

Regression Analysis and Forecasting
Prof. Shalabh
Department of Mathematics and Statistics
Indian Institute of Technology – Kanpur

Lecture 08

Testing of Hypothesis and Confidence Interval Estimation in Simple Linear Regression Model

Welcome to lecture number 8 you may recall that in the earlier lecture we had considered a model.

(Refer Slide Time: 00:20)

$y_i = \beta_0 + \beta_1 x_i + \epsilon_i \quad (i=1,2,\dots,n)$
 $H_0: \beta_1 = \beta_{10}$
 σ^2 is unknown
 $E(\epsilon_i) = 0$
 $V(\epsilon_i) = \sigma^2$

$\rightarrow \frac{SS_{\text{res}}}{\sigma^2} \sim \chi^2(n-2)$
 $\rightarrow \hat{\beta}_1 \sim N\left(\beta_1, \frac{\sigma^2}{\sum x_i^2}\right)$

independently distributed

$\frac{\hat{\beta}_1 - \beta_{10}}{\sqrt{\frac{\sigma^2}{\sum x_i^2}}} \sim t(n-2)$ under H_0

$t_1 = \frac{\hat{\beta}_1 - \beta_{10}}{\sqrt{\frac{SS_{\text{res}}}{(n-2)\sum x_i^2}}} \sim t(n-2)$ under H_0

$\hat{\sigma}^2 = \frac{SS_{\text{res}}}{n-2}$

$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$ salient i goes now 1 to n and we were trying to develop a test of hypothesis for $H_0: \beta_1 = \beta_{10}$, and we had consider the case how to test the hypothesis when sigma square is known to us. Now we are going to assume that sigma square is unknown.

Now if you try to see in practice actually this situation is more important or say more applicable than the earlier case where we assumed that sigma square is known to us. We had considered that while doing the linear regression modeling we have a only a set of data x_i and y_i and we have to infer everything on the basis of that sample of data, when I say that sigma square is known that is the value of sigma square has to be known from some past experience or from some other relevant sources.

In real life usually the value of sigma square is not known to us and we need to estimate it from the sample. We had earlier discussed that we can estimate the value of sigma square by

sum of square due to residual that $ss_{res}/n-2$. Now we are going to consider here a situation which is more applicable in real life data set, where we have a set of data on x_i and y_i and we try to estimate all the parameters only from the data set including β_1 as well as σ^2 .

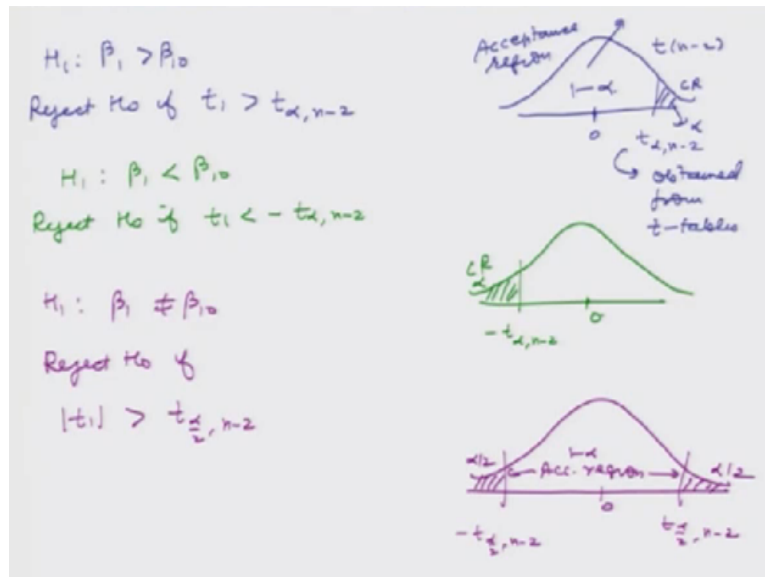
Where we assumed that accepted value $\epsilon_i=0$ and variance of $\epsilon_i=\sigma^2$. So in this case our next objective is how to find out a test statistic, so we know that ss_{res} that is the residual sum of square divided by σ^2 follows a chi-square and $n-2$ degrees of freedom and we also have seen that this $\hat{\beta}_1$ follows a normal distribution with mean β_1 and variance σ^2 upon s_{xx} , and we can also prove that these 2 are independently distributed.

Now I can use a result from the statistical inference that if x is the random variable which is following a chi-square distribution and y is another random variable which is following a normal distribution then y upon square root of x by n this follows a t distribution. So in this case I can see here that this $\hat{\beta}_1$ is a normal random variable this is following a chi-square random variable.

So, I can write down here $\hat{\beta}_1 - \beta_{10}$, which is known to us divided by square root of $\hat{\sigma}^2$ upon s_{xx} this follows a t distribution with $n-2$ degrees of freedom under H_0 , so this can be further written as $\hat{\beta}_1 - \beta_{10}$ upon square ss_{res} divided by $n-2$ s_{xx} , this follows a t distribution with $n-2$ degrees of freedom under H_0 .

Because if you try to recall we had estimated σ^2 by ss_{res} over $n-2$, now this is statistics we can denote by here say t_1 . Now we try see depending upon on the nature of alternative hypothesis we can have different types of decision rule, for example if my alternative hypothesis is $H_0: \beta_1 > \beta_{10}$.

(Refer Slide Time : 04:26)



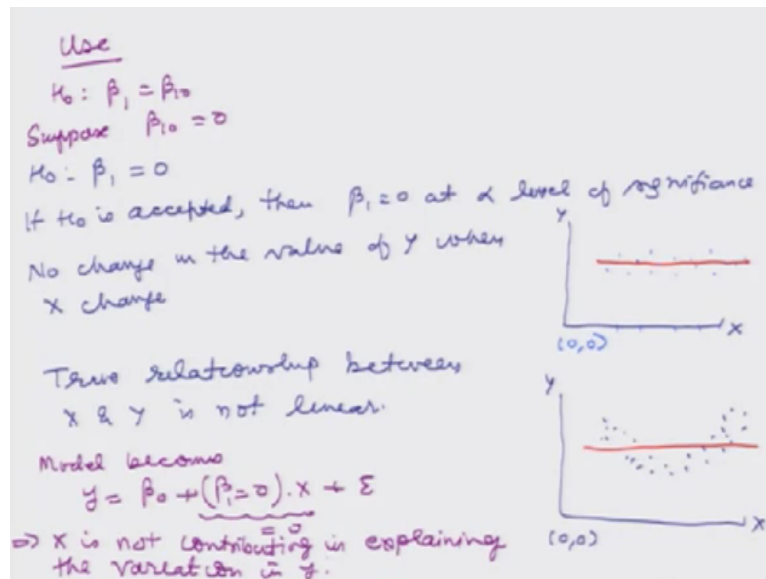
Then in this case as we have done in the earlier case this is the probably density function of t distribution with say here $n-2$, which is freedom and this is going to be the critical region of size α and this is going to be here expectance region of size $1 - \alpha$ and this is the critical value t_{α} and -2 and this value t_{α} $n-2$ because obtain from the t tables.

In this case I would say that reject H_0 if t_1 is greater than t_{α} $n-2$ on the other hand if my alternative hypothesis is H_1 , β_1 let than β_{10} where β_{10} is some known value then in this case this t distribution will be like this and the critical region will be on left hand side of size α and this value will be $-t_{\alpha}$ and -2 and somewhere here will be 0 , and in this case I would say that reject H_0 if t_1 is less than $-t_{\alpha}$ and -2 .

The third case will be of H_1 $\beta_1 \neq \beta_{10}$, in this case again the t distribution will have critical reason on both the sides of size α by 2 and α by 2 and this will be the value $-t_{\alpha}$ by 2 and $n-2$ and this will be the value t_{α} by $n-2$, so I can say here that reject H_0 , if absolute value of t_1 is greater than t_{α} by $2n-2$.

So this will again be our acceptance region of size $1 - \alpha$, so this how we proceed to test the hypothesis about the H_0 , $\beta_1 = \beta_{10}$. Next the question is what is the interpretation, and what is the utility of making this test of hypothesis, so let us try to understand the use.

(Refer Slide Time: 07:09)



Now when I am trying to do this hypotheses like $H_0: \beta_1 = \beta_0$ and suppose $\beta_0 = 0$ so essentially we are trying to now test the hypothesis $\beta_1 = 0$, now there are 2 options either our H_0 is accepted or not so in case if I say if H_0 is accepted then we say that $\beta_1 = 0$ at α level of significance on the bases of given sample of data.

What is interpretation of this thing? We had learnt earlier that β_1 is the slope of the line so when I am saying that $H_0: \beta_1 = 0$ is accepted that means the slope of the fitted line is 0 that means the fitted line is going to be something like this and all the points which are scattered around this line will look like this, so this means here what? This means the slope of the line is 0 that means x is not contributing in explaining the variation.

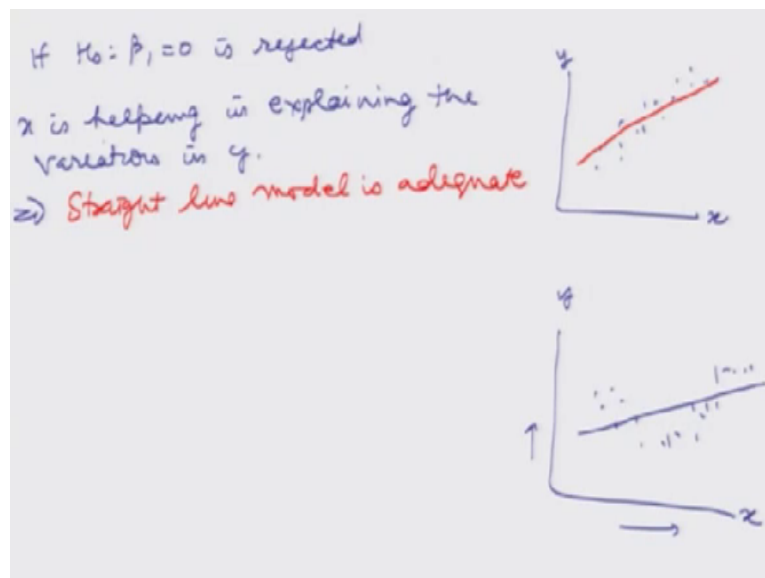
In y whatever happens here either here or here or say here or say here whatever the values x takes there is no change in the value of y , so what we can see here is the following there is no change in the value of y when values of x change. This means x is not contributing in my model building x_1 is such a variable which is not helping us in explaining the variation in the value of y .

The second situation can be it can be like this suppose we have a data like this on, which is really nonlinear and this case when we are trying to fit here a linear model this will look like this So what is this mean, that means there exists some relationship which is actually nonlinear and our x is trying to help the variation, but still it is not capable of explaining the variation.

This case I can say here in simple words the true relationship between x and y is not linear in either of the case $H_0: \beta_1=0$ is indicating the true situation, when the relationship between x and y is nonlinear why should we fit a linear model. Well, in this case in this incase if somebody want to proceed than we have the setup of polynomial regression model that can be used.

Now we can see here was $H_0: \beta_1=0$ is accepted then in that case our model now becomes say $y=\beta_0+\beta_1 x$ which takes active value here zero into $x+\epsilon$, so this quantity becomes $=0$ this implies that x is not contributing in explaining the variation in y on the other hand there can be another situation that suppose $H_0: \beta_1=0$ is rejected.

(Refer Slide Time: 11:42)



So what are we trying to say, we are simply trying to say that in this 2 dimensional graph x and y axis, yes there are some observations like this and we are trying to fit here a model like this 1 that makes that sense, that means x is helping in explaining the variation in y or in simple words this implies that a straight line model is adequate, and that is essential what we wanted.

On the other end if I try take the second case over here where we had a nonlinear relationship like this 1 and suppose if we are trying to fit here a linear model this may look like a this, even in this case we can see here that x is changing the values of y 's are changing, but that is possibly indicating a nonlinear relationship so it is suggested that by considering a polynomial regression model of say higher order like 2 or three that name work in this situation.

If p value is smaller than the significance level, which we had denoted as alpha then the null hypothesis, which we have denoted by H_0 is rejected. This gives us a very simple solution just look at the software try to look at the p value and whenever we are doing the modeling and consequently we are doing test of hypothesis we have to define the level of significant usually it is 5% or 1% depending on the situation.

So alpha =0.05 or 0.01 depending on whether we are using 5% level of significant or 1% level of significance, so you simply try to compare the p value and value of alpha and based on that we can take a correct design about the test of hypothesis. This approach we need not look into the tables of either normal probability or say t probability so that is the advantage of this approach.

We stop here now and would I request all of you to just have a close look whatever we have covered in this lecture and the earlier lecture and try to understand how we have constructed the test statistics and how we are going to use it. Now our next objective will be first of all to develop that test statistics for testing hypothesis for intercept term and for the sigma square, and after that we will concentrate on the confidence interval estimation of all these parameters till then good bye.