

第一次作业参考答案

1. 在一元线性回归方程中，在假设一到假设五成立的情况下，请推导 OLS 估计量 $\hat{\beta}_0$ 的方差

$$\text{Var}(\hat{\beta}_0|x) = \frac{(\sigma^2/n) \sum_{i=1}^n x_i^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

证明. 除了假设一到假设五，我们还需要利用

$$\begin{aligned}\hat{\beta}_0 &= \bar{y} - \hat{\beta}_1 \bar{x}, \\ \text{Var}(\hat{\beta}_0|x) &= \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}.\end{aligned}$$

首先，注意到

$$\hat{\beta}_1 = \sum_{i=1}^n w_i (y_i - \bar{y}) = \sum_{i=1}^n w_i y_i,$$

其中， $w_i = (x_i - \bar{x}) / \sum_{j=1}^n (x_j - \bar{x})^2$ 。这是因为

$$\sum_{i=1}^n w_i \bar{y} = \bar{y} \underbrace{\sum_{i=1}^n w_i}_{=0} = 0.$$

那么，

$$\text{Cov}(\bar{u}, \hat{\beta}_1|x) = \frac{1}{n} \sum_{k=1}^n \sum_{i=1}^n w_i \text{Cov}(u_k, y_i|x) = \frac{1}{n} \sum_{k=1}^n w_k \sigma^2 = \frac{\sigma^2}{n} \underbrace{\sum_{k=1}^n w_k}_{=0} = 0,$$

其中， $\text{Cov}(u_k, y_i|x) = \sigma^2 \mathbf{1}\{k = i\}$ ， $\mathbf{1}\{\cdot\}$ 是指示函数 (indicator function)，所以有

$$\sum_{i=1}^n w_i \text{Cov}(u_k, y_i|x) = w_k \sigma^2.$$

因此,

$$\begin{aligned}
\text{Var}(\hat{\beta}_0|x) &= \text{Var}(\bar{y} - \hat{\beta}_1\bar{x}|x) \\
&= \text{Var}(\beta_0 + \beta_1\bar{x} + \bar{u} - \hat{\beta}_1\bar{x}|x) \\
&= \text{Var}(\bar{u} - \hat{\beta}_1\bar{x}|x) \\
&= \text{Var}(\bar{u}|x) + \bar{x}^2 \text{Var}(\hat{\beta}_1|x) - \underbrace{2\bar{x} \text{Cov}(\bar{u}, \hat{\beta}_1|x)}_{=0} \\
&= \frac{\sigma^2}{n} + \frac{\sigma^2\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \\
&= \frac{(\sigma^2/n) \sum_{i=1}^n (x_i - \bar{x})^2 + (\sigma^2/n) n\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \\
&= \frac{(\sigma^2/n) \sum_{i=1}^n x_i^2}{\sum_{i=1}^n (x_i - \bar{x})^2}
\end{aligned}$$

□

2. 在给定二元线性回归方程 $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + u$, 以 β_1 为例, 第一种估计系数的办法是用最小二乘法直接估计: 第二种办法分成两步进行。第一步先拿 x_1 对 x_2 做回归, 得到残差, 第二步拿 y 对第一步中得到的残差做回归。请从数学上证明两种办法得到的 β_1 的估计量 $\hat{\beta}_1$ 是等价的。

证明. 在第一种方法中,

$$\min_{\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2} \sum_{i=1}^n \left[y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_{1i} + \hat{\beta}_2 x_{2i}) \right]^2$$

一阶条件:

$$\sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{1i} - \hat{\beta}_2 x_{2i}) = 0 \quad (1)$$

$$\sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{1i} - \hat{\beta}_2 x_{2i}) x_{1i} = 0 \quad (2)$$

$$\sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{1i} - \hat{\beta}_2 x_{2i}) x_{2i} = 0 \quad (3)$$

由 (1) 我们有

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}_1 - \hat{\beta}_2 \bar{x}_2$$

把上式代入 (2) 和 (3) 有

$$\begin{aligned} \sum_{i=1}^n \left[(y_i - \bar{y}) - \hat{\beta}_1(x_{1i} - \bar{x}_1) - \hat{\beta}_2(x_{2i} - \bar{x}_2) \right] x_{1i} &= 0 \\ \sum_{i=1}^n \left[(y_i - \bar{y}) - \hat{\beta}_1(x_{1i} - \bar{x}_1) - \hat{\beta}_2(x_{2i} - \bar{x}_2) \right] x_{2i} &= 0 \end{aligned}$$

因此有

$$\sum_{i=1}^n \left[(y_i - \bar{y}) - \hat{\beta}_1(x_{1i} - \bar{x}_1) - \hat{\beta}_2(x_{2i} - \bar{x}_2) \right] (x_{1i} - \bar{x}_1) = 0 \quad (4)$$

$$\sum_{i=1}^n \left[(y_i - \bar{y}) - \hat{\beta}_1(x_{1i} - \bar{x}_1) - \hat{\beta}_2(x_{2i} - \bar{x}_2) \right] (x_{2i} - \bar{x}_2) = 0 \quad (5)$$

定义 $Dy_i \equiv y_i - \bar{y}$, $Dx_{1i} \equiv x_{1i} - \bar{x}_1$, $Dx_{2i} \equiv x_{2i} - \bar{x}_2$, 式 (4) 和 (5) 可以改写为

$$\sum_{i=1}^n \left[Dy_i - \hat{\beta}_1 Dx_{1i} - \hat{\beta}_2 Dx_{2i} \right] Dx_{1i} = 0 \quad (6)$$

$$\sum_{i=1}^n \left[Dy_i - \hat{\beta}_1 Dx_{1i} - \hat{\beta}_2 Dx_{2i} \right] Dx_{2i} = 0 \quad (7)$$

(6) 和 (7) 构成了包含两个未知数、两个方程的线性方程组, 因此

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n Dy_i Dx_{1i} \sum_{i=1}^n (Dx_{2i})^2 - \sum_{i=1}^n Dy_i Dx_{2i} \sum_{i=1}^n Dx_{1i} Dx_{2i}}{\sum_{i=1}^n (Dx_{1i})^2 \sum_{i=1}^n (Dx_{2i})^2 - (\sum_{i=1}^n Dx_{1i} Dx_{2i})^2}$$

在第二种方法中, 我们在第一步考虑

$$x_1 = \alpha_0 + \alpha_1 x_2 + v$$

所以有

$$\hat{\alpha}_1 = \frac{\sum_{i=1}^n (x_{2i} - \bar{x}_2)(x_{1i} - \bar{x}_1)}{\sum_{i=1}^n (x_{2i} - \bar{x}_2)^2} = \frac{\sum_{i=1}^n Dx_{2i} Dx_{1i}}{\sum_{i=1}^n (Dx_{2i})^2} \quad (8)$$

$$\hat{\alpha}_0 = \bar{x}_1 - \hat{\alpha}_1 \bar{x}_2 \quad (9)$$

将 (9) 代入残差 $\hat{v}_i \equiv x_{1i} - \hat{\alpha}_0 - \hat{\alpha}_1 x_{2i}$, 那么

$$\hat{v}_i = Dx_{1i} - \hat{\alpha}_1 Dx_{2i} \quad (10)$$

注意到 $\bar{\hat{v}} = 0$, 因此在第二步中我们有

$$\tilde{\beta}_1 = \frac{\sum_{i=1}^n \hat{v}_i (y_i - \bar{y})}{\sum_{i=1}^n \hat{v}_i^2} = \frac{\sum_{i=1}^n \hat{v}_i Dy_i}{\sum_{i=1}^n \hat{v}_i^2} \quad (11)$$

将 (8) 和 (10) 代入 (11), 我们有

$$\begin{aligned} \tilde{\beta}_1 &= \frac{\sum_{i=1}^n (Dx_{1i} - \hat{\alpha}_1 Dx_{2i}) Dy_i}{\sum_{i=1}^n (Dx_{1i} - \hat{\alpha}_1 Dx_{2i})^2} \\ &= \frac{\sum_{i=1}^n Dx_{1i} Dy_i - \frac{\sum_{i=1}^n Dx_{2i} Dx_{1i}}{\sum_{i=1}^n (Dx_{2i})^2} \sum_{i=1}^n Dx_{2i} Dy_i}{\sum_{i=1}^n (Dx_{1i})^2 + \frac{(\sum_{i=1}^n Dx_{2i} Dx_{1i})^2}{(\sum_{i=1}^n (Dx_{2i})^2)^2} \sum_{i=1}^n (Dx_{2i})^2 - 2 \frac{\sum_{i=1}^n Dx_{2i} Dx_{1i}}{\sum_{i=1}^n (Dx_{2i})^2} \sum_{i=1}^n Dx_{1i} Dx_{2i}} \\ &= \frac{\sum_{i=1}^n Dy_i Dx_{1i} \sum_{i=1}^n (Dx_{2i})^2 - \sum_{i=1}^n Dy_i Dx_{2i} \sum_{i=1}^n Dx_{1i} Dx_{2i}}{\sum_{i=1}^n (Dx_{1i})^2 \sum_{i=1}^n (Dx_{2i})^2 - (\sum_{i=1}^n Dx_{1i} Dx_{2i})^2} \end{aligned}$$

因此 $\hat{\beta}_1$ 和 $\tilde{\beta}_1$ 是等价的。 \square

3. 子虚国政府委托乌有大学进行一个项目, 主要目的是研究劳动力市场上决定劳动收入的因素。乌有大学通过问卷调查。得到一个数据集。数据集包括如下变量: (1) *gender*: 性别, 其中 1 代表男性, 2 代表女性; (2) *birthyear*: 出生年份; (3) *marriage*: 婚姻状况, 其中 1 代表处于婚姻状态, 0 代表处于非结婚状态 (包括未婚, 离异, 丧偶等); (4) *empjob_twage*: 年总收入; (5) *schooling_yr*: 受教育年数;

在打开 Stata 后, 首先使用 `cd` 把工作目录设定为存放有数据集的目录, 再使用 `use` 打开数据集。例如,

```
cd "your_path" /* e.g. cd "/home/zhufeng/metrics" */
use "homework1_dataset.dta"
```

- (1) 创建两个新变量: (a) *male*: 1 代表男性, 0 代表女性; (b) *female*: 1 代表女性, 0 代表男性。

```
. * Question 3.1
. generate male = 1 if gender == 1
(892 missing values generated)
. replace male = 0 if gender == 2
(892 real changes made)
.
. generate female = 1 if gender == 2
(1,948 missing values generated)
. replace female = 0 if gender == 1
(1,948 real changes made)
.
```

(2) 给出以下变量的均值, 标准差, 最小值以及最大值: *female*, *male*, *birthyear*, *marriage*, *empjob_twage*, *schooling_yr*。

```
. * Question 3.2
. summarize female male birthyear marriage empjob_twage schooling_yr
```

Variable	Obs	Mean	Std. Dev.	Min	Max
female	2,840	.3140845	.464232	0	1
male	2,840	.6859155	.464232	0	1
birthyear	2,840	1974.865	11.2741	1914	1998
marriage	2,840	.6010563	.4897674	0	1
empjob_twage	2,840	4789.801	3361.39	176.1308	50725.66
schooling_yr	2,840	7.624296	2.9349	0	15

(3) 乌有大学的飘渺教授认为教育对于收入有着重要的影响, 她建议估计下面这个回归方程式:

$$\text{empjob_twage} = \beta_0 + \beta_1 \times \text{schooling_yr} + u$$

请使用给定的数据集估计这个方程, 给出回归结果。

```
. * Question 3.3
. regress empjob_twage schooling_yr
```

Source	SS	df	MS	Number of obs	=	2,840
Model	731763187	1	731763187	F(1, 2838)	=	66.25
Residual	3.1346e+10	2,838	11045077.6	Prob > F	=	0.0000
				R-squared	=	0.0228
				Adj R-squared	=	0.0225
Total	3.2078e+10	2,839	11298941	Root MSE	=	3323.4

empjob_twage	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
schooling_yr	172.9853	21.25242	8.14	0.000	131.3136 214.6571
_cons	3470.91	173.6213	19.99	0.000	3130.473 3811.347

(4) 计算 *empjob_twage*, $\widehat{\text{empjob_twage}}$, 和 \hat{u} 的均值。他们之间有什么关系。

```
. * Question 3.4
. predict yhat
(option xb assumed; fitted values)
. predict ehat, residuals
. summarize empjob_twage yhat ehat
```

Variable	Obs	Mean	Std. Dev.	Min	Max
empjob_twage	2,840	4789.801	3361.39	176.1308	50725.66
yhat	2,840	4789.801	507.6947	3470.91	6065.689
ehat	2,840	.0000148	3322.828	-5300.15	45697.88

因此, $\widehat{\text{empjob_twage}} = \text{empjob_twage} + \hat{u}$ 。