

Essentials of Data Science with R Software - 2
Sampling Theory and Linear Regression Analysis
Prof. Shalabh
Department of Mathematics and Statistics
Indian Institute of Technology Kanpur

Linear Regression Analysis
Lecture - 42
Simple Linear Regression Analysis
Properties of Least Squares Estimators

Hello, friends. Welcome to the course Essentials of Data Science with R Software – 2, where we are trying to learn the topics of Sampling Theory and Linear Regression Analysis. And, in this module on linear regression analysis, we are going to talk about the concepts of a Simple Linear Regression Analysis.

So, up to now what we have done? That we have considered general form of the simple linear regression model which has only one independent variable we have estimated the regression parameters. Of course, the regression coefficients which were β_0 and β_1 they have been estimated as b_0 and b_1 . Now, one more question arises you have estimated β_0, β_1 , by b_0 and b_1 on the values of a sample of data. If the sample changes, then what will happen?

Do you think that the same values of b_0 and b_1 will be obtained? Answer is no. These are sample dependent values. So, as the sample changes these values will be changing and why? Because they are the function of random variable, y is my random variable in the simple linear regression model. So, your b_0 and b_1 which are the estimators of the regression coefficient ordinary least square estimators of β_0 and β_1 , they are the functions of y .

Hence they are the function of a random variable. So, obviously, as the values of the variable changes these values will change. Once this b_0 and b_1 becomes a statistic, a random variable then obviously, the concept of variability comes into picture. The concept of coefficient, this confidence interval comes into picture. Just by estimating the values b_0 and b_1 will not help. But, you also need to tell that how much is the variation involve in the values of b_0 and b_1 .

Second thing is, this you have assume the random error to have variance σ^2 , but would you not like to find out what is the amount of random variation in my data? Yes, so, that can be obtained by estimating the parameter σ^2 . So, essentially we had three parameters in this case β_0 , β_1 and σ^2 .

So, we already have estimated β_0 and β_1 , now we will try to estimate σ^2 also. So, in this lecture I will try to do all the things and also I will try to show you that how these things can be obtained in the R software. Well, when I try to show you the outcome in the R software there will be many more things. But, I will try to consider only some part of that outcome and remaining part I will try to take when we try to discuss the multiple linear regression model, ok.

So, now, how to find out the variance of b_0 and b_1 ? Well, I do not know. Now how to obtain it? So obviously, first we need to do some algebra and we have to study the properties of this b_0 and b_1 . I do not know whether these are good estimator or bad estimator because we have just obtained the values.

These values can be good, these values can be bad. So, how to assure that the values which you have obtain through b_0 and b_1 they are good values and hence your model is good. So, for that we need to investigate the theoretical properties of b_0 and b_1 , ok. So, I will try to give you the details of the theoretical properties and I will try to connect it with the data sciences ok.

(Refer Slide Time: 04:27)

Properties of the direct regression estimators

Unbiasedness b_1 and b_0

$E(b_1) = \beta_1$ Thus b_1 is an unbiased estimator of β_1 .

$E(b_0) = \beta_0$ Thus b_0 is an unbiased estimator of β_0 .

$y = \beta_0 + \beta_1 x + \varepsilon$

$b_1 = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2}$

$b_0 = \bar{y} - b_1 \bar{x}$

Variances and covariance of b_1 and b_0

$Var(b_1) = \frac{\sigma^2}{s_{xx}}$

$Var(b_0) = \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{s_{xx}} \right)$

$Cov(b_0, b_1) = -\frac{\bar{x}}{s_{xx}} \sigma^2$

Find $\hat{\sigma}^2$

dependent on σ^2

statistic Random Variables

Unbiasedness

Variability

Consistency

Sufficiency

Completeness

So, let us now begin with over lecture. First I try to give you the outcome and then I will try to show you the proof, right. So, you have considered the model $y = \beta_0 + \beta_1 x + \varepsilon$, right and we have obtained this β_0 and β_1 . How, β_1 is obtained by $b_1 = \frac{s_{xy}}{s_{xx}}$ and β_0 is obtained by $b_0 = \bar{y} - b_1 \bar{x}$. So, you can see here this depends on y , this also b_1 also depends on y . So, hence they are statistic, they are estimator, they are random variables right.

So, now, so, once they are estimators we would like to see whether they have nice properties or not. So, whenever we come to any statistical estimation procedure there are several criteria, there are several properties through which we try to check whether the estimator is good or bad. First property is say unbiasedness, then we have another property variability, then we have consistency, then we have sufficiency and we have completeness.

I am not writing it in any order, right? They are these different properties are trying to give different type of information, but here I am going to consider this unbiasedness variability and consistency and sufficiency and completeness will automatically follow their, ok.

The first let me give you the final outcome. Both b_1 and b_0 , they are the unbiased estimator of their respective regression coefficient $E(b_1) = \beta_1$. So, b_1 is an unbiased estimator of β_1 and similarly, $E(b_0) = \beta_0$; that means, b_0 is an unbiased estimator of β_0 .

Now, in case if you try to find out the variance of b_0 and b_1 , then they are finally, obtained as $Var(b_1) = \frac{\sigma^2}{s_{xx}}$. $Var(b_0) = \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{s_{xx}} \right)$.

And, if you try to find out the $Cov(b_0, b_1) = -\frac{\bar{x}}{s_{xx}} \sigma^2$. So, now, you can observe hear one thing that if you want to know the amount of variation or co-variation of these estimators in your dataset you have to compute these three values.

But, they are dependent on σ^2 which is the population value and you do not know. So, if you really want to know these values for a given sample then you need to find the value of σ^2 which we can denote as a $\hat{\sigma}^2$. Once you obtain the value of $\hat{\sigma}^2$ then you can replace σ squares in these three expressions by $\hat{\sigma}^2$ and you can compute the sampling variability, ok.

(Refer Slide Time: 08:12)

Properties of the direct regression estimators

Estimate of variances and covariance of b_0 and b_1

Let the sum of squares due to residuals be $SS_{res} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$

An unbiased estimator of σ^2 is $s^2 = \frac{SS_{res}}{n-2}$

Thus $\widehat{Var}(b_0) = s^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{s_{xx}} \right)$ $\sqrt{\widehat{Var}(b_0)}$

$\widehat{Var}(b_1) = \frac{s^2}{s_{xx}}$ $\sqrt{\widehat{Var}(b_1)}$

$\widehat{Cov}(b_0, b_1) = -\frac{\bar{x}}{s_{xx}} s^2$

Standard Error of b_0 and b_1 : Positive square root of variance

So, that is what I am trying to show you here. So, this is the final visual that we are going to obtain after doing some algebra we will consider here a quantity which is called as

sum of square due to residuals. So, you remember that on the last time we had obtained the residuals e_1, e_2, \dots, e_n .

So, you simply try to find out there $\sum_{i=1}^n e_i^2$ and this is called as sum of residuals which is

here like this. So, an unbiased estimator of σ^2 is this quantity $s^2 = \frac{SS_{res}}{n-2}$, right. So, how

it is obtained that I will try to show you in the further slides.

So, now, I try to do one thing that in this expression what we have obtained here this σ^2 will be replaced by s^2 , this σ^2 will be replaced by s^2 and this σ^2 will be replaced by s^2 . So, that is what I am doing here in this slide and thus we can find out the unbiased estimator of variance of b_0 given by $\widehat{\text{var}}(b_0)$ and $\widehat{\text{var}}(b_1)$ and the estimator of covariance between b_0 and b_1 as $\widehat{\text{Cov}}(b_0, b_1)$.

So, I simply have replaced the σ^2 by here s^2 . So, based on this we can compute the standard errors of b_0 and b_1 . So, what we have to do? We simply have to take the positive square root of variance of b_0 that will give us the standard error of b_0 and then we try to take the $+\sqrt{\widehat{\text{var}}(b_1)}$ and that will give us the standard error of b_1 . So, simply try to take the $+\sqrt{\widehat{\text{var}}(b_0)}$ and $+\sqrt{\widehat{\text{var}}(b_1)}$. So, that will give us the standard errors right, ok.

(Refer Slide Time: 10:32)

Model fitting using R: Example- Estimation of σ^2 and standard errors of b_0 and b_1

Use the command

```
summary(lm(y~x))
```

There are several outcomes but we discuss here only some.

So, now the first question comes first I try to show you that how these things are obtained in R software. So, I will try to use the same data that I use in the earlier same lecture where we had obtained the 20 observations on the students on their marks and that number of hours they study.

So, but before that in order to extract the information on the variances of b_0 , b_1 means their estimate and the value of σ^2 we try to use the command here summary. So, summary of the `lm(y ~ x)` and, but this command will give you several outcomes and you have to identify that which of the outcome is going to give you the standard error of b_0 , b_1 as well as the value of σ^2 .

(Refer Slide Time: 11:32)

Model fitting using R: Example- Estimation of σ^2 and standard errors of b_0 and b_1

```
R Console
> summary(lm(y ~ x))

Call:
lm(formula = y ~ x)

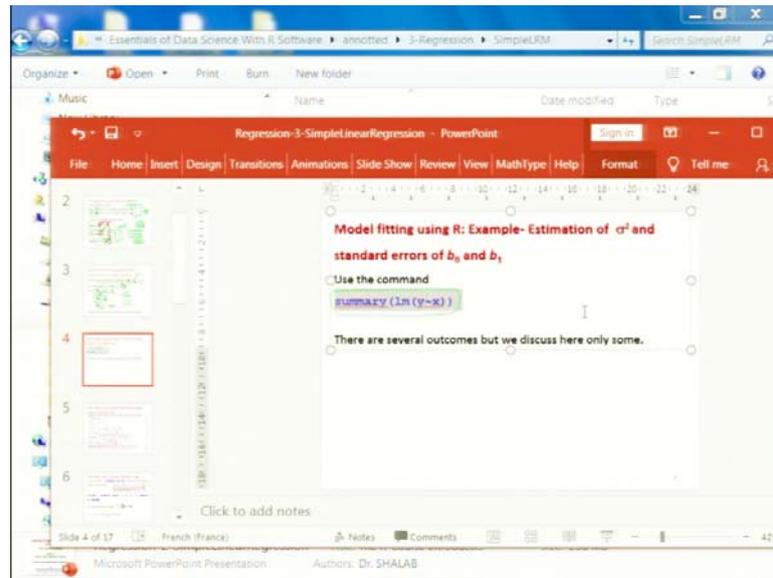
Residuals:
    Min       1Q   Median       3Q      Max
-49.195  -9.504   1.387   8.961  31.683

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) 75.7791    12.7862   5.927 1.31e-05 ***
x            3.4066     0.4379   7.778 3.65e-07 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 18.65 on 18 degrees of freedom
Multiple R-squared:  0.7707,    Adjusted R-squared:  0.758
F-statistic: 60.51 on 1 and 18 DF,  p-value: 3.647e-07
```

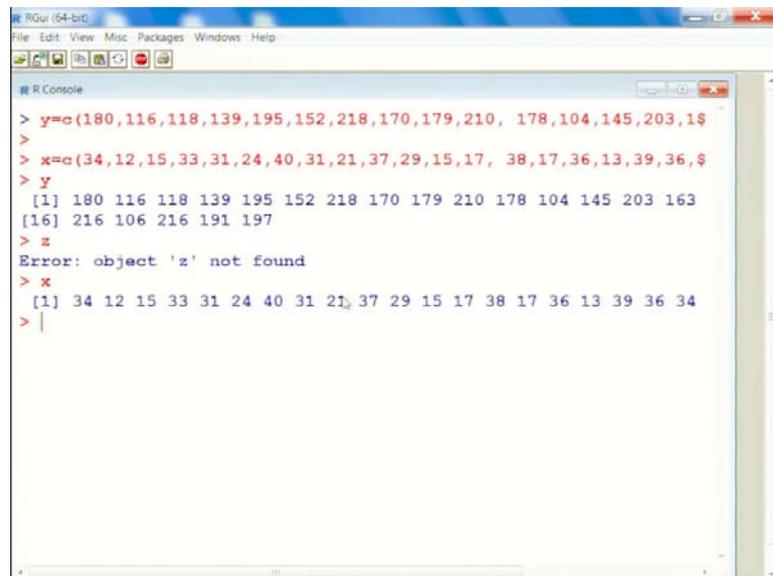
So, before I try to go to the R console I means I have taken here the screen shot to explain you in a better way, right or if you wish I can show you first on the R console so that you have a confidence on me that whatever I am doing here this is correct.

(Refer Slide Time: 11:49)



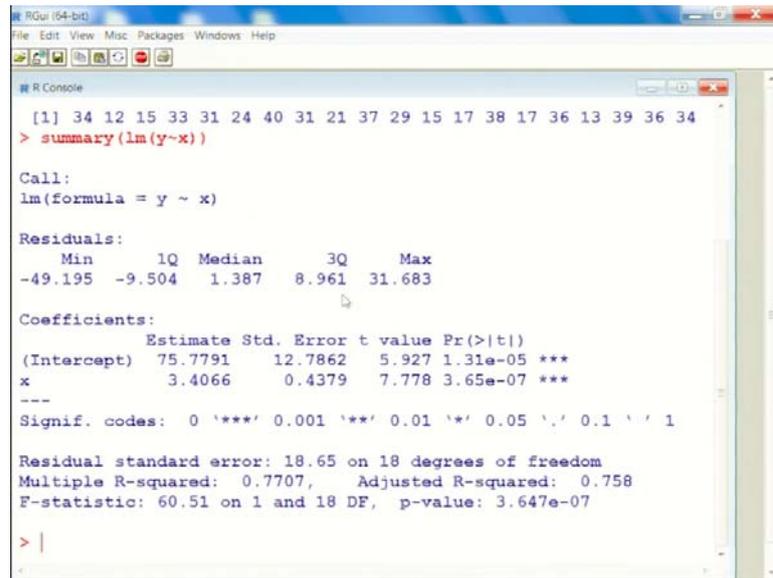
So, I try to take here the command here and I already had entered this data.

(Refer Slide Time: 12:05)



So, y is here like this, x is here like this. Sorry, x is here like this.

(Refer Slide Time: 12:14)



```
RGui [64-bit]
File Edit View Misc Packages Windows Help

# R Console

[1] 34 12 15 33 31 24 40 31 21 37 29 15 17 38 17 36 13 39 36 34
> summary(lm(y~x))

Call:
lm(formula = y ~ x)

Residuals:
    Min       1Q   Median       3Q      Max
-49.195  -9.504   1.387   8.961  31.683

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  75.7791    12.7862   5.927 1.31e-05 ***
x             3.4066     0.4379   7.778 3.65e-07 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 18.65 on 18 degrees of freedom
Multiple R-squared:  0.7707,    Adjusted R-squared:  0.758
F-statistic: 60.51 on 1 and 18 DF,  p-value: 3.647e-07

> |
```

And, now I try to find out the summary command, you can see here these are the values. So, I have just taken the screen shot of the same thing on this a slide. So, now in this output you have to just look for a minute, so that you can see there is something like here this part. You can see here this is the intercept m and this is here x and this is here estimate. So, this part is known towards this is the value of b_0 75.79 and this is the value of b_1 that 3.4, right ok.

After this there is a column standard error. So, this is actually trying to give you the value of standard error of intercept term b_0 and this is the second value which is trying to give you the standard error of the regression coefficient associated with x . And, here there is another thing which I have marked in red box, this is residual standard error. So, this is trying to give you the value of s . So, if you try to square.

(Refer Slide Time: 13:55)

Model fitting using R: Example- Estimation of σ^2

In the outcome of `summary(lm(y~x))`, observe the following marked as encircled 1

```
summary codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Residual standard error: 18.65 on 18 degrees of freedom

Multiple R-squared: 0.7707, Adjusted R-squared: 0.7411

An unbiased estimator of σ^2 is $s^2 = \frac{SS_{res}}{n-2}$

$s = \sqrt{\frac{SS_{res}}{n-2}} = 18.65$

$n - 2 = 20 - 2 = 18$ degrees of freedom

$\hat{\sigma}^2 = s^2$
↓
 $\lambda \cdot v$
↓
distribution

So, now in order to explain you these things more I have taken the snapshots of this and the they are actually here if you try to see they are here. So, first I try to look at the snapshot number 1. It is trying to give you here residual standard error is equal to 18.65 on 18 degrees of freedom.

So, what is this degrees of freedom because I will try to explain you later, but you can imagine that you are trying to estimate σ^2 by $\hat{\sigma}^2$ is equal to s^2 and this is here a random variable and this random variable will have a sampling distribution and this sampling distribution will actually turn out to be related to chi square and the chi square will have certain degrees of freedom. So, these degrees of freedom are mentioned here, right.

So, we had learnt that an unbiased estimator of σ^2 was obtained here like this. So, this is trying to give you the residual standard error means it is trying to take the positive square of this quantity and this quantity here is s which is equal to 18.65. So, this will give you an idea about the variation in the data.

So, now degrees of freedom are computed by n- 2 which are 20- 2 is equal to 18 degrees of freedom, right. So, that is the interpretation of first part. Similarly, if you come to the come to the second snap shot say 2 which is you can see here this is here, right.

(Refer Slide Time: 15:30)

Model fitting using R: Example- Standard errors of b_0 and b_1

In the outcome of `summary(lm(y~x))`, observe the following marked as encircled 2

Coefficients:

	Estimate	Std. Error	t	Pr > t
(Intercept) b_0	75.7791	12.7862	5.936	0.000166
x b_1	3.4066	0.4379	7.779	0.000166

$se(b_0) = \sqrt{\widehat{Var}(b_0)} = \sqrt{s^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{s_{xx}} \right)} = 12.7862$
 $se(b_1) = \sqrt{\widehat{Var}(b_1)} = \sqrt{\frac{s^2}{s_{xx}}} = 0.4379$

So, I am just copying and pasting it here. So, these are the estimates. So, this is the value of b_0 this is the value of b_1 and these are the standard error of b_0 and this is the standard error of b_1 . How they have been obtained? We had obtained that the estimate of variance of b_0 is given by this factor and if you want to find out the standard error you have to take the square root of this estimated variance. So, it is trying to compute it on the basis of given set of data because x and y are known to us.

So, you can compute all these quantities s^2 , \bar{x}^2 , s_{xx} and this values turns out to be 12.78. And, similarly this quantity here 0.4379 this is the standard error of the regression coefficient which is estimated by b_1 and this is the standard error here like this.

So, we had seen that the estimate of variance of b_1 was obtained by the expression s^2/s_{xx} and if you try to take the positive square this will give you the standard error of b_1 which is here coming out to be 0.4379. So, this is how you are going to interpret these 2 outcomes from the software right, ok.

(Refer Slide Time: 16:51)

Properties of the direct regression estimators

Unbiased property: Proof

Note that $b_1 = \frac{s_{xy}}{s_{xx}}$ and $b_0 = \bar{y} - b_1 \bar{x}$ are the linear combinations of $y_i (i = 1, \dots, n)$.

Therefore $b_1 = \sum_{i=1}^n k_i y_i$

where $k_i = \frac{(x_i - \bar{x})}{s_{xx}}$. Note that $\sum_{i=1}^n k_i = 0$ and $\sum_{i=1}^n k_i x_i = 1$.

$$E(b_1) = \sum_{i=1}^n k_i E(y_i)$$

$$= \sum_{i=1}^n k_i (\beta_0 + \beta_1 x_i) = \beta_0 \sum_{i=1}^n k_i + \beta_1 \sum_{i=1}^n k_i x_i$$

$$= \beta_1$$

Thus b_1 is an unbiased estimator of β_1 .

Handwritten notes on the slide:

- $b_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$
- $= \sum k_i (y_i - \bar{y}) = \sum k_i y_i - \bar{y} \sum k_i$
- $\sum k_i = \sum \frac{(x_i - \bar{x})}{s_{xx}} = \frac{\sum (x_i - \bar{x})}{s_{xx}} = \frac{\sum x_i - n\bar{x}}{s_{xx}} = \frac{n\bar{x} - n\bar{x}}{s_{xx}} = 0$
- $\sum k_i x_i = \sum \frac{(x_i - \bar{x})x_i}{s_{xx}} = \frac{\sum (x_i^2 - \bar{x}x_i)}{s_{xx}} = \frac{\sum x_i^2 - \bar{x} \sum x_i}{s_{xx}} = \frac{\sum x_i^2 - \bar{x} n\bar{x}}{s_{xx}} = \frac{\sum x_i^2 - n\bar{x}^2}{s_{xx}} = \frac{\sum (x_i - \bar{x})^2 + n\bar{x}^2 - n\bar{x}^2}{s_{xx}} = \frac{\sum (x_i - \bar{x})^2}{s_{xx}} = 1$
- $E(y_i) = \beta_0 + \beta_1 x_i$

So, now to convince you that whatever we have used, how they have been obtained and how we can assure that the properties of unbiasedness standard error etc. are correct. So, now, I am going to consider the proofs of these results. So, I will try to give you a brief description.

But my request is that unless and until you do this proof yourself using your own pen own book with your own hand it is difficult for you to understand. So, I will try to give you the steps here, but you are requested to do the same thing yourself. So, we had obtained here b_1 by s_{xy} upon s_{xx} and b_0 was y bar- b_1 x bar, right.

And both if you try to see they are the linear combinations of y_i 's. They are the function of y_i 's. So, at least if I try to consider here b_1 this I can write down as say b_1 is equal to

$$\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) / s_{xx}, \text{ right. So, in case if I try to write down this quantity here as say } k_i.$$

So, I can write down here this thing here is say k goes from $\sum_{i=1}^n k_i y_i - \bar{y}$, right and if

you try to make it more simpler this will become $\sum_{i=1}^n k_i y_i - \bar{y}$ summation i goes from 1

to n k_i right and summation k_i if you try to see this is $\sum_{i=1}^n \frac{(x_i - \bar{x})}{s_{xx}}$.

So, this is n times, ok I can write down $\frac{(x_i - \bar{x})}{s_{xx}}$ and $\sum_{i=1}^n \frac{(x_i - \bar{x})}{s_{xx}}$. So, this becomes here

0. So, that is why I am trying to write down here b_1 as a summation $\sum_{i=1}^n k_i y_i$, right. So,

and you can verify they are the $\sum_{i=1}^n k_i = 0$ and $\sum_{i=1}^n k_i x_i = 1$, ok.

So, now if you try to take the $E(b_1)$ then I can write down $\sum_{i=1}^n k_i E(y_i)$ and $E(y_i) = b_0 + b_1 x_i$

because expected value of i , i is 0 right. So, hence I can write down this quantity here as

a $\sum_{i=1}^n k_i \beta_0 + \beta_1 \sum_{i=1}^n k_i x_i$, right.

So, $\sum_{i=1}^n k_i = 0$. So, the first part becomes 0 and $\sum_{i=1}^n k_i x_i = 1$ that you can verify. So, I get

here $E(b_1) = \beta_1$ that is and that establishes that b_1 is an unbiased estimator of β_1 .

(Refer Slide Time: 20:01)

Properties of the direct regression estimators

Unbiased property: Proof

Next

$$\begin{aligned}
 E(b_0) &= E[\bar{y} - b_1 \bar{x}] \\
 &= E[\beta_0 + \beta_1 \bar{x} - b_1 \bar{x}] \\
 &= \beta_0 + \beta_1 \bar{x} - \beta_1 \bar{x} \\
 &= \beta_0.
 \end{aligned}$$

Thus b_0 is an unbiased estimator of β_0

$\bar{\varepsilon} = 0$ $E(\varepsilon) = 0$
 $\frac{1}{n} \sum_{i=1}^n \varepsilon_i = 0$
 $\bar{\varepsilon} = \frac{1}{n} \sum_{i=1}^n \varepsilon_i$
 $E(\bar{\varepsilon}) = \frac{1}{n} \sum_{i=1}^n E(\varepsilon_i)$
 $= 0$

Similarly, if you want to check the unbiasedness of b_0 then we try to consider here $E(b_0) = E[\bar{y} - b_1 \bar{x}]$. So, if you try to find out the value of \bar{y} this is your $\beta_0 + \beta_1 \bar{x}$. Do not write that that that $\bar{\varepsilon}$ is equal to 0.

Because you have to remember one thing we have made the assumption that $E(\varepsilon) = 0$ that mean this is over the whole population 1 over population size i goes from 1 to the entire population ε_i is equal to 0. But, when I am trying to write down here $\bar{\varepsilon}$ this is 1 upon n summation i goes from 1 to n ε_i this will not be equal to 0, but when you try to take the expected value of $\bar{\varepsilon}$ then this will become $\frac{1}{n} \sum_{i=1}^n E(\varepsilon_i) = 0$.

So, this is what I have done here and then- here $b_1 \bar{x}$. So, if you try to just take this quantity $b_0 + b_1 \bar{x}$ because x is the x is nonstochastic, x is fixed. So, I can take out this quantity outside the bracket and $E(b_1 \bar{x})$ this is that we already have proved that $E(b_1) = \beta_1$. So, I can write down $\beta_1 \bar{x}$. So, these two terms get cancelled out and we get here β_0 . So, this b_0 is also an unbiased estimator of β_0 , ok.

(Refer Slide Time: 21:41)

Properties of the direct regression estimators

Variance of b_1 : Proof

Using the assumption that y_i 's are independently distributed,

$$E(y_i) = \beta_0 + \beta_1 x_i$$

$$b_1 = \sum k_i y_i$$

$$Var(b_1) = \sum_{i=1}^n k_i^2 Var(y_i) + \sum_i \sum_{j \neq i} k_i k_j Cov(y_i, y_j)$$

$= 0$

$$= \sigma^2 \frac{\sum (x_i - \bar{x})^2}{s_{xx}} \quad (Cov(y_i, y_j) = 0 \text{ as } y_1, \dots, y_n \text{ are independent})$$

$$= \frac{\sigma^2 s_{xx}}{s_{xx}^2}$$

$$= \frac{\sigma^2}{s_{xx}}$$

$\sum_{i=1}^n k_i^2 \sigma^2$

10

So, after establishing the unbiasedness property of b_0 and b_1 , let us try to find out the variances. So, we are assuming here that since ϵ_i 's are IID they have got mean 0 and variance σ^2 . So, this y_i 's are also independently distributed. They will not be identically distributed remember because $E(y_i) = \beta_0 + \beta_1 x_i$ and it depends on x_i . So, if x_i is changing this mean will be changing. So, I can only assume here that y_i 's are independently distributed.

So, since I already have expressed b_1 as say $\sum_{i=1}^n k_i y_i$. So, I can write down

the $Var(b_1) = \sum_{i=1}^n k_i^2 Var(y_i) + \sum_i \sum_{j \neq i} k_i k_j Cov(y_i, y_j)$. But, since I am assuming that y_i and y_j

are independent so, this covariance will become 0 and this this part the second part will become simply 0.

So, the variance of b_1 comes out to be simply here k_i as summation k_i^2 summation I goes from 1 to n k_i^2 and $var(y_i)$ is fixed σ^2 .

So, this variance comes out to be $\sigma^2 \sum_{i=1}^n k_i^2$ and which is here like this. And, if you try to simplify it here this in $\sigma^2 s_{xx}$ this quantities nothing, but your s_{xx} . So, 1 s_{xx} gets cancelled out here and you get here the variance of b_1 as $\frac{\sigma^2}{s_{xx}}$.

(Refer Slide Time: 23:30)

Properties of the direct regression estimators

Variance of b_0 : Proof

$b_0 = \bar{y} - b_1 \bar{x}$

$$Var(b_0) = Var(\bar{y}) + \bar{x}^2 Var(b_1) - 2\bar{x}Cov(\bar{y}, b_1)$$

First we find that

$$Cov(\bar{y}, b_1) = E[\{\bar{y} - E(\bar{y})\}\{b_1 - E(b_1)\}]$$

$$= E\left[\bar{\varepsilon}\left(\sum_i k_i y_i - \beta_1\right)\right]$$

$$= \frac{1}{n} E\left[\left(\sum_i \varepsilon_i\right)\left(\beta_0 \sum_i k_i + \beta_1 \sum_i k_i x_i + \sum_i k_i \varepsilon_i - \beta_1 \sum_i \varepsilon_i\right)\right]$$

$$= \frac{1}{n} [0 + 0 + 0 + 0] = 0$$

$Var(b_0) = \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{s_{xx}}\right)$ $Var(\bar{y}) = \frac{\sigma^2}{n}$

So, you can see it is not difficult and similarly if you want to find out the variance of b_0 . So, variance of b_0 can be written here like this because $b_0 = \bar{y} - b_1 \bar{x}$. So, this will be $Var(b_0) = Var(\bar{y}) + \bar{x}^2 Var(b_1) - 2\bar{x}Cov(\bar{y}, b_1)$.

So, we already have obtained the variant of b_1 , but we have not obtained the covariance between y bar and b_1 . So, we try to find out the $Cov(\bar{y}, b_1)$ using the basic definition can be written here like this and if you try to substitute the value of \bar{y} and $E(\bar{y})$ and $b_1 - E(b_1)$.

You can simply a just simplify this expression to this thing. And, if you try to substitute here the value of your y_i , well you can obtain it directly also, but to make it more simple I will ask you to just substitute her this value of $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$ and try to simplify it here I have given you here the complete expression.

And you can see here once you take this expectation sign inside is all the terms are becoming 0. The first term will become because $\sum_i \varepsilon_i$ is 0 this term also becomes 0 and this term also becomes here 0 and you already have proved that this is equal to 0, right. So, all these terms will become 0. So, $Cov(\bar{y}, b_1) = 0$ and hence variance of b_0 comes out to be $\bar{x}^2/n + s_{xx}$.

So, hence I can write down now here $Var(b_0) = \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{s_{xx}} \right)$. So, you substitute it and you will get here the same expression which we have written earlier ok. And, one thing you have to notice here that yeah that this here bracket sign, right. So, this ε will play an important role in making all these expectations to be 0.

(Refer Slide Time: 25:52)

Properties of the direct regression estimators

Covariance between b_0 and b_1 : Proof

$$Cov(b_0, b_1) = \overset{0}{Cov(\bar{y}, b_1)} - \bar{x}Var(b_1)$$

$$= -\frac{\bar{x}}{s_{xx}} \sigma^2$$

Similarly, if you try to find out the $Cov(b_0, b_1)$, so, that will be simply covariance between $Cov(\bar{y}, b_1) - \bar{x}Var(b_1)$. So, you already have proved that this quantity is 0. So, and $Var(b_1)$ you already have obtained just substitute it here and you will get the $Cov(b_0, b_1)$.

(Refer Slide Time: 26:16)

Properties of the direct regression estimators

BLUE property

The ordinary least squares estimators b_0 and b_1 possess the minimum variance in the class of linear and unbiased estimators.

They are the Best Linear Unbiased Estimators (BLUE).

Such a property is known as the Gauss-Markov theorem.

The diagram shows a parameter β on the left. Three arrows point from β to three different unbiased estimators, labeled $\hat{\beta}_1$, $\hat{\beta}_2$, and $\hat{\beta}_3$. These three arrows are grouped by a bracket. An arrow points from this bracket to a box labeled 'Unbiased'. From the 'Unbiased' box, an arrow points to a circle labeled 'b', which is also circled. This represents the selection of the best linear unbiased estimator (BLUE) based on minimum variance.

13

So, you can see that is not difficult and after this I am just going to give you one property very important property I am not going to do the proof here. This is called the BLUE property B L U E. BLUE property and this property is also mentioned as a Gauss-Markov theorem which states that the ordinary least square estimator b_0 and b_1 possesses the minimum variance in the class of linear and unbiased estimator. What does this mean?

If you try to take any parameter say β can have different estimators – $\hat{\beta}_1, \hat{\beta}_2$, many things first of all you try to sort out all the estimators which are unbiased. So, some of them will be biased some of them will be unbiased. So, we try to choose here all the estimator which are unbiased.

Now, you try to find out the variance of each of the estimator each are unbiased then you will see out of these estimator the estimator which is based on ordinary least square estimation say b that will have the smallest variance. So, this property assures that the way you have computed the value of β_0 and β_1 using the b_0 and b_1 that is giving you an value which is unbiased and which got the minimum variance. So, you can believe on it that these are very good values.

And, we call it as if and we call it as that they are the best linear unbiased estimator. So, this here B L U and here E, BLUE, right. So, this property helps us in in assuring that the

way you have obtained this numerical value they are going to give you the good value, right.

(Refer Slide Time: 28:15)

Residual sum of squares

$$\begin{aligned}
 SS_{\text{res}} &= \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - b_0 - b_1 x_i)^2 \\
 &= \sum_{i=1}^n [(y_i - \bar{y} + b_1 \bar{x} - b_1 x_i)]^2 = \sum_{i=1}^n [(y_i - \bar{y}) - b_1 (x_i - \bar{x})]^2 \\
 &= \sum_{i=1}^n (y_i - \bar{y})^2 + b_1^2 \sum_{i=1}^n (x_i - \bar{x})^2 - 2b_1 \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \\
 &= s_{yy} + b_1^2 s_{xx} - 2b_1 s_{xy} \quad \text{where } s_{yy} = \sum_{i=1}^n (y_i - \bar{y})^2, \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i \\
 &= s_{yy} - b_1^2 s_{xx} \\
 &= s_{yy} - \left(\frac{s_{xy}}{s_{xx}} \right)^2 s_{xx} = s_{yy} - \frac{s_{xy}^2}{s_{xx}} \\
 &= s_{yy} - b_1 s_{xy}
 \end{aligned}$$

$b_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{s_{xx}}$

Now, I try to obtain the estimator of σ^2 , but before that let me try to do some homework. We try to define here the quantity residual some of squares. Residuals some of squares is denoted by SS_{res} . SS and in the subscript res this is defined as sum of squares of the residuals.

So, this is $\sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2$. The values of \hat{y}_i can be written here as $b_0 + b_1 x_i$ and if you try to write down here all these things over here and then I can simplify this quantity as $\sum_{i=1}^n [(y_i - \bar{y}) - b_1 (x_i - \bar{x})]^2$. And, if you try to open it this will come out to be the $\sum_{i=1}^n (y_i - \bar{y})^2$ + the square of second quantity.

And the cross product of first and second quantity, right. If you try to identify what are these things the first term is your s_{yy} which is $\sum_{i=1}^n (y_i - \bar{y})^2$. The second term is b_1^2 and this term here is s_{xx} and what about the third term? Third term if you remember b_1 was

$$\frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \text{ which is } \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{s_{xx}}. \text{ So, this quantity } \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \text{ can be written as}$$

$b_1 s_{xx}$.

So, this quantity here is written as say $b_1 s_{xx}$ and so, this becomes here- $2 b_1^2 s_{xx}$. So, if you try to simplify it just try to substitute the values of b_1 here I can obtain this expression over here and finally, this will come out to be like this $s_{yy} - s_{xy}^2/s_{xx}$ and if you try to write down this quantity as a s_{xy}/s_{xx} into xy . So, this can be written $s_{yy} - b_1 s_{xy}$, right. So, this is the form of SS_{res} ok.

(Refer Slide Time: 30:44)

Estimation of σ^2 : Proof

The estimator of σ^2 is obtained from residual sum of squares.

Since y_i is normally distributed, so SS_{res} has a χ^2 distribution with $(n - 2)$ degrees of freedom, so

$y_i \sim \text{Normally}$

$$\frac{SS_{res}}{\sigma^2} \sim \chi^2(n-2). \quad \text{with } y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$$

Using the result about the

$E(\text{chi-square random variable}) = \text{degrees of freedom}$,

$$E(SS_{res}) = (n-2)\sigma^2. \quad E\left(\frac{SS_{res}}{\sigma^2}\right) = n-2$$

Thus an unbiased estimator of σ^2 is

$$s^2 = \frac{SS_{res}}{n-2} \quad \sigma^2 = \frac{SS_{res}}{n-2}$$

Now, we can obtain the estimate of σ^2 using this SS_{res} , right. Since and you can see here this SS_{res} is a function of y_i 's and now we assume here that y_i are say normally distributed. They have got a normal distribution with mean $\beta_0 + \beta_1 x_i$ and variance σ^2 .

So, I am now trying to associate distribution with y_i because if you want to use the property of chi square the random variable y_i has to be normally distributed. So, assuming that y_i is normally distributed this SS_{res} has a chi square distribution with $n - 2$ degrees of freedom.

Well, the I'm not giving you here the prove because that is the part of statistics and we use it as a standard result. So, I can write down here SS_{res} / σ^2 will follow a chi square

with $n - 2$ degrees of freedom and in case if I try to find out the with the expectation of a chi square random variable, then the expectation of a chi square random variable is the same as the degrees of freedom. So, my random variable here is $SS_{\text{res}} / \sigma^2$ and which is equal to here $n - 2$.

So, I can write down here expected value of $SS_{\text{res}} = (n - 2)\sigma^2$ and from there I can write down here that $\hat{\sigma}^2$ is $SS_{\text{res}} / (n - 2)$ right. Or even I can write down here it in a much simpler way $SS_{\text{res}} / (n - 2) = \sigma^2$ and from there I can obtain that $\hat{\sigma}^2$ and $\hat{\sigma}^2$ has been indicated by s^2 . So, now, s^2 becomes an unbiased estimator of the σ^2 , right.

(Refer Slide Time: 32:52)

Estimate of variances of b_0 and b_1 : Proof

The estimators of variances of b_0 and b_1 are obtained by replacing σ^2 by $\hat{\sigma}^2 = s^2$ as follows:

$$\widehat{Var}(b_0) = s^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{s_{xx}} \right)$$

$$\widehat{Var}(b_1) = \frac{s^2}{s_{xx}}$$

= s.e.

16

So, now you can see these things are not very difficult and once you have obtained the estimate of σ^2 , now you try to consider the variance of b_0 and b_1 and try to replace σ^2 by s^2 here. So, this will give you the estimators of the variance of b_0 and variance of b_1 .

And, if you try to take the square positive square root of these quantity, so, this will give you the standard errors, right. So, it is not difficult at all. You can see a little bit thing, little bit knowledge of statistics and small algebra is only needed to obtain these things.

(Refer Slide Time: 33:25)

Some properties

- (i) $\sum_{i=1}^n x_i e_i = 0,$
- (ii) $\sum_{i=1}^n \hat{y}_i e_i = 0,$
- (iii) $\sum_{i=1}^n y_i = \sum_{i=1}^n \hat{y}_i$ and
- (iv) the fitted line always passes through (\bar{x}, \bar{y}) .

17

Well, after that just for your information possibly these properties may be useful for you I am just trying to state here simple properties of these estimator that if you try to obtain $\sum_{i=1}^n x_i e_i = 0$, if you try to obtain $\sum_{i=1}^n \hat{y}_i e_i = 0$, and $\sum_{i=1}^n y_i = \sum_{i=1}^n \hat{y}_i$. And, the fitted line which you have obtained this will always pass through the sample means \bar{x} and \bar{y} .

Well, these properties are not actually difficult to prove. If you try to look into any of the book it will just a matter of only couple of minutes. So, do not think that these are very difficult. So, that is why I am not giving you here, but my objective as I said it is essentially to motivate you to for decision sciences and for data sciences.

So, that is the reason I am not giving you the proof of these properties. So, now we come to an end for this lecture. Now, this lecture was very important at least in my opinion why because I established that whatever you are trying to compute those values have to be good values.

So, you are trying to compute β_0 and β_1 , so, you have to assure that you are computing a good value. So, for that I have given you Gauss Markov theorem and that will insure that those numerical values are good; good in the sense of having some nice statistical properties, right.

And, then I have shown you that different outcomes of R software they are not coming from sky they are simply computed and one should know that how they are computed because that will help at a later stage. How? When you are trying to work with a very huge dataset, then you will not have an opportunity to look the data physically with your eyes, but you have to look into these values. For example, if you find that the value of $\hat{\sigma}^2$ or s^2 is coming not to be very huge.

Then possibly you can look into the dataset and possibly you can do something to control the variability first because if the data has lots variability your model will not be good. That is the first principle. Sometimes whatever outcome you are getting here they may not really match with what is happening in the real life.

These types of observation will trigger inside your mind and would try to indicate that there is something wrong somewhere. Where wrong? That we do not know we have to look into the data, we have to look into these values and we have to find what wrong is happening at what point because of this the outcome is not matching with the real values.

For example, sometime you see that the coefficient of the β_1 is coming out to be negative whereas, you can see that when the values of independent variables are increasing then the response is also increasing. So, why this is happening there can be some reasons.

For example, although I am not considering here, but if you try to see such an outcome that will trigger in your mind that well there can be problem of multicollinearity into data and then you try to look into those diagnostic, those tools and try to solve the problem. As I said whenever the data comes data will not follow your rule data will have its own rules, that will have its own conditions and you are the one who should know how to control those condition and get the good outcome.

So, that is why it is very important for you that just by looking at the outcome of the software will not be helpful, but you must know that what are the values which are being computed so that if the software outcome is not matching you can look back into your analysis and try to use better tools, maybe some other type of analysis, right. So, ok so, now, you practice you think about it and I will see you in the next lecture. Till then, goodbye.