

# 回归：预测与拟合

司继春

2024年10月





# 预测问题

- **拟合 (fitting)** 以及预测是最经典的统计问题之一，而**回归 (regression)** 是解决这类问题最常用的手段。
- 如果我们观察到一系列数据 $(y_i, x_i), i = 1, \dots, N$ ，我们希望使用 $x_i$ 的线性函数： $f(x_i) = \alpha + \beta x_i$ 对 $y_i$ 进行预测，那么只要确定了其中的参数 $\alpha$ 和 $\beta$ 就确定了这个预测的函数。
- 我们称：
  - $x_i$ 为**自变量 (independent variable)** 或者**解释变量 (explanatory variable)**、**回归元 (regressor)**
  - 而 $y_i$ 为**因变量 (dependent variable)** 或者**被解释变量 (explained variable)**、**结果变量 (outcome variable)**。

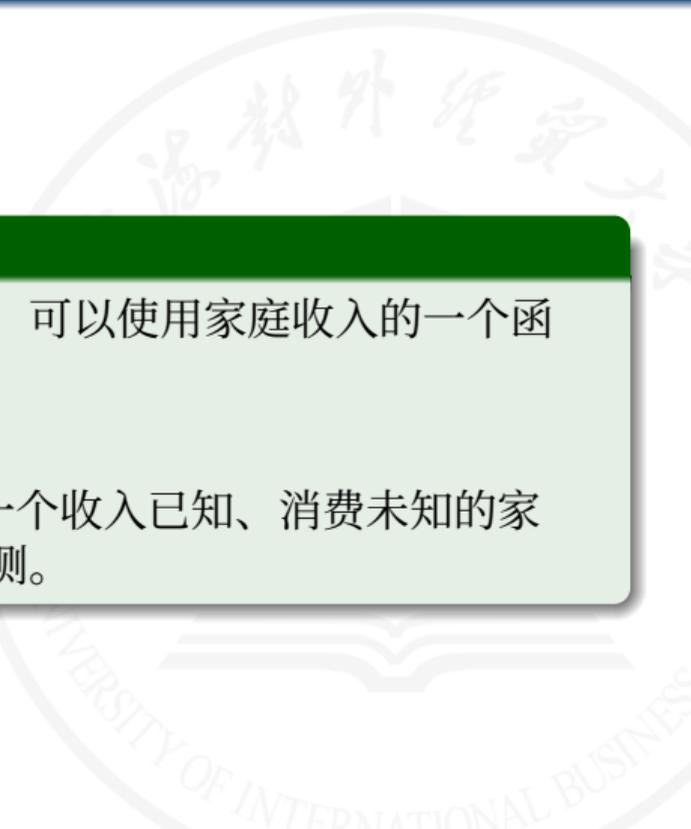
# 一元线性回归

## 一元线性回归的例子：收入与消费

一个非常经典的问题是家庭收入与消费之间的关系，可以使用家庭收入的一个函数：

$$f(I) = \alpha + \beta \cdot I$$

对家庭消费进行预测。只要我们知道了 $\alpha, \beta$ ，对于一个收入已知、消费未知的家庭，我们就可以使用以上函数对未知的消费进行预测。



# 回归与插值

我们需要区分函数关系与相关关系：

- 如果给定一个输入，有确定的输出，那么两个变量是函数关系，比如：

$$f(h) = \alpha + \beta \cdot h$$

- 如果给定一个输入，可能有不同的值与之相对应，那么是相关关系，比如身高都为170的人，体重完全可能不同

在此基础上：

- 回归需要解决的是相关关系，我们使用一个具有确定性的函数输出对不确定的 $y$ 进行预测
- 如果一个函数是未知的，为了确定这个函数，解决方案为**插值 (interpolation)** 而非回归。

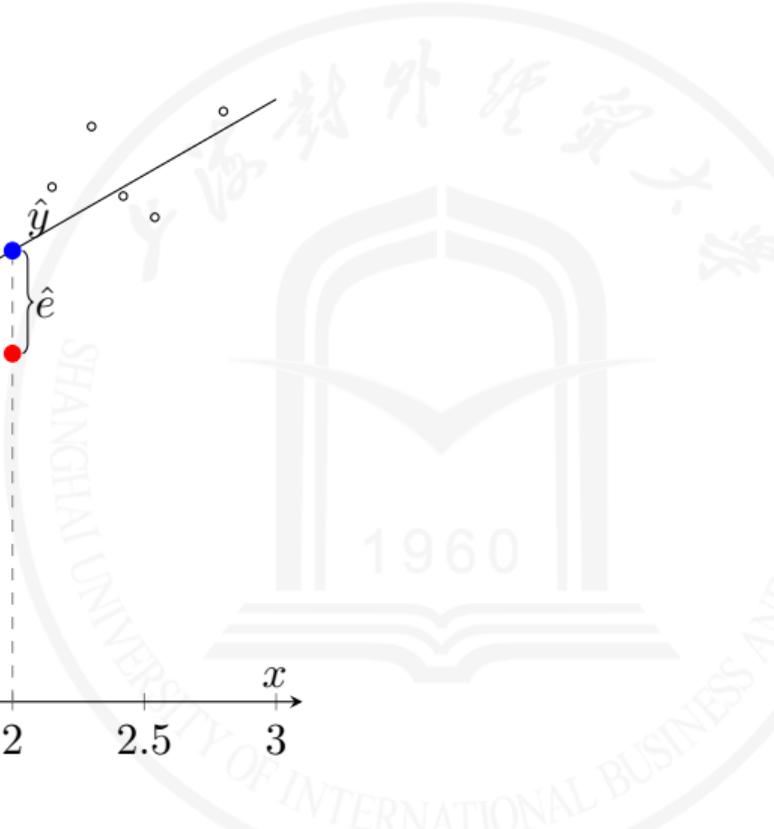
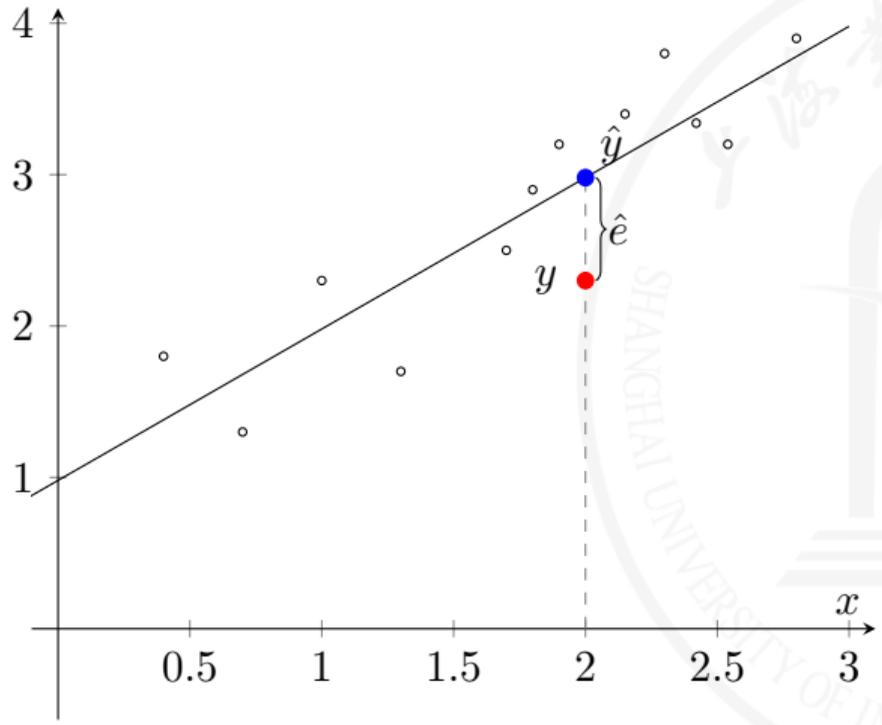
# 预测误差

如果给定一个 $\alpha$ 和 $\beta$ 的值 $(\tilde{\alpha}, \tilde{\beta})$ ，我们可以计算使用以上函数对 $y_i$ 进行预测的误差，即残差 (residuals)：

$$\tilde{e}_i = y_i - \tilde{\alpha} - \tilde{\beta}x_i$$

- 残差即使用函数 $f(x) = \tilde{\alpha} + \tilde{\beta} \cdot x$ 对个体 $i$ 的 $y$ ： $y_i$ 进行预测的预测误差
- 残差应该越“小”越好。

# 预测误差











# 一元线性回归示例

## 收入与消费

在以上程序中：

- 首先剔除了收入和消费的异常值（收入小于0的值以及收入和消费存在删失问题数据）
- 接着使用reg命令计算了消费（total\_consump）对收入（total\_income）的回归，回归结果为：

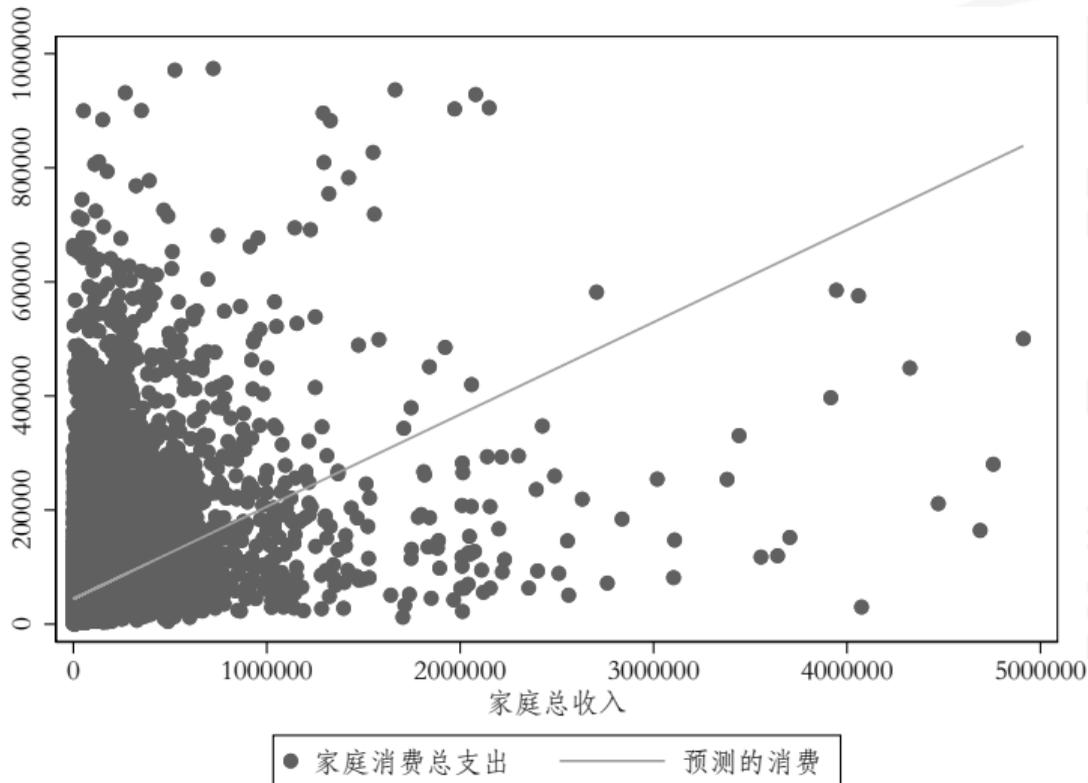
$$\hat{y} = 45696.74 + 0.15 \times x$$

意味着收入每增加1元，消费平均会增加0.15元。

- 此外，如果我们知道某个家庭收入为30万元，那么对其消费的最优预测为

$$\hat{y}_{175} = 45696.74 + 0.15 \times 300000 = 90696.74$$

# 一元线性回归示例



# 一元线性回归的三个性质

- 如果我们将 $x_i$ 的平均值 $\bar{x}$ 带入到拟合公式中，可以得到：

$$\hat{\alpha} + \hat{\beta}\bar{x} = \bar{y} - \hat{\beta}\bar{x} + \hat{\beta}\bar{x} = \bar{y}$$

因而使用最小二乘法进行预测时，在 $x_i$ 的平均值 $\bar{x}$ 处的预测即 $\bar{y}$ 。

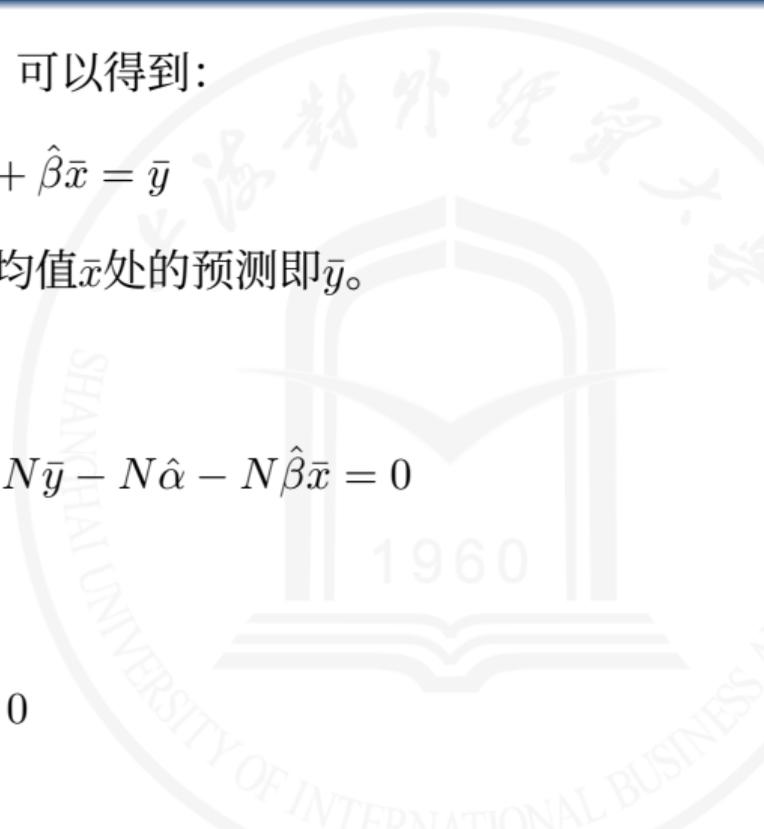
- 残差的和：

$$\sum_{i=1}^N \hat{e}_i = \sum_{i=1}^N (y_i - \hat{\alpha} - \hat{\beta}x_i) = N\bar{y} - N\hat{\alpha} - N\hat{\beta}\bar{x} = 0$$

- 残差和 $x$ 之间不相关：

$$\sum_{i=1}^N x_i \hat{e}_i = 0$$

从而残差和 $x$ 之间的样本相关系数为0。





# 0/1型变量的一元线性回归

将以上结论带入到预测方程中，可以得到：

- 当 $x_i = 0$ 时，有：

$$\hat{y} = \hat{\alpha} + \hat{\beta}x_i = \hat{\alpha} = \bar{y}_0$$

- 当 $x_i = 1$ 时，有：

$$\hat{y} = \hat{\alpha} + \hat{\beta}x_i = \hat{\alpha} + \hat{\beta} = \bar{y}_1 - \bar{y}_0 + \bar{y}_0 = \bar{y}_1$$

即，当 $x_i$ 只能取0/1两个值时，对 $y$ 的预测即分组的均值：**分组均值即最优预测。**

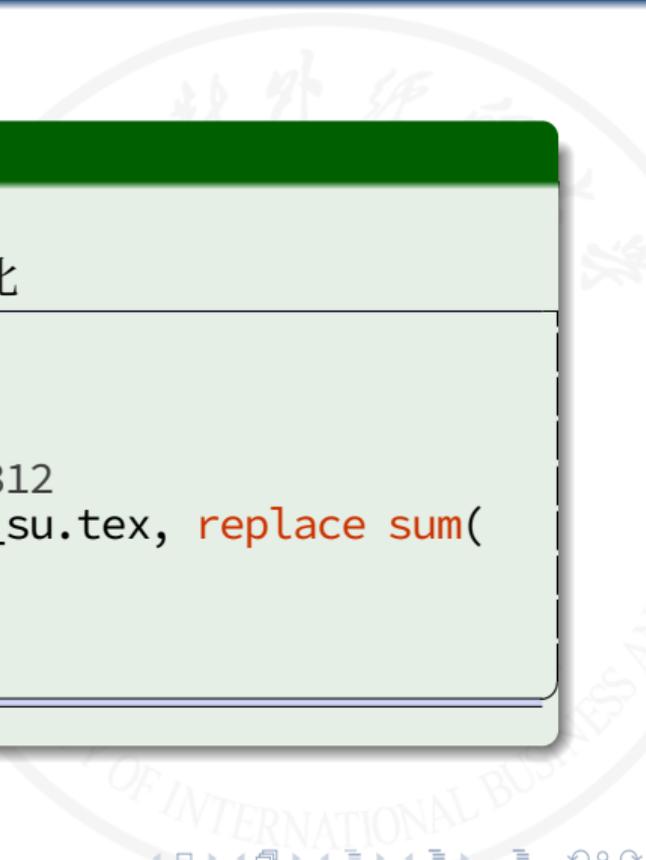


# 虚拟变量回归与均值比较

## 不同性别收入比较

代码 2: 不同性别的收入对比

```
1 // file: reg_with_dummy.do
2 use datasets/chfs2017_ind.dta, clear
3 gen p_income = a3109*12
4 gen gender = 2-a2003 // 定义为男性, 为女性a200312
5 bysort gender: outreg2 using reg_with_dummy_su.tex, replace sum(
   log) eqkeep(N mean) keep(p_income)
6 reg p_income gender
7 outreg2 using reg_with_dummy.tex, replace
```





# 虚拟变量回归与均值比较

## 不同性别收入比较

VARIABLES	(1) p_income
gender	9,297*** (464.0)
Constant	38,811*** (356.8)
Observations	33,481
R-squared	0.012
Standard errors in parentheses	
*** p<0.01, ** p<0.05, * p<0.1	

# 多元线性回归

以上讨论了一元线性回归，即使用一个解释变量 $x$ 对 $y$ 进行预测。我们还可以继续推广，即使用多个 $x$ 对 $y$ 进行预测，即使用函数：

$$f(x_i|\beta) = \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_K x_{iK}$$

其中

$$x_i = \begin{pmatrix} x_{i1} \\ \vdots \\ x_{iK} \end{pmatrix}, \beta = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_K \end{pmatrix}$$

一般而言，我们通常会保留常数项，不失一般性，我们令 $x_{i1} = 1$ 。

# 多元线性回归

在这里假设函数  $f(x_i, |\beta)$  是**参数线性的**，即不存在  $\beta_k$  之间的非线性关系。比如：

- 我们排除了如下的函数形式：

$$\hat{y}_i = f(x_i|\beta) = \beta_1 x_{1i} + \beta_1^2 x_{2i}$$

由于存在  $\beta_1$  的非线性函数，从而以上设定不是线性回归设定。

- 包含  $x_i$  的非线性函数是可以的，比如：

$$\hat{y}_i = \beta_1 + \beta_2 x_{1i} + \beta_3 x_{1i}^2$$

是完全允许的（此时不妨记  $x_{2i} = x_{1i}^2$ ）。

# 多元线性回归

为了方便起见，我们一般用向量表述上述方程：

$$f(x_i) = x_i' \beta$$

其中：

$$x_i = \begin{pmatrix} 1 \\ x_{i2} \\ \vdots \\ x_{iK} \end{pmatrix}, \beta = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_K \end{pmatrix}$$

为两个  $K$  维列向量。给定一个  $\beta$ ，我们可以使用  $x_i' \beta$  对  $y_i$  进行预测的预测值：

$$\hat{y}_i = x_i' \beta$$

以及预测的误差，即残差：

$$\hat{e}_i = y_i - \hat{y}_i = y_i - x_i' \beta$$

# 最小二乘法 (OLS)

如果我们记:

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}, X = [x_1, x_2, \dots, x_N]' = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1K} \\ x_{21} & x_{22} & \cdots & x_{2K} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \cdots & x_{NK} \end{bmatrix}$$

那么残差向量为:

$$\hat{e} = Y - X\beta = \begin{pmatrix} \hat{e}_1 \\ \hat{e}_2 \\ \vdots \\ \hat{e}_N \end{pmatrix}_{N \times 1}$$

而最小二乘法最小化:

$$\min_b \sum_{i=1}^N \hat{e}_i^2 = \min_b \hat{e}'\hat{e} = \min_b (Y - Xb)'(Y - Xb)$$

## 最小二乘法 (OLS)

对以上目标函数求导数并令其等于0, 可以得到一阶条件:

$$\frac{\partial (Y - Xb)' (Y - Xb)}{\partial b} = \frac{\partial (Y'Y - Y'Xb - b'X'Y + b'X'Xb)}{\partial b} = -X'Y - X'Y + 2X'Xb = 0$$

解以上方程可以得到:

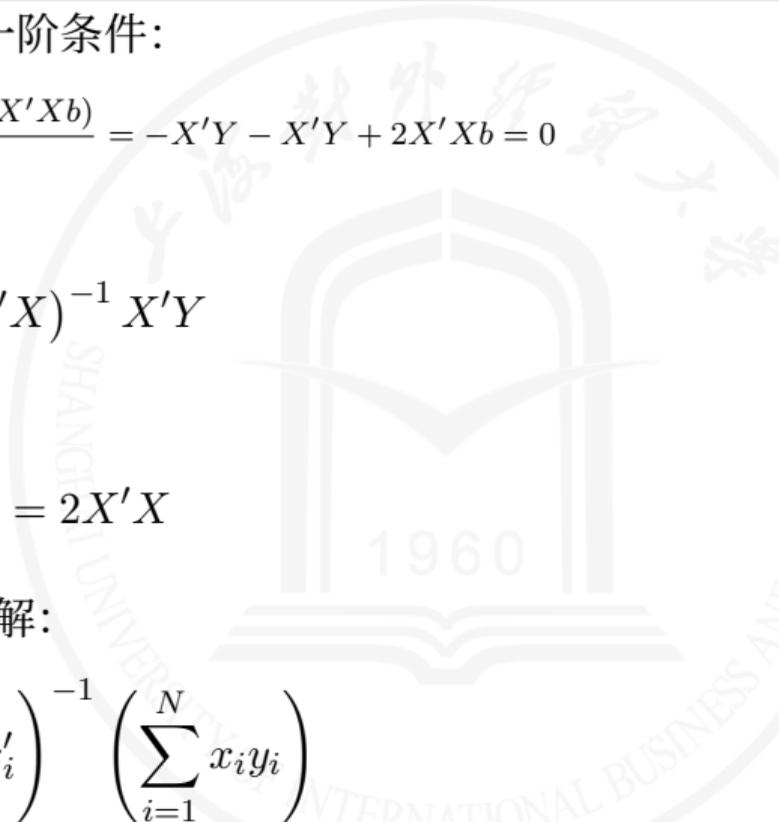
$$X'Xb = X'Y \Rightarrow \hat{\beta} = (X'X)^{-1} X'Y$$

以上最大化问题的二阶导为:

$$\frac{\partial^2 (y - X\beta)' (y - X\beta)}{\partial \beta \partial \beta'} = 2X'X$$

为一个正定矩阵, 因而以上根据一阶条件求得的解:

$$\hat{\beta} = (X'X)^{-1} X'Y = \left( \sum_{i=1}^N x_i x_i' \right)^{-1} \left( \sum_{i=1}^N x_i y_i \right)$$





# 识别条件

## 收入与消费、储蓄

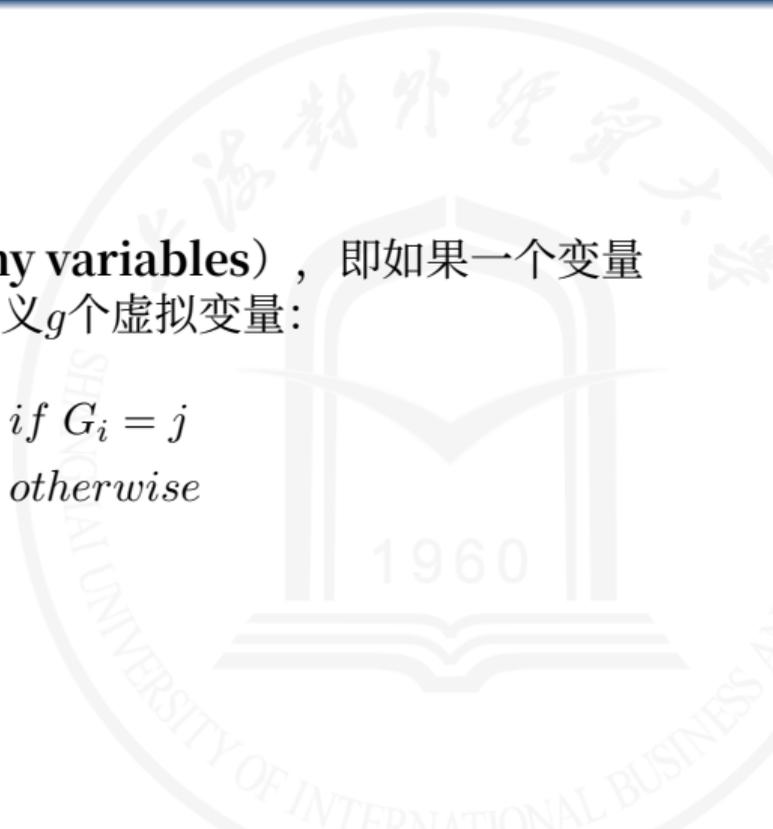
家庭收入 ( $I$ ) 等于家庭的消费  $C$  加储蓄  $S$ ,  $I = C + S$ , 那么  $I, C, S$  不能同时出现在  $X$  里面, 否则  $X$  列不满秩。但是由于  $\ln(I) \neq \ln(C) + \ln(S)$ , 因而  $X$  中同时包含  $\ln(I), \ln(C), \ln(S)$  理论上仍然是可以的。



# 虚拟变量

分类变量加入回归中需要使用**虚拟变量 (dummy variables)**，即如果一个变量的取值范围为  $G_i = 1, 2, \dots, g$ ，我们可以相应的定义  $g$  个虚拟变量：

$$d_{ij} = 1 \{G_i = j\} = \begin{cases} 1 & \text{if } G_i = j \\ 0 & \text{otherwise} \end{cases}$$













# 虚拟变量回归

## 不同教育程度的收入

- 如果一定要加入edu9这个虚拟变量，那么可以在reg命令后面加入noconstant选项，该选项即防止线性回归中包含常数项，从而我们可以包含edu9这个变量。
- 实际上，如果包含edu9而不包含常数项，那么估计的系数就是每个分组的收入的平均值，比如，edu1的系数为26115.86，意味着文化程度为文盲的平均收入为26115.86元。
- 而如果包含常数项而把edu9忽略掉，那么edu1-edu8估计的系数即每个组的收入与edu9这个组（基准组）的差异，比如edu1的系数为-100600.7，那么意味着文化程度为文盲的平均收入比文化程度为博士的平均收入低100600.7元。
- 第3列同理，如果把edu1去掉，那么edu2-edu9的系数都是与edu1组（基准组）相比的收入差异。

# 虚拟变量回归

以上的结果并非偶然。如果在回归中不加入常数项而是加入所有的分组虚拟变量，不失一般性，我们将所有的观测按照 $G_i$ 进行排序，那么 $X$ 应该是一个分块对角矩阵：

$$X = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ \vdots & 0 & \cdots & 0 \\ 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ 0 & \vdots & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \\ 0 & 0 & \cdots & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix} = \begin{bmatrix} \iota_{N_1} & 0 & \cdots & 0 \\ 0 & \iota_{N_2} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \iota_{N_g} \end{bmatrix}$$

其中 $N_j$ 为 $G_i = j$ 组的观测个数。







## 参数回归：条件期望的估计

回忆条件期望的定义：

$$\mathbb{E}(y|x) = \arg \min_{h \in L^2} \left\{ \mathbb{E} \left[ (y - h(x))^2 \right] \right\}$$

如果假设函数 $h(x) = x'\beta$ ，即 $x$ 的一个线性函数，那么上式变为：

$$\beta_0 = \arg \min_{\beta} \left\{ \mathbb{E} \left[ (y - x'\beta)^2 \right] \right\}$$

从而 $\mathbb{E}(y|x) = x'\beta_0$ ，使用样本平均代 $\left(\frac{1}{N} \sum\right)$ 替总体期望 $(\mathbb{E})$ ，得到：

$$\hat{\beta} = \arg \min_{\beta} \frac{1}{N} \sum_{i=1}^N (y_i - x_i'\beta)^2$$

即普通最小二乘估计量（ordinary least squares）。













# 最小二乘的统计性质

我们现在将最小二乘估计 $\hat{\beta}$ 看成是总体回归方程中真值 $\beta_0$ 的一个估计。在此基础上，我们继续讨论最小二乘估计 $\hat{\beta}$ 的统计性质，包括 $\hat{\beta}$ 的无偏性、一致性。为此我们引入如下假设：

## 独立同分布假设

设样本 $(x'_i, y_i)'$ ,  $i = 1, 2, \dots, N$ 独立同分布。

注意独立同分布假设与异方差并不矛盾：

- 异方差指的是条件方差 $\mathbb{V}(y_i|x_i) = \sigma^2(x_i)$ 不为常数
- 然而同分布则意味着无条件方差 $\mathbb{V}(y_i)$ 不随 $i$ 的变化而变化，两者是不矛盾的。



# 无偏性

进一步，有：

$$\begin{aligned}\mathbb{E}(\hat{\beta}) &= \mathbb{E}\left[\mathbb{E}(\hat{\beta}|X)\right] \\ &= \mathbb{E}\left[\mathbb{E}\left(\beta_0 + \left[\sum_{i=1}^N (x_i x_i')\right]^{-1} \left(\sum_{i=1}^N x_i u_i\right) \mid X\right)\right] \\ &= \beta_0 + \mathbb{E}\left[\left[\sum_{i=1}^N (x_i x_i')\right]^{-1} \mathbb{E}\left(\sum_{i=1}^N x_i u_i \mid X\right)\right] \\ &= \beta_0 + \mathbb{E}\left[\left[\sum_{i=1}^N (x_i x_i')\right]^{-1} \sum_{i=1}^N (x_i \mathbb{E}(u_i | X))\right] \\ &= \beta_0\end{aligned}$$





# 正态性假设

## 条件正态性假设

设样本  $(y_i, x_i)'$ ,  $i = 1, \dots, N$  独立同分布, 且  $y_i$  给定  $x_i$  的条件分布为同方差的正态分布, 其条件期望为线性函数, 即:

$$y_i | x_i \sim N(x_i' \beta_0, \sigma^2)$$

或者等价地:

$$Y | X \sim N(X \beta_0, \sigma^2 I)$$

以上假设等价于假设误差项  $u_i | x_i \sim N(0, \sigma^2)$ , 或者  $u | X \sim N(0, \sigma^2 I)$ 。



# 条件极大似然估计

根据以上假设，条件密度函数为：

$$f(y_i|x_i) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(y_i - x_i'\beta_0)^2}{2\sigma^2}\right\}$$

因而条件似然函数为：

$$L(\beta, \sigma|y, x) = -\frac{N}{2} \ln(2\pi) - \frac{N}{2} \ln \sigma^2 - \sum_{i=1}^N \frac{(y_i - x_i'\beta)^2}{2\sigma^2}$$

最大化以上函数，得到：

$$\begin{cases} \hat{\beta} = \left(\sum_{i=1}^N x_i x_i'\right)^{-1} \left(\sum_{i=1}^N x_i y_i\right) = (X'X)^{-1} X'Y \\ \hat{\sigma}^2 = \frac{\sum_{i=1}^N (y_i - x_i'\hat{\beta})^2}{N} = \frac{\sum_{i=1}^N \hat{u}_i^2}{N} \end{cases}$$

其中  $\hat{u}_i = y_i - x_i'\hat{\beta}$  为残差。再次，我们得到了最小二乘估计量。

## 最小二乘与条件期望

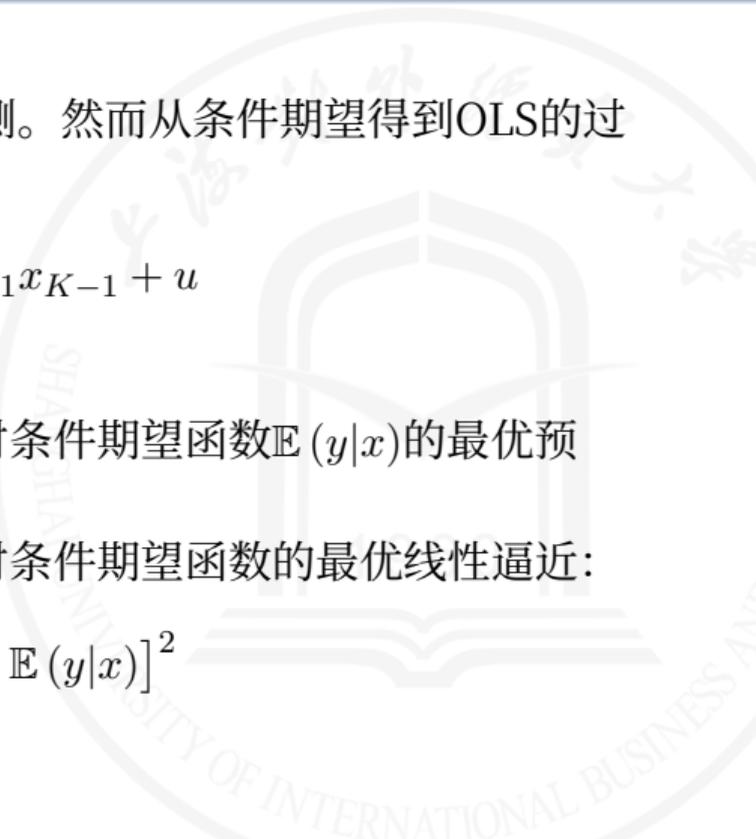
条件期望即我们使用自变量 $x$ 对因变量 $y$ 的最优预测。然而从条件期望得到OLS的过程中，我们假设了条件期望的**线性函数**形式：

$$y = \beta_0 + \beta_1x_1 + \dots + \beta_{K-1}x_{K-1} + u$$

然而这一条件未必满足：

- 如果条件期望函数的确为线性函数，OLS是对条件期望函数 $\mathbb{E}(y|x)$ 的最优预测；
- 如果条件期望函数不是线性函数，则OLS是对条件期望函数的最优线性逼近：

$$\beta_0 = \arg \min_{\beta} [x'\beta - \mathbb{E}(y|x)]^2$$





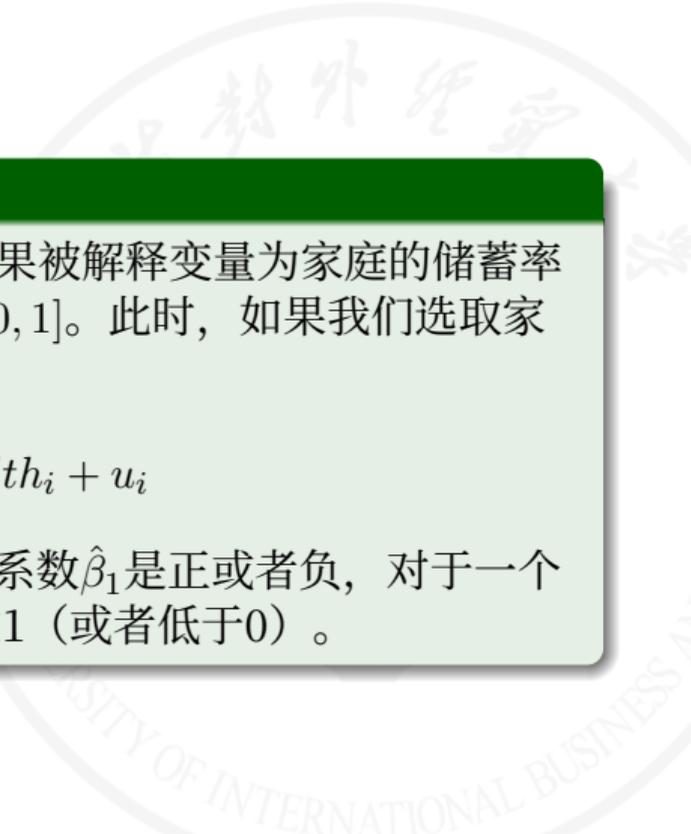
# 条件期望函数形式问题

## 支撑集问题

支撑集 (support) 即一个随机变量的取值范围。如果被解释变量为家庭的储蓄率 (saving\_rate)，我们知道  $\text{supp}(savin\_rate_i) = [0, 1]$ 。此时，如果我们选取家庭资产规模 (wealth) 作为解释变量：

$$savin\_rate_i = \beta_0 + \beta_1 \cdot wealth_i + u_i$$

由于  $\text{supp}(income_i) = [0, \infty)$ ，因而不管回归得到的系数  $\hat{\beta}_1$  是正或者负，对于一个资产规模足够大的家庭，总会使得预测的储蓄率超过1（或者低于0）。



# 条件期望函数形式问题

## 经济增长

如果令 $y_t$ 为时期 $t$ 时国家的GDP，根据索洛模型（Acemoglu, 2009, Chaper 3）， $y_t$ 满足如下关系式：

$$g_t = \beta_0 + \beta_1 \ln y_{t-1} + u_t$$

其中 $g_t = \ln y_t - \ln y_{t-1}$ 为GDP的对数增长率。根据上式，得到：

$$y_t = \exp \{ \beta_0 + (1 + \beta_1) \ln y_{t-1} + u_t \} = e^{\beta_0} y_{t-1}^{1+\beta_1} e^{u_t}$$

从而条件期望函数：

$$\mathbb{E}(y_t | y_{t-1}) = e^{\beta_0} y_{t-1}^{1+\beta_1} \mathbb{E}(e^{u_t} | y_{t-1})$$

因而条件期望函数为一个指数函数形式，而非线性函数。

# 条件期望函数形式问题

## 引力模型

在国际贸易理论中 (Head and Mayer, 2014) , 双边贸易与两个国家的GDP之间存在着被称为“引力模型”的关系, 即:

$$X_{ni} = G Y_i^a Y_n^b \phi_{ni}$$

其中下标*i*代表国家, 而*n*代表出口目的地国,  $X_{ni}$ 为两国之间的贸易额,  $G$ 为常数,  $Y$ 为国家的GDP,  $\phi_{ni}$ 则是两国之间贸易成本的函数。双边贸易额与GDP之间的关系并非简单的线性关系。

# 条件期望函数形式问题

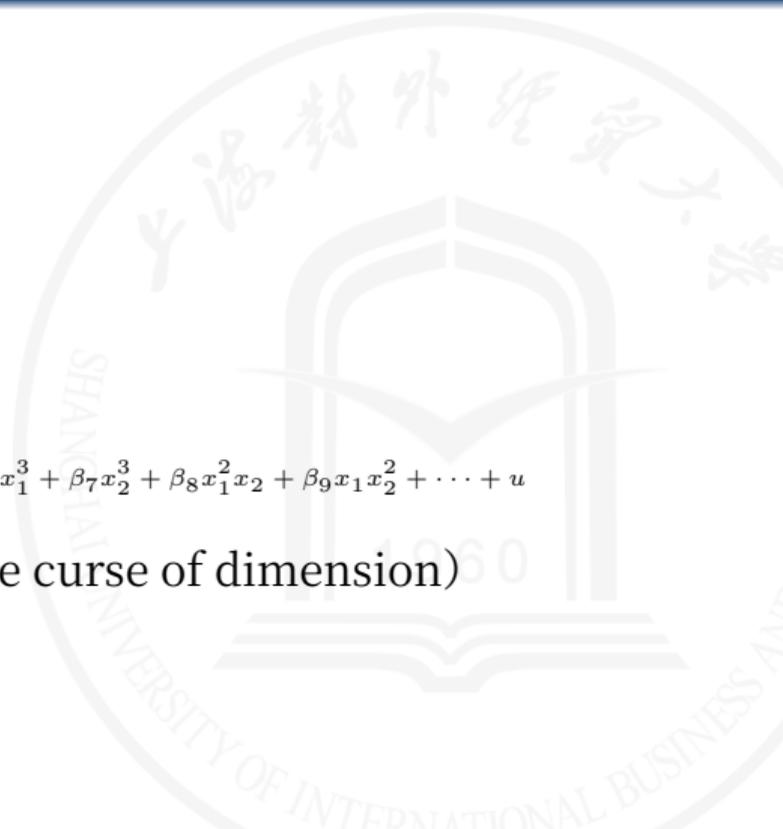
解决方案:

- 使用非参数回归 (kernel, sieve)
- 机器学习方法
- 引入多项式 (平方项、交叉项) :

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1^2 + \beta_4 x_2^2 + \beta_5 x_1 x_2 + \beta_6 x_1^3 + \beta_7 x_2^3 + \beta_8 x_1^2 x_2 + \beta_9 x_1 x_2^2 + \dots + u$$

- 以上方案可能有其缺点, 如维数的诅咒 (the curse of dimension)

更常用的解决方案——对数据进行变换:







# 对数变换

- 实际上对于一些“比例”型的数据，取对数有时也会有比较好的解释。
- 比如储蓄率例子中，储蓄率  $saving\_rate = \frac{saving}{income}$ ，如果我们将其取对数：

$$\ln saving\_rate = \ln saving - \ln income$$

从而如果将之前回归的被解释变量和解释变量取对数，即：

$$\ln saving\_rate_i = \beta_0 + \beta_1 \cdot \ln income_i + u_i$$

等价于：

$$\ln saving_i - \ln income_i = \beta_0 + \beta_1 \cdot \ln income_i + u_i$$

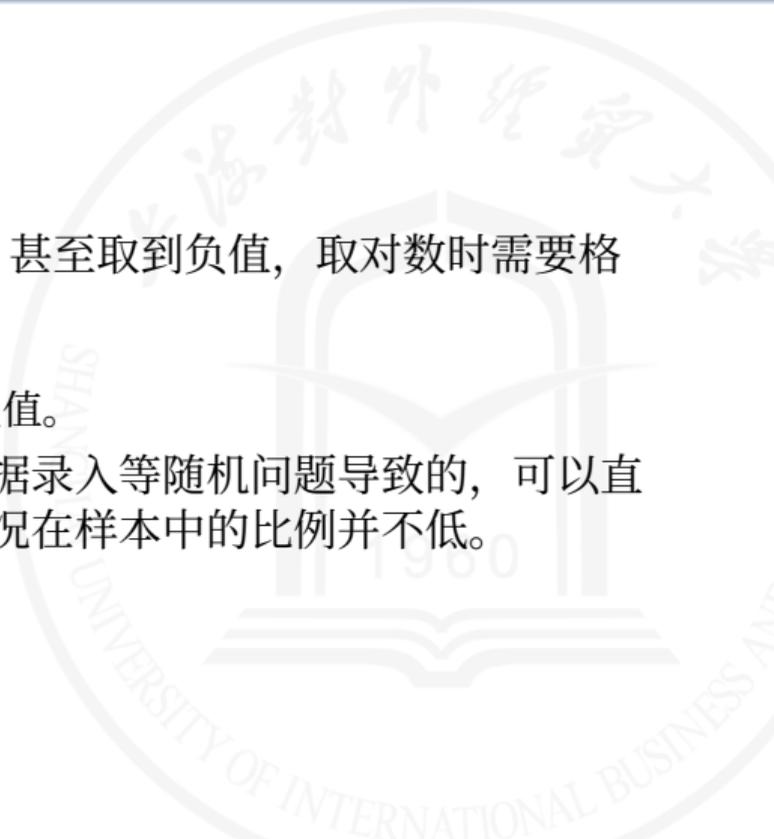
- 实际上我们可以证明（练习1.11），以上回归与以下回归是等价的：

$$\ln saving_i = \delta_0 + \delta_1 \cdot \ln income_i + u_i$$

且OLS估计量  $\hat{\beta}_0 = \hat{\delta}_0, \hat{\beta}_1 = \hat{\delta}_1 - 1$ 。

# 对数变换

- 最后需要注意的是，有些变量可能取到0值，甚至取到负值，取对数时需要格外小心。
  - 需要对收入取对数，然而很多人的收入为0；
  - 需要对净出口取对数，然而净出口可能为负值。
- 如果这些情况样本量比较少，可能是由于数据录入等随机问题导致的，可以直接忽略这些样本，然而更多的情况是这些情况在样本中的比例并不低。



# 对数变换

- 一些不太严谨的方法可以大概处理这些问题，比如：
  - 对于数据可能为0的问题，我们可以使用 $\ln(1+x)$ 取对数
    - 如此哪些 $x=0$ 的样本取“对数”之后， $\ln(1+0)=0$
    - 且这个变换是一个单调变换



# 对数变换

- 虽然上述方法被广泛应用，然而这些方法是不严谨的。
  - 原本对数变换的一个优良性质的可以“去量纲”，即不同量纲取对数只差一个常数
    - 例如收入 $x$ 如果以“万元”为单位，那么 $10000x$ 就是以“元”为单位，取对数后：

$$\ln(10000x) = \ln 10000 + \ln x$$

两者只差一个常数。

- 然而如果使用以上 $\ln(1+x)$ 的方法，该“对数”不再有此性质：

$$\ln(1+10000x) - \ln(1+x) = \ln \frac{1+10000x}{1+x} \neq C$$

- 为何是 $\ln(1+x)$ 而不是 $\ln(0.1+x)$ 或者 $\ln(10000+x)$ ?



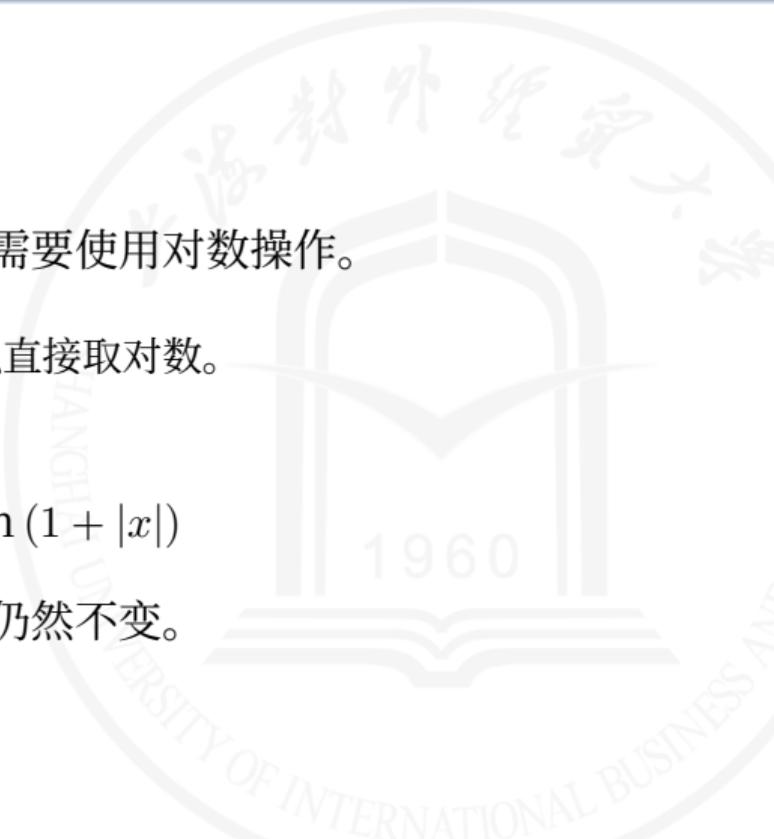
# 负数的对数变换

一些方法也许可以帮助解决这一问题

- 此外，还有些可以取到负值的变量仍然可能需要使用对数操作。
  - 比如，净出口额、人口净流入等变量
  - 取对数是一个合理的操作，然而负数不可以直接取对数。
- 对于可能为负的变量，一种方法是使用：

$$g(x) = \text{Sign}(x) \cdot \ln(1 + |x|)$$

该变换同样也是单调变换，经过变换后符号仍然不变。



# 负数的对数变换

或者，也可以使用反双曲正弦函数（如Caprettini和Voth, 2023）

- 双曲正弦函数的定义为：

$$\sinh(x) = \frac{e^x - e^{-x}}{2}$$

- 以上函数的定义域和值域都是 $\mathbb{R}$
- 当 $x \rightarrow \infty (-\infty)$ 时，以上函数趋向于 $\frac{e^x}{2} \left(-\frac{e^{-x}}{2}\right)$ ，从而当 $|x|$ 足够大时，以上函数近似于 $\frac{e^x}{2} \left(-\frac{e^{-x}}{2}\right)$
- 其反函数为：

$$\operatorname{arsinh}(x) = \ln\left(x + \sqrt{x^2 + 1}\right)$$

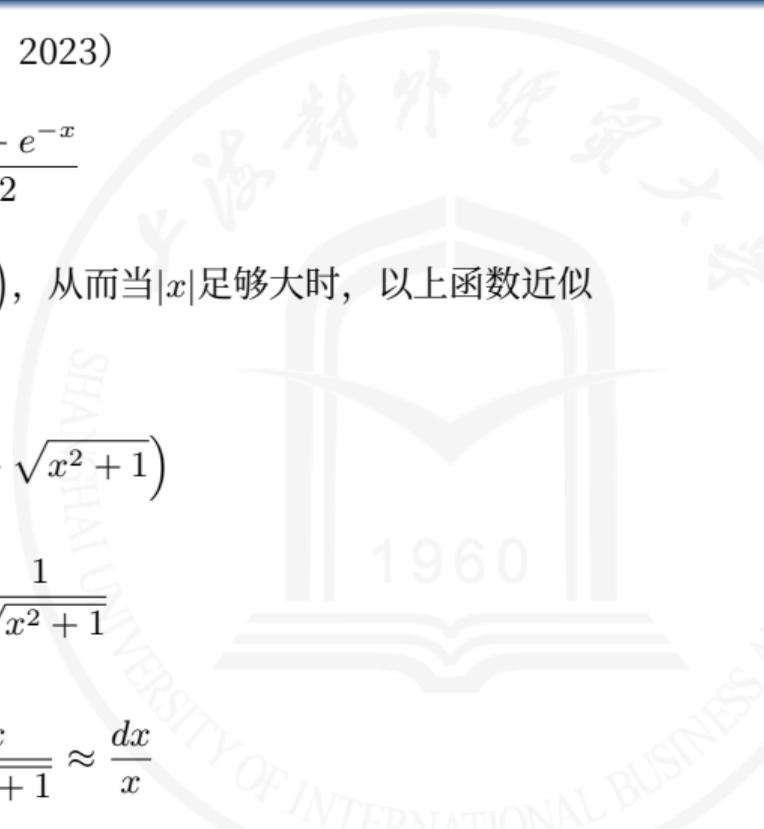
其导函数为：

$$\frac{d\operatorname{arsinh}(x)}{dx} = \frac{1}{\sqrt{x^2 + 1}}$$

当 $|x|$ 足够大时，以上导函数与 $\frac{1}{x}$ 近似，从而：

$$d\operatorname{arsinh}(x) = \frac{dx}{\sqrt{x^2 + 1}} \approx \frac{dx}{x}$$

也可以近似解释为百分比变动。



# 其他变换

其他变换:

- Box-Cox变换 (不推荐):

$$y(\lambda) = \begin{cases} \frac{y^\lambda - 1}{\lambda} & \lambda \neq 0 \\ \ln y & \lambda = 0 \end{cases}$$

- Logistic逆变换: 使用

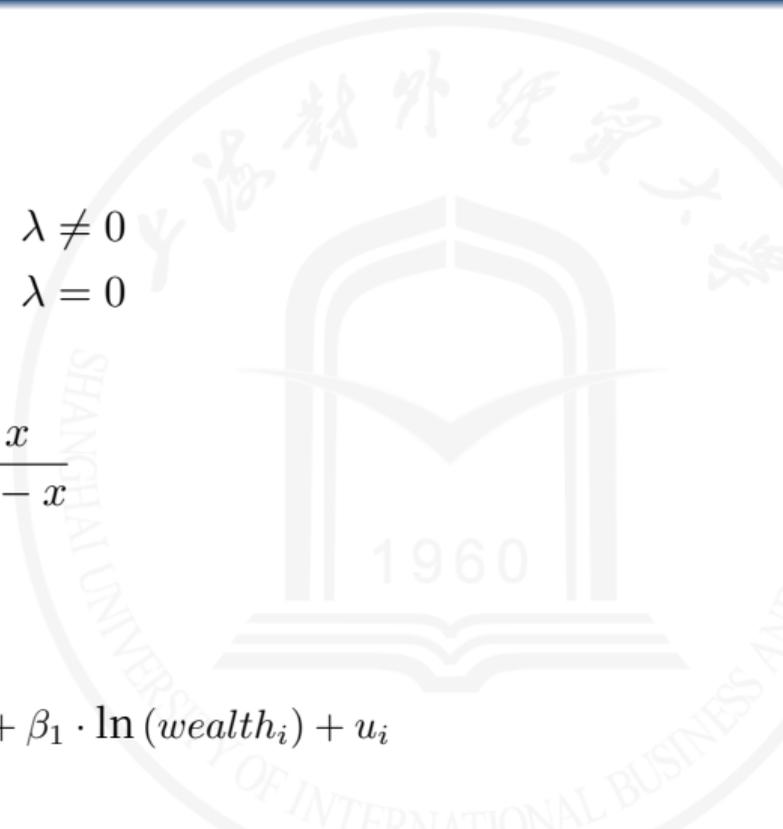
$$f(x) = \ln \frac{x}{1-x}$$

将(0, 1)区间上的实数映射到 $(-\infty, \infty)$ 上

- 比如对于储蓄率, 我们使用:

$$\ln \frac{saving\_rate_i}{1 - saving\_rate_i} = \beta_0 + \beta_1 \cdot \ln(wealth_i) + u_i$$

从而左边和右边取值范围都是 $\mathbb{R}$



# 条件期望的最优逼近

- 真实的条件期望函数我们是永远无法知道的，不过可以证明，线性回归仍然是条件期望函数的最优线性近似。
- 根据定义：

$$y_i = \mathbb{E}(y_i|x_i) + u_i$$

而最小二乘法的目标函数可以写为：

$$\begin{aligned}(y_i - x_i'\beta)^2 &= [y_i - \mathbb{E}(y_i|x_i) + \mathbb{E}(y_i|x_i) - x_i'\beta]^2 \\ &= [u_i + (\mathbb{E}(y_i|x_i) - x_i'\beta)]^2 \\ &= u_i^2 + (\mathbb{E}(y_i|x_i) - x_i'\beta)^2 + 2u_i(\mathbb{E}(y_i|x_i) - x_i'\beta)\end{aligned}$$

# 条件期望的线性逼近

由于

$$\mathbb{E} [u_i (\mathbb{E} (y_i|x_i) - x_i'\beta)] = \mathbb{E} (\mathbb{E} [u_i (\mathbb{E} (y_i|x_i) - x_i'\beta)] |x_i) = 0$$

从而:

$$\mathbb{E} [(y_i - x_i'\beta)^2] = \mathbb{E} (u_i^2) + (\mathbb{E} (y_i|x_i) - x_i'\beta)^2$$

其中第一项跟 $\beta$ 无关, 因而最小化 $\mathbb{E} [(y_i - x_i'\beta)^2]$ 等价于最小化 $(\mathbb{E} (y_i|x_i) - x_i'\beta)^2$ , 即 $x_i'\beta_0$ 是条件期望函数 $\mathbb{E} (y_i|x_i)$ 在均方误差标准下的最优线性逼近。

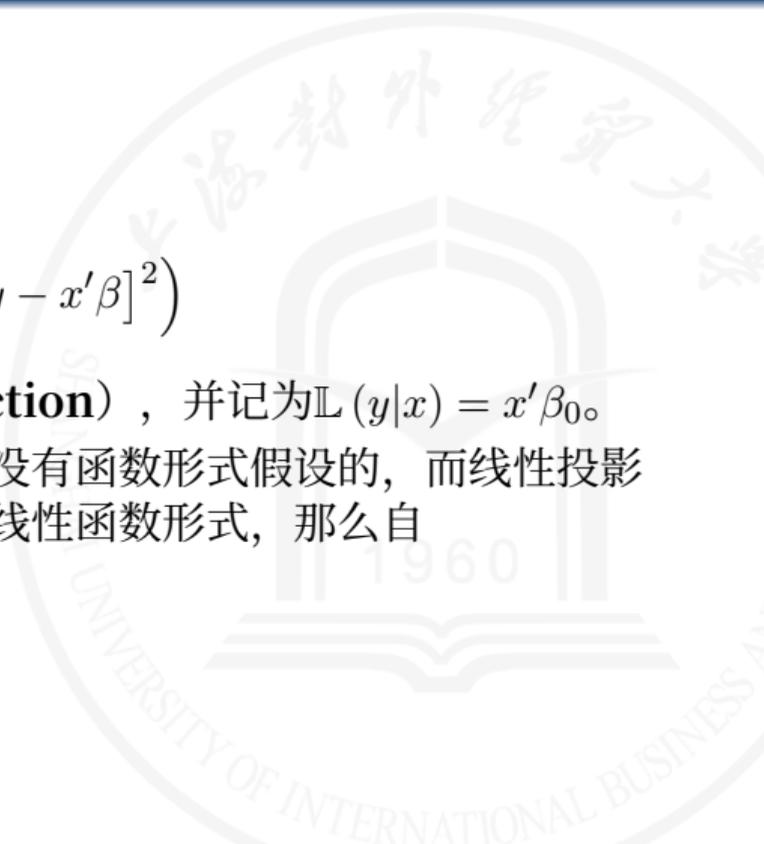
# 线性投影

- 对于最优化问题:

$$\beta_0 = \arg \min_{\beta} \mathbb{E} \left( [y - x'\beta]^2 \right)$$

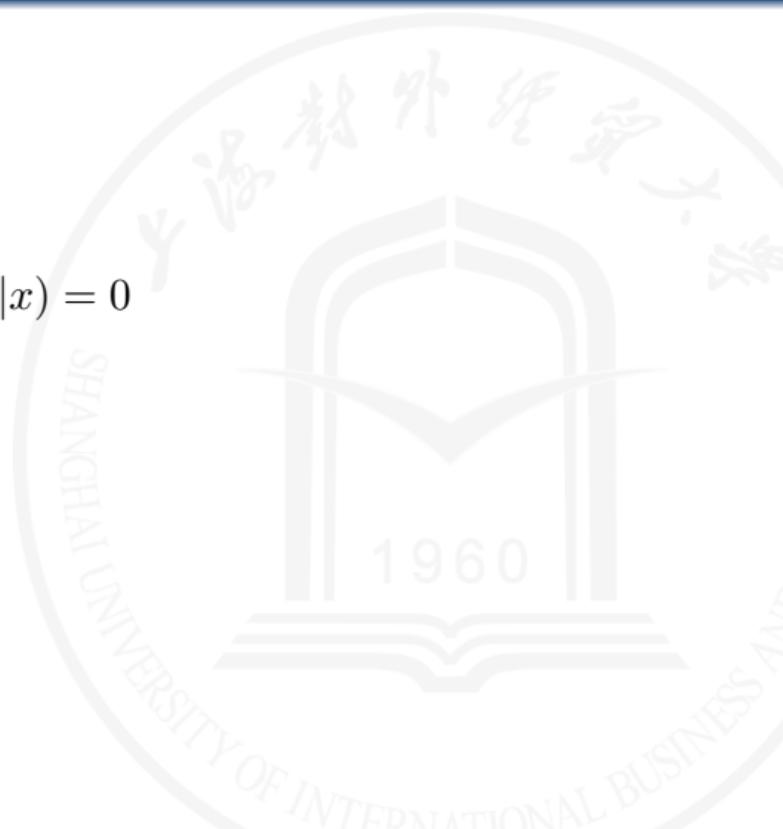
我们称 $x'_i\beta_0$ 其称为**线性投影 (linear projection)**，并记为 $\mathbb{L}(y|x) = x'\beta_0$ 。

- 线性投影区别于条件期望，因为条件期望是没有函数形式假设的，而线性投影有函数形式假设，如果真实的条件期望就是线性函数形式，那么自然 $\mathbb{E}(y|x) = \mathbb{L}(y|x)$ 。



# 线性投影的性质

- ① 令  $u = y - \mathbb{L}(y|x)$ , 那么  $\mathbb{E}(ux) = 0$ , 且  $\mathbb{L}(u|x) = 0$
- ②  $\mathbb{L}(a_1y_1 + a_2y_2|x) = a_1\mathbb{L}(y_1|x) + a_2\mathbb{L}(y_2|x)$
- ③  $\mathbb{L}(y|x) = \mathbb{L}[\mathbb{L}(y|x, w) | x]$
- ④  $\mathbb{L}(y|x) = \mathbb{L}[\mathbb{E}(y|x, w) | x]$



# 变换后的预测

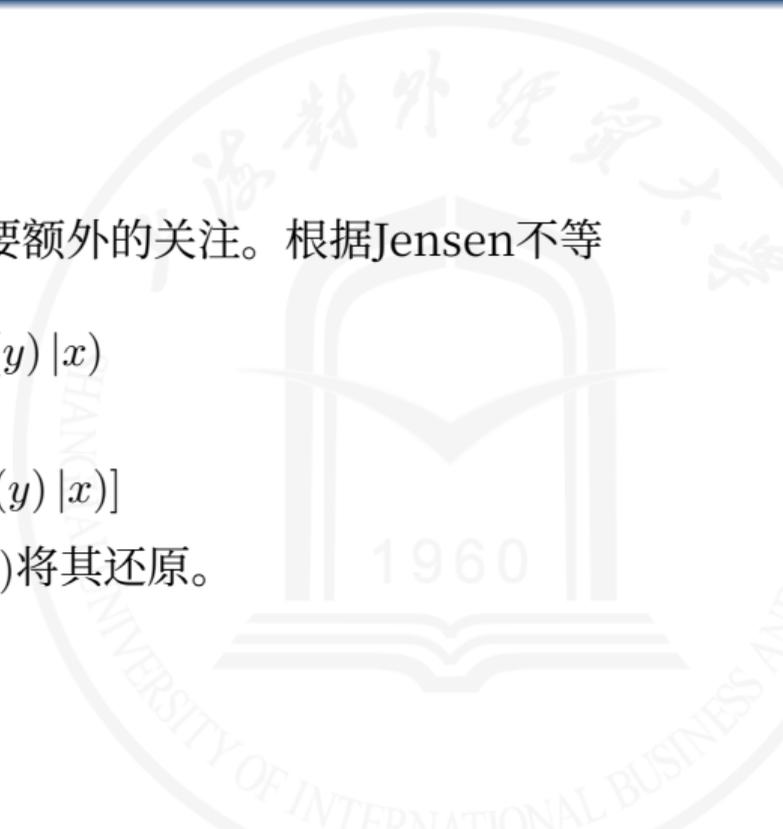
当我们使用 $y_i$ 的非线性变换时，对于 $y_i$ 的预测需要额外的关注。根据Jensen不等式，由于：

$$f(\mathbb{E}(y|x)) \neq \mathbb{E}(f(y)|x)$$

因而

$$\mathbb{E}(y|x) \neq f^{-1}[\mathbb{E}(f(y)|x)]$$

为了预测 $y$ 的值，不能先预测 $f(y)$ ，再使用 $f^{-1}(\cdot)$ 将其还原。



# 对数的预测

- 比如，如果我们设定如下方程：

$$\ln y = x'\beta + u$$

使用以上方程，我们得到的实际上是对于条件期望函数： $\mathbb{E}(\ln y|x)$ 的最优线性估计

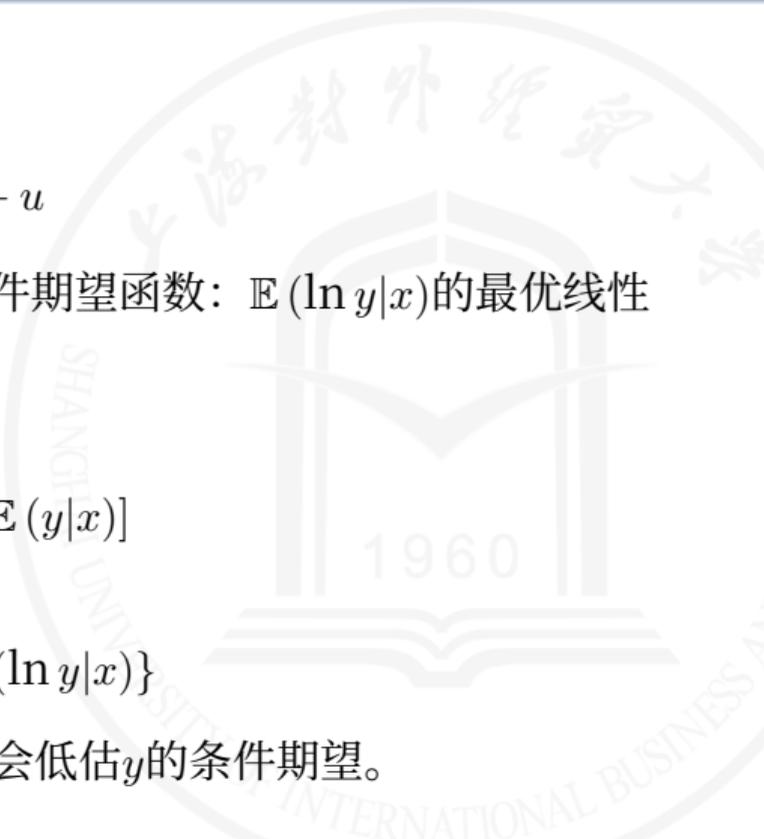
- 然而，根据Jensen不等式：

$$\mathbb{E}(\ln y|x) \leq \ln [\mathbb{E}(y|x)]$$

从而：

$$\mathbb{E}(y|x) \geq \exp \{ \mathbb{E}(\ln y|x) \}$$

因而如果我们使用 $\exp(x'\hat{\beta})$ 对 $y$ 进行预测，会低估 $y$ 的条件期望。



# 对数的预测

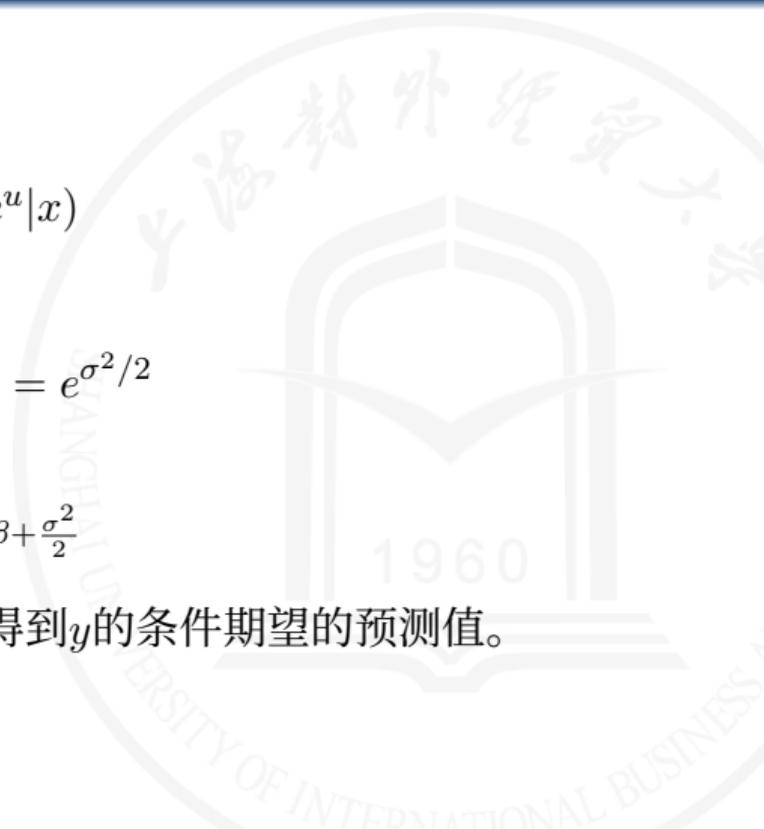
- 注意到:  $y = e^{x'\beta} e^u$  从而:  $\mathbb{E}(y|x) = e^{x'\beta} \mathbb{E}(e^u|x)$
- 如果假设  $u$  和  $x$  独立且  $u \sim N(0, \sigma^2)$ , 那么

$$\mathbb{E}(e^u|x) = \mathbb{E}(e^u) = e^{\sigma^2/2}$$

从而:

$$\mathbb{E}(y|x) = e^{x'\beta + \frac{\sigma^2}{2}}$$

将  $\beta$  和  $\sigma^2$  使用极大似然回归结果, 替代即可得到  $y$  的条件期望的预测值。







# 分步回归

- 对于总体回归:

$$y_i = x'_i \beta + u_i$$

如果我们将解释变量  $x_i$  分为两部分:  $x_i = [x'_{i1}, x'_{i2}]'$ , 那么总体回归方程可以写为:

$$y_i = x'_{i1} \beta_1 + x'_{i2} \beta_2 + u_i$$

- 如果我们将以上方程两边同时对  $x_{i2}$  求条件期望, 得到:

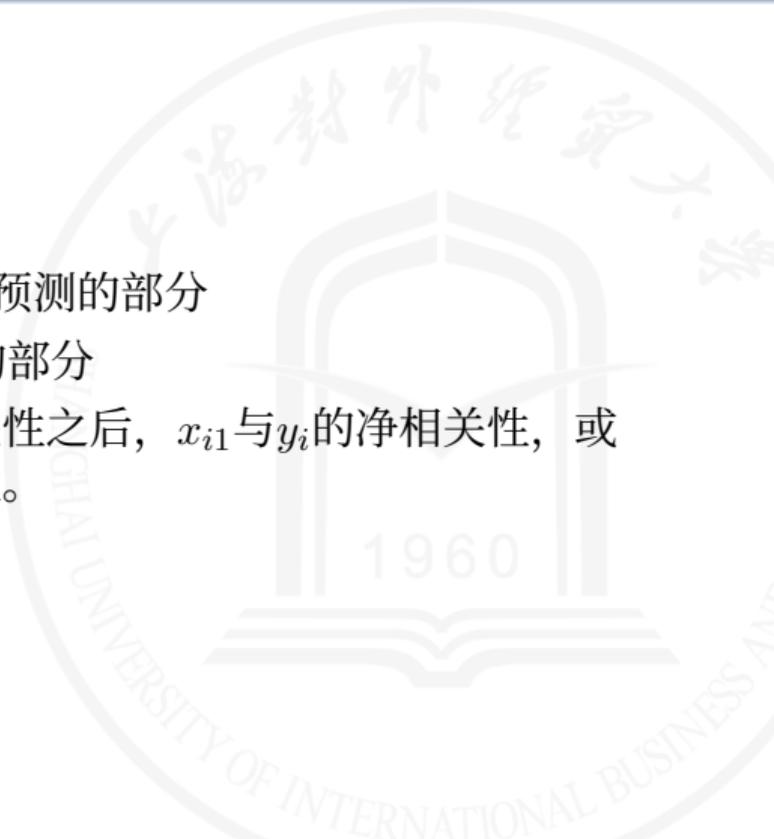
$$\begin{aligned} \mathbb{E}(y_i | x_{i2}) &= \mathbb{E}(x_{i1} | x_{i2})' \beta_1 + x'_{i2} \beta_2 + \mathbb{E}(u_i | x_{i2}) \\ &= \mathbb{E}(x_{i1} | x_{i2})' \beta_1 + x'_{i2} \beta_2 \end{aligned}$$

在总体回归方程两边同时减去上式, 得到:

$$\begin{aligned} y_i - \mathbb{E}(y_i | x_{i2}) &= x'_i \beta + u_i - \mathbb{E}(x_{i1} | x_{i2})' \beta_1 - x'_{i2} \beta_2 \\ &= [x_{i1} - \mathbb{E}(x_{i1} | x_{i2})]' \beta_1 + u_i \end{aligned}$$

# 分步回归

- 在上式中,  $y_i - \mathbb{E}(y_i|x_{i2})$ 代表 $y_i$ 中 $x_{i2}$ 所不能预测的部分
- 而 $x_{i1} - \mathbb{E}(x_{i1}|x_{i2})$ 代表 $x_{i1}$ 中 $x_{i2}$ 所不能预测的部分
- 因而系数 $\beta_1$ 代表的是排除 $x_{i2}$ 与 $x_{i1}$ 和 $y_i$ 的相关性之后,  $x_{i1}$ 与 $y_i$ 的净相关性, 或者在保持 $x_{i2}$ 不变的条件下,  $x_{i1}$ 与 $y_i$ 的相关性。



# 分步回归

- 类似的结果对于线性投影也成立:

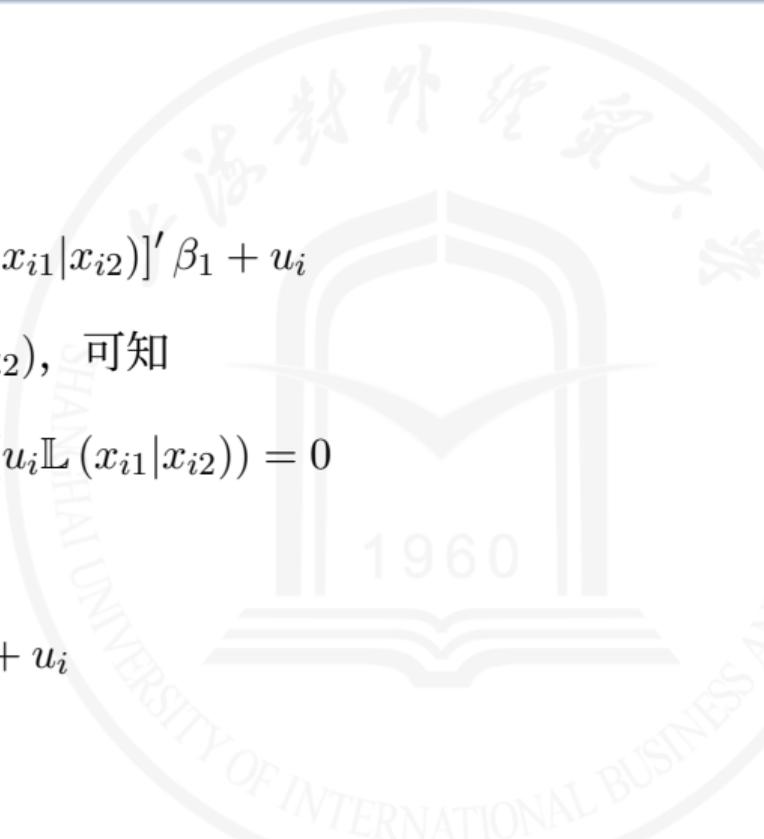
$$y_i - \mathbb{L}(y_i|x_{i2}) = [x_{i1} - \mathbb{L}(x_{i1}|x_{i2})]' \beta_1 + u_i$$

- 令  $e_{iy} = y_i - \mathbb{L}(y_i|x_{i2})$ ,  $e_{i,x_1} = x_{i1} - \mathbb{L}(x_{i1}|x_{i2})$ , 可知

$$\mathbb{E}(u_i e_{i,x_1}) = \mathbb{E}(u_i x_{i1}) - \mathbb{E}(u_i \mathbb{L}(x_{i1}|x_{i2})) = 0$$

- 从而相减后得到的式子可以写为:

$$e_{iy} = e'_{i,x_1} \beta_1 + u_i$$

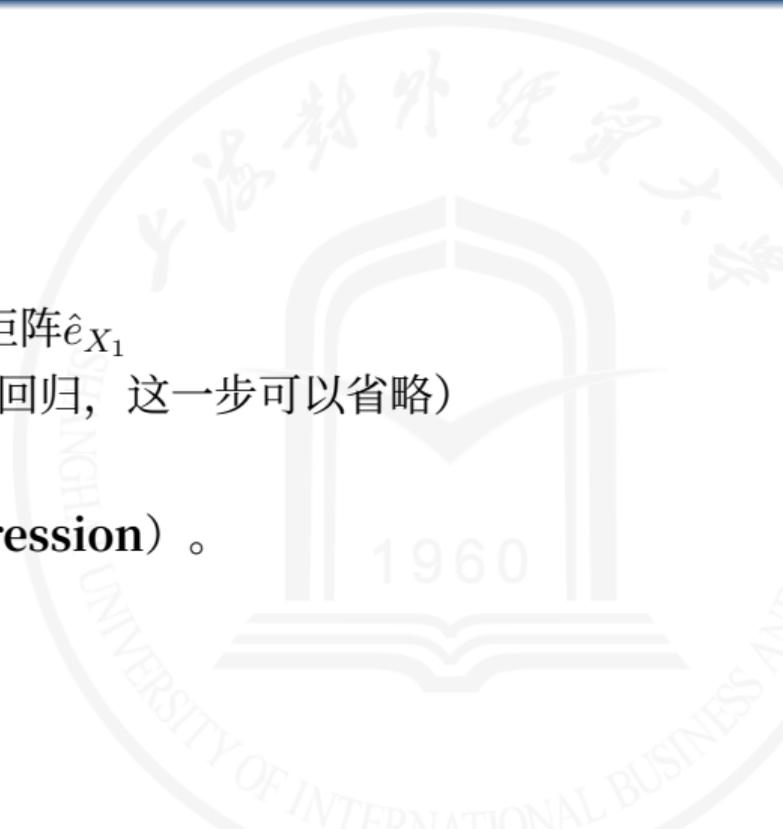


# 分步回归

为了计算 $\beta_1$ ，我们可以通过如下步骤进行计算：

- ① 使用对 $X_1$ 的每一列对 $X_2$ 做回归，得到残差矩阵 $\hat{e}_{X_1}$
- ② 使用 $Y$ 对 $X_2$ 做回归，得到残差 $\hat{e}_y$ （对于线性回归，这一步可以省略）
- ③ 使用 $\hat{e}_y$ 对 $\hat{e}_{X_1}$ 做回归，得到系数 $\hat{\beta}_1$

如上的步骤被称为**分步回归**（**partitioned regression**）。



# 分步回归

## Frisch-Waugh-Lovell定理

使用以上分步回归得到的 $\hat{\beta}_1$ 与式：

$$y_i = x'_{i1}\beta_1 + x'_{i2}\beta_2 + u_i$$

的最小二乘回归得到的 $\hat{\beta}_1$ 是完全等价的。

证明见讲义。

# 分步回归示例

## 消费与收入

partitioned\_regression.do使用对数收入和对数资产、是否农村预测了对数消费。

- 如果直接进行回归，得到的对数收入的系数为0.2088
- 而使用分步回归得到的系数同样为0.2088，两者严格相等。
- 此外，代码中还计算了 $\hat{e}_y$ 和 $\hat{e}_{x_1}$ 的相关系数，即偏相关系数，这度量了资产规模、是否农村等变量之后对数收入与对数消费之间的相关性，约为0.3288，而不排除这些的相关系数为0.5224。