

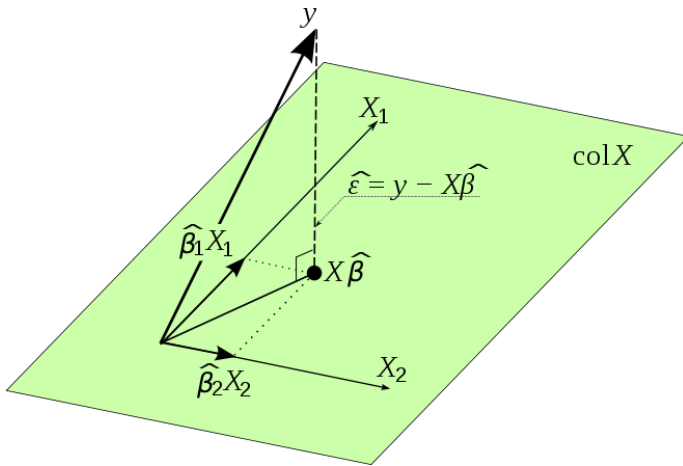
# *Lecture 4: The Classical Regression Model (II)*

*Thomas Chadeaux*

## **Contents**

1	<i>Geometric Interpretation of OLS</i>	2
2	<i>Assumptions</i>	2
3	<i>OLS is 'BEST': Finite Sample Properties of OLS estimator</i>	6
3.1	<i>Unbiasedness</i>	6
3.2	<i>The Variance of the least squares estimator</i>	8
3.3	<i>The Gauss-Markov Theorem</i>	11
4	<i>The Frisch-Waugh Theorem</i>	11

## 1 Geometric Interpretation of OLS



## 2 Assumptions

1. Linearity of the regression model: the model is linear in parameters and correctly specified. 'Linear in parameters' means that each term on the right hand side includes only  $\beta$ s and there is no built-in relationship among the  $\beta$ s. Examples of models that are NOT linear in parameters include  $Y = \beta_1 X^{\beta_2} + \varepsilon$ , but note that  $y = Ax^{\beta}e^{\varepsilon}$  IS linear in parameters after taking logs on both sides.<sup>1</sup>
2.  $X$  is an  $n \times K$  matrix with rank  $K$ . This means that our variables are linearly independent, and that there are at least  $K$  observations (i.e., at least as many observations as there are parameters to estimate). This comes from the basic identification condition from algebra.<sup>2</sup>
3. The disturbance term has zero expectation. I.e.,

$$E[u_i] = 0.$$

Sometimes that disturbance term will be positive, other times negative, but it should not have a systematic tendency in either direction. Note that this is not a big assumption as long as we can include a constant term.

4. The disturbance term is homoskedastic. That means that

$$\text{Var}[\varepsilon_i|X] = \sigma^2.$$

$$\begin{aligned} \ln(y) &= \ln(Ax^{\beta}e^{\varepsilon}) \\ &= \ln(A) + \ln(x^{\beta}) + \ln(e^{\varepsilon}) \\ &= \ln(A) + \beta \ln(x) + \varepsilon \end{aligned} \quad (1)$$

<sup>2</sup> For example, imagine that

$$y = \beta_0 + \beta_1 \text{nonlabor income} + \beta_2 \text{labor income} + \beta_3 \text{total income} + \varepsilon \quad (2)$$

Clearly total income = non-labor income + labor income, which means that there is an exact dependency in the model. Now consider instead

$$y = \beta_0 + (\beta_1 + a) \text{nonlabor income} + (\beta_2 + a) \text{labor income} + (\beta_3 - a) \text{total income} + \varepsilon \quad (3)$$

Using different coefficients, we get the same value on the LHS. I.e., there is no way to estimate  $\beta$ .

Practically, it means that the distribution of the error—loosely, our uncertainty—does NOT depend on the value of  $X$ . Consider for example a model of political voting. We expect that income affects placement on a L/R scale. But there might be a lot of variance for low-income household, as well as high income household. Or income and spending: low income household spend all their income, so there is little variance here. On the other hand some high earners spend a lot, whereas others put a lot in savings (i.e., high variance for high values of  $x$ ).

Why do we care about homoskedasticity? We care because if the variance of the error term is heteroskedastic (i.e., its variance changes with different values of  $x$ ), then OLS is no longer BLUE. It is still Unbiased, but it is no longer Best. I.e., it no longer has minimum variance among all the other linear unbiased estimators. Intuitively, this is because there is some information in the data that was not included. The information is that  $Var(\varepsilon_i|x_i) = f(x_i)$ , but I have not included that information. If I included this info, maybe I could get an estimator that would get closer to the  $y$  values more often than OLS. We'll get back to this later when we talk about the violations of the CLRM (look out for “generalised least squares”).

5. The disturbances are independently distributed. I.e.,

$$Cov[\varepsilon_i, \varepsilon_j] = 0 \forall j \neq i.$$

This means that the error terms are independent of one another. So, a shock to me should not affect you. This assumption is particularly problematic once we deal with time series, since this assumption also implies that

$$Cov[\varepsilon_t, \varepsilon_{t+1}] = 0.$$

Yet a shock at time  $t$  is likely to have repercussions at time  $t + 1$ ... Correlation between the residuals implies that our observations are not really independent, and so in a sense, that we do not really have  $N$  observations but a fraction of that. I.e., we have less information than we think. So if we estimated our standard errors using  $N$  (i.e., ignoring correlation), we will underestimate our SE and hence reject the null hypothesis more often than we should. Note that serial correlation can also happen in spatially clustered data and other situations.

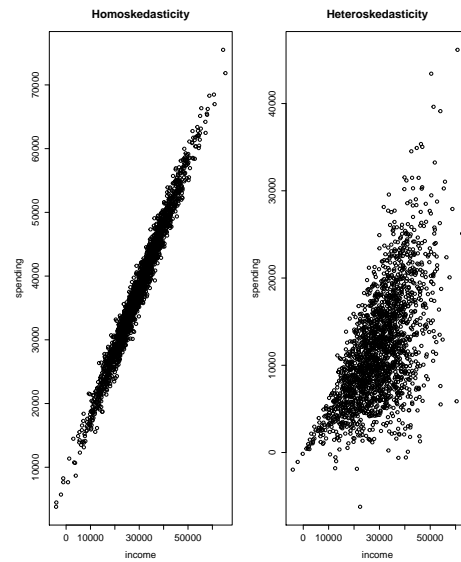


Figure 1: homoskedasticity.R

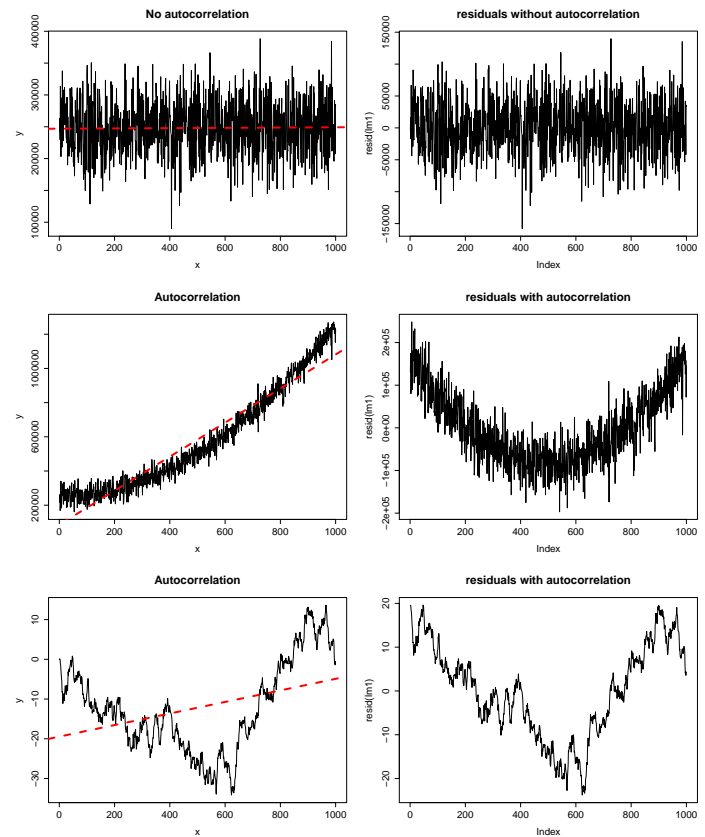


Figure 2: autocorrelation.R

Assumptions 4 and 5 can be summarised as

$$E[\varepsilon\varepsilon'] = \sigma^2 I = \begin{pmatrix} \sigma^2 & 0 & 0 & \dots & 0 \\ 0 & \sigma^2 & 0 & \dots & 0 \\ 0 & 0 & \sigma^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \ddots & 0 \\ 0 & 0 & 0 & \dots & \sigma^2 \end{pmatrix}$$

$\sigma^2$  is called the variance-covariance matrix of the residuals. The 0s on the off-diagonal show that there is no correlation between errors.

6. Normality. The disturbances are normally distributed:

$$E[\varepsilon_i | x_i] \sim N[0, \sigma^2]$$

Note that this assumption is not absolutely necessary, but will prove useful in constructing confidence intervals and test statistics.

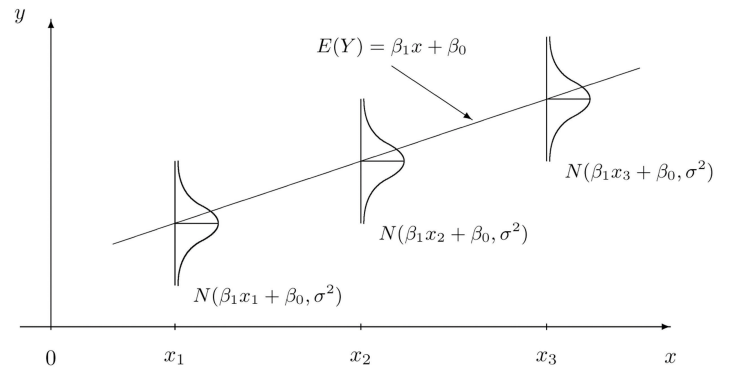
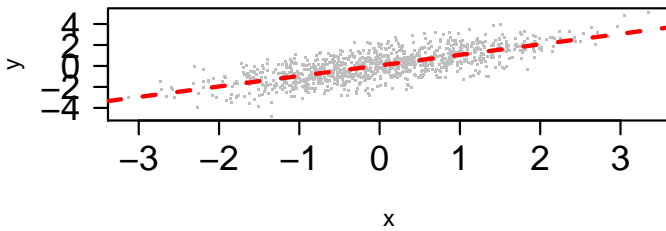
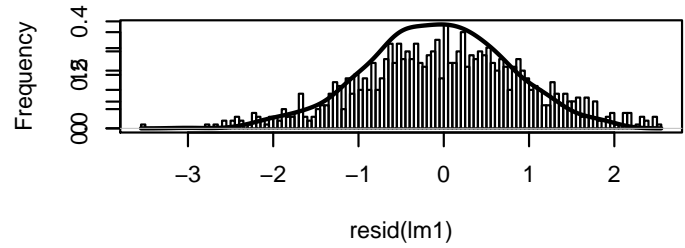


Figure 3: The normality assumption

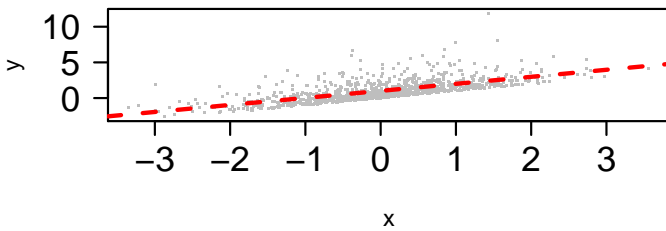
**x distributed normal  
Errors distributed normal**



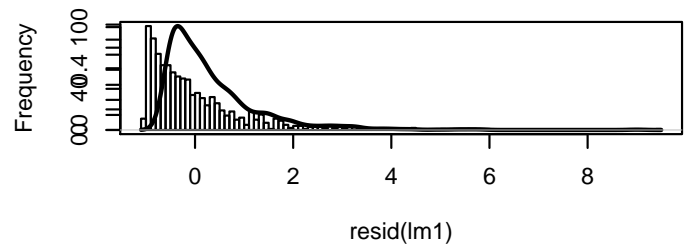
**Distribution of residuals**



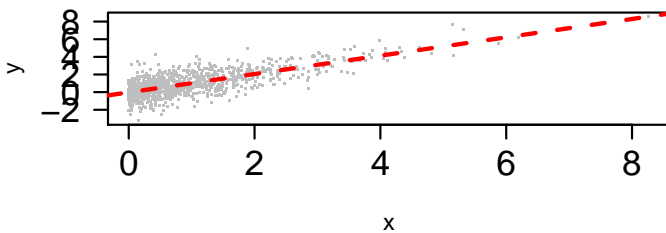
**x distributed normal  
Errors distributed exponential**



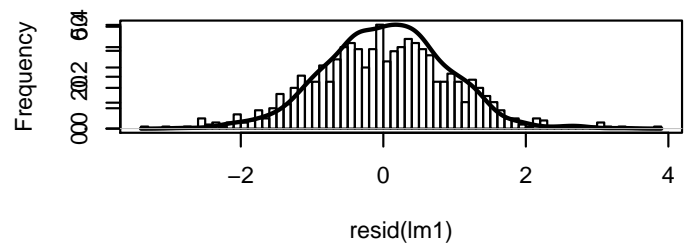
**Distribution of residuals**



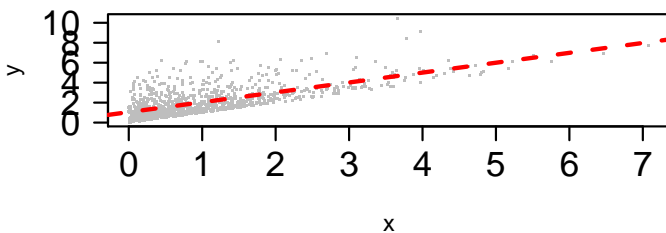
**x distributed exponential  
Errors distributed normal**



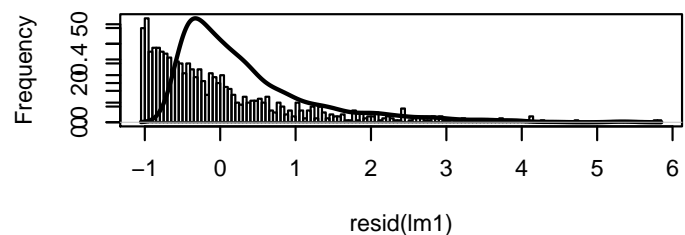
**Distribution of residuals**



**x distributed exponential  
Errors distributed exponential**



**Distribution of residuals**



nonnormality.R

### 3 OLS is 'BEST': Finite Sample Properties of OLS estimator

Why do we use OLS? After all, many other estimators are possible.<sup>3</sup> In this section we will show that OLS is the Best Linear Unbiased Estimator (BLUE). I.e., among all linear estimators and provided that the assumptions above are satisfied, then OLS has the smallest variance (it is Best) and is Unbiased .

#### 3.1 Unbiasedness

Remember that an estimator  $b$  of a true population parameter  $\beta$  is unbiased if  $E[b] = \beta$ . To find out whether  $b$  is unbiased, let's first rewrite it as:

$$\begin{aligned} b &= (X'X)^{-1}X'y \\ &= (X'X)^{-1}X'X\beta + (X'X)^{-1}X'\varepsilon \\ &= \beta + (X'X)^{-1}X'\varepsilon \end{aligned} \quad (4)$$

But remember that  $E[X'\varepsilon] = 0$  (see your homework), so that

$$E[b] = E[\beta] + E[(X'X)^{-1}X'\varepsilon] = \beta.$$

This proves that the OLS estimator  $b = (X'X)^{-1}X'y$  is an unbiased estimator of  $\beta$ . Practically, this means that if I take repeated samples of my population and apply the OLS estimator to each sample, then on average we get the true population parameter beta.

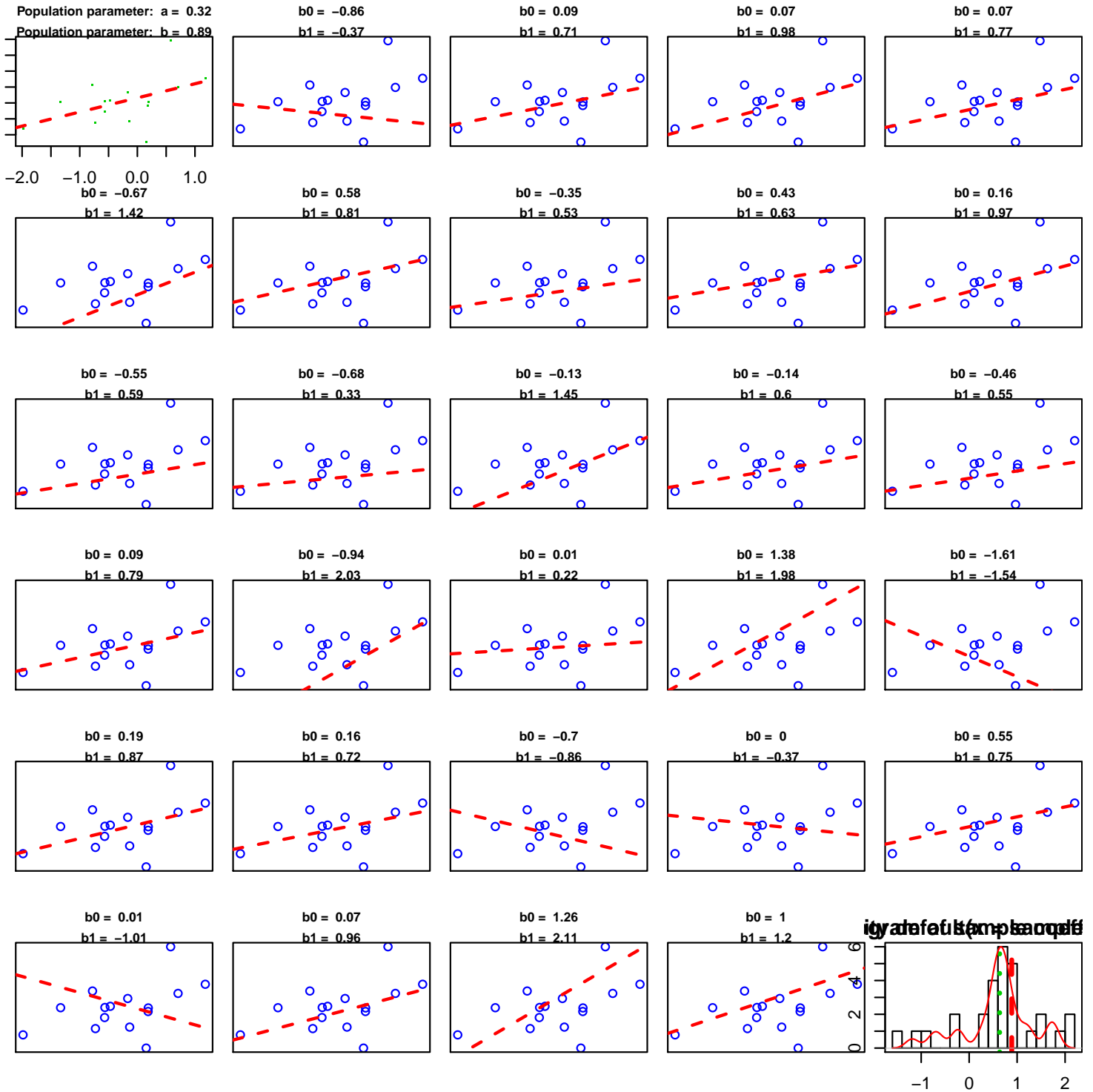
<sup>3</sup> Consider for example an estimator  $\tilde{\beta}$  that takes the first and last observation of the sample:

$$\tilde{\beta} = \frac{Y_n - Y_1}{X_n - X_1}$$

```

1 setwd('~\\Documents\\Academia\\Teaching\\TCD\\P07005_Quantitative_Methods_II/
  2016-HT/Lectures/lecture4/')
2 n <- 15
3 #unbiasedness
4 library(fields)
5 pdf('Figs/unbiasedness.pdf')
6 par(mfrow=c(6,5), mar=c(2,1,1.5,0))
7 #generate population variables
8 x <- rnorm(n)
9 y <- x + rnorm(n)
10 lm1 <- lm(y ~ x)
11 plot(x, y, cex=4, pch=19, main = paste('Population parameter: ', c('a = ',
  'b = '), round(coef(lm1),2)), col=3, cex.main=0.75)
12 abline(lm1, col=2, lty=2, lwd=2)
13 summary(lm1) # This is my true population parameters
14 true.coef <- coef(lm1)
15
16 #But I don't observe the true population. What I observe is a sample of x
  and y
17 sample.coefs <- NULL
18 for(i in 1:2000){ # generate random samples from the true population
19   samplei <- sample(x = 1:length(x), size = n/2, replace = TRUE)
20   x.sample <- x[samplei]
21   y.sample <- y[samplei]
22   lm1.sample <- lm(y.sample ~ x.sample)
23   plot(x, y,
24     main = paste(c('b0 = ', 'b1 = '), round(coef(lm1.sample),2)),

```



```

25     col=4, cex.main=0.75,
26     xaxt='n', yaxt='n', xlab='', ylab='')
27     abline(lm1.sample, col=2, lty=2, lwd=2)
28     sample.coefs <- c(sample.coefs, coef(lm1.sample)[2])
29 }
30
31 hist(sample.coefs,
32     main= 'Histogram of sample coefficients', breaks=50)
33 mean(sample.coefs)
34 xline(true.coef[2], col=2, lty=2, lwd=3)
35 xline(mean(sample.coefs), col=3, lty=3, lwd=3)
36 par(new=T)
37 plot(density(sample.coefs), col=2, xaxt='n', yaxt='n', xlab='', ylab='')
38
39 dev.off()

```

Listing 1: unbiasedness.R

Note that OLS is far from being the

### 3.2 The Variance of the least squares estimator

We care that  $b$ , our estimate of  $\beta$ , has small variance, because a small variance gives us more confidence that we have correctly estimated the true population parameter. So let us calculate the variance of  $b$ . First note that

$$\begin{aligned}
 \mathbf{b} &= (\mathbf{X}'\mathbf{X})^{-1}(\mathbf{X}'\mathbf{y}) \\
 &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'(\mathbf{X}\beta + \varepsilon) \\
 &= \beta + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\varepsilon
 \end{aligned}$$

Now remember that  $\text{Var}[X] = E[(X - E[X])^2]$ , so

$$\begin{aligned}
 \text{Var}[b] &= E[(\mathbf{b} - E[\mathbf{b}])(\mathbf{b} - E[\mathbf{b}])'] \\
 &= E[(\mathbf{b} - \beta)((\mathbf{b} - \beta)')] \\
 &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'E[\varepsilon\varepsilon']\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}
 \end{aligned} \tag{5}$$

But note that  $\varepsilon\varepsilon' = \sigma^2 I$  by assumption.<sup>4</sup> So we can rewrite (5) as

$$\begin{aligned}
 \text{Var}[b] &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\sigma^2 I\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \\
 &= \sigma^2(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \\
 &= \sigma^2 I(\mathbf{X}'\mathbf{X})^{-1} \\
 &= \sigma^2(\mathbf{X}'\mathbf{X})^{-1}
 \end{aligned} \tag{6}$$

Now we have a formula for the variance of  $\mathbf{b}$ , which is what we need to make inferences about it.

However, what is  $\sigma^2$ ? We need an estimate of it before we can proceed. We will not get into details about how this is obtained here, but loosely: a natural estimator of  $\sigma^2$  would seem to be

$$\hat{\sigma}^2 = \frac{1}{n} \sum_i e_i^2 = \frac{\mathbf{e}'\mathbf{e}}{n},$$

<sup>4</sup> Why? Because we assumed that  $\text{Cov}(\varepsilon_i, \varepsilon_j) = 0$  and  $\text{Var}(\varepsilon_i) = \sigma^2$  (see assumptions 4 and 5 in section 2 above)

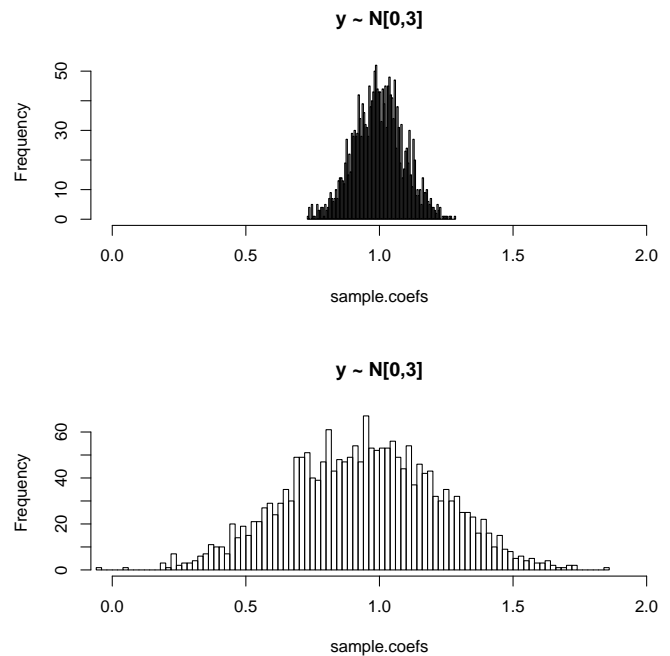


Figure 4: variance.R

```

1 # Compute variance covariance matrix of coefficients manually:
2 X <- cbind(const = rep(1,100), x= rnorm(100))
3 y <- X[,2] + rnorm(100)
4 lm1 <- lm(y ~ 1 + X[,2])
5
6 e <- resid(lm1)
7 ee <- (t(e) %*% e) / (100 - 2)
8 var.b <- as.numeric(ee) * solve(t(X)%*%X)
9 se.b <- sqrt(var.b)
10 se.b
11 summary(lm1)

```

Listing 2: vcovMatrixResid.R



but this estimator is biased for reasons similar to why the estimate of the standard deviation is biased if we divide by  $n$ . The unbiased estimator of  $\sigma^2$  is

$$s^2 = \frac{\mathbf{e}'\mathbf{e}}{n - K}.$$

The standard error of the regression coefficient, then, is  $\sqrt{s^2(\mathbf{X}'\mathbf{X})^{-1}}$ , and the square root of the  $k$ th diagonal element of this matrix is the standard error of the estimator  $b_k$ . This is crucial, because this is what will determine whether our coefficient is 'significant' or not. Equation 5 is the formula for the variance-covariance matrix of the coefficients:

$$\sigma^2(\mathbf{X}'\mathbf{X})^{-1} = \begin{pmatrix} \text{Var}(\beta_1) & \text{Cov}(\beta_1, \beta_2) & \dots & \text{Cov}(\beta_1, \beta_k) \\ \text{Cov}(\beta_2, \beta_1) & \text{Var}(\beta_2) & \ddots & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ \text{Cov}(\beta_k, \beta_1) & \dots & \dots & \text{Var}(\beta_k) \end{pmatrix}$$

Note what happens when the variance of  $X$  increases?  $(\mathbf{X}'\mathbf{X})$  gets larger (if you prefer, look at its determinant), and hence  $(\mathbf{X}'\mathbf{X})^{-1}$  gets smaller and the standard error of the coefficient becomes smaller and smaller.

We'll come back to this, but note that the standard error gives you the  $t$  value:

$$t_b = \frac{b - 0}{\text{sd}(b)}$$

### 3.3 The Gauss-Markov Theorem

We will not prove the Gauss-Markov theorem here. What it shows, however, is that given that the conditions above are satisfied, OLS is BLUE. I.e., it is unbiased and has minimum variance among all other linear estimators. The proof simply involves choosing an alternative estimator, e.g.

$$\tilde{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} + D\mathbf{y}$$

and showing that the only way for it to be unbiased is if  $D\mathbf{y} = 0$ , but even then the variance of that estimator will be greater or equal to the OLS one. Hence OLS is BLUE: it is the Best (lowest variance) of the Linear Unbiased Estimators.

We can (and should) of course look up  $t_b$  in a table of  $t$  statistics to find the associated  $p$ -value. But if you like to do things manually to understand what is happening, you could also do this. Step by step:

1. I got  $b_2 = 1.0182$  and  $SE_{b_2} = 0.5346$ .
2. From this, I infer that  $t_{b_2} = 1.0182/0.5346 = 1.905$ .
3. Either I look this up in a table, or I generate a large number of random values drawn from a  $t$ -distribution (I generated 100000), sort them, and see which value is about 1.905. In my case it turns out to be the 97030 observation. Given that we want a two sided test, this means that this is approximately a  $2970 \cdot 2 / 100000 = 0.0594$   $p$ -value (R tells me the true value is 0.0598).

```

1 # A quick look at a sample regression output
2 x <- rnorm(100)
3 y <- x + rnorm(100)
4 lm1 <- lm(y ~ x)
5 summary(lm1)

```

Console: ~/Documents/Academia/Teaching/TCD/2015-HT/PO7005\_Quantitative\_Methods\_II/Lectures/lecture4/ R Script

Call:  
lm(formula = y ~ x)

Residuals:

	Min	1Q	Median	3Q	Max
	-2.91274	-0.64145	-0.04027	0.53371	2.19282

info about residuals. Same can be obtained using: summary(resid(lm1))

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-0.03864	0.09899	-0.390	0.697
x	0.92955	0.09752	9.532	1.26e-15 ***

---  
Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.9882 on 98 degrees of freedom  
Multiple R-squared: 0.4811, Adjusted R-squared: 0.4758  
F-statistic: 90.86 on 1 and 98 DF, p-value: 1.255e-15

>

## 4 The Frisch-Waugh Theorem

Suppose we want to estimate the following model:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \varepsilon,$$

and that we are interested in  $\beta_2$ . The normal way to get  $b_2$  is to calculate

$$(X'X)^{-1}X'y$$

But there is another way:

1. Regress  $y$  on all variables but  $x_2$  and get the residuals.  
Call the residuals  $y^*$
2. Regress  $x_2$  on all other  $x$  variables and get the residuals.  
Call the residuals  $x_2^*$

3. Regress  $y^*$  on  $x_2^*$ . The coefficient on  $x_2^*$  will be the same as the one obtained by regular OLS.

Why does this work? In the first step, we ‘remove’ all variation explained by all other variables. We are left with the unexplained variance of  $y$ . Similarly, in the second step we get the variance of  $x_2$  that does not covary with other variables. In step 3, we essentially find the covariance between these two left-over variation

```

1 # Frisch-Waugh theorem
2 # Setup the data and model
3 x1 <- rnorm(100)
4 x2 <- rnorm(100)
5 b1 <- 2.1
6 b2 <- 5.3
7 e <- rnorm(100, sd=5)
8 y <- b1*x1 + b2*x2 + e
9 X <- cbind(1, x1, x2)
10
11 # estimate regression the normal way:
12 b <- solve(t(X)%*%X) %*% t(X)%*%y
13 b
14
15 # estimate it the Frisch-Waugh way:
16 #1. first, reg y~x1 and get the residuals
17 X1 <- X[,1:2]
18 resid1 <- residuals(lm(y ~ X1))
19
20 #2. regress x2 on x1
21 X2 <- X[,c(1,3)]
22 resid2 <- residuals(lm(x2 ~ x1))
23
24 coef(lm(resid1 ~ resid2))[2]
25
26 # Just to make sure it's all correct:
27 lm1 <- lm(y ~ x1 + x2)
28 summary(lm1)

```

Listing 3: frischWaugh.R