



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：

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

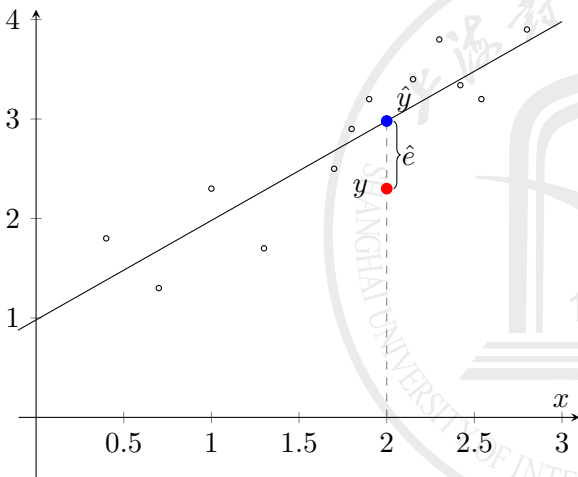
- 如果给定一个输入，可能有不同的值与之相对应，那么是相关关系，比如身高都为170的人，体重完全可能不同

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



最小二乘法

为了使得预测误差（残差）更小，一个最常见的办法是最小化均方误差（mean squared error）：

- 首先将残差计算平方，从而当预测误差为0（完美预测）时，残差的平方为0，否则不管高估还是低估，都是残差平方越小越好
- 其次对所有样本的残差平方求平均：

$$\frac{1}{N} \sum_{i=1}^N e_i^2 = \frac{1}{N} \sum_{i=1}^N (y_i - \alpha - \beta x_i)^2$$

进而，我们只要选择一个 α, β 使得以上的残差平方和最小化即可：

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

最小二乘法

- 解上述最小化问题，得到：

$$\begin{cases} \frac{\partial \sum_{i=1}^N (y_i - \alpha - \beta x_i)^2}{\partial \alpha} = -2 \sum_{i=1}^N (y_i - \alpha - \beta x_i) = 0 \\ \frac{\partial \sum_{i=1}^N (y_i - \alpha - \beta x_i)^2}{\partial \beta} = -2 \sum_{i=1}^N (y_i - \alpha - \beta x_i) x_i = 0 \end{cases}$$

- 化简上述问题，得到：

$$\begin{cases} \alpha = \bar{y} - \beta \bar{x} \\ \alpha \bar{x} = \frac{1}{N} \left(\sum_{i=1}^N x_i y_i - \beta \sum_{i=1}^N x_i^2 \right) \end{cases}$$

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

对于回归：

$$\hat{y}_i = \alpha + \beta \times x_i$$

x_i 只能取0/1两个值，此时我们称 x_i 为虚拟变量（dummy variable）

- 令 N_0 为样本中 $x_i = 0$ 的个数， N_1 为样本中 $x_i = 1$ 的个数
- 记 \bar{y}_1 为对应于 $x_i = 1$ 的 y_i 的均值，记 \bar{y}_0 为对应于 $x_i = 0$ 的 y_i 的均值，那么我们有（how?）：

$$\begin{cases} \hat{\beta} = \bar{y}_1 - \bar{y}_0 \\ \hat{\alpha} = \bar{y}_0 \end{cases}$$



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 $x_{i1} = 1$ 。



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：

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

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

注意以上我们使用了矩阵 $X'X$ 的逆矩阵，由于 $\text{rank}(X'X) = \text{rank}(X)$ ，因而 $X'X$ 可逆性要求 $\text{rank}(X) = K$ ，即要求矩阵 X 是列满秩的（同时样本量 $N \geq K$ ）。为此我们必须引入如下假设：

识别条件

矩阵 X 为列满秩矩阵，即 $\text{rank}(X) = K$ 。

矩阵 X 是列满秩的意味着：

- X 的列数小于行数，即 $K < N$ 。
- X 的任何一列不能被其他列线性表示出来——无完全共线性（perfect colinearity）



虚拟变量

- 由于在虚拟变量定义中， $\sum_{j=0}^g d_{ij} = 1$ ，即7个虚拟变量线性组合出了常数项，所以在包含常数项的回归中， d_{i1}, \dots, d_{ig} 不能同时出现。
- 解决以上问题的方法是忽略掉常数项，或者忽略掉 d_{i1}, \dots, d_{ig} 中的任何一个变量，以上两种方法都可以使得矩阵 $X'X$ 可逆，当然在现实中我们经常使用第二种方法，即抛弃其中的一个分组虚拟变量。



虚拟变量回归

不同教育程度的收入

VARIABLES	(1) p_income	(2) p_income	(3) p_income
edu1	-33,868*** (6,705)	8,211*** (1,092)	
edu2	-28,200*** (6,661)	13,879*** (781.4)	5,669*** (1,343)
edu3	-27,527*** (6,638)	14,551*** (554.9)	6,341*** (1,225)
edu4	-27,441*** (6,670)	14,638*** (850.1)	6,428*** (1,384)
edu5	-18,867*** (6,743)	23,212*** (1,306)	15,002*** (1,702)
edu6	-17,647*** (6,773)	24,432*** (1,451)	16,221*** (1,817)
o.edu7	-		
edu7		42,079*** (6,615)	33,868*** (6,705)
Constant	42,079*** (6,615)		8,211*** (1,092)
Observations	3,226	3,226	3,226
R-squared	0.042	0.383	0.042

Standard errors in parentheses

*** p<0.01, ** p<0.05, * p<0.1

条件期望与回归

重要假设：

线性函数假设

假设随机变量 y 给定 x 的条件期望 $\mathbb{E}(y|x)$ 为线性函数，

即： $h(x) = x'\beta$

那么：

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

即为真实参数，OLS估计量 $\hat{\beta}$ 为 β_0 的估计量，而条件期望的估计为：

$$\widehat{\mathbb{E}(y|x)} = x'\hat{\beta}$$



总体回归方程

- u_i 与 x_i 虽然是均值独立的, $\mathbb{E}(u_i|x_i) = 0$, 但是并没有对 $\mathbb{E}(u_i^2|x_i)$ 做任何假设, 因而 $\mathbb{E}(u_i^2|x_i)$ 或者条件方差 $\mathbb{V}(u_i|x_i)$ 可以是 x_i 的任意函数。
- 如果 $\mathbb{V}(u_i|x_i) = \mathbb{E}(u_i^2|x_i)$ 不为常数, 那么我们称 u_i 具有异方差 (heteroscedasticity) ;
- 如果 $\mathbb{V}(u_i|x_i) = \mathbb{E}(u_i^2|x_i) = \mathbb{E}(u_i^2)$, 那么我们称 u_i 具有同方差 (homoscedasticity) 性质。
- 均值独立意味着 x_i 对 u_i 没有预测能力, 而异方差的存在意味着 x_i 对 u_i 的方差仍然具有预测能力, 两者并不矛盾。
- 注意, 由于:

$$\mathbb{V}(y_i|x_i) = \mathbb{V}(x_i'\beta_0 + u_i|x_i) = \mathbb{V}(u_i|x_i)$$

即 $y_i|x_i$ 的条件方差也就是 $u_i|x_i$ 的条件方差。



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



OLS的统计性质

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

其中 $u = [u_1, \dots, u_N]'$ 。



无偏性

进一步，有：

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



一致性

如果 (x_i, y_i) 是独立同分布的, 且 $\mathbb{E}(x_{ik}^2) < \infty$ 以及 $\mathbb{E}|x_{ik}u_i| < \infty, k = 1, \dots, K$, 根据大数定律, 有:

$$\begin{cases} \frac{1}{N} \sum_{i=1}^N (x_i x_i') \xrightarrow{p} \mathbb{E}(x_i x_i') \\ \frac{1}{N} \sum_{i=1}^N (x_i u_i) \xrightarrow{p} \mathbb{E}(x_i u_i) = 0 \end{cases}$$

由于矩阵求逆为连续映射, 因而:

$$\hat{\beta} - \beta_0 = \left[\frac{1}{N} \sum_{i=1}^N (x_i x_i') \right]^{-1} \left(\frac{1}{N} \sum_{i=1}^N x_i u_i \right) \xrightarrow{p} \mathbb{E}(x_i x_i')^{-1} \mathbb{E}(x_i u_i) = 0$$

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 $u | X \sim N(0, \sigma^2 I)$ 。

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：

$$y = \beta_0 + \beta_1 x_1 + \cdots + \beta_{K-1} x_{K-1} + u$$

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

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：

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

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， ϕ_{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。

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 $saving_rate = \frac{saving}{income}$ ，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 $\hat{\beta}_0 = \hat{\delta}_0$, $\hat{\beta}_1 = \hat{\delta}_1 - 1$ 。



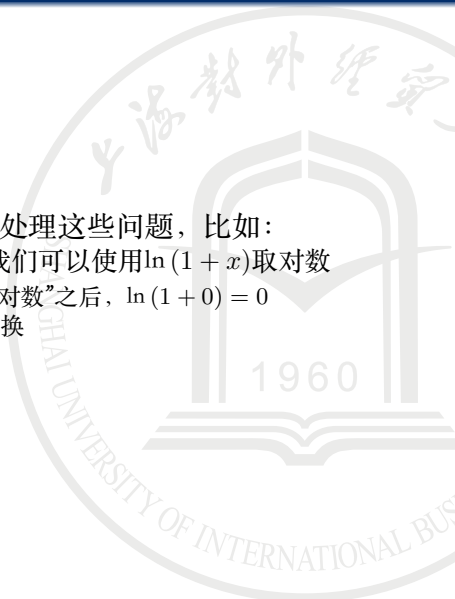
对数变换

- 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 $\ln(10000+x)$?

对数变换

一些方法也许可以帮助解决这一问题

- 比如如果需要对被解释变量 y 取对数，我们发现很多的 $y = 0$:
 - 使用Tobit一类回归，比如第I类Tobit等。
 - 在国际贸易中，Eaton和Tamura (1994) 在处理引力模型时，提出可以使用 $\ln(a + X_{ni})$ 作为被解释变量，而将 a 作为一个待估参数，即ET Tobit (Head和Mayer, 2014)。
- 如果需要取对数的变量为解释变量，一个简单的处理方法是定义两个新的变量：

$$\ln x = \begin{cases} \ln x & \text{if } x > 0 \\ 0 & \text{if } x = 0 \end{cases}$$

$$d = 1 \{x = 0\}$$

从而使用回归：

$$y_i = \beta_0 + \beta_1 \ln x_i + \beta_2 d_i + u_i$$



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。

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：

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

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 β 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。

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

- 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 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。

