


- 
- 在第2–5章中我们已经讨论了线性回归模型与OLS估计值的理论性质
 - 本章将讨论实际应用中的几个重要问题，例如：
 - (1) 数据测度单位改变会不会影响OLS标准误及相关统计量的取值？
 - (2) 怎么刻画变化的偏效应？
 - (3) 怎么选回归元？等

Chapter 6

多元回归分析： 深入专题

章节框架

- 在这一章中我们介绍回归模型实际应用中的几个重要问题：
- 首先，我们介绍数据测度单位对OLS统计量的影响
- 然后，我们对函数形式进行进一步讨论
- 之后，我们进一步探讨拟合优度和回归元选择
- 最后，我们讨论预测的构造

数据测度单位对OLS统计量的影响

- 在第二章中我们已经讨论了解释变量与被解释变量测度单位对OLS估计值的影响
- 现在我们考虑测度单位对标准误、 t 统计量、 F 统计量和置信区间的影响
- 例子：婴儿出生体重与孕妇抽烟量

$$\widehat{bwght} = \hat{\beta}_0 + \hat{\beta}_1 cigs + \hat{\beta}_2 faminc$$

其中 $bwght$ 表示以盎司为单位的婴儿出生体重， $cigs$ 表示母亲每天抽烟量， $faminc$ 表示以千美元为单位的家庭年收入

以下变化对估计值与统计量有什么影响？

(1) 如果 $bwght$ 以磅为单位（取值/16）

(2) 如果 $cigs$ 以包为单位（取值/20）

数据测度单位对OLS统计量的影响

TABLE 6.1 Effects of Data Scaling

Dependent Variable	(1) <i>bwght</i>	(2) <i>bwghtlbs</i>	(3) <i>bwght</i>
Independent Variables			
<i>cigs</i>	-.4634 (.0916)	-.0289 (.0057)	—
<i>packs</i>	—	—	-9.268 (1.832)
<i>faminc</i>	.0927 (.0292)	.0058 (.0018)	.0927 (.0292)
<i>intercept</i>	116.974 (1.049)	7.3109 (.0656)	116.974 (1.049)
Observations	1,388	1,388	1,388
<i>R</i> -Squared	.0298	.0298	.0298
SSR	557,485.51	2,177.6778	557,485.51
SER	20.063	1.2539	20.063

对函数形式的进一步讨论

- 例子：污染对住房价格的影响

$$\log(\widehat{price}) = 9.23 - 0.718\log(nox) + 0.306rooms$$

- 对数形式的讨论：
 - 方便用于百分比或弹性解释
 - 对数变量的斜率系数无关于单位变化
 - 取对数通常可以消除或缓解异常值问题
 - 取对数通常有助于确保正态性和同方差性
 - 以年等为单位变量不应取对数
 - 百分数为单位的变量也不应取对数
 - 变量有零值或负值都不能取对数


对函数形式的进一步讨论

- 为了描述递减或递增的边际效应，有时需要二次函数形式
- 例子：工资等式

$$\widehat{wage} = 3.73 + .298 \text{ exper} - .0061 \text{ exper}^2$$

(.35) (.041) (.0009)

凹函数




$$n = 526, R^2 = .093$$

- 经验的边际效应

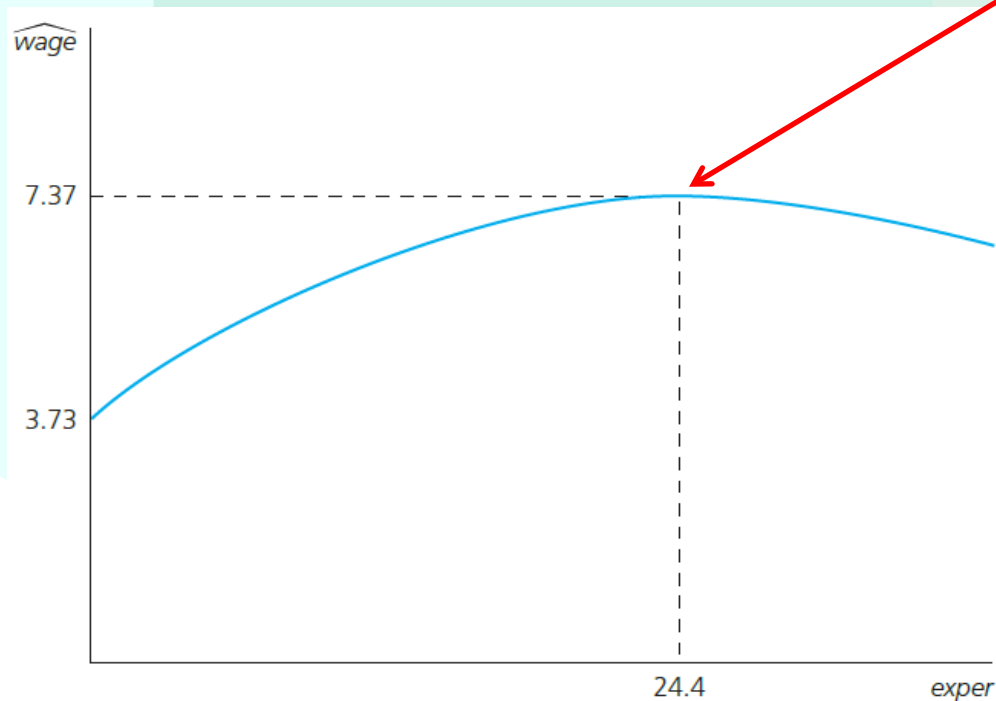
$$\frac{\Delta wage}{\Delta exper} = .298 - 2(.0061)exper$$

第一年的工作经验会增加大约0.30美元的工资，第二年会增加0.298 - 2(0.0061) = 0.29美元，以此类推。



对函数形式的进一步讨论

- 工作经验的最高工资



$$x^* = \left| \frac{\hat{\beta}_1}{2\hat{\beta}_2} \right| = \left| \frac{.298}{2(.0061)} \right| \approx 24.4$$

二次函数的代价：转折点

是否意味着24.4年后，经验的回报为负？
需要考虑以下几个情况：

(1) 取决于样本中有多少观测值位于转折点的右侧。如果很少则不需考虑（在给定的例子中，大约28%的观察值在转折点右侧）。

(2) 是否存在模型设定问题（例如忽略变量，如年龄）。

对函数形式的进一步讨论

- 例子：污染对房价的影响

空气中的氮氧化物、与就业中心的距离、
平均学生/教师比率

$$\widehat{\log(price)} = 13.39 - .902 \log(nox) - .087 \log(dist) + .545 rooms + .062 rooms^2 - .048 stratio$$

(.57) (.115) (.043) (.165) (.013) (.006)

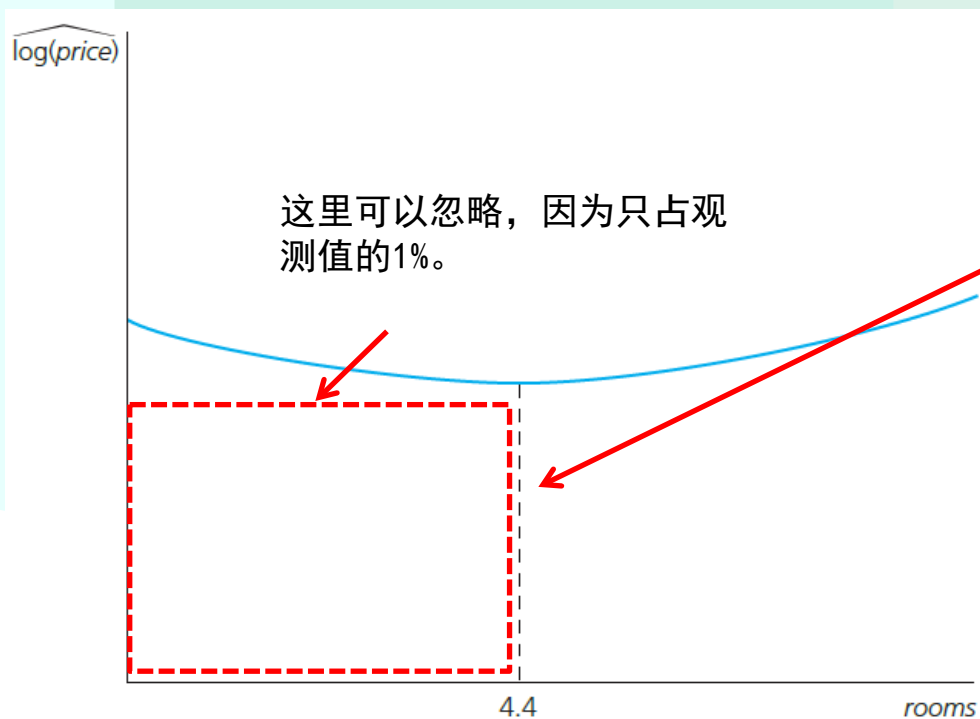
$$n = 506, R^2 = .603$$

这是否意味着，在房间数量较少的情况下，更多的房间意味着更低的价格？（少于4间的只有1%）

$$\Rightarrow \frac{\Delta \log(price)}{\Delta rooms} = \frac{\% \Delta price}{\Delta rooms} = -.545 + .124 rooms$$

对函数形式的进一步讨论

- 计算转折点



转折点：

$$x^* = \left| \frac{-0.545}{2(0.062)} \right| \approx 4.4$$

房间从5增加到6：

$$-0.545 + 0.124(5) = +7.5\% \text{ price}$$

房间从6增加到7：

$$-0.545 + 0.124(6) = +19.9\% \text{ price}$$

对函数形式的进一步讨论

- 其他可能性

$$\log(\text{price}) = \beta_0 + \beta_1 \log(\text{nox}) + \beta_2 [\log(\text{nox})]^2 \\ + \beta_3 \text{crime} + \beta_4 \text{rooms} + \beta_5 \text{rooms}^2 + \beta_6 \text{stratio} + u$$

$$\Rightarrow \frac{\Delta \log(\text{price})}{\Delta \log(\text{nox})} = \frac{\% \Delta \text{price}}{\% \Delta \text{nox}} = \beta_1 + 2\beta_2 [\log(\text{nox})]$$

- 更高阶的多项式：总成本函数

$$\text{cost} = \beta_0 + \beta_1 \text{quantity} + \beta_2 \text{quantity}^2 + \beta_3 \text{quantity}^3 + u$$

对函数形式的进一步讨论

- 含有交互项的模型

$$price = \beta_0 + \beta_1 sqrf\text{ft} + \beta_2 bdrms$$

$$+ \beta_3 sqrf\text{ft} \cdot bdrms + \beta_4 bthrms + u$$

交互项 (interaction)

$$\Rightarrow \frac{\Delta price}{\Delta bdrms} = \beta_2 + \beta_3 sqrf\text{ft}$$

卧室数量的影响取决于建筑面积的大小 (有意义)

- 相互作用效应使参数的解释复杂化

β_2 = bdrms对一套面积为零的住房的价格的影响。(没意义)

对函数形式的进一步讨论

- 例子：大学前GPA, ACT成绩和出勤率对期末考试分数的影响

出勤率的影响

$$\begin{aligned}\widehat{stndfnl} = & 2.05 - .0067 \text{ atndrte} - 1.63 \text{ priGPA} - .128 \text{ ACT} \\ & (1.36) \quad (.0102) \quad (.48) \quad (.098) \\ & + .296 \text{ priGPA}^2 + .0045 \text{ ACT}^2 + .0056 \text{ priGPA} \cdot \text{atndrte} \\ & (.101) \quad (.0022) \quad (.0043) \\ & n = 680, R^2 = .229, \bar{R}^2 = .222.\end{aligned}$$

- 相互作用效应使参数的解释复杂化

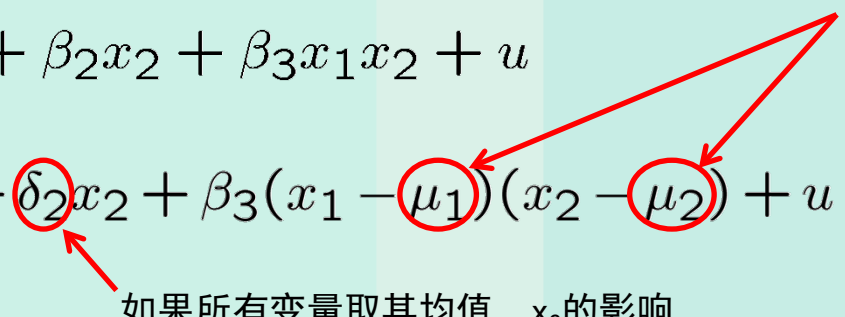
-0.0067 = priGPA为零时atndrte对期末分数的影响。（没意义）

对函数形式的进一步讨论

- 交互效应的再参数化

总体均值，由样本均值代替

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + u$$

$$y = \alpha_0 + \delta_1 x_1 + \delta_2 x_2 + \beta_3 (x_1 - \mu_1)(x_2 - \mu_2) + u$$


如果所有变量取其均值， x_2 的影响

- 再参数化的优势

- 参数的含义简单化
- 对于在均值处的偏效应的标准误可以计算

对函数形式的进一步讨论

- 在具有二次函数、交互作用和其他非线性函数形式的模型中，偏效应取决于一个或多个解释变量的值
- 例子：上个例子中atndrte对期末平均分数的偏效应：

$$-0.0067 + 0.0056priGPA$$

- 平均偏效应（average partial effects, APE）是描述因变量和每个解释变量之间关系的一种测度值
- 计算偏效应并带入估计值后，对样本中每个元素的偏效应求平均值
- 例子：上个例子中atndrte的平均偏效应：

$$APE_{atndrte} = -0.0067 + 0.0056\overline{priGPA}$$

思考：priGPA的APE是多少？

拟合优度和回归元选择的进一步探讨

- 拟合优度和回归元选择的进一步探讨
- 对R²的一般性评论
 - 高R²不意味着因果关系
 - 低R²并不妨碍精确估计偏效应
 - 普通R²测量了什么？

$$R^2 = 1 - \frac{SSR}{SST} = 1 - \frac{(SSR/n)}{(SST/n)}$$

是一个对于

$$1 - \frac{\sigma_u^2}{\sigma_y^2}$$

的估计

总体R²

- 引入新的解释变量必然会增加R², 即使该变量对被解释变量没有因果影响

拟合优度和回归元选择的进一步探讨

- 调整R2

- 一个更好的考虑自由度的估计量是

- $$\bar{R}^2 = 1 - \frac{(SSR/(n - k - 1))}{(SST/(n - 1))} = \text{adjusted } R^2$$

分子和分母的自由度

- 调整R2对于增加新的自变量进行惩罚（只有加入的新变量可以“明显地”降低SSR时，调整R2才会增加）
- 数学上可以证明：当且仅当新变量的t统计量的绝对值大于1，调整R2才会增加

- R2与调整R2的关系

$$\bar{R}^2 = 1 - (1 - R^2)(n - 1)/(n - k - 1)$$

调整R2可能为负

拟合优度和回归元选择的进一步探讨

- 利用调整R2在两个非嵌套模型中进行选择
 - 如果任一模型都不是其他模型的特例，则模型之间是非嵌套的

$$rdintens = \beta_0 + \beta_1 \log(sales) + u$$

$$R^2 = .061, \bar{R}^2 = .030$$

$$rdintens = \beta_0 + \beta_1 sales + \beta_2 sales^2 + u$$

$$R^2 = .148, \bar{R}^2 = .090$$

- 比较两个模型的R2对第一个模型不公平，因为第一个模型包含的参数较少
- 在给定的例子中，即使调整了自由度的差异，二次模型仍然是首选

拟合优度和回归元选择的进一步探讨

- 具有不同因变量模型的比较
 - R²或调整R²不能用于比较具有不同因变量的模型
- 例子：CEO薪酬与企业业绩

log(salary)比
salary的波动性
要小得多

$$\widehat{salary} = 830.63 + .0163 sales + 19.03 roe$$

(223.90) (.0089) (11.08)

$$n = 209, R^2 = .029, \bar{R}^2 = .020, SST = 391,732,982$$

$$\widehat{lsalary} = 4.36 + .275 lsales + .0179 roe$$

(0.29) (.033) (.0040)

$$n = 209, R^2 = .282, \bar{R}^2 = .275, SST = 66.72$$

拟合优度和回归元选择的进一步探讨

- 为了避免遗漏变量偏误，有时回归分析中控制因素过多
- 在某些情况下，一些变量不应被控制
 - 在啤酒税（和其他因素）导致的交通事故死亡率回归中，不应直接控制啤酒消费：

$$fatalities = \beta_0 + \beta_1 tax + \beta_2 beercons + \dots$$

- 在农民家庭健康对农药使用的回归中，不应控制看病次数
- 不同的回归可能有不同的目的
 - 在房价对房屋特征回归中，如果想研究价格评估的有效性，就应该包括房价评估；
 - 但如果只是探讨房屋特征对房价的影响，会发现不应包含房价评估（控制房价评估价值不变，讨论增加一间卧室对房屋价值的影响没有意义）

拟合优度和回归元选择的进一步探讨

- 增加自变量以减少误差方差
 - 增加自变量可能会增加多重共线性
 - 另一方面，增加自变量可以减少误差方差（残差的方差）
 - 应添加与其他自变量不相关的变量，因为它们在不增加多重共线性的情况下减少了误差方差
 - 然而，很难找
- 例子：个人啤酒消费与啤酒价格
 - 将个体特征纳入啤酒消费对啤酒价格的回归中，可以更精确地估计价格弹性

预测

- 在第三章中，我们定义了OLS预测值或拟合值：

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \cdots + \hat{\beta}_k x_k$$

具体的，给定 $x_1 = c_1, \dots, x_k = c_k$ ，我们得到预测值：

$$\hat{\theta} = \hat{\beta}_0 + \hat{\beta}_1 c_1 + \cdots + \hat{\beta}_k c_k$$

$\hat{\theta}$ 作为OLS估计值的线性组合，存在抽样波动的问题，如何得到置信区间？

$\hat{\theta}$ 可以视为预期 $\theta_0 = \beta_0 + \beta_1 c_1 + \cdots + \beta_k c_k$ 的估计，那么我们可以变换模型得到

$$y = \theta_0 + \beta_1 (x_1 - c_1) + \cdots + \beta_k (x_k - c_k) + u$$

预测

- 例子：大学GPA预测值的置信区间

$$\begin{aligned}\widehat{colgpa} &= 1.493 + .00149 sat - .01386 hsperc \\ &\quad (0.075) \quad (.00007) \quad (.00056) \\ &\quad - .06088 hsize + .00546 hsize^2 \\ &\quad (.01650) \quad (.00227) \\ n &= 4,137, R^2 = .278, \bar{R}^2 = .277, \hat{\sigma} = .560,\end{aligned}$$

- 当 $sat = 1200, hsperc = 30, hsize = 5$ 时, $colgpa$ 预测值为 2.70

预测

- 例子：大学GPA预测值的置信区间

- 模型变换：

$$sat0 = sat - 1200, hsperc0 = hsperc - 30, hsize0 = hsize - 5, hsize0^2 = hsize^2 - 25$$

- 估计结果：

$$\begin{aligned}\widehat{colgpa} &= 2.700 + .00149 sat0 - .01386 hsperc0 \\ &\quad (0.020) \quad (.00007) \quad (.00056) \\ &\quad - .06088 hsize0 + .00546 hsize0^2 \\ &\quad (.01650) \quad (.00227) \\ n &= 4,137, R^2 = .278, \bar{R}^2 = .277, \hat{\sigma} = .560.\end{aligned}$$

- 预期gpa的95%置信区间： $2.70 \pm 0.02 \times 1.96 = [2.66, 2.74]$

预测

- 当因变量为 $\log(y)$ 时对 y 的预测

$$\log(y) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_k x_k + u$$

- 给定OLS估计量，首先可以预测 $\log y$

$$\widehat{\log y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \cdots + \hat{\beta}_k x_k$$

- 之后，我们可以得到 y 的预测值：

$$\hat{y} = \exp(\widehat{\log y})$$

- 然而，这将系统性地低估 y 的预测值

预测

- 当因变量为 $\log(y)$ 时对 y 的预测

$$\log(y) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_k x_k + u$$

$$\Rightarrow y = \exp(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_k x_k) \exp(u)$$

当假设 u 独立于 x_1, \dots, x_k 时:

$$\Rightarrow E(y|x) = \exp(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_k x_k) E(\exp(u))$$

$$\Rightarrow \hat{y} = \exp(\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \cdots + \hat{\beta}_k x_k) \left(\frac{1}{n} \sum_{i=1}^n \exp(\hat{u}_i) \right)$$

对 y 的预测

小节

- 本章讨论了多元回归分析的一些重要专题
- 改变自变量或者因变量的单位，对统计量没有影响
- 对数、二次项、交互项的使用会对OLS估计值的解释产生影响
- 调整R²有些时候可以用来选择回归元
- 我们提出了预测值置信区间的构造方法并讨论了 $\log(y)$ 的预测问题