

计量经济学STAT30021

第三讲：简单线性回归分析(2)

肖志国

复旦大学管理学院

2025年9月

简单线性回归模型

基本框架: $\{(X_i, Y_i), i = 1, \dots, n\}$ 为数据, 且满足: 对于所有的 $i = 1, \dots, n$,

$$Y_i = \beta_0 + \beta_1 X_i + u_i, \quad (1)$$

以及

$$E[u_i | X_i] = 0. \quad (2)$$

- 参数 β_0, β_1 的估计:

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}, \quad \hat{\beta}_1 = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

- 最小二乘法的性质
- $\hat{\beta}_1$ 与相关系数的关系
- 模型的解释能力: R^2
- R^2 与相关系数的关系

回归模型的应用

考虑大学成绩(colgpa)与高中成绩(hsgpa)的关系。通过一组30个学生的数据，我们得到如下的回归方程：

$$\text{colgpa} = 1.41 + 0.52 \times \text{hsgpa}$$

我们可能关心的问题有如下几类：

- ① Y的变化对X的变化的敏感度：某学生A的高中成绩比学生B的高中成绩多0.5个绩点，那么A的大学成绩比B高出多少？
- ② 给定X的取值，预测Y的取值：如果某学生的高中成绩为3.5，那么其大学成绩为多少？
- ③ 模型的拟合程度：高中成绩到底能解释多少的大学成绩？

回归系数的大小

是不是回归系数 $\hat{\beta}_1$ 的值越大，就表示 X 对 Y 的影响就越大？

- $\hat{\beta}_1$: 当 X 增加一个单位时， Y 平均的变化量；或者说，是 Y 对 X 的敏感度。
- $\hat{\beta}_0$: 当 X 为零时， Y 的平均值？
- 严格说来，回归分析的结果只适用于当前数据。也就是说，他们对于新的数据来说未必成立。
- 但是我们通常假定我们通过当前数据 $\{(X_i, Y_i), i = 1, \dots, n\}$ 得到的结果能够近似的刻画 X 与 Y 的关系。
- 实际问题中， X 与 Y 的线性回归关系可能只对于特定范围取值的 (X, Y) 成立(我们常常忽视这一点)。
- 不能过度解释回归分析的结果。要对回归分析结果成立的范围有个大致的概念。

Y的测量单位的变化对回归结果的影响

假设将所有的 Y_i 都乘以一个常数 λ ，用新的 Y 来和原来的 X 做回归。回归系数会有什么变化？ R^2 呢？记新的结果为 $\hat{\beta}_{0\lambda}$, $\hat{\beta}_{1\lambda}$ 以及 R_λ^2 。

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

$$\hat{\beta}_{0\lambda} = \lambda \bar{Y} - \hat{\beta}_{1\lambda} \bar{X} = \lambda \bar{Y} - \lambda \hat{\beta}_1 \bar{X} = \lambda \hat{\beta}_0$$

$$R_\lambda^2 = r_{X(\lambda Y)}^2 = \frac{\lambda^2 \sum_{i=1}^n [(X_i - \bar{X})(Y_i - \bar{Y})]^2}{\sum_{i=1}^n (X_i - \bar{X})^2 \lambda^2 \sum_{i=1}^n (Y_i - \bar{Y})^2} = R^2$$

X的测量单位的变化对回归结果的影响

假设将所有的 X_i 都乘以一个常数 γ ，用新的 X 来和原来的 Y 做回归。回归系数会有什么变化？ R^2 呢？记新的结果为 $\hat{\beta}_{0\gamma}, \hat{\beta}_{1\gamma}$ 以及 R_γ^2 。

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

$$\hat{\beta}_{0\gamma} = \bar{Y} - \hat{\beta}_{1\gamma}(\gamma \bar{X}) = \bar{Y} - \hat{\beta}_1 \bar{X} = \hat{\beta}_0$$

$$R_\gamma^2 = r_{(\gamma X)Y}^2 = \frac{\gamma^2 \sum_{i=1}^n [(X_i - \bar{X})(Y_i - \bar{Y})]^2}{\gamma^2 \sum_{i=1}^n (X_i - \bar{X})^2 \sum_{i=1}^n (Y_i - \bar{Y})^2} = R^2$$

取log还是不取log

基本模型是

$$Y_i = \beta_0 + \beta_1 X_i + u_i,$$

这个模型描述的是X的变化对Y的变化的影响。出于以下一些原因，我们可能需要对X，或者Y，或者同时做一个非线性变换，比如，取log：

- 原模型的拟合程度不好，我们希望通过变换增加拟合程度
- 我们关心的是X的变化对Y的百分比变化的影响
- 我们关心的是X的百分比变化对Y的变化的影响
- 我们关心的是X的百分比变化对Y的百分比变化的影响，也就是说，Y对X的弹性

取log的方式

我们可以有三种方式取log，每一种方式代表一种不同的模型，从而其参数的解释也完全不同：

① 对Y取log:

$$\log(Y_i) = \beta_0 + \beta_1 X_i + u_i,$$

此模型中 β_1 表示X的单位变化对Y的百分比变化的影响

② 对X取log:

$$Y_i = \beta_0 + \beta_1 \log(X_i) + u_i,$$

此模型中 β_1 表示X的百分比变化对Y的单位变化的影响

③ 对X和Y都取log:

$$\log(Y_i) = \beta_0 + \beta_1 \log(X_i) + u_i,$$

此模型又称为一个常弹性模型，其中 β_1 表示Y对X的弹性

百分比变化

某变量 y 增长了10%，则我们称其百分比变化为+10，记作 $\% \Delta y = 10$ ；

某变量 y 减少了20%，则我们称其百分比变化为-20，记作 $\% \Delta y = -20$ 。

依定义，

$$\% \Delta y = 100 \frac{\Delta y}{y}$$

取与不取log：实际运用的原则

- log变化的效应：将大的值缩小，将小的值放大

$$\log(100) = 4.6, \log(0.02) = -3.9$$

- 一般原则：相关变量必须全部都是正数时才能取log
- 通常情况下，如果某变量的所有值都是正数，而且其跨度很大时(也就是说，方差很大时)，取log可能是一个合适的选择
- 通常会取log的变量：GDP，工资，销售额，等(我们更关心这些变量的增长率而不是其绝对数值)
- 如果某变量既有正值，又有负值，而我们又非得考虑对于它的弹性怎么办？回归原始定义，并注意可能需要去掉它的极端值：

$$\frac{\Delta Y_i}{Y_i} = \beta_0 + \beta_1 \frac{\Delta X_i}{X_i} + u_i,$$

关于模型设定

模型设定是计量经济分析的首要问题，也是一个非常难的问题。

对于模型设定，一般来说，我们无法找到真实的模型是什么。我们只能尽可能找到和真实模型接近的模型。

一个一般形式的可加模型形如：

$$\psi(Y) = \phi(X, \beta) + u, \quad E[u|X] = 0, \quad (3)$$

其中 $\psi()$, $\phi()$ 为给定函数。

当 $\psi(Y) = Y$ ， $\phi()$ 既是 X 的线性函数，又是 β 的线性函数时，这就是最常见的线性模型：

$$Y_i = \beta_0 + \beta_1 X_i + u_i.$$

关于模型设定

一个稍一般的情形是， $\phi()$ 是 X 的非线性函数，但它是 β 的线性函数。这种时候我们仍然叫线性模型，如：

$$\log(Y_i) = \beta_0 + \beta_1 X_i + u_i,$$

$$Y_i = \beta_0 + \beta_1 \log(X_i) + u_i,$$

$$\log(Y_i) = \beta_0 + \beta_1 \log(X_i) + u_i,$$

$$\sqrt{Y_i} = \beta_0 + \beta_1 \frac{1}{X_i} + u_i,$$

如果 $\phi()$ 是 β 的非线性函数，那这种模型就叫非线性模型，如：

$$Y_i = \frac{1}{\beta_0 + \beta_1 X_i} + u_i,$$

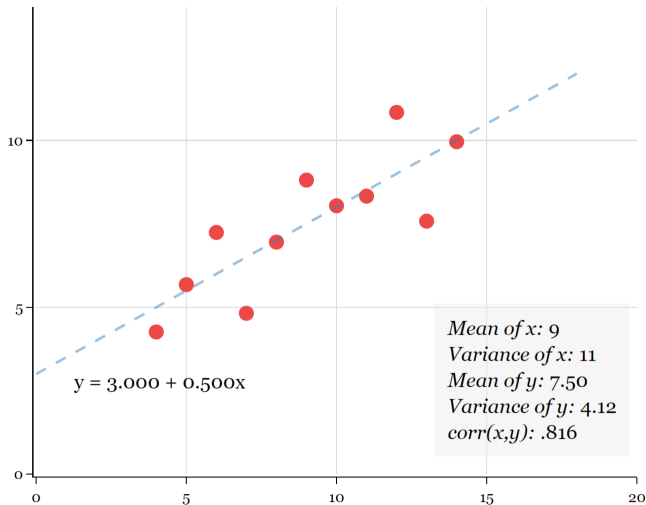
$$\log(Y_i) = (\beta_0 + \beta_1 X_i)^2 + u_i,$$

Anscombe 的回归例子

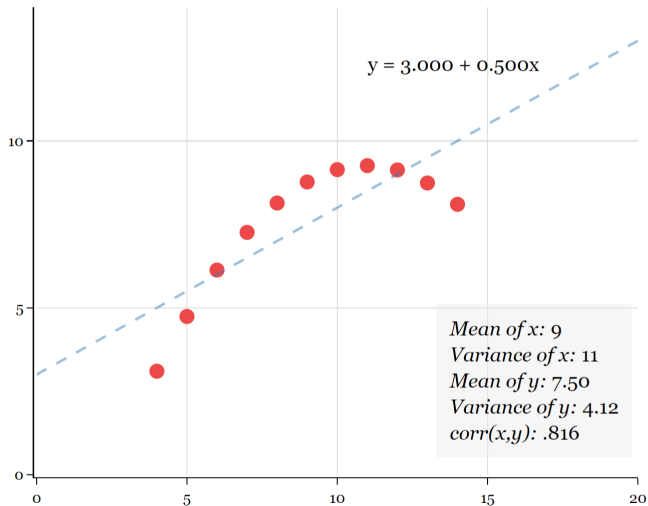
Anscombe (1973)给出了下面的一组回归数据：

X1	Y1	X2	Y2	X3	Y3	X4	Y4
10	8.04	10	9.14	10	7.46	8	6.58
8	6.95	8	8.14	8	6.77	8	5.76
13	7.58	13	8.74	13	12.74	8	7.71
9	8.81	9	8.77	9	7.11	8	8.84
11	8.33	11	9.26	11	7.81	8	8.47
14	9.96	14	8.1	14	8.84	8	7.04
6	7.24	6	6.13	6	6.08	8	5.25
4	4.26	4	3.1	4	5.39	19	12.5
12	10.84	12	9.13	12	8.15	8	5.56
7	4.82	7	7.26	7	6.42	8	7.91
5	5.68	5	4.74	5	5.73	8	6.89

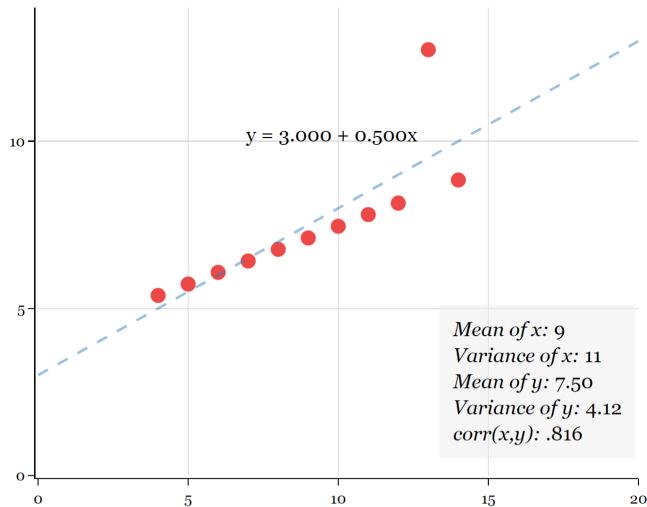
Anscombe 的回归例子: $Y_1 \sim X_1$



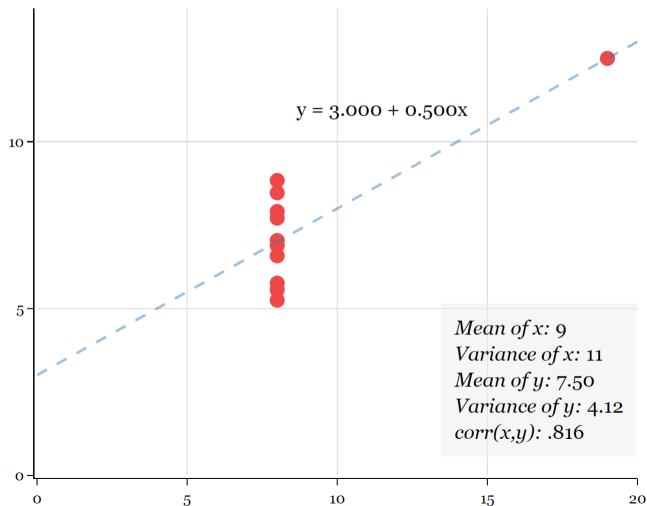
Anscombe 的回归例子: $Y_2 \sim X_2$



Anscombe 的回归例子: $Y_3 \sim X_3$



Anscombe 的回归例子: $Y_4 \sim X_4$



Anscombe 的回归例子: 讨论

四组不同的数据, 得到的回归结果, 包括各种统计量, 都是一模一样的!

Number of observations (n) = 11

Mean of the x 's (\bar{x}) = 9.0

Mean of the y 's (\bar{y}) = 7.5

Regression coefficient (b_1) of y on x = 0.5

Equation of regression line: $y = 3 + 0.5 x$

Sum of squares of $x - \bar{x}$ = 110.0

Regression sum of squares = 27.50 (1 d.f.)

Residual sum of squares of y = 13.75 (9 d.f.)

Estimated standard error of b_1 = 0.118

Multiple R^2 = 0.667

Anscombe 的回归例子: 讨论

四组不同的数据, 得到的回归结果, 包括各种统计量, 都是一模一样的!

系数 ^a							
by数据小兵			未标准化系数		标准化系数	t	显著性
group	模型		B	标准错误	Beta		
1	1	(常量)	3.000	1.125		2.667	0.026
		x	0.500	0.118	0.816	4.241	0.002
2	1	(常量)	3.001	1.125		2.667	0.026
		x	0.500	0.118	0.816	4.239	0.002
3	1	(常量)	3.002	1.124		2.670	0.026
		x	0.500	0.118	0.816	4.239	0.002
4	1	(常量)	3.002	1.124		2.671	0.026
		x	0.500	0.118	0.817	4.243	0.002

a. 因变量: y

知乎 @数据小兵

Anscombe 的回归例子: 讨论

但事实上, 只有第一组数据 $Y_1 \sim X_1$, 直接进行线性回归才是最合理的选择。

第二组数据, $Y_2 \sim X_2$, 明显存在曲线关系, 也就是说, 更好的模型可能是一个多项式形式。

第三组数据, $Y_3 \sim X_3$, 在完美的线性关系之外, 存在一个异常值。直接进行线性回归未必是最合理的选择。

第四组数据, 回归直线的斜率完全由 $x = 19$ 那一个观测值所决定。我们需要了解为什么 $x = 19$ 的这个点这么重要。同时, 我们需要确保这个观测值是可靠的。如果不可靠, 那么整个回归分析结果的可靠性就存疑了。

启示, 进行回归分析之前, 要先画图看一看, 这样可能会避免模型选择一开始就进入误区。

最小二乘估计量的合理性基础：Gauss-Markov假设

什么情况下最小二乘估计量才具有合理性质？

- ① 假设SLR.1 (线性模型): 总体变量 Y, X 满足线性模型关系

$$Y = \beta_0 + \beta_1 X + u, \quad (4)$$

- ② 假设SLR.2 (随机抽样): 样本数据 $\{(X_1, Y_1), \dots, (X_n, Y_n)\}$ 为来自总体的简单随机样本
- ③ 假设SLR.3 (解释变量的样本数据必须有波动性): 所有 X_i 的值不是恒为常数
- ④ 假设SLR.4 (条件期望为零): $E[u_i|X_i] = 0$
- ⑤ 假设SLR.5 (条件不相关性与同方差性):

$$\text{Var}(u_i|X_i) = \sigma^2, \forall i; \text{Cov}(u_i, u_j|X_i, X_j) = 0, \forall i \neq j$$

关于Gauss-Markov假设

假设SLR.1-SLR.3相对比较容易理解。假设SLR.1关注的是线性模型的合理性；假设SLR.2关注的是数据的采样方式；假设SLR.3在实际情况下都是满足的。

假设SLR.4-SLR.5相对比较复杂。也是关键的假设。这两个都是关于 u_i 关于 X_i 的条件分布的特征：假设SLR.4关注的是这个条件分布的均值（一阶矩），假设SLR.5关注的是这个条件分布的方差（二阶矩）。

这两个假设后期我们都会进行适当放松。但无论如何放松，本质上涉及的都是这两点：一个是 u_i 与 X_i 的相关性，一个是（在考虑了 X_i 以后） u_i 与 u_j 的相关性。

条件期望及其性质

设 Y 为一个随机变量, X 为一个随机向量。给定 $X = x$, Y 的条件分布的期望 $E[Y|X = x]$ 是 x 的函数。不妨记作

$$\phi(x) = E[Y|X = x] = \int y f_{y|x}(y|x) dy$$

我们称随机变量 $E[Y|X] = \phi(X)$ 为 Y 关于 X 的条件期望。
条件期望的基本性质:

- 性质1:

$$E[E[Y|X]] = E[Y].$$

- 性质2: 对任意函数 h ,

$$E[h(X)Y|X] = h(X)E[Y|X]$$

- 性质3:

$$E[h_1(X)Y_1 + h_2(X)Y_2|X] = h_1(X)E[Y_1|X] + h_2(X)E[Y_2|X]$$

条件方差及其性质

Y关于X的条件方差定义为:

$$\text{Var}(Y|X) = E\left((Y - E[Y|X])^2|X\right)$$

条件方差的基本性质:

- 性质1: 全方差公式

$$\text{Var}(Y) = E[\text{Var}(Y|X)] + \text{Var}(E[Y|X])$$

- 性质2:

$$\text{Var}(h(X)Y|X) = h(X)^2 \text{Var}(Y|X)$$

- 性质3:

$$\begin{aligned} \text{Var}(h_1(X)Y_1 + h_2(X)Y_2|X) &= h_1(X)^2 \text{Var}(Y_1|X) \\ &+ h_2(X)^2 \text{Var}(Y_2|X) + 2h_1(X)h_2(X)\text{Cov}(Y_1, Y_2|X) \end{aligned}$$

其中

$$\text{Cov}(Y_1, Y_2|X) = E\left((Y_1 - E[Y_1|X])(Y_2 - E[Y_2|X])|X\right)$$

最小二乘估计量的无偏性

记 $X = (X_1, \dots, X_n)$ 。在上述假设条件下：

- ① $\hat{\beta}_1$ 是 β_1 的无偏估计量

$$E[\hat{\beta}_1|X] = \beta_1$$

因而

$$E[\hat{\beta}_1] = \beta_1$$

- ② $\hat{\beta}_0$ 是 β_0 的无偏估计量

$$E[\hat{\beta}_0|X] = \beta_0$$

因而

$$E[\hat{\beta}_0] = \beta_0$$

最小二乘估计量的方差

在上述假设条件下：

- ① $\hat{\beta}_1$ 的条件方差与无条件方差分别为

$$\text{Var}[\hat{\beta}_1|X] = \frac{1}{\sum_{i=1}^n (X_i - \bar{X})^2} \sigma^2, \quad \text{Var}[\hat{\beta}_1] = E \left[\frac{1}{\sum_{i=1}^n (X_i - \bar{X})^2} \right] \sigma^2$$

- ② $\hat{\beta}_0$ 的条件方差与无条件方差分别为

$$\text{Var}[\hat{\beta}_0|X] = \frac{\frac{1}{n} \sum_{i=1}^n X_i^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \sigma^2, \quad \text{Var}[\hat{\beta}_0] = E \left[\frac{\frac{1}{n} \sum_{i=1}^n X_i^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right] \sigma^2$$

σ^2 的估计及性质

- σ^2 的一个估计量如下:

$$s^2 = \frac{1}{n-2} \sum_{i=1}^n \hat{u}_i^2$$

- 定理2.3: s^2 是 σ^2 的无偏估计量:

$$E[s^2|X] = \sigma^2$$

- σ 的一个估计量为:

$$s = \sqrt{\frac{1}{n-2} \sum_{i=1}^n \hat{u}_i^2}$$

- s 不是 σ 的无偏估计量

$\hat{\beta}_0, \hat{\beta}_1$ 的标准差的估计量

- $\hat{\beta}_1$ 的标准差的估计量为:

$$se(\hat{\beta}_1) = \frac{s}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2}}$$

- $\hat{\beta}_0$ 的标准差的估计量为:

$$se(\hat{\beta}_0) = s \sqrt{\frac{n^{-1} \sum_{i=1}^n X_i^2}{\sum_{i=1}^n (X_i - \bar{X})^2}}$$

注：以上标准差估计量是条件标准差（而不是无条件标准差）的估计量。为什么这样做是合理的？教材Theorem 3.2下面的文字进行了讨论。我们在大样本理论那一部分再回来讨论这个问题。

评论

简单线性回归模型的局限性：

- 现实中影响 Y 的变量有多个
- 简单线性模型没有控制住除 X 之外的变量对 Y 的影响
- 遗漏变量可能导致推断结果出现严重偏差

多变量线性回归模型：

- 引入多个影响 Y 的变量
- 其分析结果可以提供类似于“保持其他变量不变”情形下的解释