

Lecture 3: Ordinary Least Squares

Zhentao Shi

September 20, 2018

Notation: y_i is a scalar, and x_i is a $K \times 1$ vector. Y is an $n \times 1$ vector, and X is an $n \times K$ matrix.

1 Algebra of Least Squares

1.1 OLS estimator

As we have learned from the linear project model, the parameter β

$$\begin{aligned} y_i &= x_i' \beta + e_i \\ E[x_i e_i] &= 0 \end{aligned}$$

can be written as $\beta = (E[x_i x_i'])^{-1} E[x_i y_i]$.

While population is something imaginary, in reality we possess a sample of n observations. We thus replace the population mean $E[\cdot]$ by the sample mean, and the resulting estimator is

$$\hat{\beta} = \left(\frac{1}{n} \sum_{i=1}^n x_i x_i' \right)^{-1} \frac{1}{n} \sum_{i=1}^n x_i y_i = (X' X)^{-1} X' y.$$

This is one way to motivate the OLS estimator.

```
In [1]: n = 100
        beta0 = c(1.0, 1.0, 0.0)
        X = cbind(rnorm(n), rpois(n, 3) )
        e = rlogis(n) # the error term does not have to be normally distributed

        y = cbind(1, X ) %*% beta0 + e # generate data
        # in reality, we observe y and X but not e and beta0
```

Alternatively, we can derive the OLS estimator from minimizing the sum of squared residuals

$$Q(\beta) = \sum_{i=1}^n (y_i - x_i' \beta)^2 = (Y - X\beta)' (Y - X\beta).$$

By the first-order condition

$$\frac{\partial}{\partial \beta} Q(\beta) = -2X' (Y - X\beta),$$

the optimality condition gives exactly the same $\hat{\beta}$. Moreover, the second-order condition

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

shows that $Q(\beta)$ is convex in β . ($Q(\beta)$ is strictly convex in β if $X'X$ is positive definite.)

```
In [2]: reg1 = lm( y ~ X ) # OLS regression
        print(reg1)

        X1 = cbind(1, X) # the first column of X is a constant
        bhat = solve(t(X1)%*%X1, t(X1) %*% y )
        print(bhat)
```

Call:

```
lm(formula = y ~ X)
```

Coefficients:

```
(Intercept)          X1          X2
      0.84500      0.77722     -0.01699
```

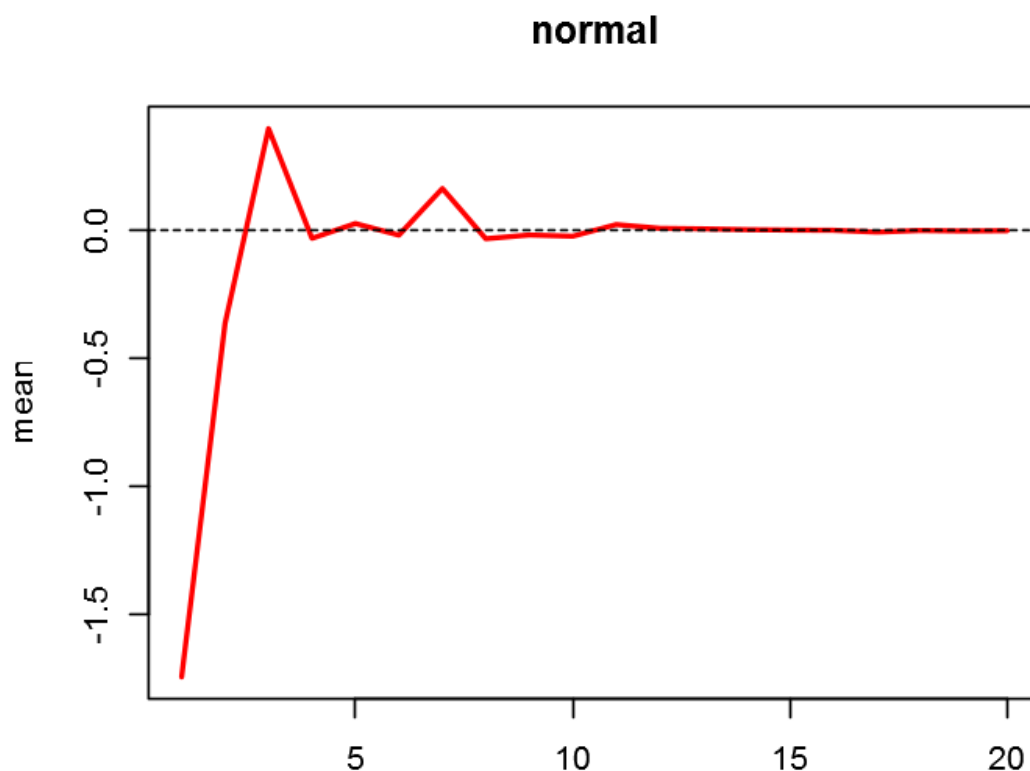
```
      [,1]
[1,] 0.84500343
[2,] 0.77722133
[3,] -0.01698568
```

Here we introduce some definitions and properties in OLS estimation.

- Fitted value: $\hat{Y} = X\hat{\beta}$.
- Projector: $P_X = X(X'X)^{-1}X'$; Annihilator: $M_X = I_n - P_X$.
- $P_X M_X = M_X P_X = 0$.
- If $AA = A$, we call it an idempotent matrix. Both P_X and M_X are idempotent.
- Residual: $\hat{e} = Y - \hat{Y} = Y - X\hat{\beta} = M_X Y = M_X(X\beta + e) = M_X e$.
- $X'\hat{e} = XM_X e = 0$.
- $\frac{1}{n} \sum_{i=1}^n \hat{e}_i = 0$ if x_i contains a constant.

```
In [3]: yhat = predict( reg1, data = X ) # predicted value from the OLS regression
        matplot( x = X[,1], y = cbind(y, yhat), pch = 1:2, xlab = "x", ylab = "y") # a graph bet

        library(repr)
        options(repr.plot.width=6, repr.plot.height=5)
        legend(x = 1.2, y = -2, pch = 1:2, col = 1:2, legend = c("y", "predicted"))
```



```
In [4]: # check the orthogonality of ehat and X1
```

```
      ehat = y - X1 %*% bhat
      print( t(X1) %*% ehat )
```

```
      [,1]
[1,] 4.218847e-15
[2,] 2.337195e-14
[3,] 0.000000e+00
```

```
In [5]: cat("The mean of the residual is ", mean(ehat), "and the sum is", sum(ehat), "\nBut the
```

```
The mean of the residual is  3.008878e-17 and the sum is 3.01148e-15
But the mean of the true error term is -0.18831
```

Real Data Example

We check the relationship between *health status* and three control variables: *the number of doctor visits*, *the number of children in the household*, and *access to health care*.

```
In [6]: library(Ecdat, quietly = TRUE, warn.conflicts = FALSE)
```

```
data(Doctor)
head(Doctor) # display the data structure
```

Attaching package: 'Ecdat'

The following object is masked from 'package:base':

sign

doctor	children	access	health
0	1	0.50	0.495
1	3	0.17	0.520
0	4	0.42	-1.227
0	2	0.33	-1.524
11	1	0.67	0.173
3	1	0.25	-0.905

```
In [7]: reg = lm(health ~ doctor + children + access, data = Doctor)
print(reg)
```

Call:

```
lm(formula = health ~ doctor + children + access, data = Doctor)
```

Coefficients:

(Intercept)	doctor	children	access
-0.02810	0.12059	0.03323	-0.63320

```
In [8]: summary(reg)
```

Call:

```
lm(formula = health ~ doctor + children + access, data = Doctor)
```

Residuals:

Min	1Q	Median	3Q	Max
-3.3370	-1.0085	-0.3261	0.6938	6.1266

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-0.02810	0.18281	-0.154	0.878
doctor	0.12059	0.01884	6.399	3.71e-10 ***

```

children    0.03323    0.04771    0.697    0.486
access     -0.63320    0.33724   -1.878    0.061 .
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

Residual standard error: 1.378 on 481 degrees of freedom
Multiple R-squared: 0.08221, Adjusted R-squared: 0.07649
F-statistic: 14.36 on 3 and 481 DF, p-value: 5.628e-09

1.2 Goodness of Fit

The so-called R-square is the most popular measure of goodness-of-fit in the linear regression. R-square is well defined only when a constant is included in the regressors. Let $M_i = I_n - \frac{1}{n} \mathbf{1}\mathbf{1}'$, where $\mathbf{1}$ is an $n \times 1$ vector of 1's. M_i is the *demeaner*, in the sense that $M_i (z_1, \dots, z_n)' = (z_1 - \bar{z}, \dots, z_n - \bar{z})'$, where $\bar{z} = \frac{1}{n} \sum_{i=1}^n z_i$. For any X , we can decompose $Y = P_X Y + M_X Y = \hat{Y} + \hat{e}$. The total variation is

$$Y' M_i Y = (\hat{Y} + \hat{e})' M_i (\hat{Y} + \hat{e}) = \hat{Y}' M_i \hat{Y} + 2\hat{Y}' M_i \hat{e} + \hat{e}' M_i \hat{e} = \hat{Y}' M_i \hat{Y} + \hat{e}' \hat{e}$$

where the last equality follows by $M_i \hat{e} = \hat{e}$ as $\frac{1}{n} \sum_{i=1}^n \hat{e}_i = 0$, and $\hat{Y}' \hat{e} = Y' P_X M_X e = 0$. R-square is defined as $\hat{Y}' M_i \hat{Y} / Y' M_i Y$.

1.3 Frish-Waugh-Lovell Theorem

The FWL theorem is an algebraic fact about the formula of a subvector of the OLS estimator. To derive the FWL theorem We need to use the inverse of partitioned matrix. For a positive definite symmetric matrix $A = \begin{pmatrix} A_{11} & A_{12} \\ A_{12}' & A_{22} \end{pmatrix}$, the inverse can be written as

$$A^{-1} = \begin{pmatrix} (A_{11} - A_{12} A_{22}^{-1} A_{12}')^{-1} & - (A_{11} - A_{12} A_{22}^{-1} A_{12}')^{-1} A_{12} A_{22}^{-1} \\ - A_{22}^{-1} A_{12}' (A_{11} - A_{12} A_{22}^{-1} A_{12}')^{-1} & (A_{22} - A_{12}' A_{11}^{-1} A_{12})^{-1} \end{pmatrix}.$$

In our context of OLS estimator,

$$\begin{aligned} \hat{\beta} &= \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} = \begin{pmatrix} X_1' X_1 & X_1' X_2 \\ X_2' X_1 & X_2' X_2 \end{pmatrix}^{-1} \begin{pmatrix} X_1 y \\ X_2 y \end{pmatrix} \\ &= \begin{pmatrix} (X_1 M_{X_2}' X_1)^{-1} & - (X_1 M_{X_2}' X_1)^{-1} X_1' X_2 (X_2' X_2)^{-1} \\ \cdot & \cdot \end{pmatrix} \begin{pmatrix} X_1 y \\ X_2 y \end{pmatrix}. \end{aligned}$$

The subvector

$$\begin{aligned} \hat{\beta}_1 &= (X_1 M_{X_2}' X_1)^{-1} X_1 y - (X_1 M_{X_2}' X_1)^{-1} X_1' X_2 (X_2' X_2)^{-1} X_2 y \\ &= (X_1 M_{X_2}' X_1)^{-1} (X_1 y - X_1' P_{X_2} y) \\ &= (X_1 M_{X_2}' X_1)^{-1} X_1 M_{X_2} y. \end{aligned}$$

Similar derivation can also be carried out in the population linear projection. See Hansen's Chapter 2.21-23.

```
In [9]: X2 = X1[:,1:2]
        PX2 = X2 %*% solve( t(X2) %*% X2) %*% t(X2)
        MX2 = diag(rep(1,n)) - PX2

        X3 = X1[:,3]

        bhat3 = solve(t(X3)%*% MX2 %*% X3, t(X3) %*% MX2 %*% y )
        print(bhat3)

        [,1]
[1,] -0.01698568
```

2 Statistical Properties of Least Squares

To talk about the statistical properties in finite sample, we impose the following assumptions.

1. The data $(y_i, x_i)_{i=1}^n$ is a random sample from the same data generating process $y_i = x_i'\beta + e_i$.
2. $e_i|x_i \sim N(0, \sigma^2)$.

2.1 Maximum Likelihood Estimation

Under the normality assumption, $y_i|x_i \sim N(x_i'\beta, \gamma)$, where $\gamma = \sigma^2$. The *conditional* likelihood of observing a sample $(y_i, x_i)_{i=1}^n$ is

$$\prod_{i=1}^n \frac{1}{\sqrt{2\pi\gamma}} \exp\left(-\frac{1}{2\gamma} (y_i - x_i'\beta)^2\right),$$

and the (conditional) log-likelihood function is

$$L(\beta, \gamma) = -\frac{n}{2} \log 2\pi - \frac{n}{2} \log \gamma - \frac{1}{2\gamma} \sum_{i=1}^n (y_i - x_i'\beta)^2.$$

Therefore, the maximum likelihood estimator (MLE) coincides with the OLS estimator, and $\hat{\gamma}_{MLE} = \hat{e}'\hat{e}/n$.

2.2 Finite Sample Distribution

We can show the finite-sample exact distribution of $\hat{\beta}$. *Finite sample distribution* means that the distribution holds for any n ; it is in contrast to *asymptotic distribution*, which is a large sample approximation to the finite sample distribution.

Since

$$\hat{\beta} = (X'X)^{-1} X'y = (X'X)^{-1} X'(X\beta + e) = \beta + (X'X)^{-1} X'e,$$

we have the estimator $\hat{\beta}|X \sim N\left(\beta, \sigma^2 (X'X)^{-1}\right)$, and

$$\hat{\beta}_k|X \sim N\left(\beta_k, \sigma^2 \eta'_k (X'X)^{-1} \eta_k\right) \sim N\left(\beta_k, \sigma^2 (X'X)^{-1}_{kk}\right),$$

where $\eta_k = (1 \{l = k\})_{l=1,\dots,K}$ is the selector of the k -th element.

In reality, σ^2 is an unknown parameter, and

$$s^2 = \hat{e}'\hat{e} / (n - K) = e' M_X e / (n - K)$$

is an unbiased estimator of σ^2 . Consider the T -statistic

$$T_k = \frac{\hat{\beta}_k - \beta_k}{\sqrt{s^2 [(X'X)^{-1}]_{kk}}} = \frac{(\hat{\beta}_k - \beta_k) / \sqrt{\sigma^2 [(X'X)^{-1}]_{kk}}}{\sqrt{\frac{e' M_X e}{\sigma^2} / (n - K)}}.$$

The numerator follows a standard normal, and the denominator follows $\frac{1}{n-K} \chi^2(n-K)$. Moreover, the numerator and the denominator are independent. As a result, $T_k \sim t(n-K)$.

2.3 Mean and Variance

Now we relax the normality assumption and statistical independence. Instead, we assume a regression model $y_i = x'_i \beta + e_i$ and

$$\begin{aligned} E[e_i|x_i] &= 0 \\ E[e_i^2|x_i] &= \sigma^2. \end{aligned}$$

where the first condition is the *mean independence* assumption, and the second condition is the *homoskedasticity* assumption.

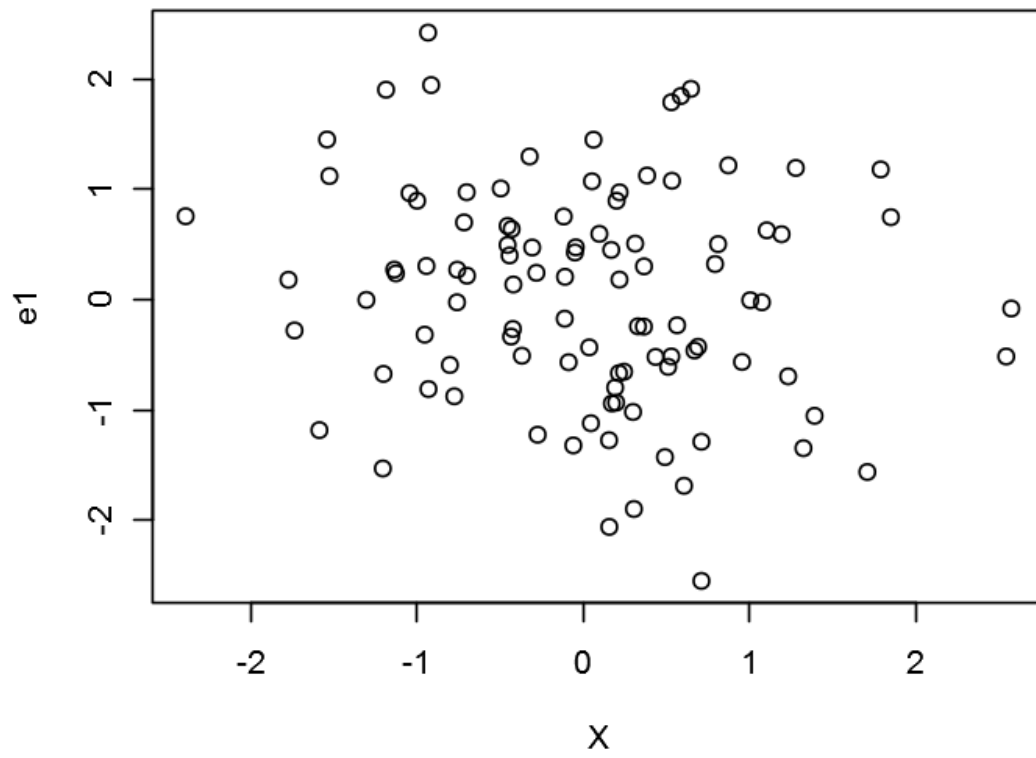
Example (Heteroskedasticity) If $e_i = x_i u_i$, where x_i is a scalar random variable, u_i is independent of x_i , $E[u_i] = 0$ and $E[u_i^2] = \sigma^2$. Then $E[e_i|x_i] = 0$ but $E[e_i^2|x_i] = \sigma_i^2 x_i^2$ is a function of x_i . We say e_i^2 is a heteroskedastic error.

```
In [10]: n = 100
         X = rnorm(n)

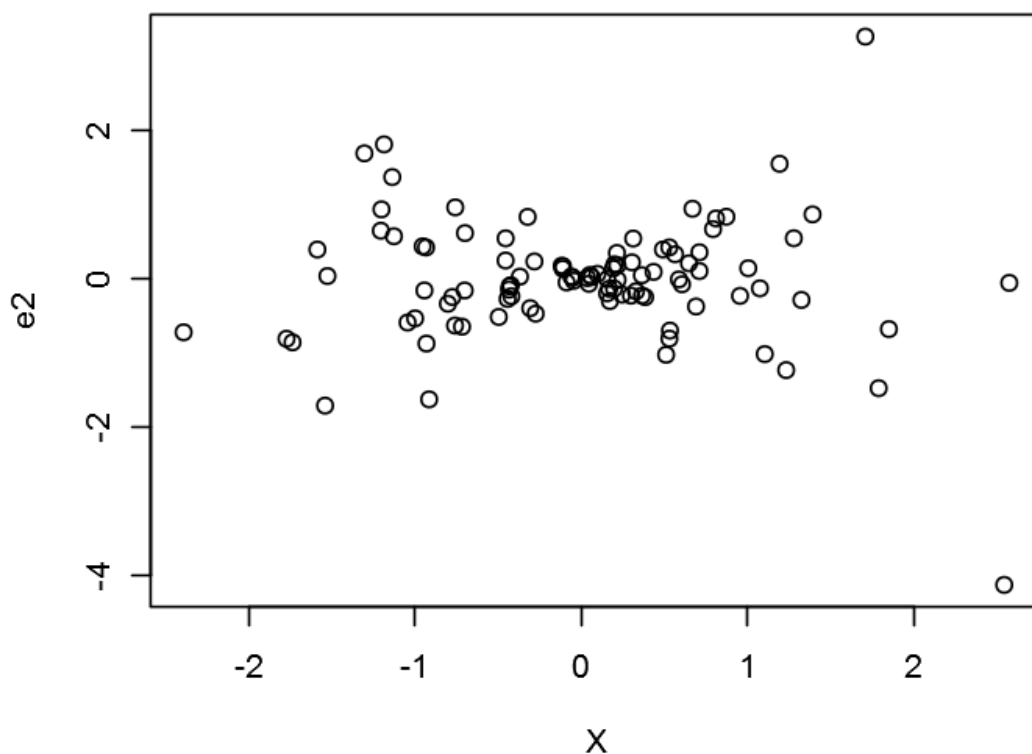
         e1 = rnorm(n)
         plot( y = e1, x = X, main = "homoskedastic")

         e2 = X * rnorm(n) # the source of heteroskedasticity
         plot( y = e2, x = X, main = "heteroskedastic")
```

homoskedastic



heteroskedastic



These assumptions are about the first and second moment of e_i conditional on x_i . Unlike the normality assumption, they do not restrict the entire distribution of e_i .

- Unbiasedness:

$$E[\hat{\beta}|X] = E[(X'X)^{-1}X'Y|X] = E[(X'X)^{-1}X(X'\beta + e)|X] = \beta.$$

Unbiasedness does not rely on homoskedasticity.

- Variance:

$$\begin{aligned} \text{var}(\hat{\beta}|X) &= E[(\hat{\beta} - E\hat{\beta})(\hat{\beta} - E\hat{\beta})'|X] \\ &= E[(\hat{\beta} - \beta)(\hat{\beta} - \beta)'|X] \\ &= E[(X'X)^{-1}X'ee'X(X'X)^{-1}|X] \\ &= (X'X)^{-1}X'E[ee'|X]X(X'X)^{-1} \\ &= (X'X)^{-1}X'(\sigma^2 I_n)X(X'X)^{-1} \\ &= \sigma^2 (X'X)^{-1}. \end{aligned}$$

2.4 Gauss-Markov Theorem

Gauss-Markov theorem justifies the OLS estimator as the efficient estimator among all linear unbiased ones. *Efficient* here means that it enjoys the smallest variance in a family of estimators.

There are numerous linearly unbiased estimators. For example, $(Z'X)^{-1}Z'y$ for $z_i = x_i^2$ is unbiased because $E[(Z'X)^{-1}Z'y] = E[(Z'X)^{-1}Z'(X\beta + e)] = \beta$.

Let $\tilde{\beta} = A'y$ be a generic linear estimator, where A is any $n \times K$ functions of X . As

$$E[A'y|X] = E[A'(X\beta + e)|X] = A'X\beta.$$

So the linearity and unbiasedness of $\tilde{\beta}$ implies $A'X = I_n$. Moreover, the variance

$$\text{var}(A'y|X) = E[(A'y - \beta)(A'y - \beta)'|X] = E[A'ee'A|X] = \sigma^2 A'A.$$

Let $C = A - X(X'X)^{-1}$.

$$\begin{aligned} A'A - (X'X)^{-1} &= (C + X(X'X)^{-1})'(C + X(X'X)^{-1}) - (X'X)^{-1} \\ &= C'C + (X'X)^{-1}X'C + C'X(X'X)^{-1} \\ &= C'C, \end{aligned}$$

where the last equality follows as

$$(X'X)^{-1}X'C = (X'X)^{-1}X'(A - X(X'X)^{-1}) = (X'X)^{-1} - (X'X)^{-1} = 0.$$

Therefore $A'A - (X'X)^{-1}$ is a positive semi-definite matrix. The variance of any $\tilde{\beta}$ is no smaller than the OLS estimator $\hat{\beta}$.

Homoskedasticity is a restrictive assumption. Under homoskedasticity, $\text{var}(\hat{\beta}) = \sigma^2(X'X)^{-1}$. Popular estimator of σ^2 is the sample mean of the residuals $\hat{\sigma}^2 = \frac{1}{n}\hat{e}'\hat{e}$ or the unbiased one $s^2 = \frac{1}{n-K}\hat{e}'\hat{e}$. Under heteroskedasticity, Gauss-Markov theorem does not apply.