

截面与面板数据分析

第一次习题课

迭代期望公式的证明与应用

2025 年 9 月 28 日

目录

1 迭代期望公式 (Law of Iterated Expectation)	5
1.1 离散随机变量证明	5
1.2 连续随机变量证明	6
1.3 积分顺序交换的数学依据	6
1.3.1 富比尼定理 (Fubini's Theorem)	6
1.3.2 托内利定理 (Tonelli's Theorem)	7
1.3.3 适用条件	7
2 经典反例：有理数/无理数随机变量	8
2.1 第 1 步：设置场景	8
2.2 第 2 步：“直觉”推导（错误的推导）	8
2.2.1 计算内层条件期望 $\mathbb{E}[X Y = y]$	8
2.2.2 计算外层期望 $\mathbb{E}[h(Y)]$	9
2.2.3 计算真实的 $\mathbb{E}[X]$	9
2.3 第 3 步：揭示矛盾	9
2.3.1 计算内层条件期望 $\mathbb{E}[Y X]$	9
2.4 第 4 步：引出教训和严格理论	10
2.4.1 问题的核心	10
2.4.2 测度论的解决方案	10
2.4.3 分析我们的例子	10

3	条件期望的信息论解释：随机变量的视角	11
3.1	条件期望作为随机变量的数学基础	11
3.1.1	传统误解与现代理解	11
3.1.2	条件期望的函数形式	11
3.2	信息论视角：信息的价值	11
3.2.1	信息的定量衡量	11
3.2.2	信息价值的经济解释	12
3.3	经济学应用：信息经济学	12
3.3.1	案例 1：股票分析师的价值	12
3.3.2	案例 2：央行的通胀预测	12
3.4	条件期望的动态更新：贝叶斯学习	13
3.4.1	信息的递归更新	13
3.4.2	信息价值的边际递减	13
3.5	信息结构的设计：机制设计理论	13
3.5.1	最优信息披露	13
3.5.2	信息租金与市场效率	14
3.6	计量经济学中的应用：工具变量	14
3.6.1	工具变量的信息论解释	14
3.6.2	弱工具变量问题的信息论解释	14
3.7	机器学习中的条件期望	15
3.7.1	预测模型的信息论基础	15
3.7.2	特征选择的信息价值	15
4	概率不等式的应用示例	17
4.1	问题设定	17
4.2	使用 Markov 不等式	17
4.3	使用切比雪夫不等式	18
4.4	假设正态分布	18
4.5	方法比较	18
5	总体参数 vs 样本统计量的重要说明	19
5.1	概率不等式中的参数	19
5.2	实际应用中的考虑	19
5.3	符号对比	19
6	协方差矩阵的正定性	19

7 正半定矩阵的等价定义	20
8 熵的性质	21
9 矩阵范数与内积	22
10 正态分布的性质	22
11 计量经济学中的条件期望：零均值假设的重要性	23
11.1 两种假设的定义与关系	23
11.1.1 假设对比	23
11.1.2 逻辑关系	24
11.2 经典反例：非线性关系	24
11.2.1 验证 $E[X_i \varepsilon_i] = 0$	24
11.2.2 验证 $E[\varepsilon_i X_i] \neq 0$	24
11.3 条件期望为零的重要性	25
11.3.1 排除所有函数关系	25
11.3.2 OLS 估计量的性质	25
11.4 误差方差估计的偏误	25
11.4.1 问题的根源	25
11.4.2 方差估计的偏误	26
11.5 数值例子验证	26
11.5.1 真实误差方差	26
11.5.2 OLS 估计的方差	26
11.6 计量经济学实践中的后果	27
11.7 检验与补救措施	27
11.7.1 检验方法	27
11.7.2 补救措施	27
11.8 两种假设的定义与关系	28
11.8.1 假设对比	28
11.8.2 逻辑关系	28
11.9 经典反例：非线性关系	29
11.9.1 验证 $E[X_i \varepsilon_i] = 0$	29
11.9.2 验证 $E[\varepsilon_i X_i] \neq 0$	29
11.10 条件期望为零的重要性	29
11.10.1 排除所有函数关系	29
11.10.2 OLS 估计量的性质	30

11.11 误差方差估计的偏误	30
11.11.1 问题的根源	30
11.11.2 方差估计的偏误	30
11.12 数值例子验证	31
11.12.1 真实误差方差	31
11.12.2 OLS 估计的方差	31
11.13 计量经济学实践中的后果	31
11.14 检验与补救措施	32
11.14.1 检验方法	32
11.14.2 补救措施	32
12 结论	33

1 迭代期望公式 (Law of Iterated Expectation)

问题 1: 迭代期望公式

证明迭代期望公式:

$$\mathbb{E}[X] = \mathbb{E}[\mathbb{E}[X | Y]] \quad (1)$$

1.1 离散随机变量证明

离散情形证明

对于离散随机变量 X 和 Y , 我们有:

$$\mathbb{E}[\mathbb{E}[X | Y]] = \sum_y \mathbb{E}[X | Y = y] \cdot P(Y = y) \quad (2)$$

$$= \sum_y \left(\sum_x x \cdot P(X = x | Y = y) \right) \cdot P(Y = y) \quad (3)$$

$$= \sum_y \sum_x x \cdot P(X = x | Y = y) \cdot P(Y = y) \quad (4)$$

$$= \sum_y \sum_x x \cdot P(X = x, Y = y) \quad (5)$$

$$= \sum_x x \cdot \left(\sum_y P(X = x, Y = y) \right) \quad (6)$$

$$= \sum_x x \cdot P(X = x) \quad (7)$$

$$= \mathbb{E}[X] \quad (8)$$

1.2 连续随机变量证明

连续情形证明

对于连续随机变量 X 和 Y ，我们有：

$$\mathbb{E}[\mathbb{E}[X | Y]] = \int_{-\infty}^{\infty} \mathbb{E}[X | Y = y] \cdot f_Y(y) dy \quad (9)$$

$$= \int_{-\infty}^{\infty} \left(\int_{-\infty}^{\infty} x \cdot f_{X|Y}(x|y) dx \right) \cdot f_Y(y) dy \quad (10)$$

$$= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x \cdot f_{X|Y}(x|y) f_Y(y) dx dy \quad (11)$$

$$= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x \cdot f_{X,Y}(x, y) dx dy \quad (12)$$

$$= \int_{-\infty}^{\infty} x \cdot \left(\int_{-\infty}^{\infty} f_{X,Y}(x, y) dy \right) dx \quad (13)$$

$$= \int_{-\infty}^{\infty} x \cdot f_X(x) dx \quad (14)$$

$$= \mathbb{E}[X] \quad (15)$$

其中：

- $f_Y(y)$: Y 的边缘概率密度函数
- $f_{X|Y}(x|y)$: 给定 $Y = y$ 时 X 的条件概率密度函数
- $f_{X,Y}(x, y)$: X 和 Y 的联合概率密度函数
- $f_X(x)$: X 的边缘概率密度函数

1.3 积分顺序交换的数学依据

理论基础

在连续随机变量的证明中，从第 (12) 行到第 (13) 行的积分顺序交换是合法的，这基于以下数学定理：

1.3.1 富比尼定理 (Fubini's Theorem)

如果函数 $|x \cdot f_{X,Y}(x, y)|$ 在 \mathbb{R}^2 上可积，那么可以交换积分的顺序。

1.3.2 托内利定理 (Tonelli's Theorem)

对于非负函数，即使不可积，也可以交换积分顺序。由于 $f_{X,Y}(x,y) \geq 0$ ，我们也可以应用托内利定理。

1.3.3 适用条件

1. $f_{X,Y}(x,y)$ 是联合概率密度函数，满足

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{X,Y}(x,y) dx dy = 1$$

2. 若 $\mathbb{E}[|X|] < \infty$ (即 X 的绝对期望存在)，则 $|x \cdot f_{X,Y}(x,y)|$ 在 \mathbb{R}^2 上可积

重要说明

因此，在 $\mathbb{E}[|X|] < \infty$ 的条件下，积分顺序的交换在概率论中是严格合法的。

2 经典反例：有理数/无理数随机变量

反例目的

虽然我们已经证明了迭代期望公式在一般条件下成立，但需要注意的是，如果不满足相应的数学条件，该公式可能不成立。以下是一个经典反例：

2.1 第 1 步：设置场景

假设我们有一个非常奇怪的随机实验。在单位区间 $[0, 1]$ 上按照均匀分布随机选取一个点 ω 。定义两个随机变量：

$$Y(\omega) = \omega \quad (\text{就是选取的点本身})$$

$$X(\omega) = \begin{cases} 1, & \text{如果 } \omega \text{ 是有理数} \\ 0, & \text{如果 } \omega \text{ 是无理数} \end{cases}$$

已知： $P(Y = y) = 0$ 对于任意 $y \in [0, 1]$ ，因为是在连续区间上取单点。

2.2 第 2 步：“直觉”推导（错误的推导）

现在，我们来计算 $\mathbb{E}[\mathbb{E}[X | Y]]$ 。

2.2.1 计算内层条件期望 $\mathbb{E}[X | Y = y]$

- “在给定 $Y = y$ 的条件下， ω 被固定为 y 。”
- 那么， X 的值也就完全确定了：
- 如果 y 是有理数，则 $X = 1$ ，所以 $\mathbb{E}[X | Y = y] = 1$ 。
- 如果 y 是无理数，则 $X = 0$ ，所以 $\mathbb{E}[X | Y = y] = 0$ 。

因此，我们可以定义一个函数：

$$h(y) = \mathbb{E}[X | Y = y] = \begin{cases} 1, & \text{如果 } y \text{ 是有理数} \\ 0, & \text{如果 } y \text{ 是无理数} \end{cases}$$

2.2.2 计算外层期望 $\mathbb{E}[h(Y)]$

- $\mathbb{E}[\mathbb{E}[X | Y]] = \mathbb{E}[h(Y)]$ 。
- 由于 Y 是 $[0, 1]$ 上的均匀分布，这个期望就是函数 $h(y)$ 在 $[0, 1]$ 上的积分。
- $h(y)$ 是一个狄利克雷函数：在有理数点上为 1，在无理数点上为 0。
- 在黎曼积分的意义下，这个函数不可积，但如果我们强行用概率的”加权平均”来想，可能会有：因为有理数测度为 0，无理数测度为 1，所以这个”平均值”应该是 $1 \times 0 + 0 \times 1 = 0$ 。
- 因此，他们可能会得出结论： $\mathbb{E}[\mathbb{E}[X | Y]] = 0$ 。

2.2.3 计算真实的 $\mathbb{E}[X]$

- X 以概率 0 取值为 1（因为有理数集是零测集），以概率 1 取值为 0。
- 所以， $\mathbb{E}[X] = 1 \times 0 + 0 \times 1 = 0$ 。

表面上的”一致”：看起来， $\mathbb{E}[\mathbb{E}[X | Y]] = 0 = \mathbb{E}[X]$ 。迭代期望定律似乎成立了？

2.3 第 3 步：揭示矛盾

现在，我们调换角色。考虑 $\mathbb{E}[\mathbb{E}[Y | X]]$ 。

2.3.1 计算内层条件期望 $\mathbb{E}[Y | X]$

- 当 $X = 1$ ：这意味着 ω 是有理数。有理数在 $[0, 1]$ 上是稠密的但可数。学生可能会困惑：”给定 ω 是有理数， Y 的期望是多少？”他们无法给出一个确定的数值，因为均匀分布在一个可数集上没有定义（不是连续分布）。
- 当 $X = 0$ ：这意味着 ω 是无理数。同样，均匀分布在无理数集上也无法用常规方法定义期望。

矛盾出现

学生发现他们根本无法用初等的方法定义 $\mathbb{E}[Y | X = x]$ ！这个表达式变得没有意义。

矛盾出现：迭代期望定律 $\mathbb{E}[\mathbb{E}[Y | X]] = \mathbb{E}[Y]$ 的右边 $\mathbb{E}[Y] = 0.5$ 是明确可知的。但左边却无法计算。这强烈地暗示了我们关于条件期望的直觉概念在此崩溃了。

2.4 第 4 步：引出教训和严格理论

2.4.1 问题的核心

当条件事件（如 $\{Y = y\}$ 或 $\{X = 1\}$ ）的概率为零时，我们熟悉的条件概率公式 $P(A | B) = P(A \cap B) / P(B)$ 完全失效。我们之前的推导是建立在沙滩上的。

2.4.2 测度论的解决方案

在测度论中，条件期望 $\mathbb{E}[X | Y]$ 不是关于点 $Y = y$ ，而是关于 Y 生成的整个 σ -代数 $\sigma(Y)$ 。

它被定义为一个新的随机变量 Z ，满足两个条件：

1. Z 是 $\sigma(Y)$ -可测的（意味着它的值只依赖于 Y ）。
2. 对于任意事件 $A \in \sigma(Y)$ ，有 $\int_A Z dP = \int_A X dP$ 。

在这个框架下，可以严格证明迭代期望定律 $\mathbb{E}[\mathbb{E}[X | \mathcal{G}]] = \mathbb{E}[X]$ 是成立的，其中 \mathcal{G} 是一个 σ -代数。

2.4.3 分析我们的例子

- 在第一个计算 $\mathbb{E}[\mathbb{E}[X | Y]]$ 中，我们歪打正着得到了正确答案，但这只是巧合。因为我们构造的函数 $h(y)$ 恰好满足测度论中条件期望的定义（可以验证它满足上述两个条件）。
- 在第二个计算 $\mathbb{E}[\mathbb{E}[Y | X]]$ 中，我们的直觉完全失败，因为它无法处理零概率事件。而测度论可以给出严格的定义。

3 条件期望的信息论解释：随机变量的视角

条件期望的本质

理解 $\mathbb{E}[X | Y]$ 不仅仅是一个数值计算，而是一个随机变量。从信息论角度探讨条件期望如何量化信息的价值和不确定性的减少。

3.1 条件期望作为随机变量的数学基础

3.1.1 传统误解与现代理解

常见误解： $\mathbb{E}[X | Y = y]$ 是给定 $Y = y$ 时 X 的期望值，是一个确定的数。

正确理解： $\mathbb{E}[X | Y]$ 是一个随机变量，它是 Y 的函数。

定理 1 (条件期望的随机变量性质). 设 (X, Y) 是定义在概率空间 (Ω, \mathcal{F}, P) 上的随机变量对。条件期望 $\mathbb{E}[X | Y]$ 是一个随机变量，满足：

1. $\mathbb{E}[X | Y]$ 是 $\sigma(Y)$ -可测的
2. 对任意 $A \in \sigma(Y)$ ，有 $\int_A \mathbb{E}[X | Y] dP = \int_A X dP$

直观理解：

- 在观察 Y 之前， $\mathbb{E}[X | Y]$ 是随机的
- 一旦观察到 $Y = y$ ， $\mathbb{E}[X | Y]$ 就确定为某个值
- 这个“某个值”就是 $\mathbb{E}[X | Y = y]$

3.1.2 条件期望的函数形式

设 Y 取离散值 $\{y_1, y_2, \dots\}$ ，则： $\mathbb{E}[X | Y] = \sum_i \mathbb{E}[X | Y = y_i] \cdot \mathbf{1}_{Y=y_i}$

其中 $\mathbf{1}_{Y=y_i}$ 是指示函数。

这清楚地表明 $\mathbb{E}[X | Y]$ 确实是 Y 的函数，因此是随机变量。

3.2 信息论视角：信息的价值

3.2.1 信息的定量衡量

从信息论角度，观察 Y 提供了关于 X 的信息。这种信息的价值可以通过以下方式量化：

无条件不确定性： $\text{Var}(X)$ **条件不确定性：** $\mathbb{E}[\text{Var}(X | Y)]$ **信息价值：** $\text{Var}(X) - \mathbb{E}[\text{Var}(X | Y)] = \text{Var}(\mathbb{E}[X | Y])$

这个等式告诉我们： $\underbrace{\text{Var}(X)}_{\text{总不确定性}} = \underbrace{\mathbb{E}[\text{Var}(X | Y)]}_{\text{条件不确定性}} + \underbrace{\text{Var}(\mathbb{E}[X | Y])}_{\text{信息价值}}$

3.2.2 信息价值的经济解释

数学表达： $\text{Var}(\mathbb{E}[X | Y])$ 衡量了 $\mathbb{E}[X | Y]$ 这个随机变量的变异性。

经济含义：

- 如果 Y 完全不相关于 X ，则 $\mathbb{E}[X | Y] = \mathbb{E}[X]$ （常数），信息价值为 0
- 如果 Y 完全决定 X ，则 $\text{Var}(X | Y) = 0$ ，信息价值最大化
- 信息价值衡量了观察 Y 能减少多少关于 X 的不确定性

3.3 经济学应用：信息经济学

3.3.1 案例 1：股票分析师的价值

设 X 为股票的真实价值， Y 为分析师报告。

投资者的决策问题： - 无信息时：基于 $\mathbb{E}[X]$ 做决策 - 有信息时：基于 $\mathbb{E}[X | Y]$ 做决策

信息的经济价值： $\text{Value of Information} = \mathbb{E}[\text{Utility}(\mathbb{E}[X | Y])] - \text{Utility}(\mathbb{E}[X])$

由于 $\mathbb{E}[X | Y]$ 是随机变量，投资者在购买信息前无法确定具体收益，但可以评估期望收益。

关键洞察：

- $\text{Var}(\mathbb{E}[X | Y])$ 越大，信息的潜在价值越高
- 但信息价值也取决于投资者的风险偏好
- 风险厌恶者可能偏好确定性更高的信息

3.3.2 案例 2：央行的通胀预测

设 X_t 为下期实际通胀率， Y_t 为当期观察到的经济指标集合。

央行的预测问题： $\pi_{t+1|t} = \mathbb{E}[X_{t+1} | Y_t]$

这里 $\pi_{t+1|t}$ 是一个随机变量，因为在时期 t 开始时， Y_t 还未完全观察到。

政策制定的不确定性：

$$\text{Var}(\pi_{t+1|t}) = \text{Var}(\mathbb{E}[X_{t+1} | Y_t]) \quad (16)$$

$$= \text{预测本身的不确定性} \quad (17)$$

这种不确定性来源于：

- 经济指标 Y_t 的随机性
- 通胀与指标之间关系的复杂性
- 模型参数的不确定性

3.4 条件期望的动态更新：贝叶斯学习

3.4.1 信息的递归更新

设投资者对资产价值有先验信念，然后逐步观察信息：

时期 0: 先验期望 $\mathbb{E}[X]$ **时期 1:** 观察 Y_1 后，更新为 $\mathbb{E}[X | Y_1]$ **时期 2:** 观察 Y_2 后，更新为 $\mathbb{E}[X | Y_1, Y_2]$

关键性质: $\mathbb{E}[\mathbb{E}[X | Y_1, Y_2] | Y_1] = \mathbb{E}[X | Y_1]$

这表明新信息的期望价值为零（在当前信息集下）。

3.4.2 信息价值的边际递减

信息的边际价值通常递减： $\text{Var}(\mathbb{E}[X | Y_1]) \geq \text{Var}(\mathbb{E}[X | Y_1, Y_2] | Y_1)$

经济含义:

- 第一份研究报告价值最高
- 后续报告的价值逐步递减
- 解释了为什么“独家信息”更昂贵

3.5 信息结构的设计：机制设计理论

3.5.1 最优信息披露

公司选择披露什么信息给投资者？

设公司知道真实价值 X ，可以选择信号结构 $Y | X$ 。

公司的优化问题: $\max_{\text{signal structure}} \mathbb{E}[\text{Stock Price}(Y)] - \text{Cost}(\text{信息精度})$

其中股价 $P(Y) = \mathbb{E}[X | Y]$ 。

关键权衡:

- 更精确的信息 \Rightarrow 更高的 $\text{Var}(\mathbb{E}[X | Y])$
- 但精确信息也可能暴露不利信息
- 最优策略取决于 X 的分布和公司的私人信息

3.5.2 信息租金与市场效率

在信息不对称市场中：

知情交易者：知道 X 的真实值 **噪音交易者：**只能观察到公开信号 Y

均衡价格： $P = \mathbb{E}[X | Y, \text{订单流}]$

这里价格本身就是一个条件期望，是随机变量。

信息租金：

$$\text{Informed Profit} = \mathbb{E}[(X - P) \cdot \text{Trading Volume}] \quad (18)$$

$$= \mathbb{E}[(X - \mathbb{E}[X | Y]) \cdot Q(X, Y)] \quad (19)$$

其中 $Q(X, Y)$ 是交易量函数。

3.6 计量经济学中的应用：工具变量

3.6.1 工具变量的信息论解释

在工具变量估计中： $Y = \beta X + \varepsilon$, $\mathbb{E}[\varepsilon | Z] = 0$

其中 Z 是工具变量。

传统解释： Z 与 ε 不相关，但与 X 相关。

信息论解释：

- $\mathbb{E}[Y | Z]$ 是 Z 的函数，是随机变量
- $\mathbb{E}[X | Z]$ 是 Z 的函数，是随机变量
- IV 估计量利用了 $\mathbb{E}[Y | Z]$ 和 $\mathbb{E}[X | Z]$ 之间的关系

IV 估计的方差分解： $\text{Var}(Y) = \text{Var}(\mathbb{E}[Y | Z]) + \mathbb{E}[\text{Var}(Y | Z)]$

IV 的有效性取决于 $\text{Var}(\mathbb{E}[Y | Z])$ 和 $\text{Var}(\mathbb{E}[X | Z])$ 的比值。

3.6.2 弱工具变量问题的信息论解释

弱工具变量： $\text{Var}(\mathbb{E}[X | Z])$ 很小

信息论含义：

- Z 提供的关于 X 的信息很少
- $\mathbb{E}[X | Z] \approx \mathbb{E}[X]$ (近似常数)
- 信噪比很低，估计不精确

解决方案：

- 寻找信息含量更高的工具变量
- 使用多个工具变量组合
- 应用正则化方法控制估计方差

3.7 机器学习中的条件期望

3.7.1 预测模型的信息论基础

在监督学习中，模型 $f(X)$ 试图逼近 $\mathbb{E}[Y | X]$ ：

理论最优： $f^*(X) = \mathbb{E}[Y | X]$

实际估计： $\hat{f}(X) \approx \mathbb{E}[Y | X]$

预测误差分解：

$$\mathbb{E}[(Y - \hat{f}(X))^2] = \mathbb{E}[(Y - \mathbb{E}[Y | X])^2] \quad (20)$$

$$+ \mathbb{E}[(\mathbb{E}[Y | X] - \hat{f}(X))^2] \quad (21)$$

其中：

- 第一项：不可减少的误差（噪音）
- 第二项：可减少的误差（偏差 + 方差）

3.7.2 特征选择的信息价值

在高维回归中，选择哪些特征 X_j ？

信息价值标准： $\text{Var}(\mathbb{E}[Y | X_j])$

边际信息价值： $\text{Var}(\mathbb{E}[Y | X_1, \dots, X_j]) - \text{Var}(\mathbb{E}[Y | X_1, \dots, X_{j-1}])$

这解释了为什么：

- 相关特征应该被包含
- 冗余特征应该被排除
- 特征选择本质上是信息价值的优化

核心洞察

将 $E[X | Y]$ 视为随机变量从根本上改变了我们对信息和不确定性的理解：

1. **信息有价值**：因为它减少了不确定性
2. **信息是随机的**：在获得前无法确定其具体内容
3. **信息可以量化**：通过 $\text{Var}(E[X | Y])$
4. **信息有成本**：获取和处理信息需要资源
5. **信息影响决策**：改变了最优行为

这种视角统一了统计学、经济学、和信息论，为理解现代数据驱动的经济分析提供了深刻的理论基础。

4 概率不等式的应用示例

应用背景

在实际问题中，我们经常需要估计随机变量取值的概率，但并不知道其确切的分布。这时，概率不等式提供了有力的工具。以下是一个考试分数估计的例子：

4.1 问题设定

假设一次考试的分数 X 是一个随机变量，其总体参数为：

$$\mu = \mathbb{E}[X] = 75, \quad \sigma = \sqrt{\text{Var}(X)} = 5$$

重要说明

- $\mu = 75$ 是**总体均值** (population mean)，即所有可能考试分数的期望值
- $\sigma = 5$ 是**总体标准差** (population standard deviation)
- 概率不等式使用的是总体参数，而非样本统计量

我们需要估计以下概率：

- 有多少人分数高于 90 分？即 $P(X \geq 90)$
- 有多少人分数少于 60 分？即 $P(X \leq 60)$

4.2 使用 Markov 不等式

Markov 不等式适用于非负随机变量：对于 $X \geq 0$ 和 $a > 0$,

$$P(X \geq a) \leq \frac{\mathbb{E}[X]}{a}$$

为了估计 $P(X \geq 90)$ ，我们考虑非负随机变量 $Y = X - 60$ （因为分数通常不低于 0 分）：

$$\mathbb{E}[Y] = \mathbb{E}[X - 60] = 75 - 60 = 15$$

$$P(X \geq 90) = P(Y \geq 30) \leq \frac{\mathbb{E}[Y]}{30} = \frac{15}{30} = 0.5 = 50\%$$

这个估计比较宽松，利用的只是一阶矩信息。

4.3 使用切比雪夫不等式

切比雪夫不等式利用总体的均值和方差信息：对于任意随机变量 X 和 $k > 0$,

$$P(|X - \mu| \geq k\sigma) \leq \frac{1}{k^2}$$

关键理解

- 这里的 μ 和 σ 是**总体参数**，不是样本统计量
- 不等式对**任何分布**都成立，无需知道具体的分布形式
- 这是一个概率上界，实际概率可能比这个值小得多

对于 $P(X \geq 90)$ ：由于 $|90 - 75| = 15 = 3 \times 5 = 3\sigma$ ，所以 $k = 3$ ：

$$P(X \geq 90) \leq P(|X - 75| \geq 15) \leq \frac{1}{3^2} = \frac{1}{9} \approx 11.1\%$$

同样，对于 $P(X \leq 60)$ ：

$$P(X \leq 60) \leq P(|X - 75| \geq 15) \leq \frac{1}{9} \approx 11.1\%$$

4.4 假设正态分布

如果我们进一步假设分数服从正态分布 $X \sim N(75, 5^2)$ ，则可以得到更精确的估计：

$$P(X \geq 90) = P\left(\frac{X - 75}{5} \geq \frac{90 - 75}{5}\right) = P(Z \geq 3)$$

查标准正态分布表得： $P(Z \geq 3) \approx 0.0013 = 0.13\%$

$$P(X \leq 60) = P\left(\frac{X - 75}{5} \leq \frac{60 - 75}{5}\right) = P(Z \leq -3) \approx 0.0013 = 0.13\%$$

4.5 方法比较

可以看出，不同方法给出的估计差别很大：

方法	$P(X \geq 90)$	特点
Markov 不等式	$\leq 50\%$	最宽松，只需一阶矩
切比雪夫不等式	$\leq 11.1\%$	更紧，利用方差信息
正态分布假设	$\approx 0.13\%$	最精确，需知道分布

5 总体参数 vs 样本统计量的重要说明

重要理论说明

理解总体参数与样本统计量的区别对于正确应用概率不等式至关重要。

5.1 概率不等式中的参数

在概率不等式中：

- 切比雪夫不等式中的 μ 和 σ 是**总体参数**
- $\mu = \mathbb{E}[X]$ 是随机变量 X 的总体期望（理论均值）
- $\sigma = \sqrt{\text{Var}(X)}$ 是随机变量 X 的总体标准差
- 这些是分布的**固有特征**，不依赖于具体的样本

5.2 实际应用中的考虑

在实践中：

- 我们通常不知道真实的总体参数
- 使用样本统计量 \bar{X} 和 S 来估计总体参数 μ 和 σ
- 当样本量足够大时，样本统计量接近总体参数
- 在应用概率不等式时，我们**假设**已知总体参数或用样本统计量近似

5.3 符号对比

总体参数（理论值）	样本统计量（观测值）
$\mu = \mathbb{E}[X]$ （总体均值）	$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ （样本均值）
$\sigma^2 = \text{Var}(X)$ （总体方差）	$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$ （样本方差）
$\sigma = \sqrt{\text{Var}(X)}$ （总体标准差）	$S = \sqrt{S^2}$ （样本标准差）

6 协方差矩阵的正定性

协方差矩阵性质

证明协方差矩阵总是正半定的（PSD）。

定理 2 (协方差矩阵的正半定性). 对于多元随机变量 Z (即向量 Z 的每个分量都是随机变量), 协方差矩阵 Σ 定义为 $\Sigma_{ij} = \text{Cov}(Z_i, Z_j)$, 简记为 $\Sigma = \mathbb{E}[(Z - \mu)(Z - \mu)^T]$, 其中 μ 是 Z 的均值向量。则 Σ 总是正半定的。

证明

对于任意向量 $v \in \mathbb{R}^n$, 我们需要证明 $v^T \Sigma v \geq 0$ 。

$$v^T \Sigma v = v^T \mathbb{E}[(Z - \mu)(Z - \mu)^T] v \quad (22)$$

$$= \sum_{i=1}^n \sum_{j=1}^n \mathbb{E}[(Z_i - \mu_i)(Z_j - \mu_j)] v_i v_j \quad (23)$$

$$= \mathbb{E} [v^T (Z - \mu)(Z - \mu)^T v] \quad (24)$$

$$= \mathbb{E} [((Z - \mu)^T v)^2] \quad (25)$$

$$\geq 0 \quad (26)$$

最后一个不等式成立是因为期望值中的项是一个平方项, 必然非负。第三个等式利用了期望的线性性。

7 正半定矩阵的等价定义

PSD 矩阵的三个等价定义

设 $A \in \mathbb{R}^{n \times n}$ 是对称矩阵, 证明以下三个关于正半定性的定义是等价的:

1. 对于所有 $x \in \mathbb{R}^n$, $x^T A x \geq 0$
2. A 的所有特征值都非负
3. 存在矩阵 $U \in \mathbb{R}^{n \times n}$ 使得 $A = U U^T$

证明

我们证明 $(1) \Rightarrow (2) \Rightarrow (3) \Rightarrow (1)$ 的循环。

$(1) \Rightarrow (2)$: 设 λ 是 A 的特征值, v 是对应的特征向量。则: $v^T A v = \lambda v^T v = \lambda \|v\|^2$
由条件 (1), $\lambda \|v\|^2 \geq 0$, 因此 $\lambda \geq 0$ 。

$(2) \Rightarrow (3)$: 考虑 A 的特征分解 $A = V \Lambda V^T$, 其中 Λ 是对角矩阵, 对角线元素为特征值 $\lambda_1, \dots, \lambda_n$ 。定义 $U := V \sqrt{\Lambda}$, 其中 $\sqrt{\Lambda}$ 的对角线元素为 $\sqrt{\lambda_1}, \dots, \sqrt{\lambda_n}$ (由于特征值非负, 平方根有定义)。显然 $A = U U^T$ 。

$(3) \Rightarrow (1)$: 设 $x \in \mathbb{R}^n$, 则: $x^T A x = x^T U U^T x = (U^T x)^T (U^T x) = \|U^T x\|^2 \geq 0$

8 熵的性质

信息熵的凹性

设 X 服从伯努利分布 $X \sim \text{Bernoulli}(p)$, $p \in (0, 1)$ 。证明熵函数 $H(X) = -p \ln p - (1-p) \ln(1-p)$ 关于 p 是凹函数。

证明

为证明 $H(X)$ 是凹函数, 我们证明其二阶导数在 $(0, 1)$ 上非正。

计算一阶导数: $H'(X) = -(\ln p + 1) + (\ln(1-p) + 1) = \ln(1-p) - \ln p$

计算二阶导数: $H''(X) = -\frac{1}{1-p} - \frac{1}{p}$

对于所有 $p \in (0, 1)$, 显然 $H''(X) < 0$, 因此 $H(X)$ 是严格凹函数。

熵的最大化

考虑有 n 个状态的离散分布。证明在所有可能的概率质量函数中, 均匀分布使熵 $H(X)$ 达到最大值。此时 $H(X) = \ln n$ 。

使用拉格朗日乘数法

设概率为 p_1, \dots, p_n , 优化问题为: $\max -\sum_{i=1}^n p_i \ln p_i \quad \text{s.t.} \quad \sum_{i=1}^n p_i = 1$

拉格朗日函数为: $L(p, \lambda) = \sum_{i=1}^n (-p_i \ln p_i) + \lambda (\sum_{i=1}^n p_i - 1)$

一阶条件: $\frac{\partial L}{\partial p_i} = -\ln p_i - 1 + \lambda = 0$

这给出 $\ln p_i = -1 + \lambda$, 即 $p_i = e^{-1+\lambda}$, 对所有 i 都是常数。

利用约束条件 $\sum_{i=1}^n p_i = 1$ 得到 $p_i = \frac{1}{n}$, 即均匀分布。

此时熵为: $H(X) = -n \cdot \frac{1}{n} \ln \frac{1}{n} = \ln n$ 。

9 矩阵范数与内积

Frobenius 内积

Frobenius 内积定义为：对于同维矩阵 $A, B \in \mathbb{R}^{m \times n}$, $\langle A, B \rangle = \text{trace}(A^T B) = \sum_{i=1}^m \sum_{j=1}^n A_{ij} B_{ij}$

证明： $x^T A y = \langle A, x y^T \rangle$, 其中 $x \in \mathbb{R}^m$, $y \in \mathbb{R}^n$, $A \in \mathbb{R}^{m \times n}$ 。

证明

$$x^T A y = \sum_{i=1}^m \sum_{j=1}^n x_i A_{ij} y_j \quad (27)$$

$$= \sum_{i=1}^m \sum_{j=1}^n A_{ij} (x y^T)_{ij} \quad (28)$$

$$= \langle A, x y^T \rangle \quad (29)$$

其中 $(x y^T)_{ij} = x_i y_j$ 。

10 正态分布的性质

正态分布的矩母函数

证明：若 $X \sim N(0, \sigma^2)$, 则 $\mathbb{E}[e^{\lambda X}] = e^{\sigma^2 \lambda^2 / 2}$ 。

证明

$$\mathbb{E}[e^{\lambda X}] = \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\infty} e^{\lambda x} e^{-x^2/2\sigma^2} dx \quad (30)$$

$$= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{\lambda \sigma z} e^{-z^2/2} dz \quad (z = x/\sigma) \quad (31)$$

$$= e^{\sigma^2 \lambda^2 / 2} \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-(z - \lambda \sigma)^2 / 2} dz \quad (32)$$

$$= e^{\sigma^2 \lambda^2 / 2} \quad (33)$$

其中第三行利用了配方： $\lambda \sigma z - z^2/2 = -\frac{1}{2}(z - \lambda \sigma)^2 + \frac{\sigma^2 \lambda^2}{2}$ 。

独立性与不相关性

设 $u, v \in \mathbb{R}^n$ 是常向量且正交（即 $\langle u, v \rangle = 0$ ）， $X = (X_1, \dots, X_n)$ 是 n 个独立同分布的标准正态随机变量。定义 $u_x = \langle u, X \rangle$ 和 $v_x = \langle v, X \rangle$ 。证明 u_x 和 v_x 独立。

证明

由于 u_x 和 v_x 是正态分布随机向量的线性变换，它们是联合正态分布。对于联合正态随机变量，独立性等价于不相关性。

计算协方差：

$$\text{Cov}(u_x, v_x) = \mathbb{E}[u_x v_x] = \mathbb{E} \left[\left(\sum_{i=1}^n u_i X_i \right) \left(\sum_{j=1}^n v_j X_j \right) \right] \quad (34)$$

$$= \sum_{i=1}^n u_i v_i \mathbb{E}[X_i^2] = \sum_{i=1}^n u_i v_i = \langle u, v \rangle = 0 \quad (35)$$

因此 u_x 和 v_x 不相关，从而独立。

11 计量经济学中的条件期望：零均值假设的重要性

核心问题

在计量经济学中，为什么条件期望假设 $\mathbb{E}[\varepsilon_i | X_i] = 0$ 比无条件期望假设 $\mathbb{E}[X_i' \varepsilon_i] = 0$ 更强且更重要？

11.1 两种假设的定义与关系

11.1.1 假设对比

1. 弱假设（正交性条件）： $\mathbb{E}[X_i' \varepsilon_i] = 0$

- 仅要求解释变量与误差项的线性关系平均为零
- 允许存在非线性关系

2. 强假设（条件期望为零）： $\mathbb{E}[\varepsilon_i | X_i] = 0$

- 要求在给定任何 X_i 值的条件下，误差项的期望都为零
- 排除了 X_i 的任何函数与 ε_i 的关联

11.1.2 逻辑关系

根据迭代期望公式: $\mathbb{E}[X_i'\varepsilon_i] = \mathbb{E}[\mathbb{E}[X_i'\varepsilon_i | X_i]] = \mathbb{E}[X_i'\mathbb{E}[\varepsilon_i | X_i]]$

因此: $\mathbb{E}[\varepsilon_i | X_i] = 0 \Rightarrow \mathbb{E}[X_i'\varepsilon_i] = 0$

但反之不成立!

11.2 经典反例: 非线性关系

例 1 (均匀分布下的二次关系). 设 X_i 服从 $[-1, 1]$ 上的均匀分布, 定义误差项为: $\varepsilon_i = X_i^2 - \frac{1}{3} + u_i$ 其中 u_i 是满足 $\mathbb{E}[u_i] = 0$ 且与 X_i 独立的随机扰动项。

11.2.1 验证 $\mathbb{E}[X_i\varepsilon_i] = 0$

$$\mathbb{E}[X_i\varepsilon_i] = \mathbb{E}\left[X_i\left(X_i^2 - \frac{1}{3} + u_i\right)\right] \quad (36)$$

$$= \mathbb{E}[X_i^3] - \frac{1}{3}\mathbb{E}[X_i] + \mathbb{E}[X_i u_i] \quad (37)$$

$$= \mathbb{E}[X_i^3] - \frac{1}{3} \cdot 0 + 0 \quad (38)$$

$$= \mathbb{E}[X_i^3] \quad (39)$$

由于 $X_i \sim \text{Uniform}[-1, 1]$, X_i^3 是奇函数, 因此: $\mathbb{E}[X_i^3] = \int_{-1}^1 x^3 \cdot \frac{1}{2} dx = \frac{1}{2} \left[\frac{x^4}{4} \right]_{-1}^1 = \frac{1}{2} \cdot 0 = 0$

所以 $\mathbb{E}[X_i\varepsilon_i] = 0$ 成立。

11.2.2 验证 $\mathbb{E}[\varepsilon_i | X_i] \neq 0$

$$\mathbb{E}[\varepsilon_i | X_i = x] = \mathbb{E}\left[x^2 - \frac{1}{3} + u_i | X_i = x\right] \quad (40)$$

$$= x^2 - \frac{1}{3} + \mathbb{E}[u_i | X_i = x] \quad (41)$$

$$= x^2 - \frac{1}{3} + \mathbb{E}[u_i] \quad (42)$$

$$= x^2 - \frac{1}{3} \quad (43)$$

显然, 当 $x \neq \pm\sqrt{1/3}$ 时, $\mathbb{E}[\varepsilon_i | X_i = x] \neq 0$ 。

关键洞察

这个例子说明：即使 $\mathbb{E}[X_i \varepsilon_i] = 0$ ，误差项 ε_i 仍然可以与解释变量 X_i 存在**系统性**的非线性关系。

11.3 条件期望为零的重要性

11.3.1 排除所有函数关系

条件期望假设 $\mathbb{E}[\varepsilon_i | X_i] = 0$ 意味着：

对于 X_i 的任何可测函数 $g(X_i)$ ，都有： $\mathbb{E}[g(X_i)\varepsilon_i] = \mathbb{E}[\mathbb{E}[g(X_i)\varepsilon_i | X_i]] = \mathbb{E}[g(X_i)\mathbb{E}[\varepsilon_i | X_i]] = 0$

这包括：

- 线性关系： $\mathbb{E}[X_i \varepsilon_i] = 0$
- 二次关系： $\mathbb{E}[X_i^2 \varepsilon_i] = 0$
- 高次多项式： $\mathbb{E}[X_i^k \varepsilon_i] = 0, \forall k$
- 任意非线性函数： $\mathbb{E}[\sin(X_i)\varepsilon_i] = 0, \mathbb{E}[e^{X_i}\varepsilon_i] = 0$ 等

11.3.2 OLS 估计量的性质

在线性回归模型 $Y_i = X_i' \beta + \varepsilon_i$ 中：

1. 仅有 $\mathbb{E}[X_i \varepsilon_i] = 0$ 时：
 - OLS 估计量 $\hat{\beta}$ 是无偏的
 - 但 $\hat{\varepsilon}_i^2$ 不是 $\sigma^2 = \text{Var}(\varepsilon_i)$ 的无偏估计
2. 有 $\mathbb{E}[\varepsilon_i | X_i] = 0$ 时：
 - OLS 估计量 $\hat{\beta}$ 是无偏的
 - 残差平方和正确地估计了误差方差
 - 标准误估计是正确的

11.4 误差方差估计的偏误

11.4.1 问题的根源

当 $\mathbb{E}[\varepsilon_i | X_i] \neq 0$ 时，设 $\mathbb{E}[\varepsilon_i | X_i] = m(X_i)$ ，则： $\varepsilon_i = m(X_i) + v_i$ 其中 $\mathbb{E}[v_i | X_i] = 0$ 。

OLS 残差实际上估计的是： $\hat{\varepsilon}_i^2 \approx [m(X_i) + v_i]^2 = m(X_i)^2 + 2m(X_i)v_i + v_i^2$

11.4.2 方差估计的偏误

$$\mathbb{E}[\varepsilon_i^2] \approx \mathbb{E}[m(X_i)^2] + 2\mathbb{E}[m(X_i)v_i] + \mathbb{E}[v_i^2] \quad (44)$$

$$= \mathbb{E}[m(X_i)^2] + 0 + \text{Var}(v_i) \quad (45)$$

$$= \mathbb{E}[m(X_i)^2] + \text{Var}(v_i) \quad (46)$$

$$> \text{Var}(v_i) = \text{Var}(\varepsilon_i | X_i) \quad (47)$$

高估问题

当 $\mathbb{E}[\varepsilon_i | X_i] \neq 0$ 时, OLS 的残差平方和**系统性地高估**了真实的误差方差, 因为它包含了:

1. 真实的随机误差方差: $\text{Var}(v_i)$
2. 模型设定误差的贡献: $\mathbb{E}[m(X_i)^2]$

11.5 数值例子验证

回到我们的例子: $\varepsilon_i = X_i^2 - \frac{1}{3} + u_i$, 其中 $X_i \sim \text{Uniform}[-1, 1]$, $u_i \sim N(0, \sigma_u^2)$ 。

11.5.1 真实误差方差

$$\text{Var}(\varepsilon_i) = \text{Var}(X_i^2 - \frac{1}{3}) + \text{Var}(u_i) = \text{Var}(X_i^2) + \sigma_u^2$$

对于 $X_i \sim \text{Uniform}[-1, 1]$:

$$\mathbb{E}[X_i^2] = \frac{1}{3} \quad (48)$$

$$\mathbb{E}[X_i^4] = \int_{-1}^1 x^4 \cdot \frac{1}{2} dx = \frac{1}{5} \quad (49)$$

$$\text{Var}(X_i^2) = \mathbb{E}[X_i^4] - (\mathbb{E}[X_i^2])^2 = \frac{1}{5} - \frac{1}{9} = \frac{4}{45} \quad (50)$$

$$\text{因此: } \text{Var}(\varepsilon_i) = \frac{4}{45} + \sigma_u^2$$

11.5.2 OLS 估计的方差

当我们用 OLS 拟合 $Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$, 但真实模型包含 X_i^2 项时, OLS 残差会包含未被捕获的 X_i^2 项的变异, 导致高估真实的误差方差。

11.6 计量经济学实践中的后果

严重后果

如果不满足 $E[\varepsilon_i | X_i] = 0$ ，会导致：

1. 模型设定错误：
 - 遗漏了重要的非线性项
 - 函数形式设定不当
2. 统计推断失效：
 - 标准误估计有偏
 - 置信区间不准确
 - 假设检验的 p 值不可靠
3. 预测精度下降：
 - 系统性预测偏误
 - 预测区间过宽或过窄
4. 政策建议错误：
 - 对因果效应的错误估计
 - 边际效应计算错误

11.7 检验与补救措施

11.7.1 检验方法

1. Ramsey RESET 检验：检验是否遗漏高次项
2. 残差图分析：绘制 $\hat{\varepsilon}_i$ 对 X_i 的散点图
3. 非参数回归：使用核回归等方法检验非线性关系

11.7.2 补救措施

1. 添加非线性项：包含 X_i^2, X_i^3 等
2. 函数形式变换：对数变换、Box-Cox 变换等

3. **非参数方法**：使用更灵活的估计方法
4. **工具变量**：当内生性严重时使用 IV 估计

核心要点

计量经济学中必须强调 $\mathbb{E}[\varepsilon_i | X_i] = 0$ 这一强假设，因为它确保了：

1. 模型正确设定
2. 统计推断有效
3. 因果解释合理
4. 预测结果可靠

仅仅满足 $\mathbb{E}[X_i \varepsilon_i] = 0$ 是远远不够的！

核心问题

在计量经济学中，为什么条件期望假设 $\mathbb{E}[\varepsilon_i | X_i] = 0$ 比无条件期望假设 $\mathbb{E}[X_i' \varepsilon_i] = 0$ 更强且更重要？

11.8 两种假设的定义与关系

11.8.1 假设对比

1. **弱假设（正交性条件）**： $\mathbb{E}[X_i' \varepsilon_i] = 0$
 - 仅要求解释变量与误差项的线性关系平均为零
 - 允许存在非线性关系
2. **强假设（条件期望为零）**： $\mathbb{E}[\varepsilon_i | X_i] = 0$
 - 要求在给定任何 X_i 值的条件下，误差项的期望都为零
 - 排除了 X_i 的任何函数与 ε_i 的关联

11.8.2 逻辑关系

根据迭代期望公式： $\mathbb{E}[X_i' \varepsilon_i] = \mathbb{E}[\mathbb{E}[X_i' \varepsilon_i | X_i]] = \mathbb{E}[X_i' \mathbb{E}[\varepsilon_i | X_i]]$

因此： $\mathbb{E}[\varepsilon_i | X_i] = 0 \Rightarrow \mathbb{E}[X_i' \varepsilon_i] = 0$

但反之不成立！

11.9 经典反例：非线性关系

例 2 (均匀分布下的二次关系). 设 X_i 服从 $[-1, 1]$ 上的均匀分布, 定义误差项为: $\varepsilon_i = X_i^2 - \frac{1}{3} + u_i$ 其中 u_i 是满足 $\mathbb{E}[u_i] = 0$ 且与 X_i 独立的随机扰动项。

11.9.1 验证 $\mathbb{E}[X_i \varepsilon_i] = 0$

$$\mathbb{E}[X_i \varepsilon_i] = \mathbb{E} \left[X_i \left(X_i^2 - \frac{1}{3} + u_i \right) \right] \quad (51)$$

$$= \mathbb{E}[X_i^3] - \frac{1}{3} \mathbb{E}[X_i] + \mathbb{E}[X_i u_i] \quad (52)$$

$$= \mathbb{E}[X_i^3] - \frac{1}{3} \cdot 0 + 0 \quad (53)$$

$$= \mathbb{E}[X_i^3] \quad (54)$$

由于 $X_i \sim \text{Uniform}[-1, 1]$, X_i^3 是奇函数, 因此: $\mathbb{E}[X_i^3] = \int_{-1}^1 x^3 \cdot \frac{1}{2} dx = \frac{1}{2} \left[\frac{x^4}{4} \right]_{-1}^1 = \frac{1}{2} \cdot 0 = 0$

所以 $\mathbb{E}[X_i \varepsilon_i] = 0$ 成立。

11.9.2 验证 $\mathbb{E}[\varepsilon_i | X_i] \neq 0$

$$\mathbb{E}[\varepsilon_i | X_i = x] = \mathbb{E} \left[x^2 - \frac{1}{3} + u_i | X_i = x \right] \quad (55)$$

$$= x^2 - \frac{1}{3} + \mathbb{E}[u_i | X_i = x] \quad (56)$$

$$= x^2 - \frac{1}{3} + \mathbb{E}[u_i] \quad (57)$$

$$= x^2 - \frac{1}{3} \quad (58)$$

显然, 当 $x \neq \pm\sqrt{1/3}$ 时, $\mathbb{E}[\varepsilon_i | X_i = x] \neq 0$ 。

关键洞察

这个例子说明: 即使 $\mathbb{E}[X_i \varepsilon_i] = 0$, 误差项 ε_i 仍然可以与解释变量 X_i 存在**系统性**的非线性关系。

11.10 条件期望为零的重要性

11.10.1 排除所有函数关系

条件期望假设 $\mathbb{E}[\varepsilon_i | X_i] = 0$ 意味着:

对于 X_i 的任何可测函数 $g(X_i)$, 都有: $\mathbb{E}[g(X_i)\varepsilon_i] = \mathbb{E}[\mathbb{E}[g(X_i)\varepsilon_i | X_i]] = \mathbb{E}[g(X_i)\mathbb{E}[\varepsilon_i | X_i]] = 0$

这包括:

- 线性关系: $\mathbb{E}[X_i\varepsilon_i] = 0$
- 二次关系: $\mathbb{E}[X_i^2\varepsilon_i] = 0$
- 高次多项式: $\mathbb{E}[X_i^k\varepsilon_i] = 0, \forall k$
- 任意非线性函数: $\mathbb{E}[\sin(X_i)\varepsilon_i] = 0, \mathbb{E}[e^{X_i}\varepsilon_i] = 0$ 等

11.10.2 OLS 估计量的性质

在线性回归模型 $Y_i = X_i'\beta + \varepsilon_i$ 中:

1. 仅有 $\mathbb{E}[X_i\varepsilon_i] = 0$ 时:
 - OLS 估计量 $\hat{\beta}$ 是无偏的
 - 但 $\hat{\varepsilon}_i^2$ 不是 $\sigma^2 = \text{Var}(\varepsilon_i)$ 的无偏估计
2. 有 $\mathbb{E}[\varepsilon_i | X_i] = 0$ 时:
 - OLS 估计量 $\hat{\beta}$ 是无偏的
 - 残差平方和正确地估计了误差方差
 - 标准误估计是正确的

11.11 误差方差估计的偏误

11.11.1 问题的根源

当 $\mathbb{E}[\varepsilon_i | X_i] \neq 0$ 时, 设 $\mathbb{E}[\varepsilon_i | X_i] = m(X_i)$, 则: $\varepsilon_i = m(X_i) + v_i$ 其中 $\mathbb{E}[v_i | X_i] = 0$ 。
OLS 残差实际上估计的是: $\hat{\varepsilon}_i^2 \approx [m(X_i) + v_i]^2 = m(X_i)^2 + 2m(X_i)v_i + v_i^2$

11.11.2 方差估计的偏误

$$\mathbb{E}[\hat{\varepsilon}_i^2] \approx \mathbb{E}[m(X_i)^2] + 2\mathbb{E}[m(X_i)v_i] + \mathbb{E}[v_i^2] \quad (59)$$

$$= \mathbb{E}[m(X_i)^2] + 0 + \text{Var}(v_i) \quad (60)$$

$$= \mathbb{E}[m(X_i)^2] + \text{Var}(v_i) \quad (61)$$

$$> \text{Var}(v_i) = \text{Var}(\varepsilon_i | X_i) \quad (62)$$

高估问题

当 $\mathbb{E}[\varepsilon_i | X_i] \neq 0$ 时，OLS 的残差平方和**系统性地高估**了真实的误差方差，因为它包含了：

1. 真实的随机误差方差： $\text{Var}(v_i)$
2. 模型设定误差的贡献： $\mathbb{E}[m(X_i)^2]$

11.12 数值例子验证

回到我们的例子： $\varepsilon_i = X_i^2 - \frac{1}{3} + u_i$ ，其中 $X_i \sim \text{Uniform}[-1, 1]$ ， $u_i \sim N(0, \sigma_u^2)$ 。

11.12.1 真实误差方差

$$\text{Var}(\varepsilon_i) = \text{Var}(X_i^2 - \frac{1}{3}) + \text{Var}(u_i) = \text{Var}(X_i^2) + \sigma_u^2$$

对于 $X_i \sim \text{Uniform}[-1, 1]$ ：

$$\mathbb{E}[X_i^2] = \frac{1}{3} \quad (63)$$

$$\mathbb{E}[X_i^4] = \int_{-1}^1 x^4 \cdot \frac{1}{2} dx = \frac{1}{5} \quad (64)$$

$$\text{Var}(X_i^2) = \mathbb{E}[X_i^4] - (\mathbb{E}[X_i^2])^2 = \frac{1}{5} - \frac{1}{9} = \frac{4}{45} \quad (65)$$

$$\text{因此：} \text{Var}(\varepsilon_i) = \frac{4}{45} + \sigma_u^2$$

11.12.2 OLS 估计的方差

OLS 会错误地估计方差为： $\mathbb{E}[\hat{\varepsilon}_i^2] \approx \mathbb{E}[(X_i^2 - \frac{1}{3})^2] + \sigma_u^2 = \frac{4}{45} + \sigma_u^2$

等等，这里似乎一样？让我们更仔细地分析...

实际上，当我们用 OLS 拟合 $Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$ ，但真实模型包含 X_i^2 项时，OLS 残差会包含未被捕获的 X_i^2 项的变异，导致高估。

11.13 计量经济学实践中的后果

严重后果

如果不满足 $\mathbb{E}[\varepsilon_i | X_i] = 0$ ，会导致：

1. 模型设定错误：
 - 遗漏了重要的非线性项

- 函数形式设定不当

2. 统计推断失效:

- 标准误估计有偏
- 置信区间不准确
- 假设检验的 p 值不可靠

3. 预测精度下降:

- 系统性预测偏误
- 预测区间过宽或过窄

4. 政策建议错误:

- 对因果效应的错误估计
- 边际效应计算错误

11.14 检验与补救措施

11.14.1 检验方法

1. Ramsey RESET 检验: 检验是否遗漏高次项
2. 残差图分析: 绘制 $\hat{\varepsilon}_i$ 对 X_i 的散点图
3. 非参数回归: 使用核回归等方法检验非线性关系

11.14.2 补救措施

1. 添加非线性项: 包含 X_i^2, X_i^3 等
2. 函数形式变换: 对数变换、Box-Cox 变换等
3. 非参数方法: 使用更灵活的估计方法
4. 工具变量: 当内生性严重时使用 IV 估计

核心要点

计量经济学中必须强调 $\mathbb{E}[\varepsilon_i | X_i] = 0$ 这一强假设，因为它确保了：

1. 模型正确设定
2. 统计推断有效
3. 因果解释合理
4. 预测结果可靠

仅仅满足 $\mathbb{E}[X_i \varepsilon_i] = 0$ 是远远不够的！

12 结论

1. 迭代期望公式在满足相应数学条件时严格成立，是概率论的基本定理之一。
2. 当涉及零概率事件时，直觉的条件期望定义可能失效，需要借助测度论的严格框架。
3. 概率不等式在不知道确切分布时提供了有用的概率上界：
 - Markov 不等式最宽松，只需要一阶矩信息
 - 切比雪夫不等式更紧，利用了方差信息
 - 当知道确切分布时（如正态分布），可以得到最精确的结果
4. 在实际应用中，选择哪种方法取决于我们对问题的了解程度和所需的精度。
5. 理解总体参数与样本统计量的区别对于正确应用统计方法至关重要。
6. **在计量经济学中**，条件期望假设 $\mathbb{E}[\varepsilon_i | X_i] = 0$ 远比正交性条件 $\mathbb{E}[X_i \varepsilon_i] = 0$ 更重要，因为它确保了模型的正确设定和统计推断的有效性。