.

设 **F 比率** (F-ratio) 为

$$F = \frac{\left((\boldsymbol{\varepsilon}^*)^T \boldsymbol{\varepsilon}^* - \boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon} \right) / r}{\boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon} / (n-K)} = \frac{(\text{RSS}^* - \text{RSS}) / r}{\text{RSS} / (n-K)}, \quad (3.69)$$

其中 $*$ 表示 CLS 的回归结果, 即 $\boldsymbol{\varepsilon}^*$ 是 CLS 估计量 $\hat{\boldsymbol{\beta}}_{\text{CLS}}$ 计算的残差向量, RSS^* 是 CLS 估计量 $\hat{\boldsymbol{\beta}}_{\text{CLS}}$ 计算的残差平方和, 而 $\boldsymbol{\varepsilon}$ 和 RSS 则是 OLS 估计量的残差向量和残差平方和. 可以证明, 式 (3.69) 与式 (3.68) 是相等的. 我们以 F 比率对假设 $H_0: \mathbf{R}\boldsymbol{\beta} = \mathbf{r}$ 的检验, 称为 F 比率的**似然比检验** (likelihood ratio test), 其与 (3.13.2) 的相同的.

回到我们引出 F 检验的问题. 特别地, 当针对总体参数 $\boldsymbol{\beta}$ 的线性约束假设为 $H_0: \boldsymbol{\beta} = \mathbf{0}$ 时, 即 \mathbf{R} 为 K 阶单位阵, \mathbf{r} 为 K 维零向量, 此时我们称 F 检验或似然比检验为模型 (3.12.1) 或方程 (3.2) 的**显著性检验** (significance test). 这时, 通过 F 检验拒绝原假设 $H_0: \boldsymbol{\beta} = \mathbf{0}$, 即是意味着在给定的显著性水平 α 下, 至少有一个回归系数 β_k 不为 0.

方程 (3.2) 的显著性检验可以与拟合优度 R^2 联系起来. 当模型 (3.12.1) 含有常数项时, 在 $H_0: \boldsymbol{\beta}_{K \times 1} = (1, 0, \dots, 0)^T$ 下注意到式 (3.69) 有

$$\begin{aligned} F &= \frac{\left((\boldsymbol{\varepsilon}^*)^T \boldsymbol{\varepsilon}^* - \boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon} \right) / r}{\boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon} / (n-K)} = \frac{\left(\frac{(\boldsymbol{\varepsilon}^*)^T \boldsymbol{\varepsilon}^*}{(\mathbf{Y}-\mathbf{1}\bar{Y})^T (\mathbf{Y}-\mathbf{1}\bar{Y})} - \frac{\boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon}}{(\mathbf{Y}-\mathbf{1}\bar{Y})^T (\mathbf{Y}-\mathbf{1}\bar{Y})} \right) (K-1)^{-1}}{\frac{\boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon}}{(\mathbf{Y}-\mathbf{1}\bar{Y})^T (\mathbf{Y}-\mathbf{1}\bar{Y})} \cdot (n-K)^{-1}} \\ &= \frac{\left[1 - \frac{\boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon}}{(\mathbf{Y}-\mathbf{1}\bar{Y})^T (\mathbf{Y}-\mathbf{1}\bar{Y})} - \left(1 - \frac{(\boldsymbol{\varepsilon}^*)^T \boldsymbol{\varepsilon}^*}{(\mathbf{Y}-\mathbf{1}\bar{Y})^T (\mathbf{Y}-\mathbf{1}\bar{Y})} \right) \right] (K-1)^{-1}}{\left[1 - \left(1 - \frac{\boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon}}{(\mathbf{Y}-\mathbf{1}\bar{Y})^T (\mathbf{Y}-\mathbf{1}\bar{Y})} \right) \right] \cdot (n-K)^{-1}} \\ &= \frac{(R^2 - (R^*)^2) (K-1)^{-1}}{(1-R^2) \cdot (n-K)^{-1}} \end{aligned}$$

其中 $(R^*)^2$ 为 CLS 估计量下的拟合优度. 而 CLS 估计量在 $H_0: \beta_{K \times 1} = (1, 0, \dots, 0)^T$ 下回归的个体方程即为 $Y_i = \beta_1 + \varepsilon_i$, 即是仅含有常数项的 OLS 估计量, 这时的 OLS 估计量的拟合值即是样本均值 \bar{Y} , 依据拟合优度的定义式 (3.35) 便可推知 $(R^*)^2 = 0$. 故在 $H_0: \beta_{K \times 1} = (1, 0, \dots, 0)^T$ 下, F 比率即为

$$F = \frac{R^2 / (K - 1)}{(1 - R^2) / (n - K)}, \quad (3.70)$$

也就是说, 对于含有常数项的模型 (3.12.1), 其方程的显著性检验的检验统计量 F 比率是关于拟合优度 R^2 的函数. 不难得知, 在 n 和 K 不变的情形下, F 是关于 R^2 的增函数.

对于式 (3.70), 其分子除以的自由度 $n - K$ 不难让我们联想至调整的拟合优度 \bar{R}^2 , 则利用式 (3.38) 对式 (3.70) 代换得

$$\begin{aligned} F &= \frac{R^2 / (K - 1)}{(1 - R^2) / (n - K)} = \left(\frac{1}{1 - R^2} - 1 \right) \cdot \frac{n - K}{K - 1} \\ &= \left(\frac{n - 1}{n - K} (1 - \bar{R}^2)^{-1} - 1 \right) \cdot \frac{n - K}{K - 1} = \frac{n - 1}{K - 1} (1 - \bar{R}^2)^{-1} - \frac{n - K}{K - 1}, \end{aligned}$$

也即有

$$\bar{R}^2 = 1 - \frac{n - 1}{(K - 1)F + (n - K)}. \quad (3.71)$$

这样, 针对在 $H_0: \beta_{K \times 1} = (1, 0, \dots, 0)^T$ 的 F 检验便通过式 (3.71) 将调整的 R^2 与 F 比率联系在一起, 同样地, 在 n 与 K 不变的情形下, \bar{R}^2 与 F 是同向变化的. 式 (3.71) 的意义在于, 由于如果检验统计量取值较大 F 比率是显著的, 则关于 F 递增的 \bar{R}^2 也是“显著的”. 这对于拟合优度 R^2 和 F 也是一样的. 在此意义而言, 对于含有常数项的模型 (3.12.1), 其方程的显著性检验亦即是针对各类拟合优度的显著性检验.

上述的分析似乎表明, 各类拟合优度越大则 F 统计量越大, 因而模型越显著——这是错误的, 因为没有考虑到样本容量 n 与解释变量的个数 K . 实际上, 只要样本容量 n 足够大, 方程的显著性总能够取得非常显著的结果. Bruce(cite) 在他的计量经济学教材中谈到:

This was a popular statistic in the early days of econometric reporting when sample sizes were very small and researchers wanted to know if there was “any explanatory power” to their regression. This is rarely an issue today as sample sizes are typically sufficiently large that this F statistic is nearly always highly significant. While there are special cases where this F statistic is useful these cases are not typical. As a general rule there is no reason to report this F statistic.

也就是说, 针对单个方程显著性的 F 检验源自于早期计量经济学发展所面临的样本有限性, 为了确保他们的回归具有解释能力, 需要汇报 F 比率, 这一方法便被传承了下来. 我们可以分析, 尽管 R^2 非常小, 但足够大的样本容量便能使得方程的显著性检验在给定 α 的显著性水平下通过. 依据定理 (3.13.2), 如果方程的显著性检验在 α 的显著性水平下拒绝 $H_0: \beta_{K \times 1} = (1, 0, \dots, 0)^T$, 则由式 (3.70) 可以推出样本容量 n 满足有

$$\left(\frac{1}{1 - R^2} - 1 \right) \cdot \frac{n - K}{K - 1} \geq F^{-1}(1 - \alpha) \Rightarrow n - K \geq \frac{1 - R^2}{R^2} \cdot (K - 1) \cdot F^{-1}(1 - \alpha).$$

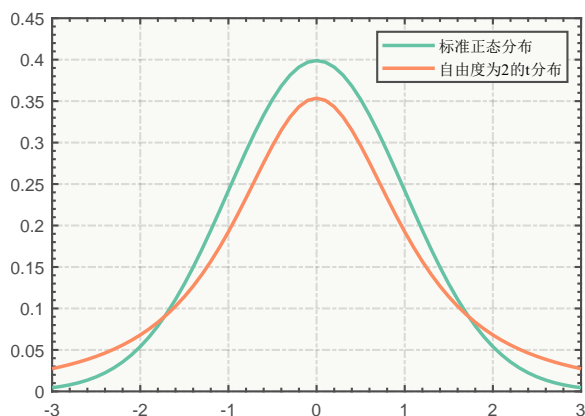
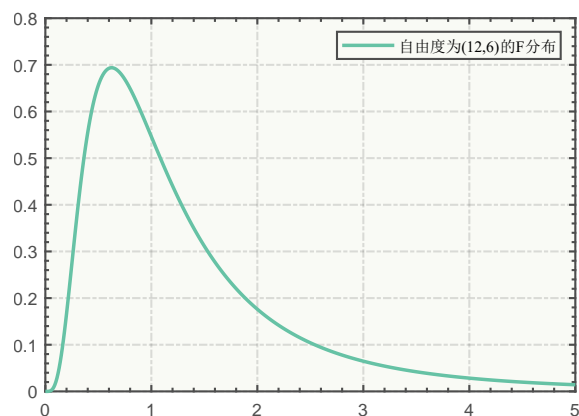
我们通过 Monte Carlo 模拟来说明这一结构. 我们给定显著性水平 $\alpha = 1\%$ 以及一个非常小的 $R^2 = 0.01$, 且解释变量共有 $K = 8$ 个, 我们编写如下的 MATLAB 代码来求解所需的样本容量 n .

```

1 clear
2 clc
3 alpha=0.01;%显著性水平
4 K=8;%变量数
5 R2=0.01;%假定的拟合优度
6 n=K+1;%样本数初始值
7 while 1
8     n=n+1;%样本数循环+1
9     F=R2/(1-R2)*(n-K)/(K-1);%F统计量
10    cv=finv(1-alpha,K-1,n-K);%计算临界值
11    if F>=cv %显著时退出
12        disp(['所需样本数:',num2str(n)])
13        break
14    end
15 end

```

代码输出的结果即是需要 $n = 1844$ 大小的样本容量. 这个数值似乎看上去很大, 但对于当代的计量经济学而言其实谈不上大样本. 中国现行的地级市共有 293 个, 如果以之作为我们研究问题的总体, 考虑 2010-2019 年的跨度为 10 年的市级年度数据, 即使是数据存在缺失, 我们也是可以轻松获得一个包含有 200 个中国地级市的跨度 10 年的面板数据——这便有了 $n = 2000$ 的样本容量. 如果是研究上市公司等问题, 则可获得的样本数据将更多. 在就机器学习的角度而言, $n = 1844$ 的样本容量可以说是有些苛刻了. 故对于当代的实证研究而言, 汇报方程的显著性检验其实不具有充分的理由. 而从这一角度来看, 即使是 $R^2 = 0.01$ 足够大的样本容量也使得回归结果也是显著的, R^2 、 \bar{R}^2 和 R_{uc^2} 的意义也仅在于度量拟合上. 需要强调, 这样不是说我们应该一味地追求“大数据”的回归分析, 好的抽样比样本总量更为重要.

图 3.5: 正态分布与 t 分布的概率密度函数图 3.6: F 分布的概率密度函数

3.14 多元线性回归模型的预测

3.15 虚拟变量的 OLS 估计量

3.16 代码

```

1  %% 正交投影与OLS估计量
2  clear
3  clc
4  Y=[-0.5;1.25;2];
5  X=[-1 1
6      1 1
7      0 1];
8  Ze=[0,0];%设定箭头位置
9  % a=0:1;%线性变换系数
10 a=[0 0 1 1;
11     0 1 1 0];%线性变换矩阵
12 Span=X*a;%张成的空间
13 Yhat=X*(X\Y);%拟合值_投影
14 e=Y-Yhat;%残差
15 selfGrootDefault(2)%绘图美化函数,没有请注释
16 figure(1)
17 quiver3(0,0,0,Y(1),Y(2),Y(3),'off','LineWidth',2.5)%绘制Y
18 hold on
19 quiver3(Ze,Ze,Ze,X(1,:),X(2,:),X(3,:),'off','LineWidth',2.5)%绘制X
20 fill3(Span(1,:),Span(2,:),Span(3,:),[0 0.4470 0.7410], ...

```

```

21     "LineStyle","none")%X张成的平面
22 quiver3(0,0,0,Yhat(1,:),Yhat(2,:),Yhat(3,:), ...
23     "off","LineWidth",2.5)%绘制投影
24 quiver3(Yhat(1,:),Yhat(2,:),Yhat(3,:),e(1,:),e(2,:),e(3,:), ...
25     "off","-.", "LineWidth",2.5)%绘制残差
26 % quiver3(-2,0,0,1,0,0,4,'k','filled','LineWidth',2);%绘制x轴
27 % quiver3(0,0,0,0,1,0,2.5,'k','filled','LineWidth',2);%绘制y轴
28 % quiver3(0,0,0,0,0,1,2.5,'k','filled','LineWidth',2);%绘制z轴
29 % xlim([-1.5,1.5])
30 % ylim([0,2.5])
31 % zlim([0 3])
32 xlabel('$x$', 'FontSize',14, 'Interpreter','latex')
33 ylabel('$y$', 'FontSize',14, 'Interpreter','latex')
34 zlabel('$z$', 'FontSize',14, 'Interpreter','latex')
35 legend('Y','X','X张成的平面','投影向量','残差向量','Location','north')
36
37 figure(2)%绘制拟合直线和散点图
38 hold on
39 xline(0,'-.','LineWidth',2)
40 yline(0,'-.','LineWidth',2)
41 scatter(X(:,1),Y,"filled","LineWidth",2)
42 y=[-5:1:5;ones(size(-5:1:5))]'*(X\Y);
43 plot(-5:1:5,y,'LineWidth',2)
44 xlim([-2.5 2.5])
45 ylim([-1 2.5])
46 xlabel('$x$', 'FontSize',14, 'Interpreter','latex')
47 ylabel('$y$', 'FontSize',14, 'Interpreter','latex')
48 legend(' ',' ','样本点','样本回归直线','Location','northwest')

```

第四章 概率渐进理论

本章的内容旨在介绍和补充一般概统课程中未能详细论述的**概率极限**知识. 概率极限对于发展计量经济学的大样本理论极为重要, 然而国内鲜有一本专门介绍这些知识的教材, 或者需要基于测度论去建立这些内容——笔者即是提供了一个较为初等的视角, 所需的数学工具不会超出之前的章节.

4.1 依概率收敛

在微积分中, 我们已经学习过关于数列极限的知识. 对于一个数列 $\{x_n\}$, 我们使用 $\varepsilon - \delta$ 语言定义了

$$\lim_{n \rightarrow \infty} x_n = a \quad \Leftrightarrow \quad \forall \varepsilon > 0, \exists N \in \mathbb{N}, \forall n > N \text{ s.t. } |x_n - a| < \varepsilon,$$

则称数列 $\{x_n\}$ 收敛于常数 $a \in \mathbb{R}$. 可以说, 极限是整个微积分的基石, 微积分的所有分析都离不开极限的影子. 而本节我们将介绍概率极限的知识, 其便是数理统计的**渐进理论** (asymptotic theorem) 的基石.

首先考虑一个序列 $\{X_n\}$, 其中的每一个 X_n 都是随机变量, 则我们称之为**随机序列** (random sequence). 我们现研究随机序列 $\{X_n\}$ 的性质, 自然, 我们可以按照数列极限的定义, 假设有 $\{X_n\}$ 收敛于常数 a , 但这样而言对于随机序列而言则是一个过强的假设, 且不能与随机变量的性质联系在一起. 故我们就随机变量提出依概率收敛.

定义 4.1.1 (随机变量依概率收敛). 对于一个随机序列 $\{X_n\}$, 对于任意的 $\varepsilon, \delta > 0$, 若存在常数 $a \in \mathbb{R}$, 使得

$$\lim_{n \rightarrow \infty} \mathbb{P}(|X_n - a| > \delta) = 0 \quad \Leftrightarrow \quad \forall \varepsilon, \delta > 0, \exists N \in \mathbb{N}, \forall n > N \text{ s.t. } \mathbb{P}(|X_n - a| > \delta) < \varepsilon,$$

则称随机序列 $\{X_n\}$ **依概率收敛** (converge in probability) 于常数 a , 同时记为

$$X_n \xrightarrow{P} a (n \rightarrow \infty) \quad \text{或} \quad \text{plim}_{n \rightarrow \infty} X_n = a.$$

常数 a 也称为随机序列 $\{X_n\}$ 的**概率极限** (probability limit).

对于事件 $(|X_n - a| > \delta)$, 其对立事件为 $(|X_n - a| \leq \delta)$, 则随机序列 $\{X_n\}$ 依概率收敛也等价于

$$\lim_{n \rightarrow \infty} \mathbb{P}(|X_n - a| \leq \delta) = 1$$

成立. 这样, 对比概率极限与数列极限的定义, 不能发现, 随机序列 $\{X_n\}$ 依概率收敛其实是概率 $\mathbb{P}(|X_n - a| \leq \delta)$ 的极限. 我们以此来认识何为概率极限.

我们知道, 如果数列 $\{x_n\}$ 收敛于 a , 则对于充分大的 n , 数列 $\{x_n\}$ 将全部落在点 a 的任意小的邻域之中——即无限接近于 a , 但不意味着等于 a . 这对于极限 $\lim_{n \rightarrow \infty} \mathbb{P}(|X_n - a| \leq \delta) = 1, \forall \delta > 0$ 而言, 即是意味着对于充分大的 n , 随机序列 $\{X_n\}$ 落在点 a 的任意小的邻域中的概率, 将可以无限的接近于 1, 但是不意味为 1. 这也就是说, 随机序列 $\{X_n\}$ 收敛于常数 a , 则无论是多么大的 n , 总有随机变量 X_n 会出现在常数 a 的给定邻域之外.

我们下面研究几个依概率收敛的具体例子.

例 4.1.1. 设随机序列 $\{X_n\}$ 满足有两点分布 $\mathbb{P}(X_n = a_n) = p_n, \mathbb{P}(X_n = 0) = 1 - p_n$, 其中 $0 \leq p_n \leq 1$. 若 p_n 收敛于 0 或 a_n 收敛于 0, 则随机序列 $\{X_n\}$ 依概率收敛于 0.

证明. 当 p_n 收敛于 0 时, 则有 $n \rightarrow \infty$ 时 $\mathbb{P}(X_n = 0) = 1 - p_n \rightarrow 1, \forall \delta > 0$, 对于事件 $|X_n - 0| = |X_n| \leq \delta$, 其包含有事件 $X_n = 0$, 故有

$$\mathbb{P}(X_n = 0) \leq \mathbb{P}(|X_n| \leq \delta) \leq 1,$$

夹逼即得 $\lim_{n \rightarrow \infty} \mathbb{P}(|X_n| \leq \delta) = 1$.

当 a_n 收敛于 0 时, $\forall \delta > 0$ 则存在充分大的 N 使得 $\forall n > N, |a_n| \leq \delta$, 而 X_n 的取值仅有 0 和 a_n , 故 $\forall n > N, |X_n| \leq \delta$, 此时

$$\begin{aligned} \mathbb{P}(|X_n| \leq \delta) &= \mathbb{P}(|a_n| \leq \delta) + \mathbb{P}(X_n = 0) \\ &= \mathbb{P}(X_n = a_n) + \mathbb{P}(X_n = 0) = p_n + 1 - p_n = 1, \end{aligned}$$

即知 $\lim_{n \rightarrow \infty} \mathbb{P}(|X_n| \leq \delta) = 1$. □

例 4.1.2. 设随机变量 X , 定义随机序列 $\{X_n\}$ 为 $X_n = b_n X$, 其中 $\{b_n\}$ 是实数列. 若数列 $\{b_n\}$ 收敛于 0, 则随机序列 $\{X_n\}$ 依概率收敛于 0.

证明. 对于任意的 $\forall \delta > 0$, 则 $\mathbb{P}(|X_n| \leq \delta) = \mathbb{P}(-\delta \leq b_n X \leq \delta)$. 随机变量 X 可能是无界的, 但无论 X 实现值 x 如何大, 由于 b_n 收敛于 0, 我们总能找到充分大的 n 使得 $-\delta \leq b_n X \leq \delta$ 成立 (序列的生成共有同一个随机变量 X), 故有 $\lim_{n \rightarrow \infty} \mathbb{P}(|X_n| \leq \delta) = 1$. □

随机序列 $\{X_n\}$ 能依概率收敛为常数, 自然我们可以考虑其能否“收敛”至随机变量. 就此而言, 一种常用的定义是, 对于随机序列 $\{X_n\}$ 以及随机变量 X , 如果 $n \rightarrow \infty$ 时有 $X_n - X \xrightarrow{P} 0$, 则称随机序列 $\{X_n\}$ 依概率收敛于随机变量 X .

依概率收敛亦可拓展至多元的情形. 设随机向量 $\mathbf{X}_n = (X_{n1}, \dots, X_{nK})^T \in \mathbb{R}^{K \times 1}$, 其每个分量都是随机变量, 则我们亦可称

$$\mathbf{X}_1 = (X_{11}, \dots, X_{1K})^T, \mathbf{X}_2 = (X_{21}, \dots, X_{2K})^T, \dots, \mathbf{X}_n = (X_{n1}, \dots, X_{nK})^T$$

构成了随机序列 $\{\mathbf{X}_n\}$.

定义 4.1.2 (随机向量依概率收敛). 对于一个 $K \times 1$ 维随机序列 $\{\mathbf{X}_n\}$, 对于任意的 $\delta > 0$, 若存在非随机的向量 $\mathbf{a} \in \mathbb{R}^{K \times 1}$, 使得

$$\lim_{n \rightarrow \infty} \mathbb{P}(\|\mathbf{X}_n - \mathbf{a}\| > \delta) = 0$$

则称随机序列 $\{\mathbf{X}_n\}$ **依概率收敛** (*converge in probability*) 于 非随机的向量 \mathbf{a} , 同时记

$$\mathbf{X}_n \xrightarrow{P} \mathbf{a} (n \rightarrow \infty) \quad \text{或} \quad \text{plim}_{n \rightarrow \infty} \mathbf{X}_n = \mathbf{a}.$$

对于随机向量 $\mathbf{X} = (X_1, \dots, X_K)^T \in \mathbb{R}^{K \times 1}$, 我们称

$$\mathbf{X}_n - \mathbf{X} \xrightarrow{P} \mathbf{0} (n \rightarrow \infty)$$

为随机序列 $\{\mathbf{X}_n\}$ 依概率收敛于 随机向量 \mathbf{x} .

可以证明, 随机序列 $\{\mathbf{X}_n\}$ 依概率收敛于非随机的 \mathbf{a} , 等价于 $\mathbf{X}_n = (X_{n1}, \dots, X_{nK})^T$ 的每个分量收敛于 $\mathbf{a} = (a_1, \dots, a_K)^T$ 的分量.

依概率收敛是通过概率的极限来定义的, 而另一种思路则是考虑极限的概率, 由此我们得到了几乎处处收敛.

定义 4.1.3 (几乎处处收敛). 对于一个 $K \times 1$ 维随机序列 $\{\mathbf{X}_n\}$, 若存在非随机的向量 $\mathbf{a} \in \mathbb{R}^{K \times 1}$ 使得

$$\mathbb{P}\left(\lim_{n \rightarrow \infty} \mathbf{X}_n = \mathbf{a}\right) = 1$$

则称随机序列 $\{\mathbf{X}_n\}$ **几乎处处收敛** (*converge almost surely*) 于 非随机的向量 \mathbf{a} , 同时记 $\mathbf{X}_n \xrightarrow{a.s.} \mathbf{a} (n \rightarrow \infty)$.

对于随机向量 $\mathbf{X} = (X_1, \dots, X_K)^T \in \mathbb{R}^{K \times 1}$, 我们称 $\mathbf{X}_n - \mathbf{X} \xrightarrow{a.s.} \mathbf{0} (n \rightarrow \infty)$ 为随机序列 $\{\mathbf{X}_n\}$ 几乎处处收敛于 随机向量 \mathbf{x} .

特别地, 当 $K = 1$ 时, 即有随机序列 $\{X_n\}$ 几乎处处收敛于常数 $a (X_n \xrightarrow{a.s.} a (n \rightarrow \infty))$ 和随机序列 $\{X_n\}$ 几乎处处收敛于随机变量 $X (X_n - X \xrightarrow{a.s.} 0 (n \rightarrow \infty))$.

在概率论中, **几乎处处** (almost surely) 指的是概率为 1 的事件——这样的事件并不是一定发生的. 例如考虑几何概型问题, 记在单位圆中选得一给定点为事件 A , 则事件

A 是概率为 0 的事件, 但可以发生; 事件的 A 的对立事件概率为 1, 但可以不发生. 几乎处处收敛便是一种将数列收敛的概率化的方法, 其比定义 (4.1.2) 的收敛性要更强, 即 $\mathbf{X}_n \xrightarrow{a.s.} \mathbf{a} (n \rightarrow \infty) \Rightarrow \mathbf{X}_n \xrightarrow{\mathbb{P}} \mathbf{a} (n \rightarrow \infty)$, 这是因为

$$\mathbb{P}\left(\lim_{n \rightarrow \infty} \mathbf{X}_n = \mathbf{a}\right) = 1 \Leftrightarrow \forall \delta > 0, \mathbb{P}(\|\mathbf{X}_n - \mathbf{a}\| \leq \delta) = 1,$$

而常数 $\mathbb{P}(\|\mathbf{X}_n - \mathbf{a}\| \leq \delta) = 1$ 的极限仍是常数.

回顾本节的例 (4.1.1) 的两点分布的随机序列, 由于 X_n 的取值为 0 或 a_n , 则可以得知随机序列 $\{X_n\}$ 满足数列收敛于 0, 故 $\{X_n\}$ 是几乎处处收敛于 0, 但 p_n 收敛于 0 时则不能保证几乎处处收敛. 几乎处处收敛的详细讨论需要涉及测度论的知识, 在此不过多介绍.

4.2 大数定律与连续映射定理

本节我们介绍大数定律与连续映射定理, 将两者结合我们便可以推知样本估计量所具有的相合性.

大数定律 (law of large numbers, 简称 LLN), 可以说是概率论与数理统计最早研究的问题之一, 其的发现伴随了数学家对于何为概率的认识, 并直接导致了概率论的频率学派的诞生. 在介绍大数定律前, 本笔记首先阐述其证明工具——Markov 不等式与 Chebyshev 不等式.

定理 4.2.1 (Markov 不等式). 对于非负的随机变量 X , 设其 $\mathbb{E}X$ 存在, 则对于 $\forall \varepsilon > 0$ 有

$$\mathbb{P}(X \geq \varepsilon) \leq \frac{\mathbb{E}X}{\varepsilon}, \quad (4.1)$$

我们称式 (4.1) 为 **Markov 不等式** (Markov's inequality).

证明. 假设非负的随机变量 X 具有概率密度函数 f , 注意到

$$\begin{aligned} \mathbb{P}(X \geq \varepsilon) &= \int_{\varepsilon}^{+\infty} f(x) dx \leq \int_0^{+\infty} f(x) dx \\ &\leq \int_0^{+\infty} \frac{x}{\varepsilon} f(x) dx = \frac{1}{\varepsilon} \int_0^{+\infty} x f(x) dx = \frac{\mathbb{E}X}{\varepsilon}, \end{aligned}$$

这里第二个 “ \leq ” 用到了事件 $X \geq \varepsilon$.

再设非负的随机变量 X 是离散的, 则同理可得

$$\begin{aligned} \mathbb{P}(X \geq \varepsilon) &= \sum_{x_i \geq \varepsilon} \mathbb{P}(X = x_i) \leq \sum_{i=1}^{\infty} \mathbb{P}(X = x_i) \\ &\leq \sum_{i=1}^{\infty} \frac{x_i}{\varepsilon} \mathbb{P}(X = x_i) = \frac{1}{\varepsilon} \sum_{i=1}^{\infty} x_i \mathbb{P}(X = x_i) = \frac{\mathbb{E}X}{\varepsilon}. \end{aligned}$$

这样, 无论非负随机变量是连续的还是离散的, 我们都证明了式 (4.1) 是成立的. \square

对于非负的随机变量 X , Markov 不等式为事件 $X \leq \varepsilon$ 的发生概率提供了一个与随机变量 X 的期望有关的上界.

定理 4.2.2 (Chebyshev 不等式). 对于随机变量 X , 设 $\mathbb{E}X$ 和 $\text{Var}(X)$ 存在, 则对于 $\forall \varepsilon > 0$ 有

$$\mathbb{P}(|X - \mathbb{E}X| \geq \varepsilon) \leq \frac{\text{Var}(X)}{\varepsilon^2} \quad (4.2)$$

我们称式 (4.2) 为 **Chebyshev 不等式** (Chebyshev's inequality).

证明. 依据定理 (4.2.1), 对于非负的随机变量 $|X - \mathbb{E}X|$ 则有

$$\mathbb{P}(|X - \mathbb{E}X| \geq \varepsilon) = \mathbb{P}((X - \mathbb{E}X)^2 \geq \varepsilon^2) \leq \frac{\mathbb{E}(X - \mathbb{E}X)^2}{\varepsilon^2} = \frac{\text{Var}(X)}{\varepsilon^2}.$$

\square

由定理 (4.2.2) 的证明不难看出, Chebyshev 不等式 (4.2) 即是 Markov 不等式 (4.1) 的一个特例. Chebyshev 不等式的意义在于, 其对于随机变量 X 出现偏离于 $\mathbb{E}X$ 的情况给出了一个概率上界, 这个上界与随机变量的方差有关. Chebyshev 不等式亦可以用于理解为何方差刻画了随机变量的离散程度.

我们下面通过计算来详细说明 Monte Carlo 不等式.

(切比雪夫不等式的意义; 蒙特卡洛模拟阐述)

例 4.2.1. 对于随机序列 $\{X_n\}$, 设 $\text{Var}(X_n) = \sigma_n^2$ 是有限的. 如果数列 $\{\sigma_n^2\}$ 收敛于 0, 则随机序列 $\{X_n\}$ 依概率收敛于 0.

证明. 依据 Chebyshev 不等式 (4.2) 则有

$$0 \leq \mathbb{P}(|X_n - \mathbb{E}X_n| \geq \varepsilon) \leq \frac{\sigma_n^2}{\varepsilon^2},$$

由于 σ_n^2 收敛于 0, 故依据夹逼定理有 $\lim_{n \rightarrow \infty} \mathbb{P}(|X_n - \mathbb{E}X_n| \geq \varepsilon) = 0$, 这样我们即知道随机序列 $\{X_n\}$ 是依概率收敛的, 但收敛于常数或是随机变量则依赖于 $\mathbb{E}X_n$. 在一些计量经济学教材中, 随机序列的方差的序列 $\{\sigma_n^2\}$ 收敛于 0 被称为**依均方收敛** (converge in mean square). \square

有了 Chebyshev 不等式 (4.2) 和 Markov 不等式 (4.1), 我们现在便可以刻画随机变量的大数定律.

定理 4.2.3 (弱大数定律, weak law of large number). 设 $X_i, i = 1 \cdots, n$ 是独立同分布

的随机变量, 且共有的总体 X 的满足 $\mathbb{E}|X| < \infty$, 则有

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{P} \mathbb{E}X, \quad (n \rightarrow \infty). \quad (4.3)$$

证明. 假设共有的总体具有有限的方差 $\text{Var}(X)$, 则依据 Chebyshev 不等式 (4.2) 可得

$$0 \leq \mathbb{P}(|\bar{X}_n - \mathbb{E}X| \geq \varepsilon) \leq \frac{\text{Var}(\bar{X}_n)}{\varepsilon^2} = \frac{\text{Var}(X)}{n\varepsilon^2},$$

当 $n \rightarrow \infty$ 时, 夹逼定理即得 $\bar{X}_n \xrightarrow{P} \mathbb{E}X$.

上面的证明为了使用 Chebyshev 不等式, 我们假定了总体具有有限的方差 $\text{Var}(X)$, 但对于定理 (4.2.3) 并不需要此条件. 由于 $\text{plim}_{n \rightarrow \infty} \bar{X}_n = \mathbb{E}X \Leftrightarrow \text{plim}_{n \rightarrow \infty} (\bar{X}_n - \mathbb{E}X) = 0$, 则令 $Y_i = X_i - \mathbb{E}X$, 定理 (4.2.3) 便等价于

$$\bar{Y}_n = \frac{1}{n} \sum_{i=1}^n Y_i \xrightarrow{P} 0, \quad (n \rightarrow \infty),$$

不需要方差假定, 我们现在证明这是成立的. 依据 Markov 不等式 (4.1), 即得

$$\mathbb{P}(|\bar{Y}_n| \geq \varepsilon) \leq \frac{\mathbb{E}|\bar{Y}_n|}{\varepsilon},$$

可以证明, 当 $n \rightarrow \infty$ 时, $\mathbb{E}|\bar{Y}_n|$ 是收敛于 0 的^[1], 这样依据夹逼定理即证当 $n \rightarrow \infty$ 时 $\bar{Y}_n \xrightarrow{P} 0$. \square

这样, 我们便证明了定理 (4.2.3), 其意义在于, 对于独立同分布的随机样本, 其样本均值将依概率收敛于总体均值. 一个非常直接且广泛的应用便是考虑两点分布, 设事件 A 发生的概率为 p , 记事件 A 发生时随机变量 $X = 1$, 不发生则 $X = 0$, 于是 X 有分布律 $\mathbb{P}(X = 1) = p, \mathbb{P}(X = 0) = 1 - p$. 考虑 n 次独立重复试验, 每次实验的结果用独立同分布的 $X_i, i = 1, \dots, n$ 表示, 则定理 (4.2.3) 意味着

$$\bar{X}_n = \frac{\text{事件}A\text{发生的次数}}{n} \xrightarrow{P} p, \quad (n \rightarrow \infty),$$

^[1]一种方法是依据 Bahr-Esseen 不等式: 对于独立同分布的随机变量 $X_i, i = 1, \dots, n$, 若 $\mathbb{E}X_i = 0$, 则对于任意的 $0 < r \leq 2$ 有

$$\mathbb{E} \left| \sum_{i=1}^n X_i \right|^r \leq 2 \sum_{i=1}^n \mathbb{E}|X_i|^r$$

成立. 对于 \bar{Y}_n 即有

$$\mathbb{E}|\bar{Y}_n| = \mathbb{E} \left| \sum_{i=1}^n \frac{Y_i}{n} \right| \leq 2 \sum_{i=1}^n \mathbb{E} \left| \frac{Y_i}{n} \right| = \frac{2}{n} \sum_{i=1}^n \mathbb{E}|Y_i| = \frac{2}{n} \sum_{i=1}^n \mathbb{E}|Y|,$$

由于 $\mathbb{E}|X| < \infty$, 则 $\mathbb{E}|Y| < \infty$, 故可以得知 $\mathbb{E}|\bar{Y}_n|$ 是收敛于 0. 对此之外, 还有多种方法可以证明 $\mathbb{E}|\bar{Y}_n|$ 的收敛性.

换句话说, n 次独立重复试验中事件 A 发生的频率将随 n 的增大依概率收敛于概率 p , 则即是概率论的频率学派的基石.

弱大数定理——加以“弱 (weak)”字修饰, 则是因为还有强大数定律.

定理 4.2.4 (强大数定律, strong law of large number). 设 $X_i, i = 1 \cdots, n$ 是独立同分布的随机变量, 且共有的总体 X 的满足 $\mathbb{E}|X| < \infty$, 则有

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{a.s.} \mathbb{E}X, \quad (n \rightarrow \infty). \quad (4.4)$$

定理 (4.2.4) 指出样本均值 \bar{X}_n 是几乎处处收敛于总体期望 $\mathbb{E}X$, 这个结论是强于定理 (4.2.3) 的依概率收敛. 自然, 强大数定律可推出弱大数定律. 对于定理 (4.2.4), 其证明需要使用测度论, 这超出本笔记所涉及的范围.

大数定理自然可以拓展至高维的情形. 设独立同分布的 k 维随机列向量 $\mathbf{X}_i = (X_{i1}, \cdots, X_{iK})^T, i = 1, \cdots, n$, 且共有的总体 $\mathbf{X} = (X_1, \cdots, X_K)^T$ 有期望向量 $\mathbb{E}\mathbf{X} = (\mathbb{E}X_1, \cdots, \mathbb{E}X_K)^T$ 和协方差矩阵

$$\begin{aligned} \text{Var}(\mathbf{X}) &= \mathbb{E} \left((\mathbf{X} - \mathbb{E}\mathbf{X})(\mathbf{X} - \mathbb{E}\mathbf{X})^T \right) \\ &= \begin{pmatrix} \text{Var}(X_1) & \text{Cov}(X_1, X_2) & \cdots & \text{Cov}(X_1, X_K) \\ \text{Cov}(X_2, X_1) & \text{Var}(X_2) & \cdots & \text{Cov}(X_2, X_K) \\ \vdots & \vdots & & \vdots \\ \text{Cov}(X_K, X_1) & \text{Cov}(X_K, X_2) & \cdots & \text{Var}(X_K) \end{pmatrix}. \end{aligned}$$

定理 4.2.5 (大数定律, law of large number). 设有 $k \times 1$ 维随机向量 $\mathbf{X}_i = (X_{i1}, \cdots, X_{iK})^T, i = 1, \cdots, n$ 是独立同分布的随机变量, 且共有的总体 \mathbf{X} 的满足 $\mathbb{E}\|\mathbf{X}\| < \infty$, 则有

$$\bar{\mathbf{X}}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \xrightarrow{P} \mathbb{E}\mathbf{X}, \quad (n \rightarrow \infty) \quad (4.5)$$

和

$$\bar{\mathbf{X}}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \xrightarrow{a.s.} \mathbb{E}\mathbf{X}, \quad (n \rightarrow \infty), \quad (4.6)$$

其中式 (4.5) 是弱大数定律 (weak law of large number), 式 (4.6) 是强大数定律 (strong law of large number).

定理 (4.2.5) 的成立条件有三, 随机向量 \mathbf{X}_i 互相独立、同分布、总体期望 $\mathbb{E}\|\mathbf{X}\|$ 有限缺一不可. 任意条件的缺失都不能保证定理 (4.2.5) 成立. 当然, 大数定理的成立条件也可以进一步的调整, 我们将在第 () 节中介绍.

依据定理 (4.2.5), 我们不难得到如下的结果.

定理 4.2.6 (弱大数定律的应用). 利用式 (4.5), 则有

1. 如果 $\mathbb{E} \|\mathbf{X}\| < \infty$, 则有 $\frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \xrightarrow{P} \mathbb{E} \mathbf{X}, n \rightarrow \infty$;
2. 如果 $\mathbb{E} \|\mathbf{X}\|^m < \infty$, 则有 $\frac{1}{n} \sum_{i=1}^n \|\mathbf{X}_i\|^m \xrightarrow{P} \mathbb{E} \|\mathbf{X}\|^m, n \rightarrow \infty$;
3. 如果 $\mathbb{E} (\mathbf{X}^T \mathbf{X}) < \infty$, 则有 $\frac{1}{n} \sum_{i=1}^n \mathbf{X}_i^T \mathbf{X}_i \xrightarrow{P} \mathbb{E} (\mathbf{X}^T \mathbf{X}), n \rightarrow \infty$;
4. 如果 $\mathbb{E} \|e^{\mathbf{X}}\| < \infty$, 则有 $\frac{1}{n} \sum_{i=1}^n e^{\mathbf{X}_i} \xrightarrow{P} \mathbb{E} e^{\mathbf{X}}, n \rightarrow \infty$;
5. 如果 $\mathbb{E} \|\mathbf{X}\| < \infty$, 则有 $\frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \xrightarrow{P} \mathbb{E} \mathbf{X}, n \rightarrow \infty$;
6. 如果 $\mathbb{E} \|\ln \mathbf{X}\| < \infty$, 则有 $\frac{1}{n} \sum_{i=1}^n \ln \mathbf{X}_i \xrightarrow{P} \mathbb{E} \ln \mathbf{X}, n \rightarrow \infty$;

上述的函数均是作用于分量. 一般地, 设映射 $\mathbf{h} : \mathbb{R}^{K \times 1} \rightarrow \mathbb{R}^{m \times 1}$, 如果 $\mathbb{E} \|\mathbf{h}(\mathbf{X})\| < \infty$, 则有

$$\frac{1}{n} \sum_{i=1}^n \mathbf{h}(\mathbf{X}_i) \xrightarrow{P} \mathbb{E} \mathbf{h}(\mathbf{X}), n \rightarrow \infty. \quad (4.7)$$

定理 (4.2.5) 只是给出了针对随机向量的收敛结果, 一个自然的问题是如果 $\mathbb{E} \|\mathbf{X}\mathbf{X}^T\| < \infty$, 那么是否有 $\frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^T \xrightarrow{P} \mathbb{E} \|\mathbf{X}\mathbf{X}^T\|, n \rightarrow \infty$ 成立? 熟知矩阵微积分的读者知道, 对于矩阵我们并不能直接的从多元微积分迁移过去, 还需要考虑不少细节问题. 在此, 答案是肯定的. $\frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^T$ 是一个 K 阶方阵, 定理 (4.2.5) 可以保证其收敛性——我们将其的证明留在的节.

我们通过 Monte Carlo 模拟来说明定理 (4.2.6).

大数定律的结论还可以进一步被加强, 现在我们介绍概率极限的连续映射定理.

定理 4.2.7 (连续映射定理). 对于随机向量的序列 $\{\mathbf{X}_n\}$, 设 $\mathbf{X}_n \in \mathbb{R}^{K \times 1}$, 映射 $\mathbf{h} : \mathbb{R}^{K \times 1} \rightarrow \mathbb{R}^{m \times 1}$ 是一个连续映射, 则当 $n \rightarrow \infty$ 时,

1. 若 $\mathbf{X} \xrightarrow{P} \mathbf{X}$, 则有 $\mathbf{h}(\mathbf{X}_n) \xrightarrow{P} \mathbf{h}(\mathbf{X})$;
2. 若 $\mathbf{X} \xrightarrow{a.s.} \mathbf{X}$, 则有 $\mathbf{h}(\mathbf{X}_n) \xrightarrow{a.s.} \mathbf{h}(\mathbf{X})$.

定理 (4.2.7) 实际上对应了一元微积分中连续函数所具有的极限性质, 只不过这里是对应的概率的收敛. 相似地, 我们对于连续映射 \mathbf{h} 的要求可以进一步缩小, 我们只需收敛的随机变量 $\mathbf{X} \in \mathbb{R}^{K \times 1}$ 是概率有界的, 即 $\exists C \subset \mathbb{R}^{K \times 1}$ 使得 $\mathbb{P}(\mathbf{X} \in C) = 1$, 而映射 \mathbf{h} 在 C 中连续, 这就可以确保定理 (4.2.7) 成立. 定理 (4.2.7) 的直接应用如下.

定理 4.2.8. 设 $K \times 1$ 维随机向量的序列 $\{\mathbf{X}_n\}$ 依概率收敛于 \mathbf{X} , $\mathbf{a} \in \mathbb{R}^{K \times 1}$ 为非随机的向量, 则有

1. 当 $n \rightarrow \infty$, 则 $\mathbf{X}_n + \mathbf{a} \xrightarrow{P} \mathbf{X} + \mathbf{a}$;
2. 当 $n \rightarrow \infty$, 则 $\mathbf{a}^T \mathbf{X}_n \xrightarrow{P} \mathbf{a}^T \mathbf{X}$;
3. 当 $n \rightarrow \infty$, 则 $\mathbf{X}_n^T \mathbf{X}_n \xrightarrow{P} \mathbf{X}^T \mathbf{X}$;
4. 若 $\|\mathbf{X}_n\| \neq 0$ 且 $\|\mathbf{X}\| \neq 0$, 当 $n \rightarrow \infty$, 则 $\frac{\|\mathbf{a}\|}{\|\mathbf{X}_n\|} \xrightarrow{P} \frac{\|\mathbf{a}\|}{\|\mathbf{X}\|}$.

结合大数定律与连续映射定理, 我们可以得到更为便捷形式的大数定律.

定理 4.2.9 (弱大数定律与连续映射定理的结合). 设 $\mathbf{X}_i, i = 1, \dots, n$ 是独立同分布的且有共有的总体 $\mathbf{X} \in \mathbb{R}^{K \times 1}$, 记 $\bar{\mathbf{X}} = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i$, 则有

1. 如果 $\mathbb{E} \|\mathbf{X}\| < \infty$, 则有 $\bar{\mathbf{X}}^2 \xrightarrow{P} (\mathbb{E} \mathbf{X})^2, n \rightarrow \infty$;
2. 如果 $\mathbb{E} \|\mathbf{X}\| < \infty$, 则有 $\bar{\mathbf{X}}^m \xrightarrow{P} (\mathbb{E} \mathbf{X})^m, n \rightarrow \infty$;
3. 如果 $\mathbb{E} \|\mathbf{X}\| < \infty$, 则有 $e^{\bar{\mathbf{X}}} \xrightarrow{P} e^{\mathbb{E} \mathbf{X}}, n \rightarrow \infty$;
4. 如果 $\mathbb{E} \|\mathbf{X}\| < \infty$, 则有 $\ln \bar{\mathbf{X}} \xrightarrow{P} \ln \mathbb{E} \mathbf{X}, n \rightarrow \infty$;
5. 如果 $\mathbb{E} (\mathbf{X}^T \mathbf{X}) < \infty$, 则有 $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i^2 - \bar{\mathbf{X}}^2 \xrightarrow{P} \sigma^2 = \mathbb{E} (\mathbf{X}^T \mathbf{X}) - (\mathbb{E} \mathbf{X})^2, n \rightarrow \infty$;
6. 如果 $\mathbb{E} (\mathbf{X}^T \mathbf{X}) < \infty$, 则有 $\hat{\sigma} = \sqrt{\hat{\sigma}^2} \xrightarrow{P} \sigma = \sqrt{\sigma^2}, n \rightarrow \infty$.

上述的函数均是作用于分量. 一般地, 设映射 $\mathbf{g} : \mathbb{R}^{K \times 1} \rightarrow \mathbb{R}^{q \times 1}$ 满足有 $\mathbb{E} \|\mathbf{g}(\mathbf{X})\| < \infty$, 则依据式 (4.7) 可知

$$\hat{\boldsymbol{\theta}} = \frac{1}{n} \sum_{i=1}^n \mathbf{g}(\mathbf{X}_i) \xrightarrow{P} \boldsymbol{\theta} = \mathbb{E} \mathbf{g}(\mathbf{X}), n \rightarrow \infty.$$

定义参数 $\boldsymbol{\beta} = \mathbf{h}(\boldsymbol{\theta})$, 其中 \mathbf{h} 为连续映射 $\mathbf{h} : \mathbb{R}^{q \times 1} \rightarrow \mathbb{R}^{m \times 1}$, 记参数 $\boldsymbol{\beta}$ 的插入估计量 (*plug-in estimator*) 为

$$\hat{\boldsymbol{\beta}} = \mathbf{h}(\hat{\boldsymbol{\theta}}), \quad (4.8)$$

则依定理 (4.2.7) 知当 $n \rightarrow \infty$, 则 $\hat{\boldsymbol{\beta}} = \mathbf{h}(\hat{\boldsymbol{\theta}}) \xrightarrow{P} \boldsymbol{\beta} = \mathbf{h}(\boldsymbol{\theta})$.

4.3 矩生成函数

按照通常的教材安排,一般是先介绍大数定律,然后介绍中心极限定理.本笔记也遵从这一顺序,但在介绍中心极限定理前,我们首先介绍一个概率论中有用的工具——矩生成函数.

定义 4.3.1 (矩生成函数). 对于随机变量 X , 我们称函数 $M(t) = \mathbb{E}e^{tX}$ 为随机变量 X 矩生成函数 (moment generating function).

依据定义 (4.3.1), 如果随机变量 X 是连续型随机变量, 则其矩生成函数为

$$M(t) = \int_{-\infty}^{+\infty} e^{tx} f(x) dx, \quad (4.9)$$

相应地, 如果随机变量 X 是离散型随机变量, 则有

$$M(t) = \sum_{i=1}^{\infty} e^{tx_i} \mathbb{P}(X = x_i). \quad (4.10)$$

我们知道, 并不是所有的随机变量其分布是存在数学期望的, 例如 Cauchy 分布的概率密度函数不能使得期望的定义积分是收敛的. 同理, 也并不是所有的随机变量都具有矩生成函数 $M(t)$ ^[2]. 矩生成函数 $M(t) = \mathbb{E}e^{tX}$ 使用了期望, 其存在性需要广义积分或无穷级数收敛, 这所有的随机变量都可以保证的. 我们本节讨论的矩生成函数都是默认其存在的.

正如其名, 矩生成函数 $M(t)$ 与随机变量的矩是紧密相关的, 我们可以使用 $M(t)$ 去计算随机变量 X 的矩.

定理 4.3.1 (矩与矩生成函数). 设随机变量 X 的矩生成函数 $M(t)$ 在 $t=0$ 的邻域内存在, 如果 X 具有有限的任意矩, 则有

$$M^{(n)}(0) = \left. \frac{d^n}{dt^n} M(t) \right|_{t=0} = \mathbb{E}X^n$$

证明. 我们仅考虑连续型随机变量 X 显然 $M(0) = \mathbb{E}e^{0 \cdot X} = \mathbb{E}1 = 1$. 当 $n=1$ 时有

$$M'(t) = \frac{d}{dt} M(t) = \frac{d}{dt} \int_{-\infty}^{+\infty} e^{tx} f(x) dx = \int_{-\infty}^{+\infty} x e^{tx} f(x) dx,$$

^[2]考虑复数域 \mathbb{C} , 若定义函数 $C(t) = \mathbb{E}e^{itX}$, 其中 i 是虚数单位, 则函数 $C(t)$ 其是一定存在的, 我们称之为随机变量 X 的特征函数 (characteristic function). 如果读者有学习过复分析, 不难发现, 实际上随机变量 X 的矩生成函数 $M(t)$ 是一个 Laplace 变换, 特征函数 $C(t)$ 是一个 Fourier 变换.

这里用到了变限积分的求导公式^[3]

$$\frac{d}{dt} \int_{-\infty}^{+\infty} F(x, t) dx = \int_{-\infty}^{+\infty} \frac{\partial}{\partial t} F(x, t) dx,$$

故有 $t = 0$ 时有 $M'(0) = \mathbb{E}X$ 成立. 容易用数学归纳法证明

$$M^{(n)}(t) = \frac{d^n}{dt^n} M(t) = \int_{-\infty}^{+\infty} x^n e^{tx} f(x) dx.$$

则有

$$M^{(n)}(0) = \int_{-\infty}^{+\infty} x^n f(x) dx = \mathbb{E}X^n.$$

□

例 4.3.1. 使用矩生成函数计算正态分布 $N(\mu, \sigma^2)$ 的原点矩.

解. 设 $X \sim N(\mu, \sigma^2)$, 则

$$M(t) = \int_{-\infty}^{+\infty} e^{tx} f(x) dx = \int_{-\infty}^{+\infty} e^{tx} \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx = \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{+\infty} e^{-\frac{(x-\mu)^2}{2\sigma^2} + tx} dx,$$

注意到

$$\begin{aligned} -\frac{(x-\mu)^2}{2\sigma^2} + tx &= -\frac{1}{2\sigma^2} [x - (\mu + \sigma^2 t)]^2 + \frac{1}{2\sigma^2} [(\mu + \sigma^2 t)^2 - \mu^2] \\ &= -\frac{1}{2\sigma^2} [(x - (\mu + \sigma^2 t))^2 + \mu^2 - (\mu + \sigma^2 t)^2], \end{aligned}$$

令 $u = x - (\mu + \sigma^2 t)$, $x = u + (\mu + \sigma^2 t)$, 则有

$$\begin{aligned} M(t) &= \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{+\infty} e^{-\frac{1}{2\sigma^2} [u^2 + \mu^2 - (\mu + \sigma^2 t)^2]} dx \\ &= \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{+\infty} e^{-\frac{1}{2\sigma^2} u^2} \cdot e^{-\frac{1}{2\sigma^2} [\mu^2 - (\mu + \sigma^2 t)^2]} dx \\ &= \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2} [\mu^2 - (\mu + \sigma^2 t)^2]} \int_{-\infty}^{+\infty} e^{-\frac{1}{2\sigma^2} u^2} dx \\ &= \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2} [\mu^2 - (\mu + \sigma^2 t)^2]} \int_{-\infty}^{+\infty} e^{-\frac{1}{2\sigma^2} u^2} du \\ &\stackrel{\text{正态分布的 pdf}}{=} \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2} [\mu^2 - (\mu + \sigma^2 t)^2]} \cdot \sigma\sqrt{2\pi} \\ &= e^{-\frac{1}{2\sigma^2} [\mu^2 - (\mu + \sigma^2 t)^2]} = e^{-\frac{1}{2\sigma^2} (\mu^2 - \mu^2 - 2\mu\sigma^2 t - \sigma^4 t^2)} = e^{\frac{1}{2}\sigma^2 t^2 + \mu t}, \end{aligned}$$

^[3]一般地, 完整的变限积分的求导公式

$$\frac{d}{dt} \int_{g(t)}^{h(t)} F(x, t) dx = F(h, t) h' - F(g, t) g' + \int_{g(t)}^{h(t)} \frac{\partial}{\partial t} F(x, t) dx$$

又被称为 Leibniz 律.

即正态分布 $N(\mu, \sigma^2)$ 的矩生成函数为 $M(t) = e^{\frac{1}{2}\sigma^2 t^2 + \mu t}$. 下面我们来求 $M(t)$ 的高阶导数. 注意到 e^x 在 $x = 0$ 处的幂级数为

$$e^x = \sum_{n=0}^{\infty} \frac{1}{n!} x^n = 1 + x + \frac{1}{2!} x^2 + \frac{1}{3!} x^3 + \cdots,$$

则可以得到 $M(t)$ 的幂级数为

$$M(t) = \sum_{n=0}^{\infty} \frac{1}{n!} \left(\frac{1}{2} \sigma^2 t^2 + \mu t \right)^n = 1 + \mu t + \frac{1}{2!} (\mu^2 + \sigma^2) t^2 + \frac{1}{3!} (\mu^3 + 3\mu\sigma^2) t^3 + \cdots, \quad (4.11)$$

由于 $M(t)$ 在 $t = 0$ 处的幂级数

$$M(t) = \sum_{n=0}^{\infty} \frac{1}{n!} M^{(n)}(0) t^n$$

是唯一的, 对比系数便可得到

$$\mathbb{E}X = \mu, \quad \mathbb{E}X^2 = \mu^2 + \sigma^2, \quad \mathbb{E}X^3 = \mu^3 + 3\mu\sigma^2, \quad \cdots$$

特别地, 由于 $Z = X - \mathbb{E}X = X - \mu \sim N(0, \sigma^2)$, 则随机变量 X 的中心矩 $\mathbb{E}(X - \mathbb{E}X)^n = \mathbb{E}Z^n$ 可以通过计算 $\mu = 0$ 的矩生成函数来得到. 令 $\mu = 0$, 式 (4.11) 即有

$$M(t) = \sum_{n=0}^{\infty} \frac{\sigma^{2n} t^{2n}}{2^n n!} = \begin{cases} 0, & k \in \text{奇数}, \\ \frac{1}{k!} \frac{k! \sigma^k}{2^{k/2} (\frac{k}{2})!} t^k, & k \in \text{偶数}, \end{cases}$$

不难得知, 正态分布 $N(\mu, \sigma^2)$ 的奇数阶的中心矩为 0, 偶数阶的中心矩可以依据上述公式计算. \square

我们将式 (4.11) 计算的正态分布 $N(\mu, \sigma^2)$ 的高阶矩整理如表 (4.1).

利用矩生成函数计算矩并不是其最重要的性质, 矩生成函数最宝贵的性质在于其唯一性.

定理 4.3.2 (矩生成函数的唯一性). 设随机变量 X 的矩生成函数 $M(t)$ 存在, 则 $M(t)$ 的形式是唯一的.

定理 (4.3.2) 的证明涉及了 Laplace 变换的反演问题, 其超出了本笔记的范围. 定理 (4.3.2) 的一个直接应用便是, 如果随机变量 X 和 Y 具有相同的矩生成函数, 那么 X 和 Y 的概率分布是相同的. 当然, 由于不是所有的随机变量 X 都存在矩生成函数^[4], 定理 (4.3.2) 在使用上便具有一定局限性.

与矩生成函数相关的另一种技巧便是累积量与累积量生成函数.

^[4]任意的随机变量 X 都具有特征函数, 与矩生成函数类似, 随机变量 X 的特征函数也是唯一的.

表 4.1: 正态分布 $N(\mu, \sigma^2)$ 的高阶矩

阶数 n	原点矩 $\mathbb{E}X^n$	中心矩 $\mathbb{E}(X - \mathbb{E}X)^n$
1	μ	0
2	$\mu^2 + \sigma^2$	σ^2
3	$\mu^3 + 3\mu\sigma^2$	0
4	$\mu^4 + 6\mu^2\sigma^2 + 3\sigma^4$	$3\sigma^4$
5	$\mu^5 + 10\mu^3\sigma^2 + 15\mu\sigma^4$	0
6	$\mu^6 + 15\mu^4\sigma^2 + 45\mu^2\sigma^4 + 15\sigma^6$	$15\sigma^6$
7	$\mu^7 + 21\mu^5\sigma^2 + 105\mu^3\sigma^4 + 105\mu\sigma^6$	0
8	$\mu^8 + 28\mu^6\sigma^2 + 210\mu^4\sigma^4 + 420\mu^2\sigma^6 + 105\sigma^8$	$105\sigma^8$

定义 4.3.2 (累积量与累积量生成函数). 设随机变量 X 存在有矩生成函数 $M(t) = \mathbb{E}e^{tX}$, 我们称函数 $K(t) = \ln M(t)$ 为随机变量 X 的**累积量生成函数** (*cumulant generating function*). 我们称累积量生成函数 $K(t)$ 的在 $t=0$ 处的 n 阶导数 $\kappa_n = K^{(n)}(0)$ 为随机变量 X 的 n 阶**累积量** (*cumulant*).

累积量 κ_n 使用希腊字母 κ 表示, 英语单词写作 Kappa. 由于累积量生成函数 $K(t) = \ln M(t)$, 显然 $K(0) = \ln 1 = 0$. 利用复合函数求导公式不难得到

$$K'(t) = \frac{M'(t)}{M(t)}, \quad K''(t) = \frac{M''(t)}{M(t)} - \left(\frac{M'(t)}{M(t)} \right)^2,$$

对于正态分布 $X \sim N(\mu, \sigma^2)$, 则对应的累积量 κ_n 为

$$\kappa_1 = \frac{M'(0)}{M(0)} = \mu, \quad \kappa_2 = \frac{M''(0)}{M(0)} - \left(\frac{M'(0)}{M(0)} \right)^2 = \mu^2 + \sigma^2 - \mu^2 = \sigma^2,$$

即我们 X 的 1 阶累积量、2 阶累积量分别是 X 的期望与方差. 对于更为高阶的累积量, 我们可以利用复合函数的高阶导数公式 (又称为 Faà di Bruno 公式) 去计算. 可以证明, 前三个累积量 κ_n 等于随机变量的对应阶的中心矩, 但更高阶的累积量 κ_n 是则是随机变量中心矩的多项式函数. 依据定义 (4.3.2), 我们可以得到累积量生成函数在 $t=0$ 处的幂级数为

$$K(t) = \sum_{n=0}^{\infty} \frac{\kappa_n}{n!} t^n = 0 + \kappa_1 t + \frac{\kappa_2}{2!} t^2 + \frac{\kappa_3}{3!} t^3 + \cdots \quad (4.12)$$

我们将矩、累积量与常用的随机变量的数值特征总结为表 (4.2).

本节的最后我们研究一下样本均值 \bar{X}_n 的生成函数. 设随机变量 $X_i, i=1, \dots, n$ 是独立同分布的且具有共有的总体 X , X 的期望和方差记有 $\mu = \mathbb{E}X$ 和 $\sigma^2 = \mathbb{E}(X - \mathbb{E}X)^2$.

表 4.2: 矩、累积量与随机变量的数值特征

阶数	矩			累积量	
	原点矩	中心矩	标准化	自身	归一化
1	均值	0	0	均值	\
2	\	方差	1	方差	1
3	\	\	偏度	\	\
4	\	\	峰度	\	超峰度

我们知道, 样本均值 $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ 的期望和方差分别有 $\mathbb{E}\bar{X}_n = \mu$ 和 $\text{Var}(\bar{X}_n) = n^{-1}\sigma^2$.

令

$$Z_n = \sqrt{n}(\bar{X}_n - \mu) = \frac{1}{\sqrt{n}} \sum_{i=1}^n (X_i - \mu), \quad (4.13)$$

这样随机变量 Z_n 便具有零均值和 σ^2 的方差. 我们现在计算 Z_n 的矩生成函数 $M_{Z_n}(t)$ 以及累积量生成函数 $K_{Z_n}(t)$. 记总体 X 的矩生成函数 $M_X(t) = \mathbb{E}e^{tX}$ 和累积量生成函数 $K_X(t) = \ln M_X(t)$, 则 Z_n 的矩生成函数为

$$\begin{aligned} M_{Z_n}(t) &= \mathbb{E}e^{tZ_n} = \mathbb{E}e^{t \frac{1}{\sqrt{n}} \sum_{i=1}^n (X_i - \mu)} = \mathbb{E} \left(\prod_{i=1}^n e^{t \frac{X_i - \mu}{\sqrt{n}}} \right) \stackrel{\text{i.i.d.}}{=} \prod_{i=1}^n \mathbb{E} \left(e^{t \frac{X_i - \mu}{\sqrt{n}}} \right) \\ &= \prod_{i=1}^n \mathbb{E} \left(e^{t \frac{X_i}{\sqrt{n}}} \cdot e^{t \frac{-\mu}{\sqrt{n}}} \right) = \prod_{i=1}^n e^{t \frac{-\mu}{\sqrt{n}}} \mathbb{E} e^{t \frac{X_i}{\sqrt{n}}} = \prod_{i=1}^n e^{t \frac{-\mu}{\sqrt{n}}} M_X \left(\frac{t}{\sqrt{n}} \right) = e^{-\mu\sqrt{n}t} \left(M_X \left(\frac{t}{\sqrt{n}} \right) \right)^n, \end{aligned}$$

相应地, Z_n 的累积量生成函数为

$$\begin{aligned} K_{Z_n}(t) &= \ln M_{Z_n}(t) = \ln \left[e^{-\mu\sqrt{n}t} \left(M_X \left(\frac{t}{\sqrt{n}} \right) \right)^n \right] \\ &= n \ln M_X \left(\frac{t}{\sqrt{n}} \right) - \mu\sqrt{n}t = nK_X \left(\frac{t}{\sqrt{n}} \right) - \mu\sqrt{n}t. \end{aligned}$$

对于总体 X , 利用式 (4.12) 以及 $\kappa_1 = \mu, \kappa_2 = \sigma^2$ 便得到 Z_n 的累积量生成函数有

$$\begin{aligned} K_{Z_n}(t) &= n \sum_{m=0}^{\infty} \frac{\kappa_m}{m!} n^{-\frac{m}{2}} t^m - \mu\sqrt{n}t = \sum_{m=0}^{\infty} \frac{\kappa_m}{m!} n^{1-\frac{m}{2}} t^m - \mu\sqrt{n}t \\ &= -\mu\sqrt{n}t + \left(\kappa_1 n^{\frac{1}{2}}t + \frac{\kappa_2}{2!}t^2 + \sum_{m=3}^{\infty} \frac{\kappa_m}{m!} n^{1-\frac{m}{2}} t^m \right) \\ &\stackrel{\kappa_1=\mu, \kappa_2=\sigma^2}{=} -\mu\sqrt{n}t + \mu n^{\frac{1}{2}}t + \frac{\sigma^2}{2!}t^2 + \sum_{m=3}^{\infty} \frac{\kappa_m}{m!} n^{1-\frac{m}{2}} t^m \\ &= \boxed{\frac{\sigma^2}{2!}t^2 + \sum_{m=3}^{\infty} \frac{\kappa_m}{m!} n^{1-\frac{m}{2}} t^m}, \end{aligned}$$

亦即有

$$M_{Z_n}(t) = \exp(K_{Z_n}(t)) = \exp\left(\frac{\sigma^2}{2!}t^2 + \sum_{m=3}^{\infty} \frac{\kappa_m}{m!}n^{1-\frac{m}{2}}t^m\right).$$

4.4 渐进分布与中心极限定理

本节我们将介绍渐进分布的内容, 与概率极限相比, 其具有更为广泛的应用场景.

定义 4.4.1 (依分布收敛). 设 $K \times 1$ 维的随机向量的序列 $\{\mathbf{X}_n\}$, $\mathbf{X}_n = (X_{n1}, \dots, X_{nK})^T$ 的累计密度函数为 $F_n(\mathbf{x}) = \mathbb{P}(X_{n1} < x_1, \dots, X_{nK} < x_1)$. 若存在 $K \times 1$ 维随机变量 \mathbf{X} , 其具有连续的累计密度函数为 $F(\mathbf{x})$, 若极限

$$\lim_{n \rightarrow \infty} F_n(\mathbf{x}) = F(\mathbf{x})$$

成立, 则我们称随机序列 $\{\mathbf{X}_n\}$ **依分布收敛** (*converges in distribution*)于 \mathbf{X} , 记作 $\mathbf{X}_n \xrightarrow{d} \mathbf{X}$. 我们也称 F_n 弱收敛, *converge weakly*于 F . 对于随机序列 \mathbf{X} 和其分布 F , 我们称之为随机向量的序列 $\{\mathbf{X}_n\}$ 的**渐进分布** (*asymptotic distribution*)、**大样本分布** (*large sample distribution*)或**极限分布** (*limit distribution*).

依据定义 (4.4.1), 如果 F 对应是一个非随机的向量 \mathbf{a} 的累计密度函数函数, 即有 $\mathbb{P}(\mathbf{X} = \mathbf{a}) = 1$, 这时我们称 F 是退化的 (degenerate). 此时当 $n \rightarrow \infty$ 时, $\mathbf{X}_n \xrightarrow{d} \mathbf{a}$ 便也有 $\mathbb{P}(\mathbf{X}_n = \mathbf{a}) = 1$, 则可以推出 $\mathbf{X}_n \xrightarrow{P} \mathbf{a}$.

如果随机序列 $\{\mathbf{X}_n\}$ 依概率收敛于随机向量 \mathbf{X} , 则一定意味着依分布收敛于 \mathbf{X} . 这是因为, 对于 $\varepsilon > 0$, 我们有不等式^[5]

$$\begin{aligned} & \mathbb{P}(X_1 < x_1 - \varepsilon, \dots, X_K < x_1 - \varepsilon) - \mathbb{P}(\|\mathbf{X}_n - \mathbf{X}\| > \varepsilon) \\ & \leq \mathbb{P}(X_{n1} < x_1, \dots, X_{nK} < x_1) \\ & \leq \mathbb{P}(X_1 < x_1 + \varepsilon, \dots, X_K < x_1 + \varepsilon) + \mathbb{P}(\|\mathbf{X}_n - \mathbf{X}\| > \varepsilon) \end{aligned}$$

成立, 亦即

$$F(\mathbf{x} - \mathbf{1}\varepsilon) - \mathbb{P}(\|\mathbf{X}_n - \mathbf{X}\| > \varepsilon) \leq F_n(\mathbf{x}) \leq F(\mathbf{x} + \mathbf{1}\varepsilon) + \mathbb{P}(\|\mathbf{X}_n - \mathbf{X}\| > \varepsilon). \quad (4.14)$$

由于 $\mathbf{X}_n \xrightarrow{P} \mathbf{X}$, 则 $n \rightarrow \infty$ 时得式 (4.14) 便有

$$F(\mathbf{x} - \mathbf{1}\varepsilon) \leq \lim_{n \rightarrow \infty} F_n(\mathbf{x}) \leq F(\mathbf{x} + \mathbf{1}\varepsilon),$$

由于 ε 具有任意性, 当 $\varepsilon \rightarrow 0$ 时, 即有 F_n 收敛于 F . 我们将依概率收敛与依分布收敛收敛的关系总结如下.

^[5]这个不等式也不显然啊 XD...

定理 4.4.1. 对于的随机向量的序列 $\{X_n\}$,

1. 当 $n \rightarrow \infty$ 时, 如果 $X_n \xrightarrow{P} X$, 则有 $X_n \xrightarrow{d} X$;
2. 当 $n \rightarrow \infty$ 时, 如果 $X_n \xrightarrow{d} a$, 其中 a 为非随机的向量, 则有 $X_n \xrightarrow{P} a$.

通过定理 (4.4.1), 我们可以知道依概率收敛是比依分布收敛更为强的收敛. 而大数定律告诉我们样本均值 \bar{X}_n 依概率收敛于总体期望 $\mathbb{E}X$, 自然的结论便是样本均值 \bar{X}_n 依分布收敛于 $\mathbb{E}X$. 但 $\mathbb{E}X$ 是非随机的, 我们希望能够得到 \bar{X}_n 更为精确的渐进性质. 我们在微积分中曾学习过估计极限的阶 (order) 的方法, 例如极限

$$\lim_{n \rightarrow 0} \frac{\sin n - n}{n} = 0$$

即意味着 $\sin -n$ 收敛于 0 与 n 收敛于 0 是同阶的, 这也被认为 $\sin -n$ 与 n 有相同的“速度”收敛至 0, 我们使用小 o 表示法记有当 $n \rightarrow 0$ 时 $\sin n - n = o(n)$. 在概率极限中, 我们也有类似概率估阶以及概率小 o 表示法, 但我们不打算展开介绍这些内容.

回顾式 (4.13) 标准化后期望为 0、方差为 σ^2 的 Z_n , 则当 $n \rightarrow \infty$ 时 Z_n 与 $1/\sqrt{n}$ 是同阶. 我们现在研究 Z_n 的概率分布, 首先我们指出如下定理.

定理 4.4.2 (Lévy 连续性定理). 对于随机变量的序列 $\{X_n\}$, 若存在矩生成函数 $M_{X_n}(t) = \mathbb{E}e^{tX_n}$ 满足有

$$\lim_{n \rightarrow \infty} M_{X_n}(t) = \lim_{n \rightarrow \infty} M_X(t) = \mathbb{E}e^{tX}, \forall t \in \mathbb{R},$$

则有 $n \rightarrow \infty$ 时 $X_n \xrightarrow{d} X$.

4.5 随机过程的渐进性质

第五章 OLS 估计量的渐进性质

下面我们通过 Monte Carlo 模拟来说明定理 OLS 估计量的相合性. 考虑数据生成过程 $Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$, 其中 ε_i 独立同分布且具有共有的 $N(0, 50^2)$. 我们使用 Monte Carlo 模拟生成一万个样本个体, 并逐个计算 OLS 估计量 $\hat{\beta}_0$ 和 $\hat{\beta}_1$, 绘制了如图 (5.1) 和图 (5.2) 所示的半对数坐标的 OLS 估计量的折线图.

图 (5.1) 和图 (5.2) 的横轴是 10 为底数的对数, 这样对于 $n = 10000$ 个回归系数的估计值, 图 (5.1) 和图 (5.2) 可以较为直观地展示不同数量级下 OLS 估计量与真实值的差异. 首先, 对于 10^0 的数量级, 不能发现在最初样本容量很小时偏差很大. 这是因为, $n = 1$ 时起始点的 $\hat{\beta}_1 = 0$, 这是因为只有一个样本点的一元 OLS 即等价于直线 $Y = \bar{Y}$, 当 $n = 2$ 时则是一元 OLS 便是平面中两点确定的一条直线, 信息有限 OLS 估计量不能剔除回归误差 ε_i 的影响——尽管我们这次模拟中 β_0 的估计非常准确. 随着 n 的增加, 在 10^1 的数量级的样本容量内 β_1 的估计值没有最初那样的大的偏离, 但依旧与真实值差距不小. 对于 10^1 的数量级, 图 (5.2) 中可以直观感受 n 的增加缩小了 $\hat{\beta}_1$ 与 β_1 差距, 尤其是较 10^0 的数量级时. 当样本容量 n 破百后, 在 10^2 的数量级下 $\hat{\beta}_1$ 的偏差经历了一个小幅的增长过程, 则也说明了依概率收敛、几乎处处收敛都不能保证有充分大的 n 使得估计量与真实值有任意小的偏差. 当 n 的数量级达到 10^3 的数量级, 不难发现 $\hat{\beta}_1$ 与真实值 $\beta_1 = 25$ 的偏差在图中很难体现, 这样便是说 10^3 的数量级下 OLS 估计量要较更小的数量级更为精确.

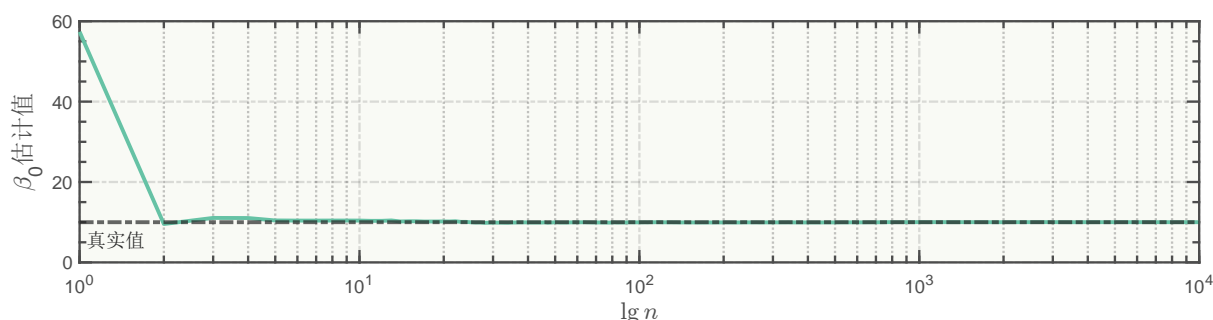
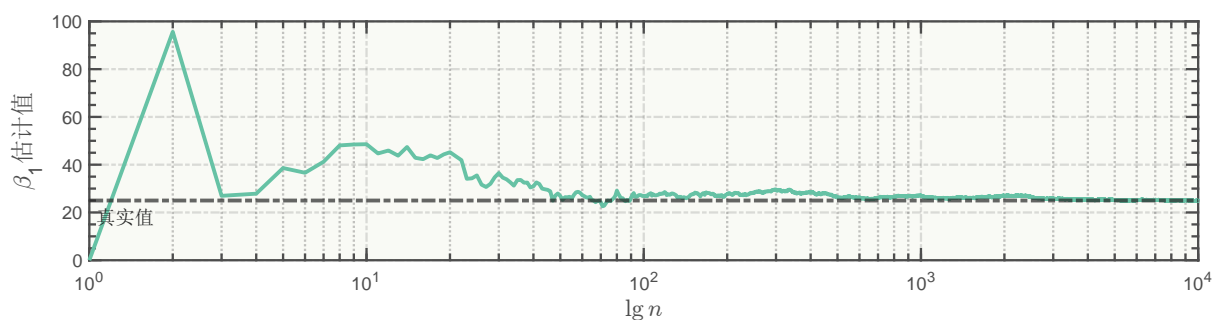


图 5.1: 样本容量增大时的 $\hat{\beta}_0$

图 5.2: 样本容量增大时的 $\hat{\beta}_1$

5.1 遍历平稳下 OLS 估计量的渐进性质

5.2 代码

```

1 %% OLS估计量的相合性
2 clear
3 clc
4 s=RandStream('mcg16807','Seed',114514);%设定复现的随机流以及种子
5 %设定E(Y|X)=10x+25,e~N(0,50^2)
6 n=10000;%样本容量
7 X=[ones(n,1),randsample(s,-100:100,n,true)'];%生成样本X
8 e=50*randn(s,n,1);%生成扰动项
9 Y=X*[25;10]+e;%数据生成过程
10 Beta=zeros(n,2);%初始化存储每次回归结果的Beta
11 for i=1:n
12     beta=X(1:i,:)\Y(1:i,:);%OLS估计量
13     Beta(i,:)=beta;%保存OLS估计量
14 end
15 selfGrootDefault(2)%绘图美化函数,没有请注释
16 %绘制估计值均值变化
17 figure(1)%普通坐标
18 set(gcf,'unit','centimeters',...
19     'Position',[0 0 21*1.25 29.7*0.35*1.25]);%设置figure为A4纸宽度
20 tiledlayout(2,1);%创建2*1画布
21 nexttile
22 plot(Beta(:,1),'LineWidth',2)
23 xlabel('\fontname{宋体}样本容量\fontname{Times New ...
24     Roman}n','FontSize',12.5)
25 ylabel('\beta_1\fontname{宋体}估计值','FontSize',12.5);

```

```

25 yline(25, '-.', '\fontname{宋体}真实值', 'LineWidth', 2, 'LabelHorizontalAlignment' ...
    ...
26     , 'left', 'LabelVerticalAlignment', 'bottom')
27
28 nexttile
29 plot(Beta(:,2), 'LineWidth', 2)
30 xlabel('\fontname{宋体}样本容量\fontname{Times New ...
    Roman}\n', 'FontSize', 12.5)
31 ylabel('\beta_0\fontname{宋体}估计值', 'FontSize', 12.5);
32 yline(10, '-.', '\fontname{宋体}真实值', 'LineWidth', 2, 'LabelHorizontalAlignment' ...
    ...
33     , 'left', 'LabelVerticalAlignment', 'bottom')
34
35 figure(2)%对数坐标
36 set(gcf, 'unit', 'centimeters', ...
37     'Position', [0 0 21*1.25 29.7*0.35*1.25]); %设置figure为A4纸宽度
38 tiledlayout(2,1); %创建2*1画布
39 nexttile
40 semilogx(Beta(:,1), 'LineWidth', 2)
41 xlabel('$\ln n$', 'FontSize', 12.5, 'Interpreter', 'latex')
42 ylabel('\beta_1\fontname{宋体}估计值', 'FontSize', 12.5);
43 yline(25, '-.', '\fontname{宋体}真实值', 'LineWidth', 2, 'LabelHorizontalAlignment' ...
    ...
44     , 'left', 'LabelVerticalAlignment', 'bottom')
45
46 nexttile
47 semilogx(Beta(:,2), 'LineWidth', 2)
48 xlabel('$\ln n$', 'FontSize', 12.5, 'Interpreter', 'latex')
49 ylabel('\beta_0\fontname{宋体}估计值', 'FontSize', 12.5);
50 yline(10, '-.', '\fontname{宋体}真实值', 'LineWidth', 2, 'LabelHorizontalAlignment' ...
    ...
51     , 'left', 'LabelVerticalAlignment', 'bottom')

```

概念索引

符号

\bar{R}^2	80
\tilde{R}^2	80

B

BLUE	32
BUE	32
备择假设 alternative hypothesis	91
被解释变量 explained variable	6
标准化残差 standardized residual	44, 76
标准误	49, 85
标准误差 standard error	49, 85

C

Chebyshev 不等式 Chebyshev's inequality	108
残差 residual	35, 75
残差平方和 residual sum of square	39
插入估计量 plug-in estimator	112
差一法 leave one out	76
常数项 constant term	12
超平面 hyperplane	69
抽样误差 sampling error	47, 66

D

demeaned regressor	12
大数定律 law of large number	110
大样本分布 large sample distribution	118

概念索引	124
第二类错误 type II error	92
第一类错误 type I error	92
调整的 R^2 adjusted R^2	80
迭代期望定律 law of iterated expectations	5
独立同分布 independent and identically distributed	17
E	
Eicker-White 协方差矩阵估计量 Eicker-White covariance matrix estimator	83
F	
F 比率 F-ratio	99
F 检验 F test	98
方差分析 analysis-of-variance	38
非中心化 R^2 uncentered R^2	79
G	
Gram 矩阵	68
概率极限 probability limit	104
杠杆值 leverage value	41, 75
功效函数 power function	91
归零矩阵 annihilator matrix	37, 73
H	
回归标准差	9
回归方差 regression variance	9
回归方程 regression equation	6
回归函数	3
回归模型 regression model	6
回归平方和 regression sum of square	39
回归误差 regression error	6, 24
回归系数 regression coefficient	12
回归元 regressor	6
回归子 regressand	6
J	
极限分布 limit distribution	118

概念索引	125
------	-----

几乎处处 almost surely	106
几乎处处收敛 converge almost surely	106
假设 hypothesis	90
假设检验 hypothesis test	91
检验 test	91
检验统计量 test statistic	91
简单迭代期望定律 simple law of iterated expectations	4
简单随机样本 simple random sample	17
渐进分布 asymptotic distribution	118
渐进理论 asymptotic theorem	104
接受域 acceptance region	91
截距项 intercept term	12
解释变量 explanatory variable	6
解释平方和 explained sum of square	38
经典线性回归模型 classic linear regression model	46
矩估计 moment estimation	18
矩生成函数 moment generating function	113
矩阵 \mathbf{A} 的平方根 square root of the matrix \mathbf{A}	74
拒绝域 rejection region	91
均方误差 mean square error	8
均值独立 mean independence	7

K

可决系数 coefficient of determination	39, 79
-----------------------------------	--------

L

LLO 残差 LLO residual	77
LLO 回归 LLO regression	77
LLO 预测值 LLO prediction value	77
LS 估计量 LS estimator	19
累积量 cumulant	116
累积量生成函数 cumulant generating function	116
离差平方和 sum of squared deviations	38
临界值 critical value	91

M

Markov 不等式 Markov's inequality	107
MVU 估计	8

N

内生性 endogeneity	58
拟合优度 goodness of fit	39, 79
拟合值 fitted value	34, 75

O

OLS 估计量 OLS estimator	19, 62
-----------------------------	--------

P

p 值 p-value	94
偏差校正估计量 bias-corrected estimator	43, 82
平方和分解	38
普通最小二乘估计量 ordinary least squares estimator	19
谱分解 spectral decomposition	74

Q

强大数定律 strong law of large number	110
球形误差方差 spherical error variance	56

R

R^2	39
弱大数定律 weak law of large number	108, 110
弱收敛 converge weakly	118

S

似然比检验 likelihood ratio test	99
受约束的最小二乘估计量 constrained least squares estimator	99
枢轴变量 pivotal variable	49
数据矩阵 data matrix	55
数据生成过程 data generating process	17
数据向量 data vector	55

随机扰动项 stochastic disturbance	6
随机误差项 stochastic error	6
随机序列 random sequence	104

T

t 检验 t test	93
t 统计量 t-statistic	49, 89
特征函数 characteristic function	113
条件标准差 conditional standard deviation	9
条件定理 conditioning theorem	5
条件方差 conditional variance	9
条件期望函数 conditional expectation function	3
条件期望函数误差 CEF error	6
条件同方差 conditional homoskedasticity	56
条件协方差矩阵 conditional covariance matrix	58
条件异方差 conditional heteroskedasticity	59
同方差 homoskedasticity	9
同方差的 Wald 统计量 homoskedastic Wald statistic	97
同方差线性 CEF 模型 homoskedastic linear CEF model	11
同方差线性回归模型 homoskedastic linear regression model	11, 25
投影 projection	69
投影矩阵 projection matrix	36, 70
投影误差 projection error	13

W

Wald 检验 Wald test	93
Wald 统计量 Wald statistic	97
White 协方差矩阵估计量 White covariance matrix estimator	83
稳健的协方差矩阵估计量 robust covariance matrix estimator	84
无多重共线性 no multicollinearity	56
误差平方和 sum of squared errors	19

X

显著性检验 significance test	94, 99
显著性水平 significance level	92

概念索引	128
线性 CEF 模型 linear CEF model	11
线性回归方程 linear regression equation	12
线性回归模型 linear regression model	11, 24, 66
线性假设 linearity assumption	56
线性条件期望函数 linear CEF	11
线性投影 linear projection	16
线性投影方程 linear projection equation	13
线性投影模型 linear projection model	16, 62
序列相关 serial correlation	57
Y	
严格外生性 strict exogeneity	56
厌恶参数 nuisance parameter	48, 89
样本原点矩 sample origin moment	18
样本中心矩 sample center moment	18
依分布收敛 converges in distribution	118
依概率收敛 converge in probability	104, 106
依均方收敛 converge in mean square	108
异方差 heteroskedasticity	9, 59
异方差稳健的协方差矩阵估计量 heteroskedasticity-robust covariance matrix estimator	84
异方差一致的协方差矩阵估计量 heteroskedasticity-consistent covariance matrix estimator	84
因变量 dependent variable	6
预测残差 prediction residual	77
预测误差 prediction error	77
原假设 null hypothesis	91
Z	
张成 span	70
正规方程组 normal equations	15
正交矩阵 orthogonal matrix	73
正交投影 orthogonal projection	69
正态回归模型 normal regression model	45, 87
正态假设 normality assumption	45

置信区间 confidence interval	49
置信水平 confidence level	49
自变量 independent variable	6
自相关 autocorrelation	57
子样本 sub-sample	76
总体回归函数 population regression function	3
总体平方和 total sum of square	38
组间方差 across group variance	10
组内方差 within group variance	10
最小二乘估计量 least squares estimator	19
最小方差无偏估计 minimum variance unbiased estimate	8
最优无偏估计量 best unbiased estimator	32, 69
最优线性无偏估计量 best linear unbiased estimator	32, 68
最优线性预测量 best linear predictor	13, 60
最优预测量 best predictor	7

模型索引

模型 2.1.1 (回归模型)	6
模型 2.2.1 (线性 CEF 模型)	11
模型 2.2.2 (同方差线性 CEF 模型)	11
模型 2.2.3 (一元线性投影模型)	16
模型 2.4.1 (一元线性回归模型)	23
模型 2.4.2 (同方差一元线性回归模型)	25
模型 2.7.1 (一元正态回归模型)	45
模型 2.7.2 (一元经典线性回归模型)	46
模型 3.3.1 (多元线性投影模型)	62
模型 3.3.2 (多元线性投影模型 (含有截距项))	62
模型 3.6.1 (多元线性回归模型)	66
模型 3.9.1 (leave-one-out 回归)	77
模型 3.12.1 (多元正态回归模型)	86