**Quality Control and Improvement with MINITAB**
**Prof. Indrajit Mukherjee**
**Shailesh J. Mehta School of Management**
**Indian Institute of Technology, Bombay**

**Lecture - 27**
**Multicollinearity, Best Subset Regression, Multiple Regression, Basics on Design of Experiment**

Hello and welcome to our session 27 on Quality Control and Improvement with MINITAB. I am Professor Indrajit Mukherjee from Shailesh J Mehta School of Management, IIT Bombay. So, previous session what we are doing is that, we are discussing about the multiple regression and multicollinearity problems like that. So, we will take another examples to again emphasize on the importance of how to tackle Multicollinearity like that ok.

(Refer Slide Time: 00:45)



So multicollinearity, what we have explained is that relationship between X variables that exist, and that can distort the relationship between what is and we are not able to generalize the equations which can be used in reality for prediction like that ok. So, we need to deal with this and variation inflation factor that we have to calculate and try to see that it is not more than 5.

**Quality Control and Improvement using MINITAB**

1. Remove some of the highly correlated independent variables

2. Attempt Stepwise Regression for variable selection

3. Attempt Best Subset Regression for variable selection

4. Partial least squares regression

5. Principal components analysis

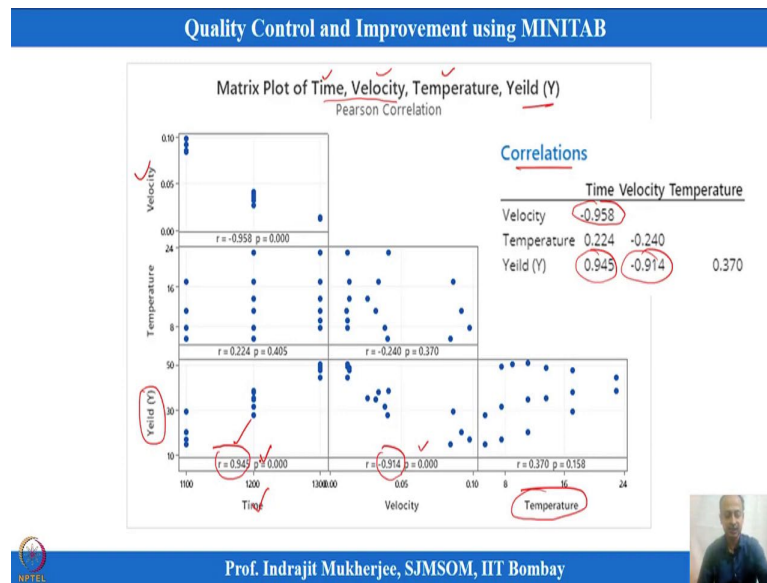Prof. Indrajit Mukherjee, SJMSOM, IIT Bombay

And, if it is more than 5 then we have to adopt some other means to deal with that. And how to deal with that? That we have mentioned over here. That means, highly correlated variables, one of them will we can take and the other ones can be eliminated. And, we can attempt stepwise regression method which takes care of which can suggest that which are the variables to be taken and we can reconfirm the multicollinearity exist or not in the final models like that ok.

And, we can also see best subset regression where multiple options exist because we do not want to stop only considering what stepwise regression gives. Because, some sometimes what happens say, alternate model is easy to control in real life; that means, those variables are easy to control which is suggested.

And if it is fine, if we have to sacrifice some amount of R square values or something like that, that is not a constraint for us in production processes, but variables which are very difficult to control and stepwise regression keeps those variables in the final equation then it becomes difficult for us to change that one, if we adopt that one, so in that case it becomes difficult. So, best subset regression gives you some more options to select the variables like that.

Then, appropriate methods like partial least square regression, principal component regression all these things can be adopted ok.

So, we will take one more examples where we have variables like, yield is the Y characteristics and this depends on time, velocity and temperature. These are the variable X that was considered to be significant or potential variables and we want to get a model out of this regression model out of this. All are continuous variable over here, so in this case.

And when we did the correlation analysis of this correlation analysis, what we observe is that velocity is highly correlated with time variables over here and then we can see that yield is highly correlated with time, velocity over here. So, but this P-Value will indicate whether which of the variable is highly correlated.

So, yield is one of the variables. So this seems to be significant over here, velocity is also significant over here, but temperature does not seem to be significantly influencing the yields over here ok. And, what we can see is that temperature with time is not significant, temperature with velocity is also not significant ok.

So, velocity with time is highly correlated, this with time velocity and time is highly correlated over here. So, this correlation matrix gives you some preliminary information, what is expected means when we run the regression analysis, what is expected in all these variables. So what we can see is that time and velocity are two variables which are highly correlated, so one we have to adopt over here. So how do we select which one to adopt over here? So, maybe when we do a trial and error basis.

If we have to select the variables the policy may be that, because yield is the correlation coefficient is 945 for time and for velocity is 914 maybe this variable we should select like that ok. Time may be the only variable and replace the velocity with the time like that.

So, yield has to be regressed with time and, if temperature is significant then we will include that one. So anyhow, this is the suggested one of the guidelines. So, which is highly correlated with the Y variable, we can select that one. And, we want to see that what happens if we select stepwise regression what happens like that.

(Refer Slide Time: 04:38)



So, this is the data set C11 to C20. This is the data set that we are having and we want to; we want to implement this.

(Refer Slide Time: 04:48)



So, what we have done is that basic statistics over here we have gone to correlation.

(Refer Slide Time: 04:52)



And then we have identified which are the variables over here. We want to understand the relationship. So, yield is the first variable time, velocity and temperature.

(Refer Slide Time: 05:02)



(Refer Slide Time: 05:07)



And, in options what we have given, Pearson correlation we want to see. And then in graphs what we want is that correlation with P-Values.

(Refer Slide Time: 05:11)



And if I click ok, and results what we have mentioned is we can also see pairwise correlation table.

(Refer Slide Time: 05:15)

And if you click ok, this is the; this is the diagram that I have shown in the PPT slides like that. So, this shows relationship between yield and other variables. Yield is highly related with time, the P-Value is point near to 0 and velocity is also highly related with yield over here, but temperature does not seem to be significantly influencing over here.

So, this is the relationship and time is having a high time is not having a correlation with temperature over here ok, and velocity is having no correlation with temperature, but velocity is having a correlation with time over here. So that is reflected over here, temperature and time. So, this is highly 0.958 is the minus 0.958 is the.

(Refer Slide Time: 06:03)



So, this is so what we can do is that we can go to stepwise regression. So I will use stepwise regression over here.

(Refer Slide Time: 06:09)

(Refer Slide Time: 06:11)



So, fit regression model, and here only thing I have to do is that I have to introduce stepwise regression over here, and click ok.

(Refer Slide Time: 06:18)



And then validation we can pair wise validation as usual.

(Refer Slide Time: 06:21)



And also storage we can save the residual standard as residual over here, and I click ok.

(Refer Slide Time: 06:29)



And let us try to see which are the variables will go in the model and which will go out of the model ok. So, what happened is that, it has identified time and temperature.

(Refer Slide Time: 06:35)



So, initially we are thinking that it is not so highly correlated, but this is not so significant 0.05, just above the cross line that is 0.05 that is the cut off over here. So, it is the analysis over here says that, we have to take into consideration time and temperature, that is the best model that is coming out of this.

(Refer Slide Time: 06:55)



And, the R square adjusted value is 90 point something like that and 10 fold cross-validation more or less close.

(Refer Slide Time: 07:03)



So this model seems to be adequate, if I consider that alpha level of significance is 0.1, in that case we can also consider temperature as significant variable we may retain this one.

(Refer Slide Time: 07:13)
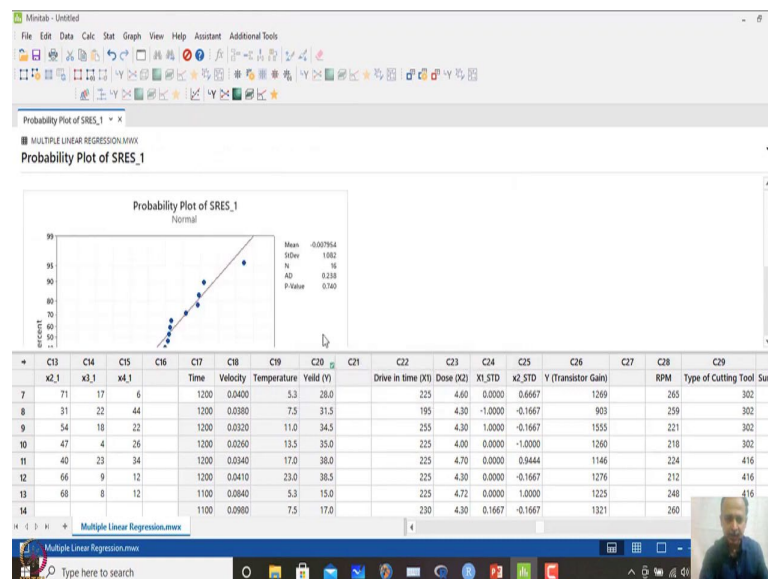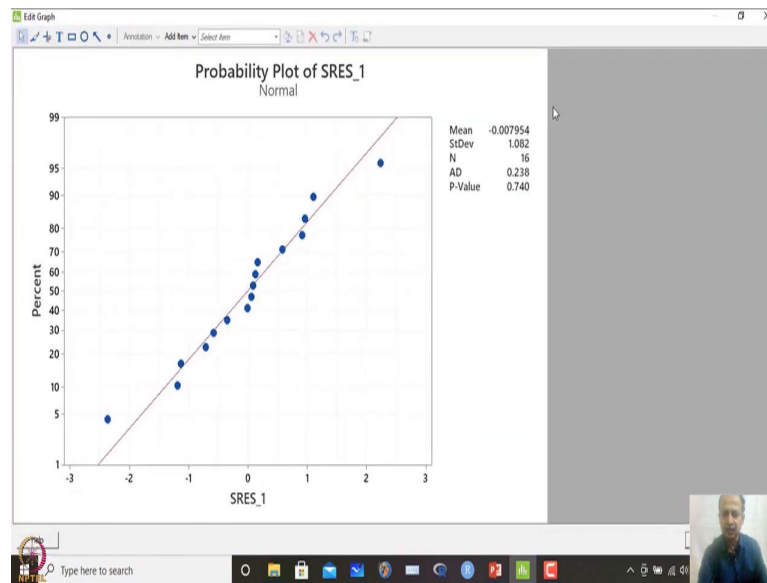
(Refer Slide Time: 07:16)



So, let us try to see what about the normality distributions of this residual over here.
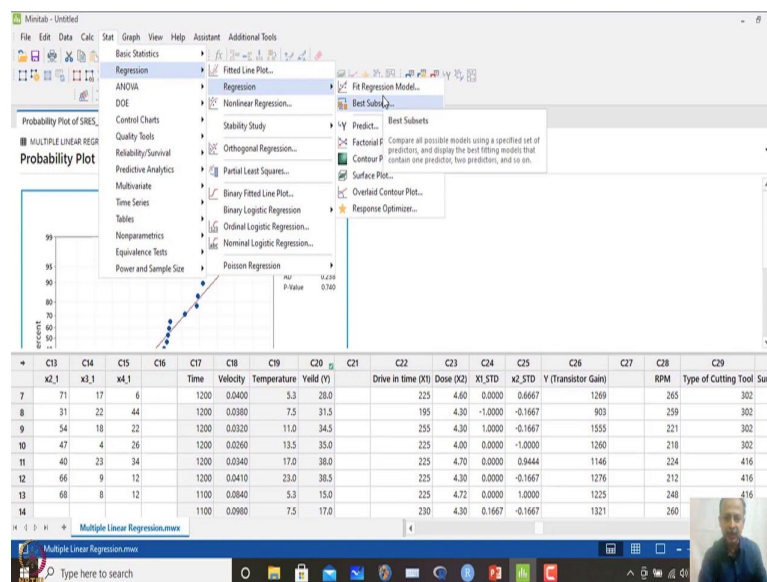
(Refer Slide Time: 07:20)
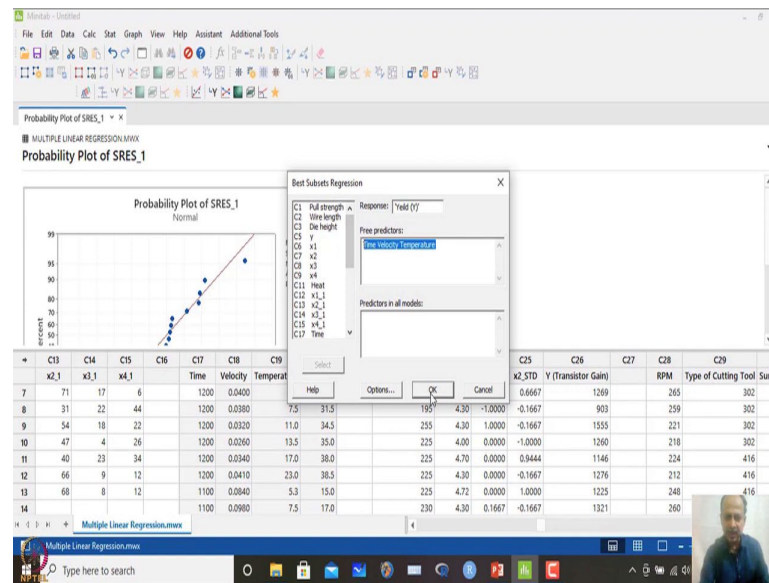
(Refer Slide Time: 07:23)



So, last residual will be normality we can check, and in this case 0.74, so there is not much problem with the normality assumptions over here. And, we do not have any problem over here.
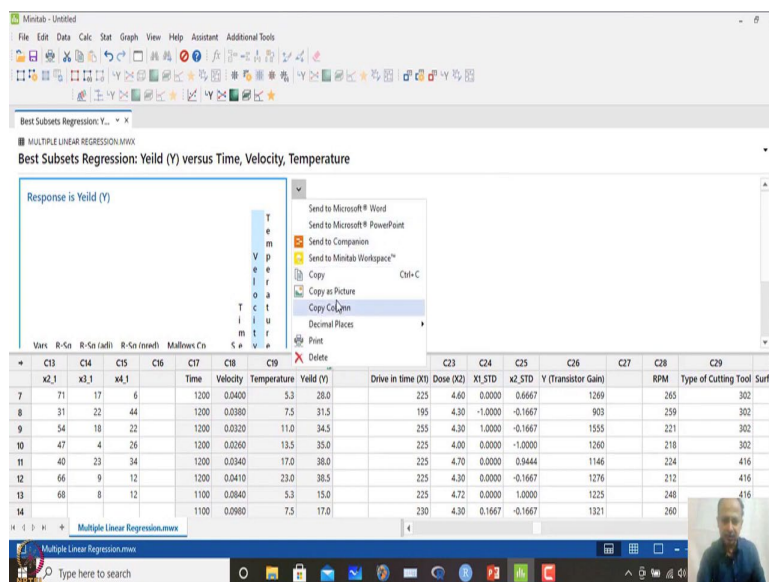
(Refer Slide Time: 07:31)

So, in this case normality is not a constraint. So, in this case we can also see, when I use best subset regression which is the best model it is giving, so I have taken yield over here and we have taken all the three variables.

And if I click ok, then what happens let us try to see. And we can just paste it copy as a picture over here so then we can paste it in excel and try to see what happens ok. So, we can adopt, we will try to see what happens if we select, so excel we are opening one excel sheet ok. So, let us do that one and we can paste it over here.
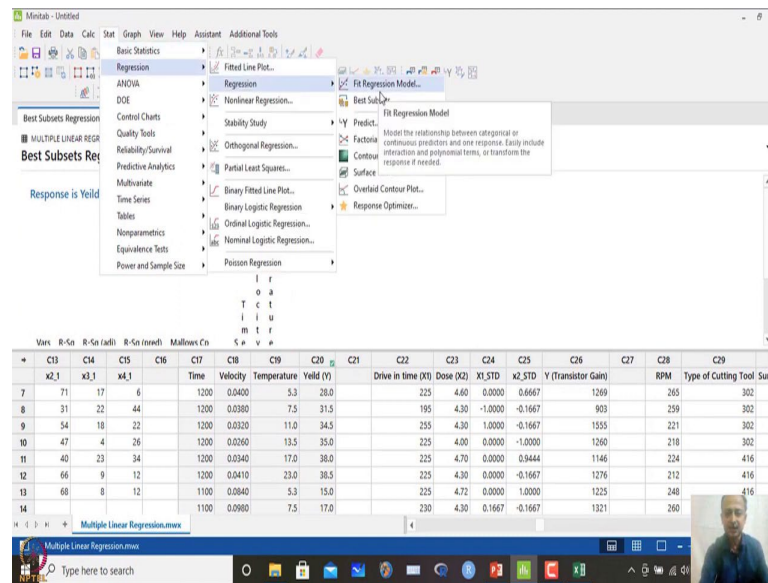
(Refer Slide Time: 08:14)



So, when we are pasting it over here. What, observations we have that we want to see. So, over here this is time variable over here, this is temperature and this is the velocity over here. With one variable what we are seeing mallow Cp is about 4, R square adjusted is 88.5.

When I take consider only time as the variables like that ok, but mallow Cp is higher than the number of variables plus 1. So this is not recommended. Second one is also not recommended with one variable, third one what we are seeing, time and temperature it is mallow Cp is less than the number of variables plus 1, so 3 is 2 is less than 3. So, this seems to be adequate, this one of the suggested model over here.

And we are getting R square value of 92 which is higher than the earlier one. And, what we are getting is that R square adjusted is also high over here. And, prediction is also more or less same what we are getting over here, only thing is that when we adopt this one tenfold cross-validation also we can check ok. 5.9 is quite high, so this goes away and last one is considering all models, we want to deduce this one.
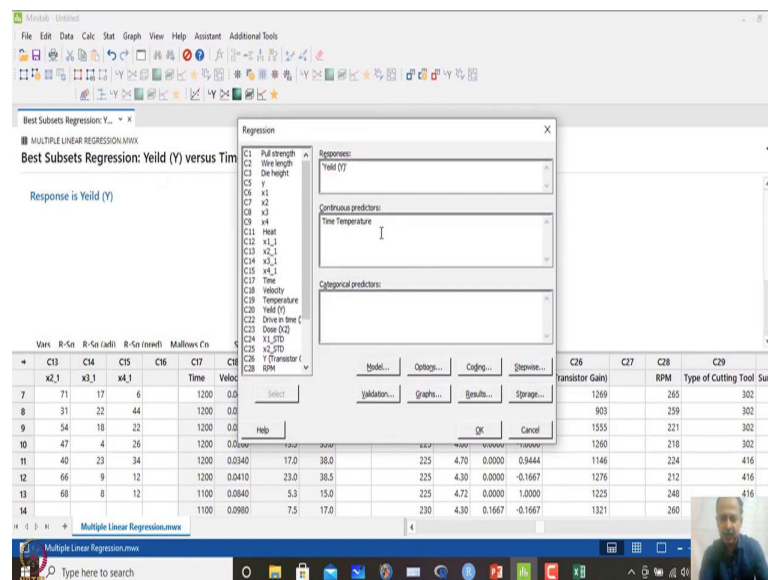
So last one does not, we do not see the last one. So, mallow Cp suggest that this is the one when time and temperature can be considered like that. So, so that is also suggested by stepwise regression what we have seen like that.
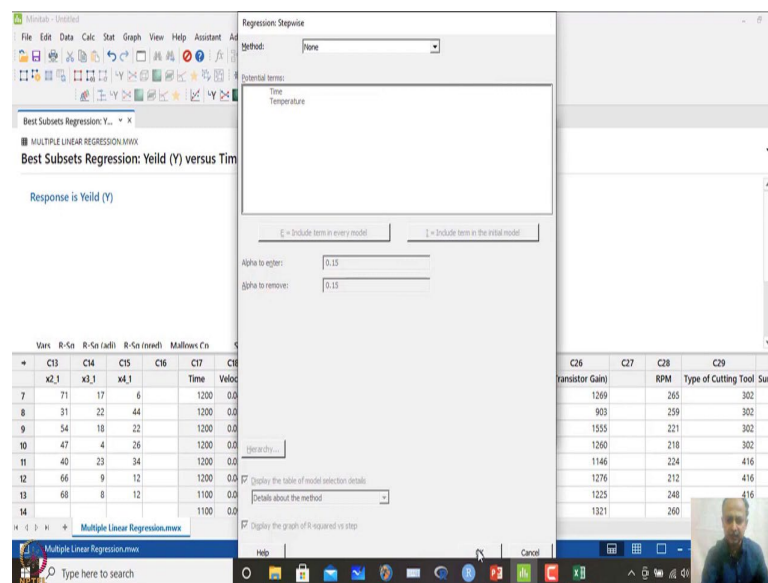
(Refer Slide Time: 09:33)



So in this case, again I am doing this. So, if we select this one fit regression model over here.
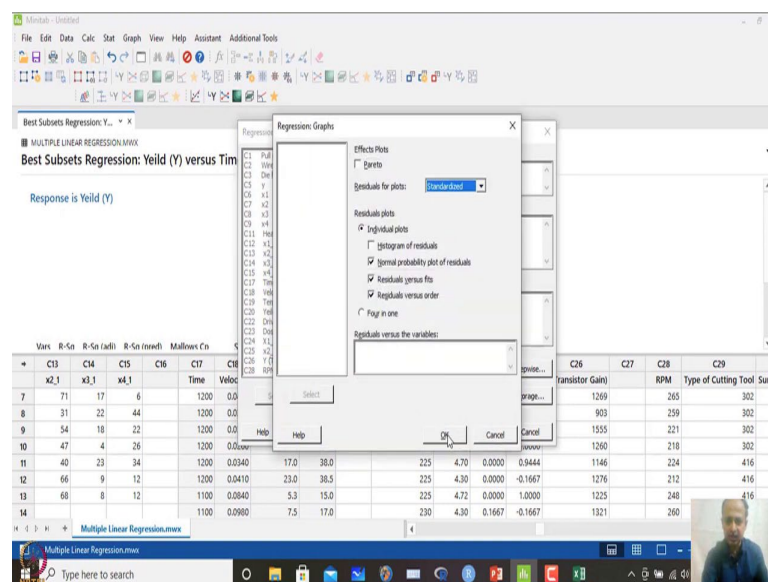
(Refer Slide Time: 09:39)



And in this case, time and temperature. So, I am removing this one and these are the two variables.
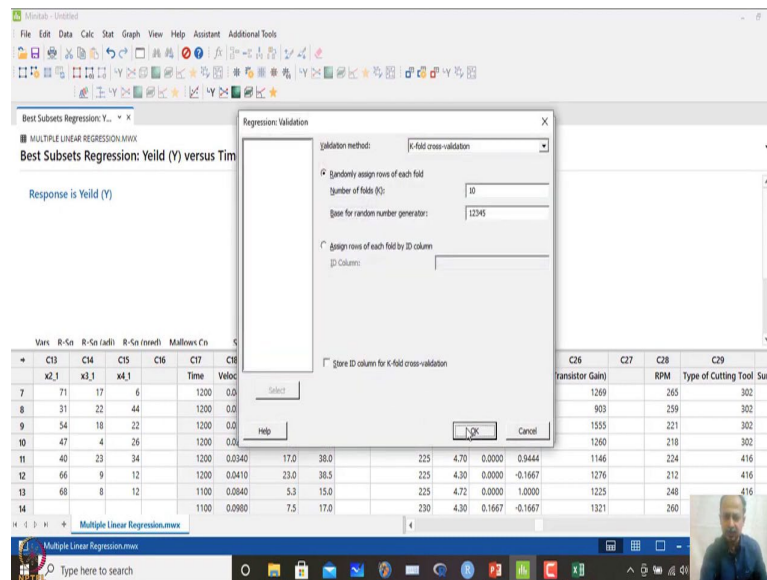
(Refer Slide Time: 09:44)



And stepwise regression, I will not use now. So, in this case I have already identified the variables.
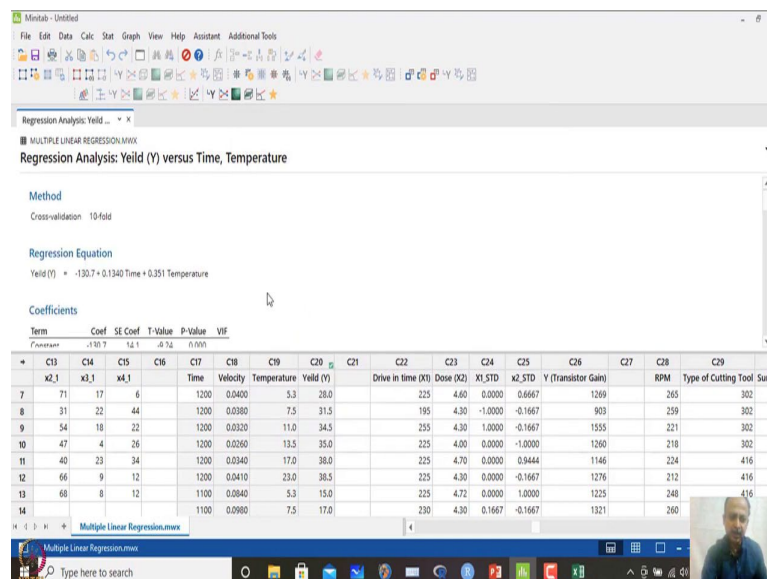
(Refer Slide Time: 09:51)



And, I want to see all graphs, whether it is and this residual can be a standardized residual over here, so if I click ok.

(Refer Slide Time: 10:00)



And, validation we have already careful cross-validation we have given.
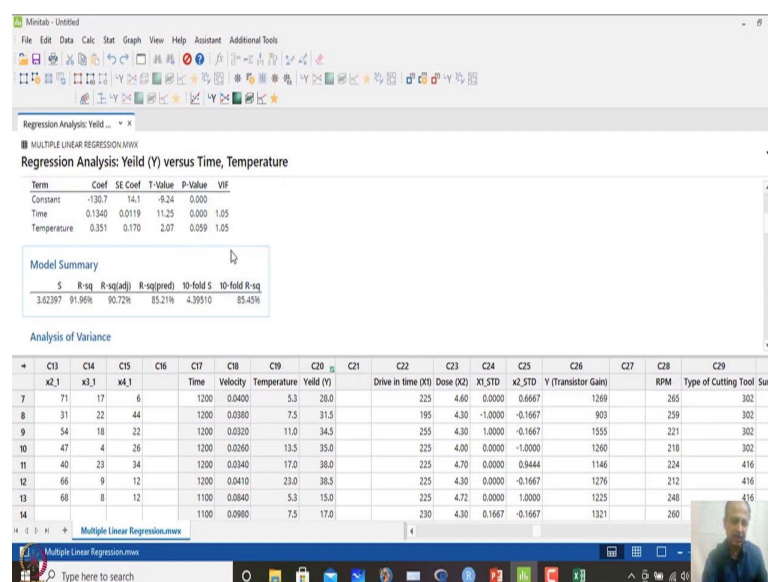
(Refer Slide Time: 10:05)

(Refer Slide Time: 10:05)



And if you click ok, what happens is that, this is the final equation minus 130.7 and plus time it is positively correlated and second one temperature is also positively correlated like that.
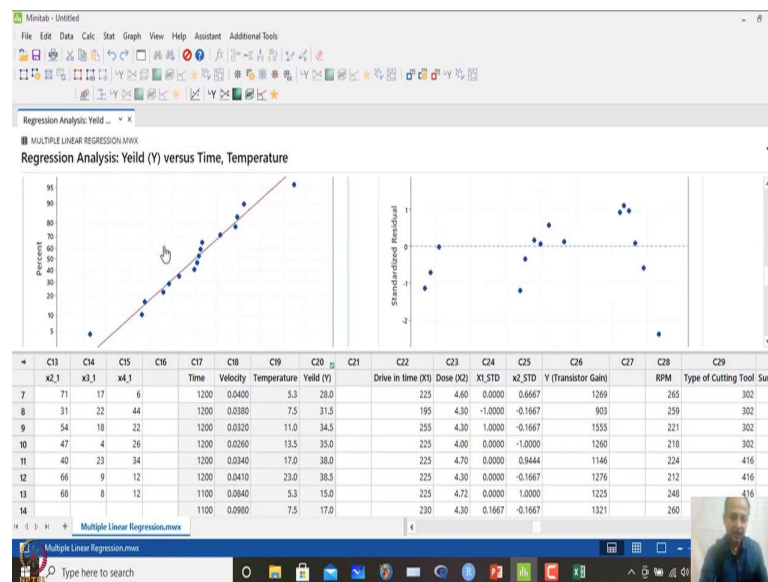
(Refer Slide Time: 10:14)



So now, temperature is retained over here, although the P-Value is not so significant over here. We can also just eliminate and see whether that model performs in any way, because we have a variable to enter and exit we have given a alpha value of 0.15, so that

is why this model, it has come in the model when we have done stepwise regression. But we can retain this one, because this is very close to 0.05.

(Refer Slide Time: 10:40)
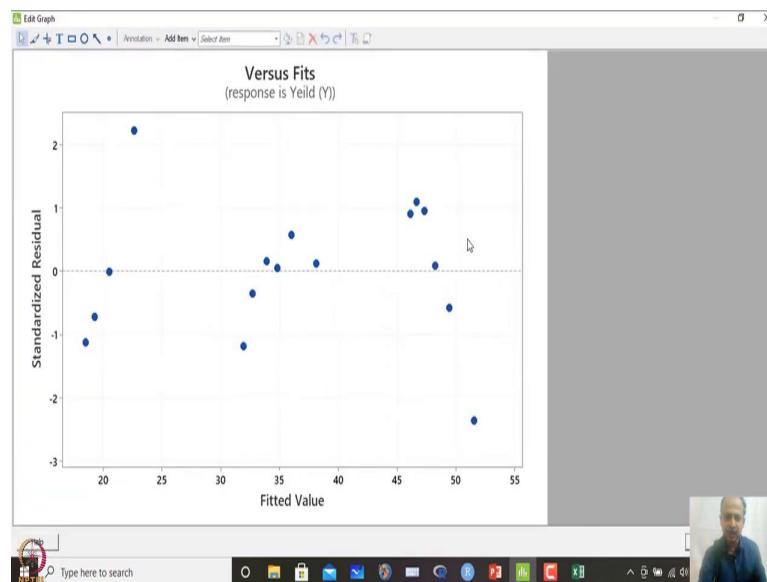


(Refer Slide Time: 10:42)



And, we can see that normal distribution assumptions over here that seems to be satisfactory as all the points.
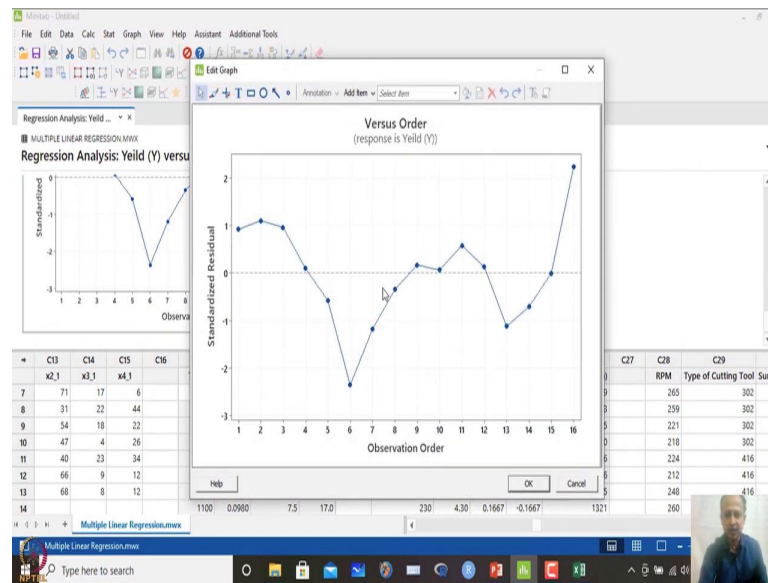
So this graphically when we see seems to be satisfactory points on the line. So, but we can cross check that one.
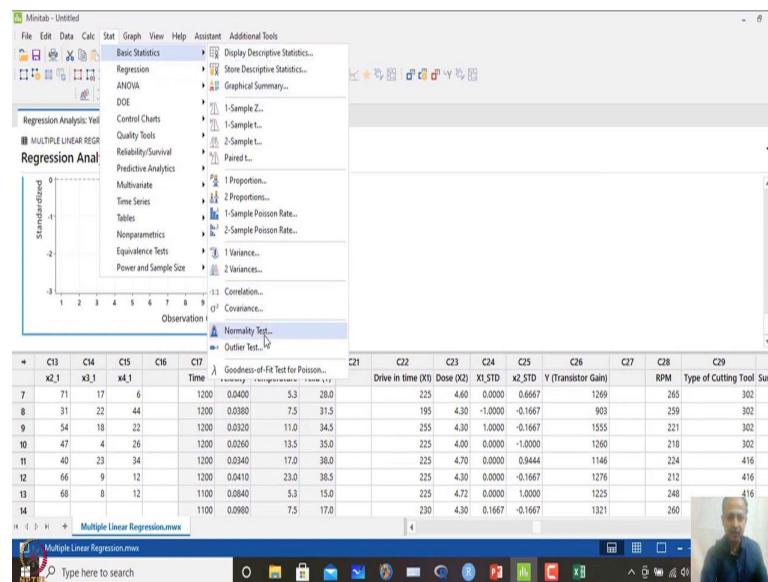
Even the residual versus fit plot also does not show any abnormality or patterns like that, so maybe Breusch Pagan test will also confirm this one.
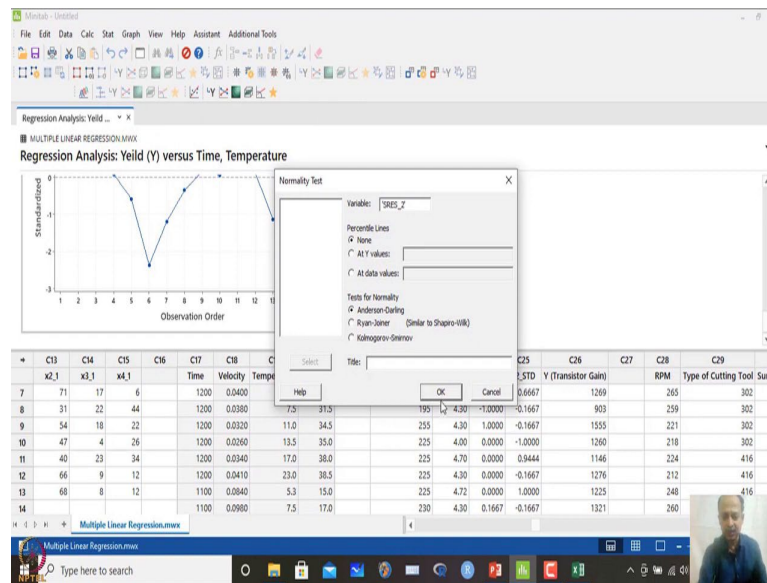
(Refer Slide Time: 11:01)



And also there is no as such abnormalities in auto correlation, what we observe some trend or something is not observed it is on both side of the 000 point like that.
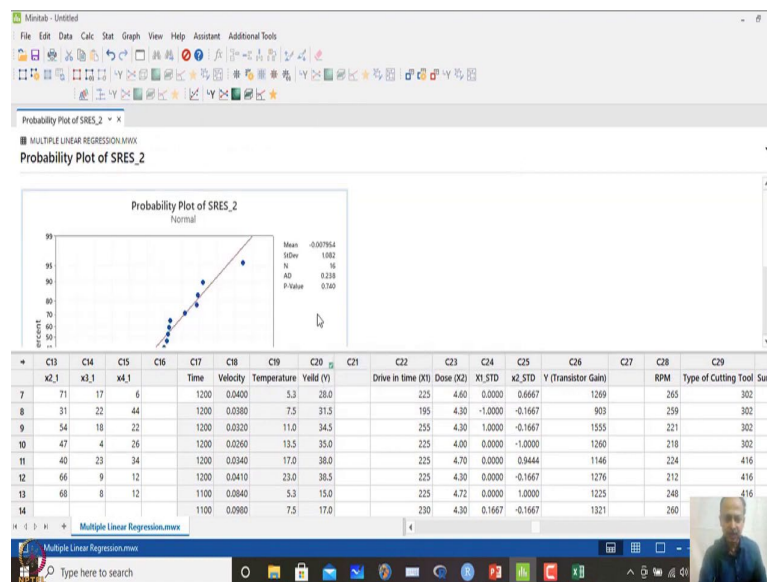
(Refer Slide Time: 11:14)



So in this case, what is expected is that it should confirm the normality and other assumptions over here.
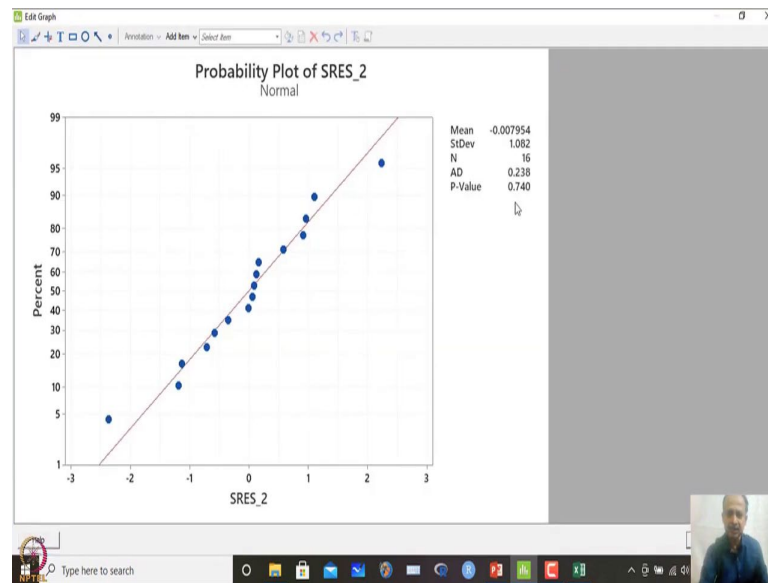
(Refer Slide Time: 11:17)



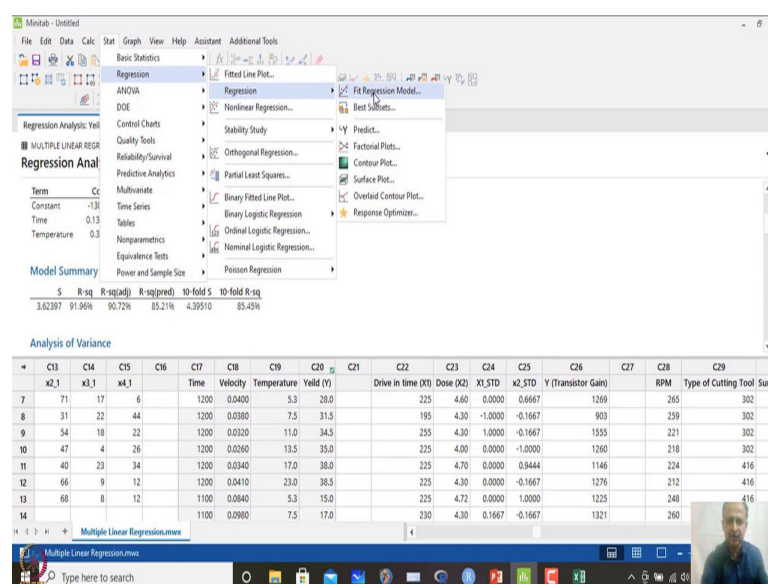So we can just check this one and the last row will be the last column will be the variables residuals that we have to check.

(Refer Slide Time: 11:24)

So, in this case it is coming out to be 0.7 seems to be satisfactory ok. So, in this case you have to practically also think that whether to retain the last variables or to remove, but R square adjusted as has improved what was observed.

So, R square adjusted has improved. If I consider only one variable, maybe R square adjusted will be low. So in case, I consider one variable as time as the only variable, so in this case we can do that.

(Refer Slide Time: 11:52)



So, we can just consider time as the only variable, so fit regression model instead of temperature I will remove this one. So, I just keep a note of this R square adjusted is around 90, this is around 85 R square predicted and tenfold cross-validation earlier was around 85.
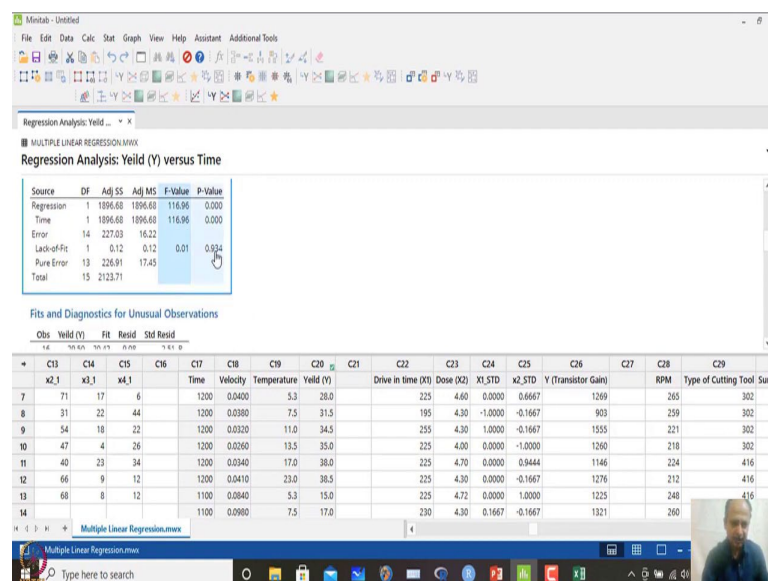
(Refer Slide Time: 12:07)



So, 90 and 85 approximately that is the range what we are getting.
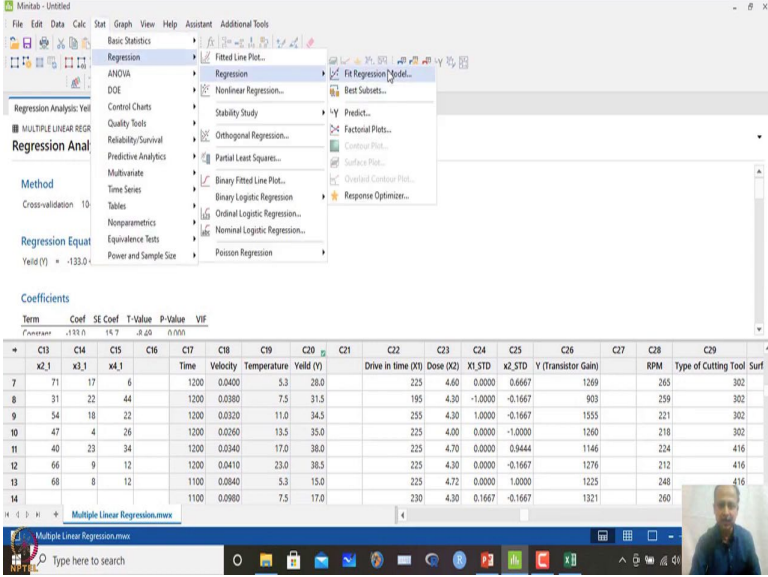
(Refer Slide Time: 12:09)



And, when we do this here you see R square adjusted is also low lower than the earlier one. So, it depends on the practical sense of your, whether to adopt the variable or whether to remove the variable which is not significant like that that depends on the process engineers or somebody who is knowledgeable about that whether to retain that one or to eliminate that one. And, there is no lack of fit as such that is observed over here.

(Refer Slide Time: 12:31)

So this is lack of fit. So, linear model seems to be adequate so we can adopt this one. So, you have to think from practical aspects like that. So, multicollinearity problem when we are considering, both the variables over here, so if I consider both the variables, and in that case we are not getting any multicollinearity issue over here.
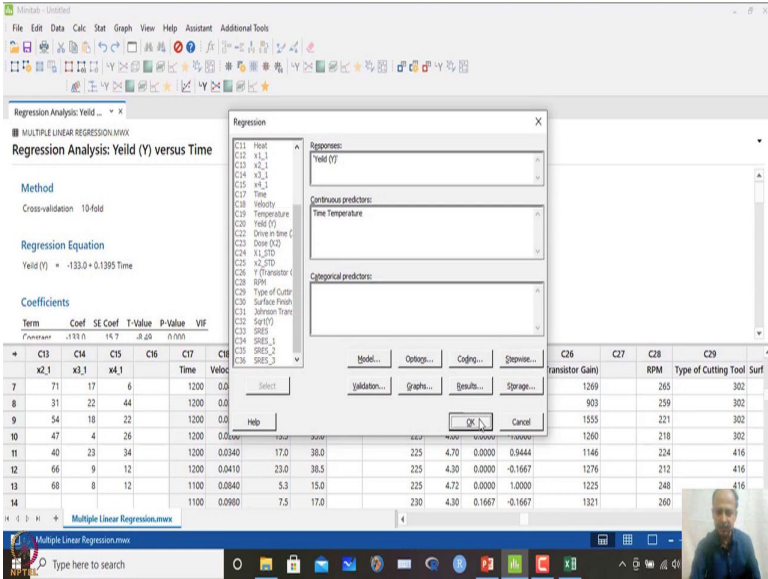
(Refer Slide Time: 12:43)



So, this time and temperature if we consider both of them.

(Refer Slide Time: 12:46)



And variation inflation factor is around 1.05, so which is quite satisfactory.

(Refer Slide Time: 12:56)



(Refer Slide Time: 12:56)



And, so whether to include both the variables or not to include both the variables, that is the judgment that you have to take. But, I will suggest to retain this one because this is improving R square adjusted and also tenfold cross-validation and R square predicted is also improving so we can retain this one ok.

So, but there is no black and white scenario using regression like that. So it depends on the process engineers and then see the predictive behaviors and then we try to adopt whichever model is very close ok. Suggested model over here is, we can include time

and temperature both the variables. But if you go by significance, in that case temperature may be dropped and we stick to time only that is the only variable ok.

(Refer Slide Time: 13:46)



So, this is one scenario. And, we can have in regression some other scenarios like that we can place this is the example, I want to place over here. Maybe, we will delete this one residual over here. So, this is one example that we are taking, so RPN type of cutting tool and surface finish this is the variable surface finish is Y and RPM is the continuous variable and there can be a categorical variable, that means, or type of the cutting tool.

So, here the number is categorical. It can take two values, 302 and 416 like that. So this is the two types of tool. So this is categorical variable basically, so this is no I cannot say 302 is greater than 416 like that. So this cannot be arranged in ascending order, descending order. So this is like color, different types of color. So, this is categorical variable we have to treat, and regression has an option to deal with categorical variable also.

(Refer Slide Time: 14:34)



So, surface finish is Y over here. So what we can do is that stat go to stat regression, fit regression and fit regression model over here.

(Refer Slide Time: 14:38)



And, instead of this variables over here I can just use continuous variable and categorical predictor also we can include in the regression model over here. So in this case, what we will do is that, we will we will incorporate the response as surface finish which is the actual variable and then continuous predictor over here is RPN and then categorical predictor over here is will be type of cutting tool over here.

(Refer Slide Time: 15:07)



So, in this case stepwise regression we can use and we can see which variables will go in and which will go out. So same significance level we have used.

(Refer Slide Time: 15:14)

(Refer Slide Time: 15:18)



And, in models we have included both the variables and included the constant term also, validation also we have taken, cross-validations over here.

(Refer Slide Time: 15:21)



And results, what we can do is that we can see all results over here. And here, there is a options of Durbin Watson statistics which has to be compared with tabulated value, then we can use this one. So, I am not using because I can convert into R I can go to R and model that one and see Breusch Pagan test and also the DW test and corresponding P-Values.

(Refer Slide Time: 15:41)



What I will do is that, standardized residual I will save over here again.

(Refer Slide Time: 15:45)



And then graphically what we want to see, normal plot, residual plot, and order of the data. Plots like that.

(Refer Slide Time: 15:52)



(Refer Slide Time: 15:53)



So, if you click ok what happens is that, it its suggest selection process and then what we have selected like that by default. And, these are the two for 302 types of cutting tool, this is the surface finish is related to RPM, this is the equation, and for 416 it will give a different equation.

(Refer Slide Time: 16:15)



So, categorical variables, different levels that we have selected, each levels we will have a different equation with the continuous variable. So, this is shown over here.

(Refer Slide Time: 16:30)



And, the coefficients are also given over here, so based on which we have developed the regression equation. And, VIC, variation inflation factors is approximated, so that is not an issue. So, these two variables type of cutting tool and also RPN both are significant over here. And, about R square adjusted value is 97.5 and cross-validation is 97.02, so very close. So, this model seems to be very accurate.

(Refer Slide Time: 16:38)



And, what we are getting is that no lack of fit is observed over here.

(Refer Slide Time: 16:42)

(Refer Slide Time: 16:44)



Only thing is that, we have to check whether normality assumptions is violated.

(Refer Slide Time: 16:47).



So, here what you see, is maybe there is a problem with normality over here. One point is over here, some has gone outside like that, so we have to test that one.

(Refer Slide Time: 16:55)



And, heteroscedasticity does not seem to be a problem because it seems to be random.

(Refer Slide Time: 17:01)



And, also we can see that this one may not be auto correlation is not so significant. So, these two checks can be done, but what we want to see normality, so this C 33 column will tell me whether normality assumptions is violated or not.

(Refer Slide Time: 17:15)



(Refer Slide Time: 17:17)



So what I will do is that, basic stat and normality test and what we will do is that, normality residual we can check and try to see what happens.

(Refer Slide Time: 17:22)



(Refer Slide Time: 17:24)



Normality test and what we observe over here for the residuals is that P-Value is less than 0.05. That indicates that there is a problem of normality issue over here of the residuals.

So, in this case, directly the Y characteristics needs some transformation and we have Box-Cox transformation and also Johnson transformation, both the options are there.

(Refer Slide Time: 17:46)



So, in this case what we will do is that, first we will see whether Box-Cox transformation works for this, for the Y variable. So what we will do, control chart Box-Cox transformation.

(Refer Slide Time: 17:51)



And, for this surface finish we have taken this variable over here. Subgroup size is 1 selected.

(Refer Slide Time: 17:57)



And in options optimal or rounded value of alpha just mention that one.

(Refer Slide Time: 18:00)

(Refer Slide Time: 18:01)



And when we when you just click that option, what happens is that it suggest you that rounded value is 0.5; that means, lambda is approximately we can take as 0.5, although the actual value is 0.72, but we can take a rounded value of which is lying within the confidence interval over here.

So, 0.5 we are selecting over here as rounded value. So, I have taken a square root transformation over here. So, 0.5 means Y to the power 0.5 means, square root transformation that is suggested ok.

(Refer Slide Time: 18:30)

(Refer Slide Time: 18:33)



With square root transformation, then what have what we have to do is, that regress variables fit. Instead of surface finish what we will do is that, we will mention that let us go for square root transformation of Y that values and do all the analysis same analysis over here. So, regression analysis over here.

(Refer Slide Time: 18:43)

(Refer Slide Time: 18:46)



And, here also we see that R square adjusted 97.

(Refer Slide Time: 18:50)



And so it has regressed square root of Y with the variable.

(Refer Slide Time: 18:53)



And only one with the variable RPN over here. So this both the variables are significant.

(Refer Slide Time: 18:58)



And, the equations is given, lack of fit is not prominent.

(Refer Slide Time: 19:01)



So in this case, one out layer is recorded over here as it is minus 3 point more than plus or minus 2.

(Refer Slide Time: 19:09)



So in this case, what we have to see is that, again there is a after this transformation there is a residual that is generated over here.

(Refer Slide Time: 19:15)



Let us try to check whether the correction has happened with Box-Cox transformation or this is not adequate. So, in this case what we will do is that, we will see the residual over here and try to see ok.

(Refer Slide Time: 19:32)



So, let me just check this is stat basic stat, normality test over here, so I will go to the last variable, that is recorded.

(Refer Slide Time: 19:36)



(Refer Slide Time: 19:37)



And then I see what is the value of the P-Value. So, P-Value is again less than 0.05. So again there is a problem and that the problem is not resolved, so error is not coming out to be normal, so in this case again it is not white noise. So, in this case what we have to do is that we have to so then what I have done is that I have gone for Johnson's transformation, family of transformation so Johnson's transformation.

(Refer Slide Time: 20:00)



So, what I have done is that basic stat, sorry this is quality tools, and in that case Johnson's transformation is there.

(Refer Slide Time: 20:06)



And in this case, I have reported that place it into single column data are arranged where the data is. So, I will say surface finish is the data and store in which column. So, I have mentioned over here as C31.

So, when you click that one C31 over here, so in this case what happens is that, if I click ok.

(Refer Slide Time: 20:25)



And options what we have given is that 0.1 is the value to select the best fit, so in this case ok.

(Refer Slide Time: 20:30)

(Refer Slide Time: 20:32)



So, it will give you the transformation that is required. So, this is the transformed equation transform function over here that is given L n what you see the last over here and this is transforming the variable. So initially P-Value is less than point the original data is less than 0.05 and after transformation P-Value is coming out to be 0.906; that means it has done a rightful transformations.

Now we have to only confirm. And, this is saved over here Johnson's transformation is saved over here in C31 column. Now, what we will do is that, with this column we will regress with RPN and type of over here.

(Refer Slide Time: 21:02)



(Refer Slide Time: 21:05)



So, what we will do is that we will go to regression, regression analysis, fit regression, instead of this square root what we will do is that we will use Johnson's transform variable and other things remain same. So, I will click ok.

(Refer Slide Time: 21:14)



(Refer Slide Time: 21:17)



And we will get a residual over here and we will get the equation. So, after Johnston's transformed variable with RPN, this is given for 302 and 416 types of tools like that.

(Refer Slide Time: 21:24)



And, what we observe is that 91 percent R square adjusted 88, so very close.

(Refer Slide Time: 21:28)



So, this is quite acceptable and in this case there is no lack of fit and both the variables are significant that is observed over here. Categorical as well as RPM over here.

(Refer Slide Time: 21:39)



So, in this case residual we can check, whether the correction has happens some positive things has happened over here.

(Refer Slide Time: 21:42)



So, I am going to normality test over here, and what I do is that I go to the last residual values and I click ok.

(Refer Slide Time: 21:48)



(Refer Slide Time: 21:50)



And, what I observe is that, when I have done this transformation suddenly this problem has gone so; that means the normality problem of the error residual has gone. So, when I am using Johnson's transformation this has happened. When I have used Johnson's transformation this is giving results like that.

So, which one Box-Cox or Johnson you have to try out and figure out that whichever gives you error as white noise, so that has to be adopted like that. So, this is one of the

example, when categorical variable is also considered in the model and we are able to address that one.

(Refer Slide Time: 22:22)



So, categorical variable and how to address in case there is a non normal situation in multiple regression, how it is to be addressed that we can see ok.

And, how to select the variables in case we are in dilemma which variable. So we have talked about stepwise regression. So, this regression is an important aspects which can suggest that which are the variables or potential variables can be considered in experimentation.

(Refer Slide Time: 22:46)



So, but we have to understand that regression does not say we cannot extrapolate regression equation. So, whatever regression that experiment is happening; that means, this is y and this is the range of x domain within which we have got the information; that means, data are collected over here and these are the observations over here.

So, I can only restrict to this region and predict. So, for a given value of x, at a given value of x what is the predicted value of expected value of y. So, expectation of y for a given x can be only calculated within the domain or within the range of x over here. So, x we can say upper bound and x lower bound within this, so no extrapolation is generally preferred in regression equation ok.

And, association does not mean, that means, y is related with x function of x, does not mean that it is a real variable that impacts y. So, certain scenarios can be, there is a relationship between two variables which are not physically any way connected but there can be high correlation.

So, many examples can be cited like that. So, it does not mean that there is a causality. So regression does not prove causality. So, design of experiment is the only way to define causalities like that. So, to understand causalities between the variables we have to intentionally induce variation by changing, the factors changing the factors or variables that we are interested into and try to understand what is the functional relationship.

So, the best appropriate way to develop the functional relationship is by design of experiments. So, there is no other alternative. So, regression with historical data does not give you proper association or causality cannot be proved based on that ok.

So, then what is the option, then what is the option for quality improvements. Then the options we have identified, few variables based on regressions. So there are variables X, X matrix over here, and I have collected also y variable information. So, how do I connect these two and this connection and develop the function. So, I need a function over here, that will explain the variability of y over here. So, what is required is that systematic way of variation is required, systematic way.

(Refer Slide Time: 24:40)

**Quality Control and Improvement using MINITAB**

**Improvement Phase:**

**Design of Experiments (DOE)**

Prof. Indrajit Mukherjee, SJMSOM, IIT Bombay

And then what comes is statistical experimentation which is the improvement phase. Until and unless I have a right function I cannot improve. I cannot reach to the global optimal point like that or setting points like that ok. So, design of experiment is the stepping stone for improvement phase like that.

So, basic things that our understanding over here is that, there are types of variability types of variation we have already spoken about it and some basic information I am providing over here which we can think of as recap. So, there are common cause variability in the process, there are special cause abnormal variability which are easy to detect based on SPCprocess control chart.

And, then what we do is that, to understand the functional relationship between y and x, we have to do it a systematic change in we have to induce variability we have to induce variability by changing the factor, by changing the factor x over here. And, if we can do that systematically what happens is that we can generate a function and the function can be optimized.

That function, we know this is the real equation between the in the process. This is the empirical relationship that exists and this is based on systematic theoretically strong models that we have developed over here. Then we can use optimization technique just to reach to the global optimal point like that. So, we are interested into and for that design of experiment is most preferred like that ok.

So, design of experiment what we are seeing over here is basically reducing the common cause variability. So, this can be applied when the process is in statistical control. So, when the process is in stable process, so whenever it is stable process we say and it is in statistical control then only we can go for design of experiment, because we want to reduce the common cause variability further, because more and more I reduce the variability, so accuracy and precision both improves over here.

So, we are interested in being and variance both aspects like that, what we have seen in capability analysis also. So, for that what is required reduced for the variability requires design of experiments. So, for that we need to, so that helps, statistical design of experiment helps us to bring the mean to target value and also the reduced variability near to 0.

So, we want minimum variability over here. So, various design of experiment techniques are proposed and huge amount of resource is available for design of experiments. So you can see. But, preliminary book that you can see is Montgomery's book which is which can be a good resource of learning initial steps like that, but there are other books also you can see, so there are many other books which you can see like that. Amitava Mitra is one of the books where design of experiment is detailed ok.

And, so this is one of the techniques. And also, I want to recap that one of the concepts that if you have variability, that means, causes and this is impacting my y and causes are interlinked also somewhat they are interlinked with each other and they together can impact the y like this. So this is the scenario like this.

This is one of the variable X1 this is one of the variable X2, X3 and this is X4 over here, and they are interrelated with their in this what you are seeing and they are not completely independent over here what is observed over here.

So, if they are completely independent in that case, we can understand that what is the level of X that will deliver the optimal y over here. So, this is possible, but this is not possible when we are having a complex relationship between multiple X's over here and together they can impact the y, that can be scenario.

So, it can be scenarios which is not possible to identify. So, by simple other analysis like, design failure mode and effect analysis or process failure mode and effect analysis that is not possible by that ok. So, complex relationship and developing the function, when we are developing the functions over here only design of experiments can help.

So, design of experiments can help. And, before we go into details of design of experiment this process p information is again I am recapping over here. There will be certain control variables which is in the hand of experimenters. So, experimenter will change this control variable basically.

So when you go to a process what you observe is that, you will see the operator is changing some of the variables, not all variables, some of the key variables which is possible to control. So, these are the variables X1 to Xp which is known as controllable variables.

Design of experiment is all about controllable variable, most of the time we deal with controllable variables. There will be some inputs and I also talked about covariates which can influence the outcomes over here, and that can also be considered in the model when we are doing experimentation we can also consider the covariates as a variable which will vary.

We cannot control that one, but influence of that has to be considered when we are developing the mathematical models or something like that. Or we can deal with that, so covariates also we can see ok. And, there will be some uncontrollable factor; that means, it is uneconomical to control.

And sometimes we do not have any information enough information about the some of the variables, because you cannot build a perfect model like that y is a function of all Xp variables X1 to Xps like that. There will be some amount of error in the model because of this noise variables over here, because of this noise variable. So, we do not get a perfect function.

Because all the variables that is impacting the process is some of the variables are unknown to us. So, in that case, these unknown as noise variables or and the influence of that is very less because we have identified most of this. So, these are the variables and some of them are controllable, some of them are covariates over here, so we have knowledge about that.

So, and there will be some noise variable which we do not have any control or we do not have enough information about these variables like that ok. So, we want to; we want to get a setting so in design of experiment, what we are trying to do is that this is control variable, so it is in our hand to control. So, we want to get a combination of this combination of this that will optimize my Y CTQ.

And y will be the average value of y should be close to the target value, and the variability of y that is that we will generate is very close to 0 values like that. That means there is no as such variability we want to develop. So, our objective of design of experiment is twofold objective; one is which are the variables impacting y, and then if these are the variables which is significantly impacting y, what can be the combination or setting conditions of this Xp variables that is impacting over here, which I can control.

So that I get the y exactly to the target value defined by the designer and also the variability near the target value is near to 0 like that. So, that is the overall objective of design of experiments. So, and another objective is to develop the functional relationship so that we can optimize like that.

So, whenever I have a function, I can optimize. So, I need a function. So, function to need a function I need what are the Xp variables which is controllable which is significantly contributing to the variability in y; that means, this is and there will be some error that can be because of this noise variables which are not considered and which is difficult to control and uneconomical to control.

This can also be impact with these errors may also be impacted by these covariates which are influencing the , but we do not have any control on that. So, I can only control this Xp variables, I do not have any control on the inputs over here and the uncontrollable variables over here. So, most of the time we try to get the best combination of X1 to Xp that in presence of this noise and in presence of covariates or inputs variability that we are experiencing in the process like that ok.

So this is all about what we will try to understand in design of experiments and how to do the experiments that we will try to understand over here. There are different ways of doing experimentation, different designs are available, and we will see some very few of them. So that that is our objective in experimentation, so that we will cover in our next session ok.

Thank you for listening.