

虚拟变量回归

不同教育程度的收入

同样使用2017年CHFS数据，对不同教育程度的收入进行分解。在数据集中，变量a2012代表教育程度，比如a2012=0时表示文盲，=1代表小学，=9代表博士等。我们使用如下程序计算分组差异或者分组平均：

条件极大似然估计

根据以上假设，条件密度函数为：

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

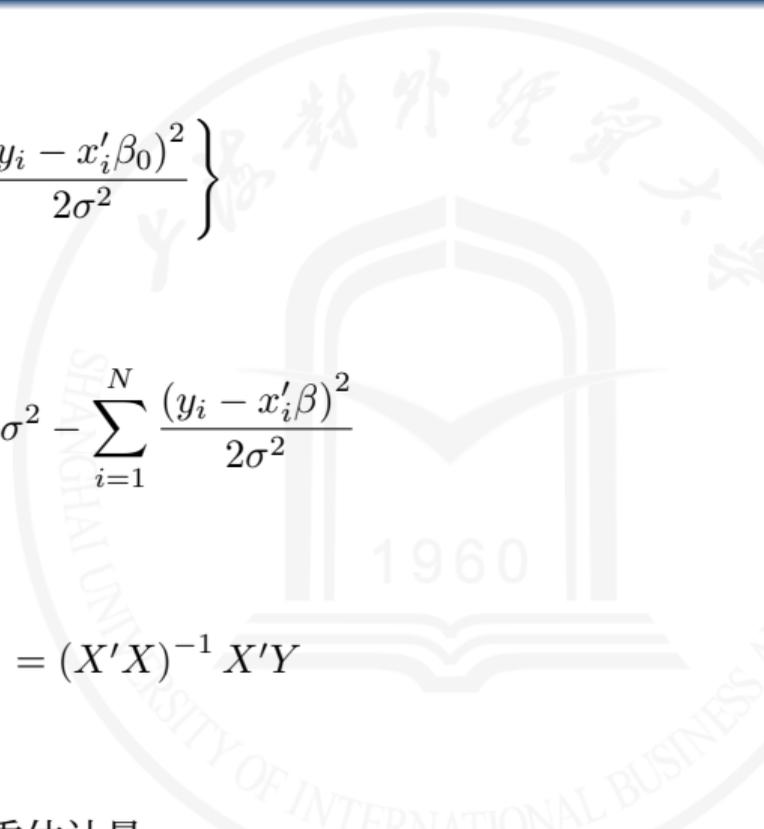
因而条件似然函数为：

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

最大化以上函数，得到：

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

其中 $\hat{e}_i = y_i - x'_i\hat{\beta}$ 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

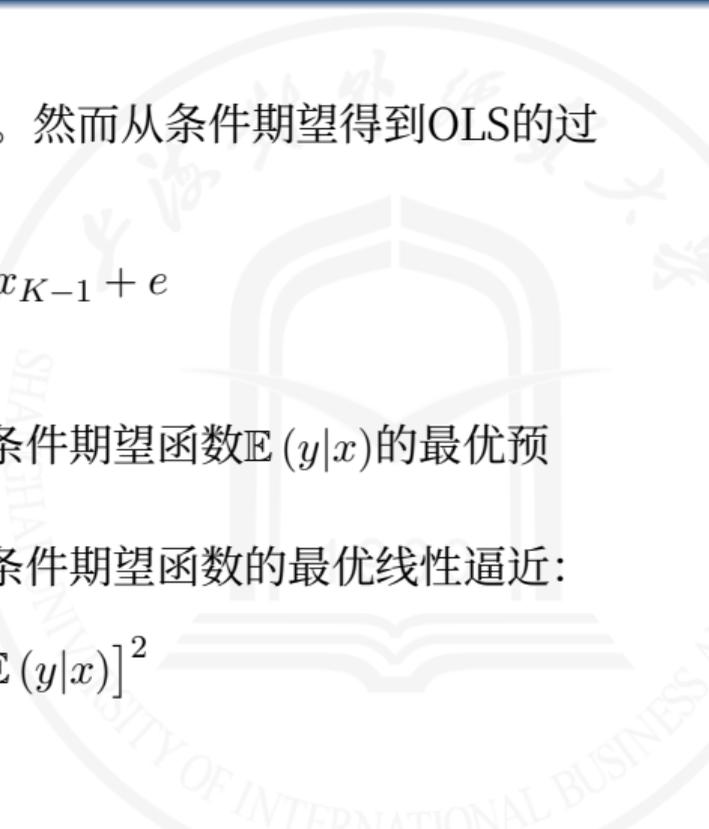
条件期望即我们使用自变量 x 对因变量 y 的最优预测。然而从条件期望得到OLS的过程中，我们假设了条件期望的线性函数形式：

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_{K-1} x_{K-1} + e$$

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：

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



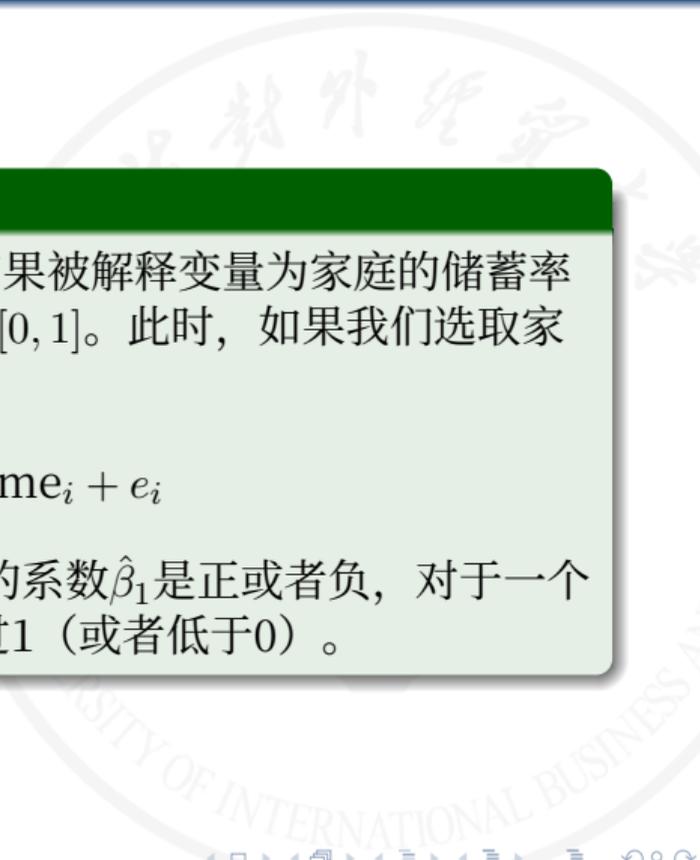
条件期望函数形式问题

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 $\text{supp}(\text{savin_rate}_i) = [0, 1]$ 。此时，如果我们选取家庭资产规模 (wealth) 作为解释变量：

$$\text{savin_rate}_i = \beta_0 + \beta_1 \cdot \text{income}_i + e_i$$

由于 $\text{supp}(\text{income}_i) = [0, \infty)$ ，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：

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

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：

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

从而条件期望函数：

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

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

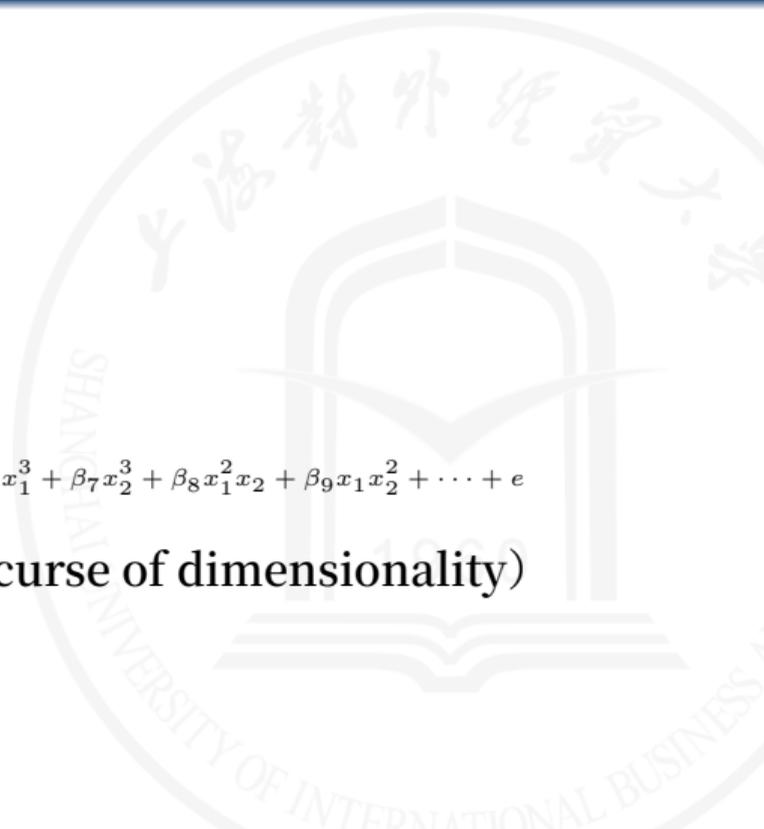
解决方案:

- 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) :

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

- 以上方案可能有其缺点, 如维数诅咒 (the curse of dimensionality)

更常用的解决方案——对数据进行变换:



对数变换

对数变换的最常用的变换：

$$x \rightarrow \ln(x)$$

进行对数变换的理由：

- 将 $(0, +\infty)$ 变换到 $(-\infty, \infty)$ 上，更符合取值范围的逻辑
- 有偏的分布，如收入、财富等右偏分布，取对数之后可以得到一个近似对称的分布
 - 解释变量和被解释变量都是对称的更符合直觉
 - 一些右偏的分布比较难以线性组合出对称的分布
- 理论预期。如上例中的GDP和出口，变量取对数后都可以变成线性函数关系，这是经济学理论预期的。

对数变换

- 具有弹性 (elasticity) 解释, 经过取对数后, 其变化可以解释为百分比变化:

$$d \ln y = \frac{dy}{y}$$

人口、GDP等具有比较平稳的增长率, 取对数更容易与其他变量之间满足线性关系。

- 比如对于GDP: $y_t \propto (1 + \beta_1)^t y_0$, 取对数后:

$$\ln y_t = C + t \log(1 + \beta_1) + \log y_0$$

更容易与其他变量形成线性关系

- 如果GDP是指数增长, 那么:

$$X_{nit} \propto (1 + \beta_{i1})^{ta} y_{i0}^a (1 + \beta_{n1})^{tb} y_{n0}^b$$

从而出口也类似, 取对数后容易与其他变量形成线性关系

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 $\text{saving_rate} = \frac{\text{saving}}{\text{income}}$ ，如果我们将其取对数：

$$\ln \text{saving_rate} = \ln \text{saving} - \ln \text{income}$$

从而如果将之前回归的被解释变量和解释变量取对数，即：

$$\ln \text{saving_rate}_i = \beta_0 + \beta_1 \cdot \ln \text{income}_i + e_i$$

等价于：

$$\ln \text{saving}_i - \ln \text{income}_i = \beta_0 + \beta_1 \cdot \ln \text{income}_i + e_i$$

- 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：

$$\ln \text{saving}_i = \delta_0 + \delta_1 \cdot \ln \text{income}_i + e_i$$

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

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：

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

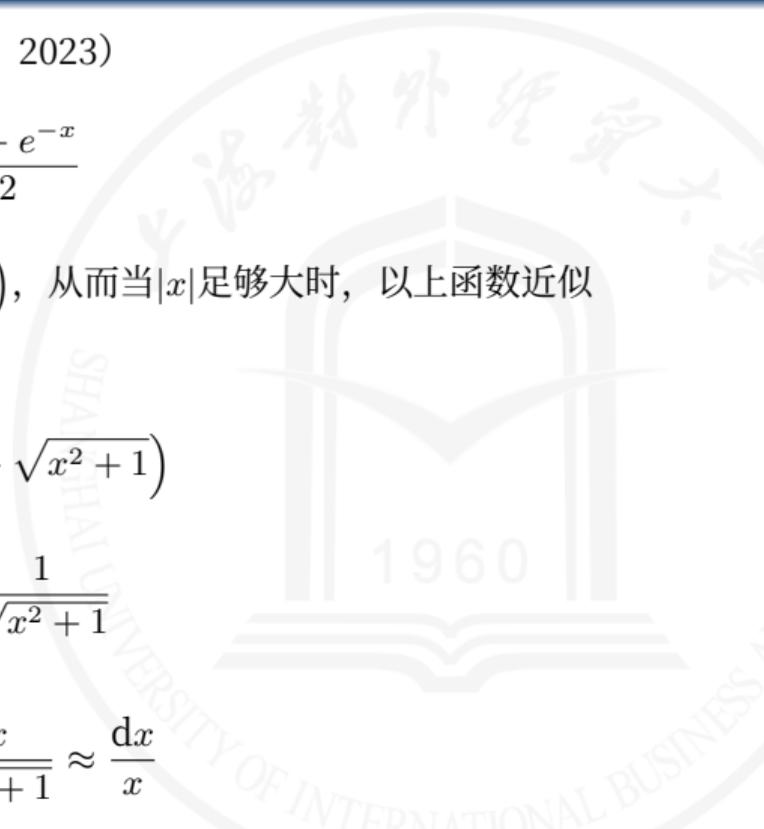
其导函数为：

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

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

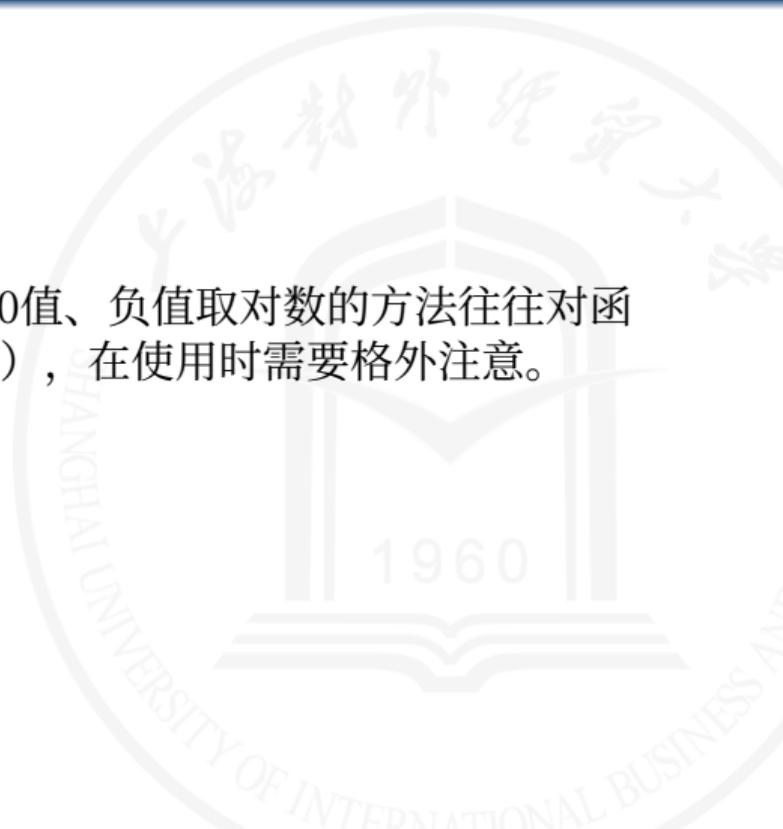
$$d\operatorname{arsinh}(x) = \frac{dx}{\sqrt{x^2 + 1}} \approx \frac{dx}{x}$$

也可以近似解释为百分比变动。



Log with Zeros

- 最后，仍然需要再次提示的是，以上的解决0值、负值取对数的方法往往对函数形式有很强的假设（Chen和Roth, 2023），在使用时需要格外注意。
- 其他方法：
 - 拟泊松回归



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

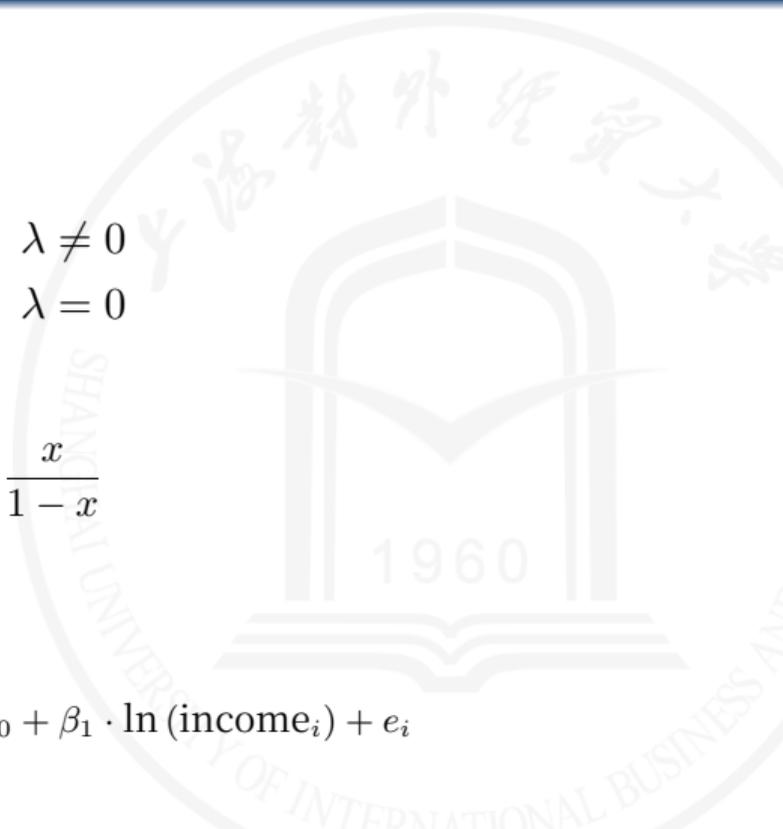
$$f(x) = \ln \frac{x}{1-x}$$

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:

$$\ln \frac{\text{saving_rate}_i}{1 - \text{saving_rate}_i} = \beta_0 + \beta_1 \cdot \ln(\text{income}_i) + e_i$$

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：

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

而最小二乘法的目标函数可以写为：

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

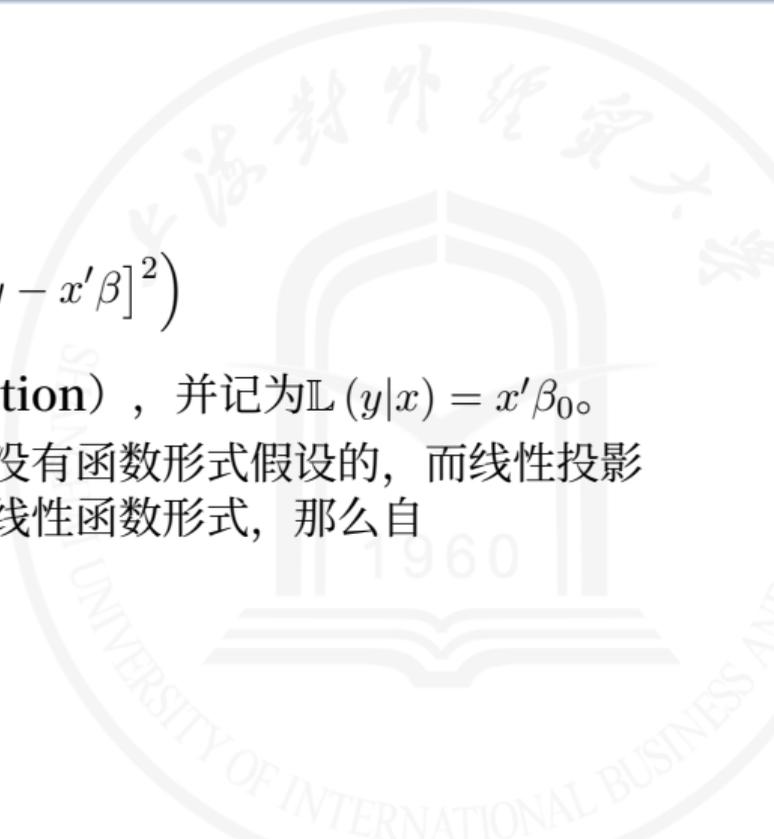
线性投影

- 对于最优化问题:

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

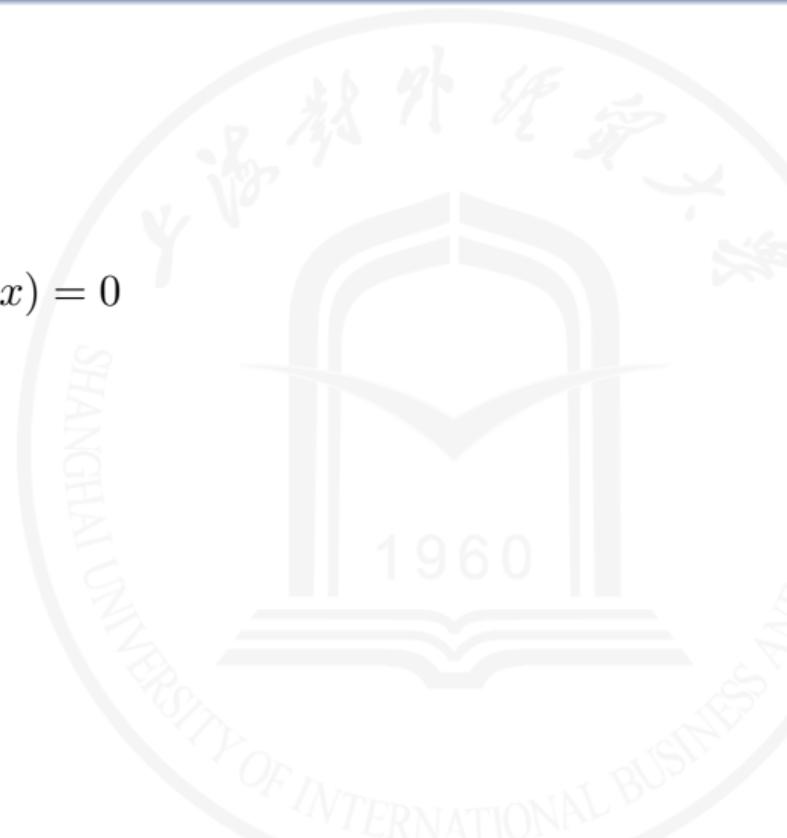
我们称 $x'_i\beta_0$ 其称为**线性投影 (linear projection)**，并记为 $\mathbb{L}(y|x) = x'\beta_0$ 。

- 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 $e = y - \mathbb{L}(y|x)$, 那么 $\mathbb{E}(ex) = 0$, 且 $\mathbb{L}(e|x) = 0$
- 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 $\mathbb{L}(y|x) = \mathbb{L}[\mathbb{E}(y|x, w) | x]$



对数的预测

- 注意到: $y = \exp(x'\beta) \cdot \exp(e)$, 从而

$$\mathbb{E}(y|x) = \exp(x'\beta) \mathbb{E}(\exp(e)|x)$$

- 如果假设 e 和 x 独立且 $e \sim \mathcal{N}(0, \sigma^2)$, 那么

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

从而:

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

将 β 和 σ^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。

