Exploring Survey Data on Health Care Prof. Pratap C. Mohanty Department of Humanities and Social Sciences Indian Institute of Technology, Roorkee

Lecture - 26 Regression Models of Quantitative Healthcare Variables

Welcome friends to this NPTEL MOOC module on Exploring Health Care Survey Data. We are on the particular week explaining n influential statistics in health care. The first lecture is on understanding regression models specifically on the quantitative and qualitative variables that are used in models. In this particular lecture, we will be emphasizing only the quantitative variables, in our next lecture will be on qualitative variables.

Again, you might be having some sort of problem with whether it is a quantitative dependent variable or quantitative independent variable or qualitative dependent variable, or qualitative independent variable. We will give all sorts of clarifications to you so that your understanding in this regard will be much crystallized.

Without discussing the background of it let us move on and clarify what is all about. And starting from the very first aspect of regression as we all know that this is quite common, in different statistical techniques as researcher used to answer say several questions.

The regression results and their analysis have various uses, especially since it gives certain predictions about the value of the dependent variable based on at least one independent variable.

(Refer Slide Time: 02:07)



This explains the impact of changes in an independent variable on the dependent variable. This helps in determining the strengths of the relationship between these two variables X and Y, also it helps in forecasting the value of Y based on the control variables or X variables.

(Refer Slide Time: 02:28)



This gives cause and effect relations though the causal influences cannot be drawn, certain logical understanding through the theoretical model can be mapped through the regression coefficients. And regression analysis assumes a linear relationship between the predictor and outcome variable.

(Refer Slide Time: 02:46)

	Correlation vs Reg	ression
BASIS FOR COMPARISON	CORRELATION	REGRESSION
Definition	Determines co-relationship or association of two variables.	Describes how an independent variable is numerically related to the dependent variable(cause-effect relationship)
Usages	To represent linear relationship between two variables.	To fit a best line and estimate one variable on the basis of another variable
Dependent and Independent variables	No difference	Both variable different
Indicates	The extent to which two variables move together	Regression indicates the impact of a unit change in the known variable(x) on the estimated variable(y).
Objective	To find a numerical value expressing the relationship between variables.	To estimate values of random variable on the basis of the values of the fixed variable.
🏐 swayam 🔮		4

We need to first establish the relationship between correlation and coefficient by their definition usage which kind of variable they consider what are the indicators of the objective etc. Starting with the definition of correlation we know that there exists a co-relationship or association between two variables. Whereas, in the case of regression this describes how an independent variable is numerically related to the dependent variable.

With a certain logical relationship between these two. Now, you might be confused about whether it gives cause and effect, in reality, it does not give perfect cause and relations or effect. It gives some understanding related to why it has happened and how it has happened, how it is relevant? But the exact causal conclusions are not derived through regression.

But it is for sure how much it is affecting, what is the quantification of it can be derived from the regression. Correlation and its uses are like it gives you linear relationship whereas, a perfect fit line or trend line (best-fit trend line) could be derived, based on the variables and their relationship with another variable.

Then in correlation, the perfect trend line is not derived. Regarding dependent and independent variables, correlation does not differentiate these two. Whereas, in the case of regression we have to differentiate then only we can find out the impact of one on another one. Containing the indicators of correlation and regression the extent to which two variables move together is called correlation.

Whereas in the case of regression, this indicates the impact of a unit change in the known variable on the estimated variable. The objective of correlation is to find a numerical value expressing the relationship between the variables whereas, in the case of the regression analysis this gives estimates and their values of random variable on the basis of values of the fixed variable.

(Refer Slide Time: 05:14)

Linear Regression Analysis
Most commonly used method to examine the relationship between a quantitative outcome and one or more explanatory variables.
The dependent variable must be a continuous variable.
The independent variables can take any form: continuous or discrete or indicator variable (dummy variable).
The linear regression model describes how the dependent variable is related to the independent variable(s).
Relationship between the mean of the response variable and the level of the explanatory variable assumed to be approximately linear (straight line).
S Swayani 🧕 S

We are now attaching information about linear regressions and their analysis of how these are read. This is the most commonly used method to understand the outcome variable and the explanatory variable. Then the dependent variable must be a continuous one. If it is not a continuous one, then we will discuss some forms of non-linearity or some transformation to a linear model.

That we will discuss in our successive models. At this moment we are saying that the dependent variable must be continuous. The independent variable can take any form and may be continuous-discrete indicator type like in dummy variables etc. The linear regression model describes how the dependent variable is related to the independent variable.

The relationship between the mean of the response variable and the level of the explanatory variable is assumed to be approximately linear. So, that is why it is called linear regression analysis.

(Refer Slide Time: 06:18)



There are different types of linear regression models; one is called simple linear regression another is called multiple linear regression model. Simple where we also call it as bivariate regression model other one is called multivariate regression model. Variate means the extent of variability due to the number of factors including a number of variable factors including, independent factors included.

At this moment there are only two factors in fact of X and Y. But here there are so many factors. So, that is why even if we are saying X is in fact a vector of information of control variables. So, $\beta_1 X$ the equation which we have mentioned here is a kind of vector of dependent and independent variables. And this equation is best explained in Greene's book of Econometrics.

Where the matrix method how is explained, how different vectors are explained in the shaping the equation like this are well explained you can have a look and find out the exact reasoning behind writing on $\beta_0 + \beta_1 X + \epsilon$.

Epsilon is the error term. We know that the X cannot explain everything about Y. There must be some other component called some fixed component as explained by β_0 . And there will be still some unexplained aspects captured through the error term.

In the case of multiple linear regression, we have to discuss so many other independent variables along with the first one and we can be able to estimate each of their coefficients differently and that will be helpful in interpreting the result correctly.



(Refer Slide Time: 08:34)

Let us understand the structure of the model. As we know that Y represents the dependent variable, and β_0 is our intercept or the fixed coefficient, which is independent of the change in the covariates. And likewise, we can give an example of consumption as a function of income.

And other factors if I write down C as a function of income when income tends to 0, C may not be tending towards 0, there would be some positive consumption even if income boils downs to 0 for a particular consumer. The particular consumer or household may not get income in that particular month where they are supposed to consume something by borrowing or by other means of getting the food.

Therefore, consumption there must be a fixed component. Similarly, almost all the variables in the world that we are trying to project must follow this kind of structure. Other details we have already highlighted, and I am sure this will help you. Epsilon is called the random error component, which is also called the unexplained or random or error component.

This term is called error term and or also called disturbance term. Sometimes denoted by U and that is in fact distributed normally with 0 and 1 standard deviation. This shows the

normality of the distribution, that all sorts of projections in the model can be made. We will explain some of these things in this particular lecture.



(Refer Slide Time: 10:31)

Population linear regression model how is explained in this diagram you can see a trend line it is fitted with that is simply called the expected Y value is equal to β_0 plus β_1 Xi. And β_0 is the intercept and any dots you could see in this particular diagram, dots like here some of the dots are deviating from the trend line.

And the trend line is explained in the expected value given they are control variables. Now our purpose here is to identify the extent of error, and the differences of each value from the standard line or from the trend line, or from the fitted line. So, it is explained here, and these are also all dots that explain the exact information, or they are also called observed value, we have also highlighted that.

How the observed value is actually deviating from the predicted value of Y for Xi? if our observed value is almost on the trend line or on the predicted line; that means, our estimation is actually expected to be very good in that case. Your R square is expected to be quite good, so it will be close to 1.

But in most the cases like in survey data R square value as we already explained earlier in our previous lectures that it is usually lesser. You need not worry much by we have given all reasons behind R square and its range which one could read how it should be interpreted.

Though this is what we already explained in our previous weeks' lecture. You have to follow and find out the differences.

(Refer Slide Time: 12:39)



Now here we are presenting a difference between the population and sample regression model. A population where we have written the equation Yi is equal to $\beta_0+\beta_1Xi+\epsilon$ from here a random sample is taken. That random table is going to give a sample series of data.

That series we can estimate, and here is our estimated β values, and that is going to project the population β values and if it is correctly projecting then our sample is in fact giving the right result in terms of predicting the population information. This is what we usually follow in our sampling design and the rest of the details are very much clarified from the diagram itself. (Refer Slide Time: 13:38)



Now, we are explaining sample linear regression we have already explained the population linear regression in this chart and this table here is sample linear regression. Therefore, the population now you can see how on sampled observations, we have we are just highlighting sample. The sample linear regression line provides an estimate of the population regression line and these are all going with the estimated value.

These are in fact called estimated values, one that is written; that means, they are the estimated values based on the sample data we have collected. So, those hat is in fact useful in predicting the population values. So, like we take so many varieties of samples those have different distributions.

In each case, we will establish a number of tests, whether the sample is actually catering to the or in fact representing the population. So, whichever is estimating quite well, then our regression model in fact predicts the population very correctly.

Sometimes we take the number of a different sets of samples if each of the estimated values are not varying much; that means, our estimation is in fact robust. Those techniques we do. I am not emphasizing much here. Now, coming to the method of estimation there are different methods of estimation sometimes use maximum likelihood estimation as well, but at this hour we are explaining about ordinary least square method OLS.

(Refer Slide Time: 15:39)



One of the methods of estimation through OLS; OLS is everywhere deviated and is used frequently by researchers. This is in fact as I already said the most commonly used method. And estimated values represented with they are hat beta β_0 or β_1 hat. They are determined based on the error least square error.

So, now you can see the estimated value is presented here and this is in fact each of the observed values, the observed value minus the estimated value the difference is in fact called an error. So, that is Ei ε_i . ε_i and its square are taken wherever we can minimize the errors. So, that eventually indicates our estimation is going to be fitting with the trend line. Therefore, minimizing that error is important.

Now, if you will not have taken it might be the case that the minimum is fine, but if you do not take this square, it is very difficult to minimize the error. Because sometimes some values of this might be positive some, if it is not having square term then the some might be negative, the average might boil down to be zero. So, then in that case estimating or minimizing the value nearest to the estimated value is in fact very difficult.

Squaring is going to give us non-negative values. Then we have options to minimize it as close to the 0-value as possible. That is why it is called squaring the error term. That is why it is called the least square. We are trying to understand which one is in fact giving us the least error. So, that is why this is called the sum of ε_i^2 and then these are estimated like this.

(Refer Slide Time: 17:58)



The graphical illustration is given in this particular picture. There are different $\varepsilon_i \varepsilon_{i1}$ hat ε_{i2} hat ε_{i3} if all these are in fact estimated and derived with the lowest levels, minimized level with their square term is called least square error.

(Refer Slide Time: 18:19)

Coefficient Equations	
Predicted equation: $\widehat{Y}_i = \widehat{\beta}_0 + \widehat{\beta}_1 X_i$ $E(\underline{w}_i - \overline{w}) = \widehat{D}_i (\underline{w}_i - \overline{v}) = 0$ Sample Slope: $\widehat{\beta}_1 = \frac{SS_{xy}}{SS_{xx}} = \frac{\Sigma (X_i - \overline{X})(Y_i - \overline{Y})}{\Sigma (X_i - \overline{X})^2}$ Sample Y-intercept: $\widehat{\beta}_0 = \overline{Y} - \widehat{\beta}_1 \overline{X}$	
S swapan 🧕	13

The next one is coefficient equations. What do we mean by this coefficient? First is the estimated value of ε_i basically is estimated. So, the estimated value of the error term is missing here because the expected value of ui minus u bar is equal to 0 and then you can say

that ε_i minus ε bar is equal to 0 because of the assumption. The assumption is that our distribution is normal; since we have already said this is distributed N with $(0, \sigma^2)$.

Therefore, this is what we all already said. When this assumption is fulfilled; that means, the error term is estimated to be 0. Whatever is left is in fact the estimated value in terms of β_0 and β_1 hat. now, based on this we can estimate finally, the value of β_1 hat and β_0 hat. you can just see further on β , how β values fulfilling BLUE properties or not, whether best linear unbiased estimation matters or not.

So, this is estimated as β_1 hat is derived as the sum of the square relationship between X and Y that is the covariance divided by the variance of X. So, covariance basically is the relationship between Y and X divided by it is the sum of squares of it is own variance. Then once the β hat is determined then we can estimate the β_0 hat as well.

(Refer Slide Time: 20:08)



The assumptions of the linear regression model are very important, and you must take note of them. Each of it is very important to validate your model of OLS whether it is in fact applicable or can be applied in many contexts or not. So, the first assumption is that the linear regression makes the following assumptions. The first one is called linearity the relationship between the dependent and independent variable must be linear.

Linear by its coefficient not by its variables and linear by its parameter not by its variables. So, variables could be like this, Yi if I said, so it should be if it is β_0 then β_1 it could be of Xi square and ε if I take this as a modified equation because linearity stills hold with the parameters.

We are estimating, this is explained here already and if it is Yi is equal to β_0 plus β_1^3 then Xi and ε this is in fact not a linear model. This is rejected at this moment, we are not estimating a non-linear model. The next one is called non-autocorrelation. Autocorrelation the word itself says it is having a covariance between pairs of random error. So, error inferior period over the time if there has been some continuity of the data. The errors are actually having a certain path of progress.

In one period there are set of errors, in the second period there are again sets of distribution of error. There are those errors, if errors i and j, given the control variables are X if it is equal to 0. That means there is no covariance. So that means, the error terms are considered to be and have no autocorrelation, if you are having autocorrelation this will be obviously non-zero.

So, the error term has a population mean of 0, the error term should be normally distributed as we already said, especially its population mean i.e., for the entire population. The expected value of the error term given its control variable should be equal 0, all regressors are correlated with the term with the error terms. So, the error term should be uncorrelated this is what we wanted to emphasize here.

(Refer Slide Time: 23:14)



The next assumption is on homoscedasticity, which simply says that your variance of the error term which we have been emphasizing those variances must be constant that must not be different every time. So, it should be constant this is what is explained with sigma square.

Once it is constant; that means, we are assuming we have assured ourselves that our distribution is not going to be changed much. So, our prediction is going to be better. So, as against homoscedasticity, in that case, your variance is going to be explained as sigma i square.

Therefore, i stands for some variability and it is different in different scatter plots. if you take several samples out of your population your variance is going to be different every time. If it is not varying much near about the same value; that means, that the data is in fact following OLS assumption or is more or less normal.

What about multicollinearity? Multicollinearity is basically the word that says collinearity. Therefore, how linear collinear; that means, more than one variable the explanatory variable. If there are multiple linearity relationships; that means, we simply say multicollinearity. In other words it says that there should not be any perfect linear relationship among X variables.

They are the among the explanatory variable if there exist some linear relationship or perfect linear relationship, then we are actually violating. OLS is not going to be an efficient because from one β value you wanted to predict something about the Y, but in fact that β value is not independent, it is explained by another independent variable.

So, your model is not correctly estimated. Therefore, we have to assume the model with no multicollinearity. Next is about the normality of the error term where the error term is in fact distributed like this that I have just explained. So, it should have 0 mean and constant variance and the next assumption is on a number of observations must be greater than the number of parameters to be estimated.

Like beta 1, beta 2, beta 3, beta naught etcetera number of parameters must be less than that of the observations we have considered. Otherwise, the wireless estimation cannot able to predict, predict correctly. Next to guidance is on illustration of linear regression in health care data.

(Refer Slide Time: 26:25)



For this purpose, we are using NSS 75th round data on health care for the year 2017-18. We are interested in examining the factors affecting the total health expenditure in the case of element or element persons. Some households and individual-level variables are taken as an explanatory variables for the analysis.

(Refer Slide Time: 26:45)



One note is here for you that the purpose of this lecture is to show how to do regression analysis using health care variables, it does not give the exact model used for research purposes and the required testing after regression. Therefore, we are just giving some sample set of explanation and practical results that you may not directly use for writing papers.

Therefore, you need to understand the concept and develop your own model. Now here the sample variables we have filtered from the 75th round of the national sample survey. The dependent variable we have considered for using the ordinary least square regression model is medical expenditure.

(Refer Slide Time: 27:34)



And medical expenditure we have transformed to a log medical expenditure the log transformation is taken, to make it more rationalizing. Because usually, expenditure as a variable is more skewed. Therefore, expenditure is overstated whereas, income is usually understated.

Therefore, when expenditure is overstated a transformation log transformation is going to nullify and simplify the data to make a better distribution. Similarly some of the relevant independent variables we have noted down like sectored dummy rural with any so urban categories then social groups, age square we have taken age into age, household size, household usual consumer expenditure and nature of the element, medical insurance in their premium, nature of treatment etcetera.

Even you can refer to some of these variables in one of the models we published this is our paper published this year just recently in the journal on the repository of a science direct. So,

clinical epidemiology and global health journal. You can just follow and I have highlighted this on the screen just click on that link and I am sure you will get the exact paper for your reference.

(Refer Slide Time: 29:01)



Further this will also explain about the non-linear qualitative dependent variable models in the next class and you will understand things better. Now, we are explaining the steps for running the regression model. Steps are here like we will examine first of all they are statistics the basics descriptive statistics then we will understand them graphically and also understand their correlation.

Then we will draw and interpret some regression results, then regression assumptions to be tested whether the regression was actually drawn correctly or not. These four important steps all must be mandatorily followed to write an article research paper. Examination of descriptive statistics first is on descriptive statistics. I will simultaneously open the database in front of you and then operate it.

(Refer Slide Time: 30:14)



The rest of the result I will leave at the end I will simply leave to you and I am sure you going to enjoy it, but let me now open a database in front of you.

(Refer Slide Time: 30:22)



I am just going to open the sample data.

(Refer Slide Time: 30:24)

The Edit Data Graphics Statistics See Edit Data Graphics Statistics Services T € X	i Uler Window Help 3 🔝 O - O							- Viribles	7 Tax
A Commentation 0								A Lowenship has	
A Kannad	ニテンテン							A construction of the set	
• Commana //c		Copyright	1905-2017 StateCorp 11	c				name Laber	
There are no dems to show.	statistica/sata Analysis	4905 Lakew	wy Drive					There are no items to s	how.
	Special Edition	College St	ation, Texas 77045 USA						
		E00-STATA-	PC http://www.a	tata.com					
		979-696-46	00 stataĝetata.	008					
		979-196-41	(
			🚺 Open				×		
	Read and reading the second second		6	W - Parities -			u B. O. Sawa Durature		
	terminate or 10 formula		4 I	C / Bessilp /			· · ·		
			Organise • New folder				II + □ 0		
	Notes:		A	Nerra	Datamonthad	hea	Site		
	1. Unicode is supported; see b	help unloade	> CrieDrive			11-			
	Maximum number of variables	s is set to	🗇 🌰 QeeDrive	Convocation 2021	04 09 2021 12:33	File folder			
				Dr.S.K. Gupta	02-01-2021 16-67	Lie folder			
	l'		V INSPC	Foundation Day FF to	10-11-2021 10:11	FOCTODO La latita			
) J SD Objects	Inches 202	01.11.2021 11.52	Els folder			
			V Desktop	Waren lablet 6.1.19-1	09-10-2021 15-39	filefolder			
			Convection 2	WRD & M	16-09-2121 18-02	File folder			
			Dr. S.K. Gupta					Properties	ų ×
			Foundation Da						
			🔿 📒 IT Mandi					* variables	
			> 🗧 January 2022					Label	
			WacomTablet					5.04	
			> WRD & M					Tornat	
			> B Documents					Value label	
) & Downloads					Notes 6 Defe	
			Date				- All Onda Eller	Limane	
								Label	
							Open Cancel	Notes	
							4	Variables 0	
								Memory 64M	
	Command							4	
C1Uses/Studio A/Desktog/January 2022/Dr.	Pratap Mohanty/Stata15							CUP	NUM OIR
11 O. Tune here to courch	0 15 🧥	- A	Date: 1	角 NETT Lachard S 🖓 I	hannadhaine Gida C. 📰 Burtar	F 151	A 200 Ounter A 6	FD \$ 4 .4 Dec 1635	
- ogye nere to search	0 H 💡	6		Anna and a fully	Hanna Hanna		🚽 zric Goldy in 🦉	2-11-2-11-2	21 🔨

(Refer Slide Time: 30:28)

τax									 Variables 	
commands have	(0)								1 (ter variables here	
mand rc		Conversions	1995-2017 StateCorp.	112					Name Lab	el
e are no items to show.	Statistics/Data Analysis	StateCorp							There are no iter	ens to sh
		4905 Lakew	wy Drive							
	Special Edition	College 3t	Ation, Texas 17045 1 PC http://www	GA						
		979-696-46	00 state@stat	8.008						
		979-696-46	01 (fex)							
			🔯 Open					×		
	Annual and a distance		← → + + → +	his PC > Dealtop > January 2022 >			~ 0	D Search Sensary 2022		
			Organise • New fol	der				⊫• □ 0		
	Notes:		OreDrive *	Name	Date modified	lype	Size			
	 Unicode is supported; see Maximum number of variable 	help unicode	-	Dx. B.R. Gurier	27-11-2021 19-52	File Folder				
			 OreDrive 	Dr. Harshit Lakra	17-11-2021 10-39	l de tolder				
	·		- 🛄 This PC	Dx. M V Sund	26-10-2021 15:03	Fiefolder				
			🗊 3D Objecta	Dr. Praty Mohanty	29-11-2021 15:59	File folder				
			Cesktop	Dr. Rannesh	25-10-2021 14:38	Fiefulde				
			Convocation 2	Dr. Rachmi Gaur	26-11-2021 15:27	Lie folder				
			Dr. S. K. Gupta	Dr. Santosh Kangnekai	25-09-2121 18/01	Frie fulder			Properties	
			Foundation Da	Dr. Valuendra	27.11.2021 15.15	Eletable			A 19 10	
			ITMandi						 Variables 	
			January 2022						Label	
			WacomTablet						3,94	
			WRD & M						Format	
			Decuments						Value label	
			L Develoads						+ Data	
			Re	NITE				All Stata Files 🗸	Liename	
								from from	Label	
								open Carce		
									Sat 0	
									Memory 64	4M
									Sorted by	
	Command									
		_								

(Refer Slide Time: 30:30)

	0.0							10-110-0	
A local data and a	(8)							 Variables 	
Liber commands have								Filter variables here	
# Command _rc	1 1 15.1	Copyright	1905-2017 StateCorp	LLC				Name Label	
There are no items to show.	Statistics/Data Analysis	statecorp 4905 Lakew	av Driva					There are no items	to show.
	Special Edition	College St	ation, Texas 77045 U	54					
		E00-STATA-	PC http://www						
		979-696-46	00 stata@stat	a.cos					
		979-196-41	VA (E&R)						
			🔯 Open				×		
	formation and the second second		6	No.PC 1 Depleters 1 January 2022 1 Dr. Prates 1	Mohanhy 1		A D South Dr Puter Mohern		
	to make an of the state								
			Organise • New fold	le			🗈 • 🖬 \tag		
	Notes:		A	Name	Date modified	1,pe	507		
	 Unicode is supported; see 	help unicode		1	10.00.00010.00	The Desider			
	Neximum number of variable.	s is set to	 OreDrive 	Let 10	12.15.2021 10.24	r ve folder L (a toldar			
			Dis FC	ler II		File faibles			
			20.06	Ler 12		f de tolder			
			Destroe	Lec 13	29-10-2121 17:47	Filefulder			
			Constant a	Lecture,17,syntax,and,dofile	05-11-2021 16-21	f de tolder			
			Convecation 2	Lecture_18_dictionary/ile	48-11-2021 17:48	File folder			
			Dr. S. K. Gupta	Lecture, 30, independence, of inviewent, a.	20-11-2021 17-22	f de folder		Properties	
			Foundation Da	Lecture, I7, Propensity, Scene, Metching	26-11-2021 16:38	Filefulder		+ Variables	
			IT Mande	Lecture28, provession, model on healten.	29-11-2021 15:59	File folder		Name	
			January 2022	NPTEL_Lecture_25_measure_of_goodness.	20-11-2021 16:53	File folder		Label	
			WacomTablet	NPTIL,Lecture19,20, Data,Browsing,Bas.	05-11-2021 17:40	l de folder		3.04	
			WRD & M	NPIE, GESHAL, MANYE, SA, MAR, DA.	2011/2021 1940	For fuller		Formal Value John	
			Decuments	State 15	27.10.2021 16.14	Einfahles		Notes	
			👃 Downloads 👋	AND .		10010000	~	+ Data	
			Ret	une			 All Stata Files 	Liename	
							Oran Cancel	Label	
							- the case	Note: Verables 0	
								Memory 640	4
								Sorted by	
	Command								
When Statis Uffectual Income 2022 f	h: Postan Mehant/Stata15								COLUMN STREET, NO.

Which we have kept for this purpose.

(Refer Slide Time: 30:33)

Kevev • • • ×				 Variables 	т а х
1 Iter commands have	(R)			1 (ter variables here	
# Command _rc		Conversions	1981-1017 StateCore 110	Name Label	
There are no items to show.	Statistics/Data Analysis	StateCorp		There are no items t	to show.
		4905 Lakew	ey Drive		
	Special Edition	College St DOD-STATA-	ation, Texas TTO45 UDA RC http://www.atata.com		
		979-696-46	00 stata@stata.com		
		979-696-46	01 (fex)		
			🚺 Open X		
	from agent formation		← i · · ↑		
			Onnies Buchtler		
			Vigence • Heavitater		
	1. Unicode is supported; see	help unicode	OraDina Name Date modified Ape Soe		
	Meximum number of variable	a is set to	OveDrive Milek1,2,3,4,5,8,9(a88),ata 28-11-2021 1427 070-File 14,416-13		
			💭 Tha PC		
			3 10 Objects		
			Centrop		
			Convocation 2		
			E S.K. Gupta	Properties	ų ×
			E Foundation Da		
			17 Mendi	Name	
			January 2022	Label	
			WaccenTablet	5,94	
			WRD & M	Format Voluminated	
			b) Documents	Notes	
			& Devricads	* Data	
			File name Islaud (2,3,4,3,4,0).dbis da 🗸 🗸 🗸 🖌 🗸 🗸 Vice Files	Elename	
			Open Cancel		
				Variables 0	
				Obsenations 0	
				Sax 0	
	Commano				

And I am and we will experiment with this some of the variables derive results.

(Refer Slide Time: 30:35)

eer T # x	(i)		in Variables	T & J
Command rc use "C1Uses/Studio A/Des		Cepyright 1995-2017 StataCorp LLC StateCorp	Name FSU	Label
describe Sector	Special Edition	4405 Lakeway Drive College Station, Teams 77045 USA 00-USTALS: https://www.state.com	Segment SSS	
		979-454-4600 stategetata.com 979-656-4600 (fax)	Househeld	
			Filter Centre_Round	
	transfer of the state		Schedule Samale	
	Notes:		Sector	
	 Unicode is supported; see Maximum number of variable 	elp unicode_advice. 1# met to 5000; mee help met_maxwar.	District Stratum	
	. use "C:\Users\Studio A\Desktop\Jan	ary 20210r. Pratap Mohanty(Lecture26regression_model on healtoare variables(block1_2_3_4_5_5_9(skb).dta*	Subjitratum SubjRound	
	- describe Sector		Sub_sample FCO_Sub_Re	
	storage display v variable name type format 1-	lue Bel variable label	Streteries	
	Sector byte \$10.00 5	erot		
			Name	Sector
			Label	
			1,94	
			Format	510.09
			Note upor	SELLEY
			4 Date	
			Liename	No.817.145.11
			Label Notes	
			Variables	136
			Observations	41,240
			Sax	13.80M
			Memory	ыM
	Command			

Now what I will do? I will first go through explaining describe; describe command that will give you the basic description of the variable. I will write down des or describe in full that will give us the correct result. Any variables I can just write it down.

Sector I have taken here alright. So, now, the describe I have entered with that variable. Now, this has given me information about what is the storage type of this particular variable and what is the label of this, since label are not defined correctly what categories it has. So, label it is not showing in detail that we will show it later. But at this moment I will explain it differently.

Similarly coming to another command like; so, like in the describe command once again let me explain that this list various features of the variables to be used in the linear regression. So, by typing describe all your set of variables may be starting with dependent independent it will give you the complete list of which is the description.

(Refer Slide Time: 32:03)



Like here if you have two variables one is dependent that log medical expenditure then the second one is reimbursement of medical expenditure. Then it has given which kind of storage it is and how it is presented. Variable time and in format columns indicate that all the data are in this case numeric or not it is since we have not in destring it, but at this moment it is not required for you to destring even if the data is in string this describe table will give you result.

But in case of summary, it is not going to give you, we are just coming to it. The next important estimation is called the summary table, I will straight away go to this particular software. And now I will summarize the sum; the only a sum is fine or a complete summary of the variables.

(Refer Slide Time: 33:09)

nninn T B X Genraund hers Command //c use "Cliben/Studie ADes describe Sector"	I I I II III III III IIII IIII IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII	Sayraya (M-Sull'Sauray MC Sauray 485 Galaya Yana Galaya Katai, Sauray Katai Galaya Katai, Sauray Katai M-Gal-Galay (Katai M-Gal-Galay) (Katai	 Verdation Ing. Ing. research. Ing.medical, 	
	Notesi 1. Unicode is supported; see 1 2. Maximum number of variable:	ig wante office. In or to 1001 on help we parter.		
	. use "C:\Deere\Studio A\Desktop\Zeon . describe Sector storage display va	ny 1979ar, finda Maanykeetereks, regreening, moli oo baltaare versikarikkeeksjojityksjonek dae. M		
	Genter byte kib.dg Se	5.6	Bit Packate Range right Range right	nedical, org
	Command			

If I just take one variable for your reference this is going to give you the result.

(Refer Slide Time: 33:15)

Verv T 4 x Effer commands have Command Command Converted Coloribe Station Sociale Sta	(a) 	Segregar 1916-1917 Section 1927 Best-Server Street 1916 Bill Johnson (Street Herman 1916) Die Methall (Street Herman) Die Methall (Street Herman) Die Methall (Street Herman)	Variables Vare	Label
	Motee: 1. Unicode im supported; see he	g unanymus.		
	 Haximum number of Variables use "C:\Users\Stmfio A\Desktop\Janua describe Sector 	n në ti svin in dag nëjatur. 19 META: Freig Mahayileturele, repressio, koki o haltare virikles(31043.3.3.3.5.3.3.1048).46*		
	atorna display sal			
	storape display val variable name type format lab	en el versable [stel]	Properties	
	storage display val variable name type format lak Sector byte %10.0g Sec	er el seriable iabel Mar	Properties	
	etorape display val variable name type format lak Sector byte %10.0g Sec . emm log_medical_exp	se 11. vrijalje jelej 101	Properties M + + + Variables Name	log, medical, en
	storage display val variable name type format lab factor byte %10.0g Sec . em log_metical_exp variable 7 cm Man	ar a versative later. And The Man	Properties	log, nedsal, o
	etorape dirpley val variable name type dormat lan fector byte tilo.0g Sec . em log_metical_emp Variable Ote Hean log_metica-p 37,160 5.353011	of second to labor. 197 197, Sec. No. No. No. 1 (1994) 8 13 (1934)	Properties	log,medical_p Rost NADg
	111100 60(1) 111 11101 100 100 100 64400 130 640 100 64400 130 640 100 10100 100 640 100 10100 100 640 100 10100 27,340 5,3500 1	9 9(1) 17 17 18 19 19 19 19 19 19 19 19 19 19 19 19 19	Poperies B + + + Model I det I de	10 g.metrol.co Reat Skilly North 2, 1, 4, 5 10 41,240 11,961 840 840
	netropy degli y netropy degli y lanter byte thicky to netropy degli y very degli y ver	of seconds (and)	Properties	log,melvid, or Rust Sklog Novi 1, 2, 1, 4, 5 13 4, 24 13, 50 13, 50 13, 50 13, 50 13, 50 14,
	ericipi energi e	9 911. 1941 1. 1947 0 1. 1949 1. 194910 0 1. 19391	Pagentes in + + - Solidar - Solidar - Sonor - Sonor	1 Rog, medical, os Road, 34 Og Barcel (2, 1, 4, 5, 41, 20 13, 20 13, 20 13, 20 13, 20 13, 20 13, 20 13, 20 13, 20 14, 20,

Now, this has given us information that log medical expenditure which we wanted to understand. The summary table is going to give us information about how many observations are there, what is the mean value of the expenditure, what is the standard deviation, and what is the mean and maximum value.

Now, just for your clarity, I am going to give you another summary which is string data, let it be the string data of let it be the sector the same sector variable which we have already done it.

(Refer Slide Time: 33:51)

1.1.2	<		 Veriables 	
commands have	⁽⁰⁾		A 10	
mmand ,rc		Copyright 1995-2017 StateCorp 110	Name	Label
"C1Uses/Studio A/Des	Statistics/Data Analysis	StateCorp	Geprert	
uibe Sector		4905 Lakeway Duive	Household	
log_medical_exp	Special Edition	College Station, Texas 77045 USA	Sector	
	-	000-STATA-PC http://www.stata.com	Response, C	4
		979-696-4600 stata@stata.com	Romarks_eb	d
		8/8-680-4601 (E8K)	HouseheldS	ice .
	In control to an an annual sector		childbirthEx	4-
	former applies discovered		Household	,L.
	summer or 10 books		Type_of_late	i
			Hew, many,	
			OutbreakCo	
	1. Delocia in supportadi ana	ale velocia advica.	Househeld	9.
	2. Maximum number of variable	is set to 5000; see help set maximum.	no_of_times	<u>-</u>
			cont provi	E.,
	. use "C:\Dsers\Studio A\Desktop\Jan	ary 2022\Dr. Pratap Mohanty\Lecture26_ regression_model on healtcare variables\block1_2_3_4_5_8_9(a&b).	dta" communica	e
			weeke no	
	. describe Sector		tree, resp	
			Wathrest	6
	variable name type format 1.	bel variable label	A	
			Properties	
	Sector ofte 410-00 S	101	# Variables	
	. sum log medical exp		Name	Sector
			Label	
	Variable Obs Nea	Std. Dev. Nin Max	har	byte
			Format	
	log_medica-p 37,100 5.95301	1.250563 0 13.12236	Value label	
			4 Deta	
			Linowe	
			Label	
			Notes	
			Variables	136
			Observations	41,240
			Sax	13.58M
			Memory	64M
	Command		,	
	sum Reador			

Now we will summarize the sector variable here now you will find the difference between string variable and nonstring variable.

(Refer Slide Time: 34:05)

🚦 Stata/SE 15.1 - C.(User/Studio N/Deskto	(January 2022/Dr. Pistap Mohanty/Lecture26, regression, n	nodel on healtrare variables/bloch1,2,3,4,5,3,5(adm).dm		- 0 ×
File Edit Data Graphics Statistics	User Window Help			
S ■ S U = · ⊥ · Ø · Ø :	0.0			
Review T # X			Variables	тах
1 itse commands have			3 m	
A Command IV			Name	Label A
1 use 'Cillien/Statio #/Des		opyright 1985-2017 StateCorp LLC	Second	
2 describe Sector	4	And the second sec	Household	
3 sum log medical exp	Special Edition C	ollege Station, Texas 77045 USA	Sector	
4 sum Sector		00-STATA-PC http://www.stata.com	Response_Co	
	5	79-696-4600 stata@stata.com	Remarks, ebe	
	5	79-696-4601 (EAR)	HouseholdSize	
			childhithEsp	
			Household_t	
			Spejeljati	
			Hex_many	
			OutbreakCo	
	Notes:	n unione efficie	Household_U.,	
	 Maximum number of variables i 	a set to 5000; see help set magner.	no_el_times	
			child_bith_e	
	. use "C:\Users\Studio A\Desktop\Javaar	y 2022/Dr. Fratap Mohasty/Lecture26regression_model on healtcare variables/block1_2_3_4_5_8_9(a&b).dta*	communicab	
			in time break	
	. describe Sector		Doot service	
	storage display valu		Whathwadic	~
	variable name type format labe	1 variable label	Properties	9 ×
	August 10 10 10 10 10		8 * *	
	secor ofte thoug seco	ot	 Variables 	
	. sum log medical exp		Name	times/hospitalised
			Label	
	Variable Obs Mean	Std. Dev. Nin Max	3,94	and .
			Formul	2.91
	100_8e310a-p 37,100 5.953011	1.250563 0 13.12236	Value label	
	. sum Sector		+ Data	
				Not12145190
	Variable Obs Hean	Std. Dev. Min Max	Label	
			Notes	
	Sector 43,219 1.492509	-4997025 1 2	Variables	136
			Observations	41,240
				13 SGM
			forence y	5451
	Command		,	
	sum times_hospitalised			
C1Uses/Studio A/Desktog/January 2022/Dr.	Piatap Mohanty/Stata15			CAP NUM OIR
P Type here to search	o = 🂿 🛪	💉 💽 📙 Statul S 🥼 NPTEL, Lacture 26, 🥐 Provedhiet State S 🔠 Statu SE 1531–CASL	4× JA 046	1629
				or maker 🕚

(Refer Slide Time: 34:18)

w T # X		Copyright 1985-2017 StataCopp LLC	A Variables T
Command _rc	statistics/sata Assiyata	405 Lakeway Drive	Nave Label
use "C1Uses/Studio A/Des	Special Edition	College Station, Texas 77045 USA	Segment
lescribe Sector		800-STATA-PC http://www.stata.com	Household
an log, medical, exp		\$75-656-4600 state@state.com	Sector
m Sector		179-696-4601 (fax)	Response_Ca
n times, hospitalised			Remarks_ebe
			HouseholdSize
	comment of the boundary		childhirthExp
			Househeld_t
			Type_of_lats
	Noteel		How, many
	 Unicode is supported; see 	elp unicode_aduice.	OutbreakCo
	Maximum number of variable	is set to 5000; see help set_maxwar.	Household_U
			re_ef_times
	. use "C: (Users(studio A(Desktop)/so)	ary 2022/br. Histob Woushth/Peoplesie Lediessien Bogel on Designate Aslightes/Diodx1 5 3 4 2 8 A(990) 404.	child_bith_t
	describe Sector		communicab
			Whether_hos
	storage display w	lue	up.a. times, hospit
	variable name type format 1	bel variable label	not_seeking
			Whethreads
	Sector byte 410.0g 5	250F	Properties
	and has ended and and		A **
	. aum tod mentcar esb		 Variables
	Variable Des New	Sed Save. Non. Max	Name times/septal
			Label
	log_medica-p 37,100 5.95301	1.250563 0 13.12236	Type and
			Format Sile
	. sum Sector		Value label
			Noter
	Variable CD-8 Mea	Std. Dev. Min Max	* Data
			Liename block1_J_1_4
			Laber
	. sum times hospitalised		Versibles 116
			Dhansheer Al Sal
	Variable Cbs Mea	Std. Dev. Min Max	Cor ITelas
			Memory 64M
	times_hosp-d 0		Sorted by
	P		v.
	Command		
	br times hospitaliand		*
	at times maprications		

Now enter you will find that this also is going to give us values because this has already been destringed. Will give you another one let it be this is time spent in hospitals. Since this is a destring variable this is a string variable, the summarize table is going to give us 0. How do that this is a string variable? Just check browse and this variable br and times hospital.

(Refer Slide Time: 34:42)

times, hospitalised;									
time_hosp-d								Variables	
L								A Fite variable	
1								92 Name	Label
								Directory	
								0 150	
								() Segment	
								0 555	
								() Household	
								[] Level	
								[] file	
								Centre,Round	
								[] Round	
								[] Schedule	
								C) Strate	
					Ν.			C Mit Kerne	
					4			C fided	
								C) Statum	
								() Sub stratum	
								C) Sub-Round	
								() Sub_sample	
								0 100,5x8,8x	
								Variables Snapsk	Auto
								rropertes	
								* Variables	
								Label	
								test	
								format	5.01
								Value lubel	
								Notes	
								* Data	
								Limana	block1.7,
								Label Notes	
								Valuables	116
									41,240
								544	13 SDM
								Mannary	5454
								Serted by	
							Ψ.		

The font of this one must be in red or there are no entries at all or I will just do one thing there is no observation here.

(Refer Slide Time: 34:51)

times_hospitalised[1]											
time_hosp-d										 Variables 	
375										A Filter variable	
974										S Name	Label
976										@ times.hospt	
974										C) FSU	
177										Segment	
170										S22	
379										Household	
100										🗋 Level	
101										() filler	
02										Centre, Round	
10										U Round	
04										C) Schedule	
15										C) Sample	
N										C MS Reside	
17										O Dated	
										O Suba	
										0.54 mmm	
										Sub Round	
10										() Sub_sample	
21										0 100 Sub Re	
52										Variables Seams	helt
99											
14										Properties	
15										* Variables	
м										Paperne	
97										Line	
10										- Lograt	2.2
10										Value label	
0										Notes	
1										# Data	
3										Elmane	block1,2,
3										Label	
4										Notes	
5											135
4											11.004
,											MM
1										a Sorted by	
										R.	
0											
										-	
										×	

(Refer Slide Time: 34:53)

8 II 3 · L · B · N	1 O O		
Tex	Statistics/Data Analysis	StateCorp	 Veriables
		4905 Lakeway Drive	4
Company and	Special Edition	College Station, Texas 77045 USA	
mmand _rc		800-STATA-PC http://www.stata.com	Name Label
"C10ses/Studio A/Des		979-696-4600 stata@stata.com	Segment
cribe Sector		979-690-4601 (EAX)	Royshold
ind medications	Manager Annu Sala annual Sala		5600
s Sector	terry again transmission		Report Co
iner houridated	and a second second second		Monadal Con
ind indpande	1		childhithEan
			Household 1
	Notes:	hele verende edular	Sor of Inti-
	2. Maximum number of variabl	a is set to 1000; see bala set mayor.	How many -
			OutbreakCo
	. use "C:\Users\Studio A\Desktop\Ja	mary 2022\Dr. Fratap Mohanty\Lecture26 regression model on healtcare variables\block1 2 3 4 5 8 9(a&b).dts"	Household U.
			no of times
	. describe Sector		child bith e
			 communicable disease
	storage display	alue	Whether hos
	terrente name cipe tormet	arei isilanie isisi	time, lespt.
	Sector byte \$10.0g	ieccor	not, seeking
			Whetheredic
	. sum log_medical_exp		Properties
			8.11
	Variable Cbs No	n 318. Dev. Rin Max	* Variables
	log medica-p 37,100 5,9530	1 1.250563 0 13.12236	Nava temp
			Label
	. sum Sector		ited and
			Format Silv
	Variable Gbs Ne	n Std. Dev. Hin Hea	Value label
			hister
	546501 63,214 1.492	9 .4997825 1 2	a Data
	sum times hospitalized		Land Land
			Notes
	Variable Obs Ne	n Std. Dev. Min Max	Variables 136
			Obsensibles 41.740
	times_bisp-d 0		Sax 13.800
	be along benefitived		Memory 64M
	. er times_sospitations		Sorted by
			·
	Command		
	24		
- Contraction - MARA	Access Million and Access 11		

Just browse we will just see which are the string variable browse only browse br enter.

(Refer Slide Time: 35:00)



So now, you can see which variables are in fact your string end which are nonstring or numeric, or labeled numeric. I have already clarified to you that the values or the variables which are in red font red color are string variables.

(Refer Slide Time: 35:23)



Those values are in black font.

(Refer Slide Time: 35:34)



They are all black font they are in fact like the black font they are called numeric values.

(Refer Slide Time: 35:46)



Those who are in blue color they are called labeled numeric they are not exact numeric values they are in fact labeled. So, if I just do one piece of information like at the last these are all our numeric values. From the beginning I will just do one thing, I will just try to summarize one piece of information on your screen; religion let it be. So, now, I will go to explain religion here in red font.

(Refer Slide Time: 36:15)

nand res Silben/Studio A/Des	Securit Antenne Cellarge Francis, Fees Ville Tall Electrical Antenne Tellar Tall Electrical Antenne Ville Tall Electrical Antenne Ville Tall Electrical Antenne Ville Tall Electrical Antenne Ville Tall	A Variables Vari
inbe Sector i Tog, medical, exp i Sector i times, hospitalised imes: hospitalised	[2] Control and a subject of control for the subject of the sub	Ne Relation_to_h
	 Dittoi is supported see help uscode whice. Rashan namber of variables is set to 10000 see help set marray. 	
	. use "C:\Derre\Sindio A\Deskip\January 2022\Dr. Pratap Mohanty\Lecture26_ regression_model on healtoare variables\blockl_2_3_4_5_8_9(ab).dta"	
	. describe Sector	
	storage diaplay value variable name type format label variable label	
		_
	. em log_metical_emp	
	Variable Obs Nean Std. Dev. Nin Nax	Presenties
	log medica-p 37,100 5.453011 1.250563 0 13.12236	A
		* Variables
	, and Sector	Label
	Variable Gbs Hean Std. Dev. Hin Max	Type shi1
	Sector 43.319 1.492599 .492525 1 2	Format 1001 Value label
		Notes
	. sum times_hospitalised	* Data
	Varisble Obe Nean Std. Sev. Min Max	Label
		Notec
	clnes_hosp-d 0	Variables 136
	. br times_bospitalised	Sax 1139M
		Memory 54M
		Sorted by
		v.
	Command	

So, now, summarize religion.

(Refer Slide Time: 36:23)

🔯 Sata ST 33- C. Dever Studio N. Destrop Sanary 2022 Dr. Fratag Mohardy Liercowski, respective, no odel on hashcare vaniliteitation 2, 3, 4, 5, 9, 90 Millionta	- 0 ×
File Edit Data Graphics Rotatics Unar Window Help	
◎ 第 ● □ 〒・1 □ - □ □ □ ○ ●	
Deview T # X	 Variables T # >
A little constant has 0	A et
# Command jrc	Name Label
ue "Clubes/Studio Albes	Religion
2 decode Setter 1. Unicode is augorited; see delp Unicode Avrice. 1 unita medica no 2. Maximum sublect of variables is used to 5000 see belp set maxwar.	Kdation_to_h
4 um Settor	
5 un tem_hopkind	
6 brinnsjouptiond . describe Sector	
8 sum Religion atorage diaplay value	
variable name type format label variable label	
Sector byte 410.0g Sector	
. was log_medical_emp	
Variable Obe Mean Std. Dev. Min Max	
log_medica-p 37,100 5.853011 1.250363 0 13.12236	
. was Sector	
Variable Oba Mean St.S. Dev. Min Maa	Properties 8.3
Sector 43,219 1.482599 .4942525 1 2	8
. emm times hospitalised	Name Religion
Variable Obs Nean Std. Dev. Nin Nau	Label Type stri
	Format %N
Class_Disp-3 0	Notes
. br times_hospitalised	+ Data
	filmane block1,2,3,4,5,1,9
	Notes
. ema Religion	Variables 136
Variable Obs Mean Std. Dev. Nin Max	Obsenations 41240
	Memory 6M
881.55100 U	Sorted by
	v
Command	
tab Heligion	
Cilibard Studie #Deskteplanuary 2020: Protag Maham/Stata1	CAP NUM ON

Now, this is going to give us 0 values. So, 0 values, because this is in this, is in a string value. Now, if I destring and run the same command it will give me the right result.

(Refer Slide Time: 36:39)



I can destring it quickly and then I can save it. I think I have already discussed that and we will also use it several times later. So, now, we will go to our PPT and stick to the requirement.

(Refer Slide Time: 36:52)



Similarly next is our understanding frequencies; tab, tabulate, there are tables, they are different commands, I know you will get it very correctly from my modules. So, tab is going to give frequencies of a particular variable tab command you can give it and then the variable

name. But if you like here if I say t a b tab the only t a then any variable here let it be same religion enter.

(Refer Slide Time: 37:29)

🔤 Stata/SE 15.1 - C1/User/Studio N/Desite	pilanuary 2022/Dx. Px	tap Mohanty Lectu	neth, regression, i	model on healtcare	variable/block	00.000	δh).du		- 0 ×
File Edit Data Graphics Statistic	User Window	Help							
🗃 🖬 🕲 🗟 - 🖬 - 🛃 - 🛃 - 🛃 - 🛃 -	3 🗍 0 - 0								
Review T # X	Sector	byte 41	0.0g Sect	tor			×	Aviables	T # X
A Litter commands have 0	ann log ned	ical emp						A ni	
# Command ,rc								Name I	abel
1 use "Cillises/Studio A/Des	Variable	08.4	Nean	Std. Dev.	Nin	Max		Religion	
3 sum log medical, exp	log_medica-p	37,100	5.953011	1.250563	0	13.12236		Keason, o, K.	
4 sum Sector									
5 sum times, hespitalised									
7 br	Variable	08.8	Nean	Std. Dev.	Nin	Max			
8 sum Religion	Sector	43,219	1.402509	. 4997025	1	2			
r an ranger	. em times_h	ospitalised							
	Variable	08.8	Nean	Std. Dev.	Nin	Маж			
	times_hosp-d	0							
	. br times_ho	spitalised							
	. br								
	. em Religio								
	Variable	06.8	Nean	Std. Dev.	Nin	Маж		Properties	# ×
	Patron of							* Variables	
		, ,						Narra	Région
	. tab Religio							Label	(d)
	Religion	Freq.	Percent	Cum.				Format	Mr.
								Value label	
	2	7,321	16.94	90.00				+ Data	
	3	2,495	5.77	95.17				Liename	No.117,174,5150
	1	1,177	2.72	90.49				Label	
		215	0.60	\$9.57					
	,	13	0.03	99.60				Character	41.343
	, ,	171	0.40	100.00					11304
	Total	43,219	100.00					Memory Sorted by	ым
	Command								
	tab Sector Rel	lgios							
C1Uses/Studio A/Desktop/January 2022.Dr	Pratag Mohanty/Stat	15	-	_					CAP NUM OIR
P Type here to search		0	H 🌖 -	< 💽 🛛	Stati		🕒 MPRL Lecture 28 👔 Received Side S. 🗮 Statu SE 531-656. 🌰 24°C. Cloudy \land 📀 🎵 🚦	4× d ⁶ 046	29-11-2021

This is going to give me what? This is going to give me frequency and its cumulative frequencies. And tabulate is the command is also possible to operate in case of the string variable. But since it is not calculating anything once it is estimating anything the string variable is not going to give you any result.

So, now if in the place of double tab and two variables that will be giving you the values of two that is frequencies cross-tabulation of the variables. So, tab then two variables let us take tab and two variables sector and religion.

(Refer Slide Time: 38:25)

Norm Norm <th< th=""><th></th><th>Variable</th><th>06.4</th><th>Neas</th><th>Std. Den</th><th>. Mir</th><th>n Ma</th><th></th><th></th><th></th><th></th><th></th><th></th><th>A Variables</th><th>14</th></th<>		Variable	06.4	Neas	Std. Den	. Mir	n Ma							A Variables	14
art Standbark are Stan	the commands have	Sector	43,219	1.492599	.4997025			,						N nl	1 Mart
Notice of the second	e "Cillion/Studio AlDes	. em times_	ospitalised											Religion	CREW
Note (note) 0 0 (note), hadded (monophild) 10 0 10 10 (note), hadded (monophild) 10 10 10 10 10 (note), hadded (monophild) 10 10 10 10	escribe Sector um log_medical_exp	Variable	084	Nean	Std. Der	. N14	n Me							Relation_to_h	
41 Use_Merit Uses	m Sector m times, hospitalised	times_hosp-d	0					-					N		
Instrume	times, hospitalised	. br times_ht	spitalised										4		
Angene and an and an angene and an angene and angene and angene a	m Religion	. he													
Vitikalit Ref R	o Faligion di Sector Religion	. em Religio													
bittom 0 state Margine Con- bittom Fractor Con- 1 39,501 30.46 33.46 2 39,501 30.46 33.46 3 10,501 57.97 80.40 40.46 4 10,77 70.98 40.46 10.47 5 10.40 57.97 80.40 40.46 4 10.46 80.46 40.46 40.46 4 10.46 80.46 40.46 40.46 5 10.46 90.46 90.46 40.46 5 10.46 90.46 90.46 40.46 6 10.46 10.46 10.46 10.46 10.46 101 2.43 10.46 10.46 10.46 10.46 10.46 101 2.43 1.41 13.4 12.40 11.40 12.40 10.46 10.46 101 1.43 1.43 1.44 11.40 12.40		Variable	06.0	Nean	Std. Den	. Nir	n Me								
Autory France		Religion	0					-							
Initial Frank Frank Conc 1 1,534 13,64 13,64 14,64 16,64 1 1,534 13,64 13,64 16,64 16,64 16,64 1 1,534 13,64 16,64 <td< td=""><td></td><td>. tab Religio</td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td></td<>		. tab Religio													
1 17.94 13.64 9.64 2 7.93 5.79 40.71 3 7.93 5.79 40.71 4 17.91 19.84 40.71 4 17.92 4.83 40.71 6 107 4.84 40.71 7 10 4.93 5.97 10 4.98 40.71 101 4.88 40.71 101 4.88 40.71 101 4.88 40.71 101 4.88 40.71 101 4.88 10.11 101 4.88 10.11 101 4.88 10.11 102 1.33 3.84 10.11 101 1.33 3.48 10.11 10.19 1101 1.33 3.48 1.47 10.31 10.41 1101 1.33 3.48 1.47 10.31 10.31 10.11 1101 10.11 <t< td=""><td></td><td>Religion</td><td>freq.</td><td>Percent</td><td>Cum.</td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td></t<>		Religion	freq.	Percent	Cum.										
2 7,301 14.40 90.60 4 2,401 37.90 40.60 5 1,303 4.00 90.70 6 103 4.00 90.70 7 10 4.00 100.70 10 4.00 100.70 100.70 10 4.00 100.70 10 4.00 100.70 10 4.00 100.70 10 4.00 100.70 10 10.70 100.70 10 2 3 4 5 10 2 3 4 100.70 10 10.70 100.70 100.70 100.70 10 10.70 100.70 100.70 100.70 100 10.70 100.70 100.70 100.70		1	31,574	73.05	73.06										
4 1,177 2,79 40,40 40		2	7,321 2,495	16.94	90.00 95.77									Properties	4
1 133 5.00 10.91 10.00<			1,177	2.72	90.49										
1 103 6.0 10.13 6.0 10.13 6.0 10.14 10.00		5	173	0.40	50.85									* Variables	
1 10 4.00 84.00 100.00		6	235	0.69	99.57									Paperse	reigen
Image: Second		?	13	0.03	99.60									Laber	
Strate <th colspan="</td> <td></td> <td>,</td> <td>171</td> <td>0.40</td> <td>100.00</td> <td></td> <td></td> <td></td> <td></td> <td></td> <td></td> <td></td> <td></td> <td>194</td> <td></td>		,	171	0.40	100.00									194	
		Total	43,219	100.00										Value label	
Interview 1 2 3 4 5 4 7 9 71141 10000 1000		. tab Sector	Religion											+ Data	
Series 1 2 3 6 5 6 7 5 7041 1011 <th< td=""><td></td><td></td><td></td><td></td><td></td><td>B-1-1-1-1-1-1-1-1-1-1-1-1-1-1-1-1-1-1-1</td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td>Liename</td><td>block1,7,1,4,3,</td></th<>						B-1-1-1-1-1-1-1-1-1-1-1-1-1-1-1-1-1-1-1								Liename	block1,7,1,4,3,
inter i <td></td> <td>541714</td> <td></td> <td>,</td> <td>,</td> <td>##13@100</td> <td>· .</td> <td></td> <td>,</td> <td></td> <td>Torel</td> <td></td> <td></td> <td>Label</td> <td></td>		541714		,	,	##13@100	· .		,		Torel			Label	
rttl: 16.30 3.40 1.40 10 14 19.20 10.10 10.40 10 obse 15.33 3.40 1.60 30 145 15 12 15 10.40 10.40 10 10.40 1		.40107		*			,			,				Noter	
othe 15.33 5.493 1.41 240 145 145 12 24 26.01 156.01<		rurel	16.341	3,430	1.476	794		190	1	114	22.362			Vanables	135
Testal 31,394 3,301 3,485 1,317 132 245 13 171 43,219 Mathematical		uzban	15,233	3,083	1,019	383	165	105	12	57	20,857			Ubsenations Cost	11,257
Tetral 31, 574 7, 301 2, 405 1, 177 173 285 13 171 43, 289 Sector 17														Mamory	ALM.
		Total	31,574	7,321	2,495	1,177	173	592	13	171	43,219			Sorted by	
•														*	
Commond .		A													

Then religion gives us information about the cross-tabulation of how rural and urban sectors, are actually distributed across religion, and its clearly understood. Now there are further sophisticated commands possible and we have already explained I have explained them all in detail in the review of command section in my previous module as well you can just follow and find it out.

We will do some new addition in this module on understanding some graphical representation correlation.

(Refer Slide Time: 39:02)



So, in order to understand the correlation or to check the relationship between variables scatter plot helps in understanding this graphically. Now what we will do? We will copy this command and simply run it. So, and just to show you quickly how it follows and that that is the reason why we have kept it on your screen and I am sure you can operate on your own as well.

So, we will copy this scatter we could have entered manually, but keeping the time constant.

(Refer Slide Time: 39:43)



I am only directly operating it in the screen.

(Refer Slide Time: 39:45)

T # X	. sum times he	epitalised											 Variables 	
ommands have	Variable	0.4	Nean	Std. Dev	. Nin	Ма							🔨 ni	
and re-							-						Name	Label
Sector	Cines_Disp-3	0											Relation_to_h	
medical_exp	. br times_hos	pitalised												
or 15 hospitalised	. br													
hospitalised	. sum Religion													
oine.														
jon .	Variable	08.8	Nean	Std. Dev	. Kin	Ма	-							
u Religion	Religion	0												
og medical josp ro 111	. tab Religion													
	8+114100	F1+1.	Percent	Cun.										
		iii qi												
	1	31,574	23.06	73.06										
	5	2,495	5.77	95.77										
	•	1.177	2.72	90.49										
	5	173	0.40	50.05										
		295	0.69	99.57									Properties	
		13	0.03	99.60									8.11	
	,	174	0.40	100.00									 Variables 	
	Total	43.219	100.00										Nava	Edigion
													Label	
	. tab Sector B	lel1gion											3,04	10/1
													Format	5.01
					Religion								Value label	
	Sector	1	2	3	•	5	6	7	,	Total			higher	
	TUTA1	14.141	1.418	1.476	194		140		114	22 162			* Data	March 1
	urben	15,233	3,083	1,019	303	165	105	12	57	20,057			DM	
													Notes	
	Totel	31,574	7,321	2,495	1,177	173	295	13	171	43,219			Variables	136
													Observations	41,240
	. scatter log	medical exp	responsed a	by med lasts									Sax	13.90M
		arrest by last	Cumu tor 1	1.04114									Memory	64M
													Sorted by	
	_		_											
	Command													

And this is going to give us the scatter plot.

(Refer Slide Time: 39:50)



So, variable symbols are not available we can just change them.

(Refer Slide Time: 40:05)



Now it has the name has been changed you can just check, how we also find that there is some errors in the variable. The variable name has been changed in my our data, which is why the model the command was not running. You have to give the exact name by searching through the window.

Now, you can find out how the two variables are related and where they are concentrating what is in fact the trend line all those things we are going to explain, in this model. Then next I am going to operate with we can also operate two-way table as well two ways scatter plate as well, then the average could be also run we are not going to do it at this moment. We will simply skip it and that is kept for your reference only.

(Refer Slide Time: 40:55)

	average reimbursed medical expenditure to log medical expendit	
🎯 swayani 🤶		25

(Refer Slide Time: 40:56)

We can combine scatter with lfit to show a scatterplot with fitted values. graph twoway (scatter log_medical_exp reimbursed_by_med_insur)(lfit log_medical_exp reimbursed_by_med_insur), title("average reimbursed medical expenditure to log medical expenditure")

Like here you can do like you can combine scatter with ifit to show the scatter plate with fitted values. So, the fitted values can also be explained, so that command we have already given for your reference.

(Refer Slide Time: 41:16)



I am sure the trend line the result will look like this and the trend line basically the fitted line command if you attach it will give you this kind of result. So, fitted line in this case it seems that the fitted line is not that correct. So, you have to include multiple variables and the relevant variables to draw the fitted line.

(Refer Slide Time: 41:39)



Similarly, correlation can also be checked through this command where I am just going to run it as per this one. So, it is here. I am just copying it once again and we will go to that

particular page. And we have kept all those things for your reference. Yes, this is the one. We are just copying it and running it on the window.

(Refer Slide Time: 42:12)

	3 0 0 0												
w 193	Variable	CB-8	Nean	Std. De	r. Mil	n Me	*					 Variables 	1
iter commands have	times hisp-d						-					* m	
Command _rc	the other											Name	Label
use "C1Uses/Studio A/Des	. br times_bo	spitalised										reimbursed,	
describe Sector													
um log_medical_exp													
um timus houritalised	. sum Religio	e.											
times hospitalised													
	Variable	06.8	Nean	i Std. De	r. Rij	n Ke	*						
un Religion	Religion	0											
b Religion													
ab Sector Religion	. tab Religio												
atterlog,medical,exp.rol., 111		F		0.0									
cabe log_medical_exp re	, service	riet.	Percent	0.081									
	1	31,574	23.05	73.06									
	2	7,321	16.94	\$0.00									
	3	2,495	5.77	95.11									
		173	0.40	50.05									
	6	235	0.68	99.57									
	,	13	0.03	99.60									
	,	131	0.40	100.00								Properties	
	Total	43,219	100.00									A	
												 Variables 	
	. tab Sector	Religion										Taprie Label	removied_b
	1 1				Religio							Tat	long
	Sector	1	2	3	4	· .	6	7	+	Totel		format	5.10.09
												Value label	
	rurel	16,341	3,438	1,476	794		190	1	114	22,362		Notes	
	01045	15, 651	3,043	1.019	202	165	145	14	57	20,457		* Data	
	Total	31,574	7,321	2,495	1,177	173	235	13	171	43,219		Label	
												Notec	
	. scatter log	medical_exp	reisbursed	by med_insu								Variables	135
	r(111);	poreed by sec	Creater cos	0.000								Observatione	41,240
												548	13.55M
	. scatter log	medical_exp	reisbursed	by_medical_	10925							Sorted by	
	P											v	
	_												
	Command												
		tion) and pair	inburned by	medical ins	or, star(0.	05) pig							
	pecorr log_med	treat_eap re-											
	pwcorr log_med	inter_opy in											

The window is here alright.

(Refer Slide Time: 42:15)

📾 🗍 🖉 - 🖬 - 🖾 - 🖾 i	i User Window H	44p											
	3 0 0 0												
111	Religion	0										A Variables	'
ter commands here													
command jrc	. tab Religion											Nave	Label
lescribe Sector	Religion	Freq.	Percent	Cum.								intervente,	
un log_medical_exp													
um Sector	1	31,574	73.05	73.06									
um times, hospitalised		7,321	16.94	90.00									
etimes, hospitaloed	4	1,177	2.72	90.49									
e de la companya de la	5	173	0.40	50.05									
um Keligion	6	235	0.69	99.57									
ub Santos Rabaino	1 1	10	0.03	99.60									
catter log medical expires. 111			0.00	100.00									
catter log_reedical_expresi	Total	43,219	100.00										
week log_medical_exp_m_													
	. Cab sector Re	singson											
	1 I				Religion								
	Sector	1	2	3	4	5	6	7	*	Total			
	UTAL	15,233	3,438	1,476	194	145	190	12	57	22,362			
												Properties	
	Total	31,574	7,321	2,495	1,177	173	295	13	171	43,219		8.1.1	
	energian loss a	netical em	rainburged (to net losse								+ Variables	
	variable reinho	orsed by med	insur not	found								Name	reinbused_b
	(610)											Label	
												1.04 Locat	
	. scatter log_s	Bearcer_exb a	responsed	by_menical_i	UNIE							Value label	
	. pwcorr log me	edical_exp :	reinbursed	by medical i	netr, star(0	05) #1g						Notes	
												# Data	
		jod we-b tet	tabu-r									Liename	bick1,7,3,4,
	los medica-n	1.0000										Label	
												Variables	136
												Observations	41,240
		0.1110+ 1	1.0003									Sar	13.90M
	14192-11840-1											Memory	64M
	reinburee-r	0.0000										Sorted by	
	telbpoteeq-t	0.000											
	teimbursed-r	0.000											
	telmpoteeq-t	0.000										-	
	reisbureet-r	0.0000											
	resaburee6-r	0.0000			- h							4	
	resaburee6-r	0.0000			ħ							4	
	resabureerr	0.000			h,							4	

So, the correlation and the significance value basically the power of the correlation we have estimated. So, the star mark represents the level of significance at the 5 percent level 0.045

command is attached to this star. So, wherever the relationship is established at a 5 percent significance level it indicates the star mark.

Therefore, star mark you can just see, that I am just explaining once again. Similarly multiple linear regression linear in case of for the multiple linear regression we can also understand the power of the correlation basically the level of significance. You need to just attach this part.

In addition to that comma followed by comma then star with the bracket at 0.5 percent etcetera if you give it then it will give you the right result. So, the plus and minus sign indicates the value and the relationship whether positive related or negative related.

(Refer Slide Time: 43:26)



Higher value of rho represents stronger relationship or weaker relationship.

(Refer Slide Time: 43:27)

Normality of Y-Variable	
 Normality of Y variable can be shown graphically in two ways <u>Histogram:</u> <i>Command:</i> histogram log_medical_exp We can use the normal option to superimpose a normal curve on this graph and the bin(20) option to use 20 bins. Command: histogram log_medical_exp, normal bin(20) 	;: ve
swym 🧕	30

The next check is on the normality of the Y variable, the dependent variable how it is normal, and what extent we emphasize the normality of the distribution. The command we usually follow to go for a histogram how the histogram looks like, just to check histogram log a medical expenditure which the variable we have taken that we can check to do it once again it is here. So, we will run this command histogram to log medical expenditure that is that has been already taken.

(Refer Slide Time: 44:10)



So, the histogram will be displayed on your screen.

(Refer Slide Time: 44:11)



Now, in this one, you can easily see what the trend line looks like. The line we can also draw through the next command is the 1. So, histogram, then we will specify it.

(Refer Slide Time: 44:37)

	. tab Religio	ie.												A TRUNCH	
r commands have														3 mi	
mand r	Religion	Freq.	Percent	Cum.										Name	Label
Cillian/Stufe #Des		11.574	11.04	13.44										minhuned	case
uibe Sector	2	7,321	16.94	90.00											
ina medical ma	3	2,495	5.77	95.77											
Sector	4	1,177	2.72	50.49											
imes_hespitalised	5	173	0.40	90.89											
es hespitalised	6	235	0.69	99.57											
	1	10	0.03	99.40											
Religion			0.40	100.00											
eligion extor Religion	Total	43,219	100.00												
n log medical copirei 111	. tab Sector	Religion													
in log medical exp. it					Religion										
gram log_medical_exp	Sector	1	2	3	+	5	6	7	1	Total					
		14.340		1.474	104		1.0.0								
	EUL+1	14, 341	3,438	1,476	194		190	1	114	22,342					
	1175.45	15.233	3 633	1.019	10.1	165	105	12	52	20.057					
	urben	15,233	3,093	1,019	383	165	105	12	57	20,857					
	Totel	15,233	3,083	2,495	383	165	295	12	171	20,857					
	Total	15,233 31,574 [medical_exp	3,083 7,321 reisbureet_3	2,495 2,495 39_866_10400	1,177	163	295	12	171	43,219				Properties	
	Totel . scatter log variable reis	15,233 31,574 [medical_exp_ shoreed_by_med]	3,083 7,321 reinbursed_3	2,495 2,495 by_med_tasts found	383	143	295	12	171	43,239				Properties	
	Total Total . scatter log variable reis r(11));	15,233 31,574 [medical_exp shoreed_by_med]	3,083 7,321 reinbursed_3	2,495 2,495 ty_med_inate found	1,177	145	295	12	171	20,857				Properties +	
	Total Total . scatter log variable reis r(111); . scatter log	15,233 31,574 [medical_exp dorsed_by_med	3,083 7,321 reinbursed_3 _insur_not_1 reinbursed_3	2,495 2,495 ty_med_insta found	1,177	163	295	12	171	20,057				Properties	rembused
	Total Total . scatter log variable reis r(11); . scatter log	15,233 31,574 [medical_exp protect_by_med [medical_exp	3,003 7,321 residureed_3 _ineur_not_1 residureed_3	1,019 2,495 by_med_instations found by_medical_1	383 1,177 0.655	143	235	12	171	20,057				Properties + v = Name Label	rembursed
	Total Total - scatter log variable reis r(11); - scatter log - pwcorr log	15,233 31,574 [medical_exp - [medical_exp - medical_exp -	3,003 7,321 resolution 3 (insur cot) resolutioned 3 resolutioned 3	2,495 2,495 ty_med_instant tound ty_medical_t	383 1,377 	165 173	295	12	171	20,857				Properties • + + • Variables Name Latel 3.94	rembused long
	Total . scatter log variable reis r(11); . scatter log . pwcorr log	15,233 31,574 [medical_exp: pedical_exp: 	3,003 7,321 retabursed J insur not 1 retabursed J retabursed J	1,019 2,495 sy_med_instat iound sy_medical_1 sy_medical_1	383 1,377 neze neze, star(0	145 173	295	12	171	20, 857				Properties • • • • • • • • • • • • • • • •	rembursed long %10.0g
	Total Total . enatter log variable reis r(11); . enatter log . pwcorr log	15,233 31,574 [medical_exp: medical_exp: medical_exp: medical_exp: log_me-p.re	3,003 7,321 retabursed 3 inter not 1 retabursed 3 retabursed 3 into-r	1,019 2,495 sy_med_instat iound sy_medical_1 sy_medical_1	383 1,377 	145 173	295	12	57	20, 857				Properties • • • • • • • •	reinbursed long 5/10.5g
	Total . scatter log variable reim r(111): . scatter log . pwcorr log	15,233 31,374 [medical_exp : medical_exp : medical_exp : log_me-p_re	3,003 7,321 retabursed 3 [insur not 1 retabursed 3 retabursed 3 insu-r	1,019 2,495 ty_med_loates tound ty_medical_1	383 1,377 	145 173	295	12	57	20,857				Properties • • • • Varlables Name Label Tomat Varba label Nature Nature Nature 1 about Nature 1 about 1 abo	hambursel long %10.5g
	Total . scatter log variable reis r(11); . scatter log . pwcorr log log_medica-p	15,233 31,574 [medical_exp medical_exp = medical_exp = log_me-p_re 1.0000	3,003 7,321 reisbureed_i ineur_not_i reisbureed_i reisbureed_i	1,019 2,495 by med_inser found by medical_1 by medical_1	383 1,177 netr netr, star(5	145 173	295	12	57	20,857				Properties Properties Farme Kane Label Spe Formut Label Notes Dute Label Notes Dute Label Notes Notes Notes Notes Notes Notes Notes Notes Notes Note	hembursed long Michig
	UIDen Total - scatter log variable reis r(11)): - scatter log - pwcorr log log_medica-p	15,233 31,374 medical_exp - medical_exp - medical_exp - log_me-p_re 1.0000	3,803 7,321 resubureed 3 (near not 1 resubureed) resubureed) inbure	1,019 2,495 Sy_med_insectiond Sy_medical_1 Sy_medical_1	383 1,377 neter neter, star(0	145 173	295	12	37	20,857				Properties	Hamburga Jong SVDDg Heckl./, J
	UIDEN Totel . scatter log variable reim r(111); . scatter log . pwcorr log log_medica-p reimbursed-r	15,233 31,574 gmedical_exp medical_exp medical_exp log_me-p_re 1.0000 0.1110*	3,883 7,321 reisbureed_i [insur not i reisbureed_i isburr 1.5000	1,019 2,495 Sy_med_instant cond y_medical_i	383 1,377 	145 173	235	12	171	20,857				Properties	hordwood long 5005g block1,7,1
	Total . enatter log variable rele r(11); . enatter log . pwcorr log log_medics-p relabursed-r	15,233 31,574 gedical_exp gedical_exp gedical_exp 1.000 0.1110 0.000	3,833 7,321 retabureed_1 [team not 1 retabureed_1 retabureed_1 inbu-r 1.0000	1,019 2,495 ny_med_inser found ny_medical_1 ny_medical_1	383 1,377	145 173	295	12	57	20,857				Propertus	nomburied long %200g block1,2,3,
	utben Total . seatter log variable ress r(111); . seatter log . pwcorr log log_medica-p reimbursed-r	15,233 31,574 medical_exp_ medical_exp_ medical_exp_ 1.0000 0.1110* 0.0000	3,883 7,321 reisbureed_l interr not 1 reisbureed_l isbure 1.0000	1,019 2,495 by_med_insus found y_medical_1 y_medical_1	383 1,377 nemer nemer, star(0	145 173	295	12	37	20, 857				Properties * Modules Nores Liber Topit Forme Liber Nores Nores Liber Nores	rembused long 3203g Biocl1,2,3, 135 41340
	utben Total . scatter log verskle fels fill) . scatter log . perorr log log_medica-p reimbursed-r	15,233 31,574 [medical_wap medical_wap medical_wap 1.0008 0.1118* 0.0009	3,883 7,321 reisbursed_ [indur not 1 reisbursed_] reisbursed_1 1.0000	1,019 2,495 by medical_insur- found by_medical_i y_medical_i	383 1,377 newr, star(0	145 173	295	12	37	20, 857				Posperies W = = = Valables Latel Vors	Homburstel Jong Nobij Bieck1, Z, 1 136 41, 20 13.934
	utben Total - eastter log variable relative (illi) - preorr log netabursed-r - histogram l	15,233 31,574 gendical_wap periodical_wap gendical_wap 100_mer-pice 1.0000 0.1116* 0.0000 0.1116* 0.0000	3,883 7,321 restabureed_3 issurr not 1 restabureed_3 issu-r 1.8000 1.8000	1,019 2,495 by_med_insec loand yy_medical_1	383 1,377 neter neter, star(0	145 173	295	13	171	20,857				Popertas i + + Vadates Nore Liter Son Format Note Liter Note Liter Note Liter Note Son Note Son Liter Note Son Liter Note Son Note	10mborad long 5.05g block1,2,3 13 13 13 13 13 13 13 13 13 13 13 13 13
	utben Total - seatter log variable reis r(11): - seatter log - preorr log reisbursed-r - histogram 1 (blo45, stal	15,233 31,374 guedical_axp guedical_exp medical_exp 1.0000 0.1110+ 0.0000 0.1110+ 0.0000 0.1110+ 0.0000	3,883 7,321 retibureed_1 _insur not i retibureed_1 insured_1 insured_1 1.5000 p \$160007)	1,019 2,495 by_med_insco found by_medical_1 by_medical_1	383 1,177 netr, star(0	145 173	285	13	171	20, 857				Properties •	rembursed horg 5.00g block1,2,1 136 43,20 13.90M bM

So, it is here.

(Refer Slide Time: 44:39)



And it will specify the graph correctly. The graph now will you can identify whether the graph is approximating a normal distribution or not. This seems that this is a kind of leptokurtic diagram and then more it is symmetric though not mesokurtic. So, it is leptokurtic, but it has normality assumptions fulfilled. So, then coming to the next one this is what we have already explained.

(Refer Slide Time: 45:15)



So, this is what we have derived.

(Refer Slide Time: 45:17)



Then another one is called kernel density plot the through that you can also check the normality of the distribution. This is an alternative way to check the plots. [FL] this also gives approximate values of the probability density function of the particular variable.

So, here we need to specify kdensity the variable name then normal the word we need to I mention. Kernel density plots have the advantage of being smooth and of being independent of the choice of origin like histograms.

(Refer Slide Time: 45:57)



So, this will look like this. I am not going to draw it here once again. So, you can see how the normal distribution plot is derived kernel distribution estimation is highlighted with the blue line. This is almost overlapping with the normal density function. Here you can easily say that yes your data is more normal.

(Refer Slide Time: 46:20)

Model Estimation	
Beginning with simple linear regression in who one predictor variable.	nich we only have
□ In Stata, the dependent variable is listed in the regress command followed by one or variables.	mmediately after more predictor
regress log_medical_exp reimbursed_by_medica	al_insur
swayam 🧕	

The next one is about model estimation beginning with sample simple linear regression in which we only have one predicted variable. In Stata the dependent variable is listed immediately after the regress command, regress is the command we give. So, the dependent variable will quickly follow the regress command.

Then one or more predictor variables are specified, is like this. This is one, this is our dependent variable; then all sorts of independent variables could be mentioned alright. So, this is what is going to give us the estimation.

(Refer Slide Time: 47:01)



Therefore, first of all, we will copy this then if I remember this variable we can simply click that log medical expenditure then reimbursement then medical insurance etcetera. Therefore, will straight away go to Stata we will type regress.

(Refer Slide Time: 47:17)

🜉 StatarSI 15.1 - C.Usevidhadio Aldreatri	opiJanuary 2022/Dx. P	iatap Mohanty'Lect	turežt, regression	, model on heat	care variables/,bl	w0,0,0,0	i(aðit).ita										- 0 ×
The Edit Data Graphics Statistic	the Window	Help															
	0.0.0																
$\label{eq:response} \begin{array}{ c c c c c c c c c c c c c c c c c c c$	1 2 3 4 5 6 7 9 Total . tab Sector rural urban Total	31, 574 7, 321 2, 485 1, 137 123 285 13 131 43, 219 8e11gten 1 16, 341 15, 233 31, 574	73.06 16.94 5.77 2.72 0.60 0.69 0.69 0.63 0.40 100.00 2 2 3,438 3,083 7,321	73.66 90.00 95.37 99.49 99.99 99.57 99.60 100.00 100.00 3 1,476 1,019 2,495	Religio 4 764 303 1,177	n 5 8 145 173	6 190 105 295	7 1 12 13	5 114 57 171	Tot#1 22,362 20,857 43,239						A Variadars 3 m Name Televisured	T & X
	. scatter lo variable reis r(111): . scatter lo	g_medical_exp g_medical_exp	reinbursed	found by_medical_	r 10827											Properties a + -	ũ ×
	log_medica-p	_medical_exp log_me-p_tr 1.0000	reinbursed_	by_medical_	inser, star	(0.05) #1g					6	r				Name Label Tgar Format Value label	reminaned by medi hosp tracky
	relaburee5-r	0.1110* 0.0000	1.0000													Noter	Net17,133,196
	. histogram : ibin=45, sta . histogram : ibin=20, sta	log_medical_ec rt=0, width=.: log_medical_ec rt=0, width=.:	XD 29160007) XD. normal b 65611015)	5n(20)												Variables Observation Sax Memory Sorted by	135 H 41,240 12,300 64M
C Ubert Studio & Destrapilencery 2020 P type here to search	Command regress Log_m	edical_exp re tal5	inbaraed_by	medical_ins	er Januts		20 MILLION	n3,- 👔	PowePoint	State 8.	State 19.1 - 1	за. <u>Н</u> а	sph - Graph	•	23°C ^ ()	7 17 17 17 18 4 17	CLP NUM OIR 1453 29-11-2021

Therefore, then all the variables are listed here log medical expenditure we are taking at the dependent variable, then we will take reimbursement.

(Refer Slide Time: 47:32)

iter constrands have		 Variables
	0	3.0
Command or	rurel 16,341 3,438 1,476 794 0 180 1 114 22,369	Notes Label
ce Vilkeniškolo kiđer	urban 15,233 3,083 1,019 383 165 105 12 57 20,857	name (nee
ine Contratory Statement Witholds	Total 31.574 7.321 2.495 1.177 173 295 13 171 43.219	terrester
m log medical exp		
um Sector	. statter log medical_exp_relabored_by_med_instr	
an times, hespitalised	variable reimbursed by med insur not found	
ctimes,hospitalised	4(193)	
	. scatter log medical exp reimburged by medical insur	
um Religion		
ub Patigion	. pwcorr log_medical_exp_reimbursed_by_medical_insur, star(0.05) sig	
an sector mengron	The main statistics	
and regimental (op rol., 111	1 And the second s	
accertos metical pra st	log_medica-p 1.6009	
stoaram loa medical esa		
istegram log medical ex		
egres log medical exp rei	reisbursed-r 0.1114 1.0003	
	. histogram log_medical_exp	
	(b):#45, #tat(+0, v):01:#-29140807)	Properties
	(b)0.445, statted, statted, visto-zyskosz) . histogram log medical egy, normal bin(20)	Properties B = = =
	(14)**5, #44(*6, *45)*.79(4007) . http://pail.org.mol.ml.html/20 (14)**2, #44(*6, *45)*.5(11))	Properties B = - + Variables
	Lineda, mared, and Andre Jordan Lineda, mared, and and Shifti	Properties
	Lines, Annue, Annue, Annue, Annue, Annue Annues, Annue, Annue, Annue Annues Annues, Annue Annues, Annues Annues, Annues,	Properties
	Lines, Lines, Lines, Juno Juno J. Junigen Ling Junio J., Second J. Juno Lines, Lines, S. 1998 - State J. Junio J. Juno - report Lines J. 1998 - Marco J. 1998	Properties W = = = None is removed by total Type long format Staller
	Linets, Nations, States, 2010 States in States, 2010 Linets, Nations, States, 2011 States, 2011 Linets, Nations, States, 2011 States, 2011	Properties B = = = Voidate Norse rendure(by) Label La
	Interest, Spaces, Spaces, Zupper, Spaces, Space	Properties © + + + 2 Modules Later Later Grand Analy Variabil Note
	Linker, Statuter, Statu	Properties
	Ballets, Rations, States, State	Augustina a
	Linest, States, States, Junes, Jung,	Paperies B P Mode I and I
	Lines, States,	Paperins Paperins Paperins Rome embody, Latin ber Volay Volation Rome both J.K.N. Latin Rome both J.K.N.
	Linest, States, States, Junit, States, Junit, States, S	Paper Image: Control Product Product Product Product Star Image: Control Star Image: Control Star Image: Control Torret Image: Control Torret Image: Control Star Image: Control
	Lines, Antonio, Santonio, Antonio, 2010007 Lines, Nanova, Santonio, Santonio, Santonio, Santonio, Lines - respecto Lagunational, que constante da gundical, Lines - Santonio 1976 (2016) (2017) (2	Paperin
	100000 100000 100000 100000 100000 100000 100000 100000 100000 100000 100000 100000 100000 100000 100000 100000 100000 100000 100000 100000 1000000 1000000 1000000 1000000 1000000 10000000 10000000 10000000 10000000 10000000 100000000 100000000 100000000 10000000 100000000 1000000000000000000000000000000000000	Paperin
	Linker, Statuter, Statu	Paperine
	100000 100000 100000 100000 100000 100000 100000 100000 100000 100000 100000 100000 100000 100000 100000 100000 100000 100000 100000 100000 100000 100000 100000 100000 100000 100000 1000000 1000000 1000000 1000000 1000000 10000000 10000000 10000000 10000000 1000000000000000000000000000000000000	Appendix A Modeline A Modeli
	Lines, States,	Paperine Paperi
	Lines, Jacobs, Lines, Jacobs,	Paperine

Then related to reimbursement of insurance reimbursement for the medical expenditure. Now, on the screen you have got your regression result the ordinary least square regression model. This suggests that how the f coefficient beta value and beta naught and beta 1 the constant is in fact is your beta naught that is 6.01 9169 and your beta 1 is 0.0002226.

Most importantly we need to understand whether the data the coefficient is actually significantly estimating our model or not. Now, the P value are very effective P value if it is less than 0.05; that means, it is significant at the 5 percent level if it is less than 0.01; that means, it is significant at the 1 percent level.

Therefore, here it seems that your beta value is significant at even 1 percent level as well. Therefore, I am just explaining through the diagram through the PPT, here are the results of many interpretations you can do it. First, you need to check about the entire model and whether it is significant or not through the F test by default it considers the F test comparing the differences of residuals.

Therefore, the F test is defined, and the degrees of freedom is more important. So, F test values and its degrees of freedom are numerator and denominator are 1 and 14942, this is this given. And because a number of numbers in a number of observations is 14944 and accordingly others are understood. Now, the fitted line for log medical expenditure is in fact that we can derive, now at this moment let me just clarify its P values which I have just said is given here.

The model says that regarding whether it has a perfect fitted line or not R square value which we explained earlier R square and adjusted R square value is important. R square suggested that around 1.2 percent of the variance in medical expenditure is explained by the reimbursement by medical insurance. So, similarly adjusted R square where the degrees of freedom due to the parameter is changed. So, the number of our estimated parameters is changed.

So, accordingly, since you are only one parameter, we have taken one independent variable you have taken. So, a number of parameters is the same as for the entire model. The R square and adjusted R square are the same. It is not different When you increase the number of more than 1 independent variable your degrees of freedom will be different, and the adjusted R square value will be lesser than that of the R square value.

Adjuste R square value could be also negative which is all those things we already explained. Another piece of information in this chart around about linear or null in linear regression is on root means the sum of the square that is root means square error. root means square error is in fact the standard deviation of the regression, this is the closer to 0 better the fit.

And from all those things we can estimate a fitted line. So, fitted line for log medical expenditure is equal to antilog is 6.02 of the constant variable plus 0.00 0002 times the variable that is reimbursed by medical insurance. If I do this calculation I can get a fitted line, this is what your Excel or even your software usually follows to derive the fitted trend line.

(Refer Slide Time: 52:01)

reg squ	ress log_ uare_age	medical_	exp rei	imburse	ed_by_	medi	cal_insu	r	
•	Source	55	df	MS	Number of (F(2, 14939)	obs =	14,942		
	Model Residual	308.215609 22826.454	2 1 14,939 1	54.107804 .52797738	Prob > F R-squared Adj R-squa:	= = red =	0.0000		
	Total	23134.6697	14,941 1	.54840169	Root MSE		1.2361 -		
	lc	g_medical_exp	Coef	. Std. Er:	1. (T	₽> t	[95% Conf.	Interval]	
	reimbursed_by_	medical_insur Square_age _cons	.000221 .000020 5.97108	5 .0000163 2 5.22e-00 1 .0159963	3 13.59 5 3.87 1 373.28	0.000 0.000 0.000	.0001896 9.99e-06 5.939726	.0002535 .0000305 6.002435	

Similarly, we can do a multiple linear regression by adding another variable here we have added in addition to this insurance we have added now square age. So, with the same regress command, we have got the information, all those are a very important number of observations F values whether the model is significant or not R square adjusted R square root mean square we have already all explained.

Therefore, similarly residual sum of the square is total sum of the square. So, you can also estimate t values and then P values. So, and their coefficient everything is presented, so I have already explained I think I need not spend more time here.

(Refer Slide Time: 52:52)



The above result including only continuous variables as an explanatory variables we can also include indicator variables into the model with a little caution. If there are certain indicator variables or dummy variables the easiest way to include a set of dummies in regression is by using the prefix I.

If there are set of dummies is included or an indicator variables is included, the i dot command is more suggested that will be comparing your categorical variable from a base category. By default, the first category is considered the reference category. In that case, some conditioning should be given. So, putting xi should be there, xi command along with the regress must be there and that will give you better information related to the interpretation of the dummy variables.

(Refer Slide Time: 53:45)



Therefore, this is what the command looks like. In our next module, we will explain all those things in detail, but at this moment we are not drawing. I will interpret the result to some extent. If I give Xi along with regress I am doing it because I have some categorical variables like sector, so i dot is taken. Then another categorical variable is here social groups, education etcetera similarly here as well.

Therefore, this will compare its base category like if sector we have rural and urban; rural is our best category and urban is another category. So, the coefficient whatever we have received rural by default is the base category you can also change the base category your coefficient is going to interpret as, how urban is comparatively better off or incurring expenditure as compared to rural on the medical expenses.

Therefore, in that case, another interesting aspect is to control for heteroscedasticity in the model, you can also include a robust option. So, the robust option is going to control your heteroscedasticity, it will give one by standard errors of OLS coefficients under heteroscedasticity.

(Refer Slide Time: 55:03)

Standardized Coefficient	
Regression equation for unstandardized regression coefficients	
$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon_i$	
\succ Here β_1 , β_2 are unstandardized coefficient.	
Regression equation for standardized regression coefficients	
$\frac{Y - \overline{Y}}{SD(Y)} = \begin{pmatrix} \beta_1^* & \gamma_1 - \overline{X_1} \\ \beta_2 & D(X_1) \end{pmatrix} + \begin{pmatrix} \beta_2^* & \gamma_2 - \overline{X_2} \\ \beta_2 & D(X_2) \end{pmatrix}$	
> Here β_1^* , β_2^* are standardized coefficient.	
\succ There is no $oldsymbol{eta}_{0}$.	
$\succ \beta_1^* = \beta_1 \frac{SD(X_1)}{SD(Y)}$	
💿 swayan 🧕	39

So, there are some other details like standardized coefficient, unstandardized coefficient since we have already run for huge minutes in this lecture. We are not spending much time. I am just wrapping up all those things detail in our next episode maybe in the next year we will come up with detailed lecture of all those things.

At this moment I am just clarifying this concept called what is what you mean by standardized coefficient, what do you mean by unstandardized coefficient. Regression equation for standardized regression coefficients it is simply unstandardized regression coefficient is simply taking its variable as it is X1, X2 etcetera. So, beta 1, beta 2 are called unstandardized coefficients.

Regression equation for standardized regression coefficients if anything any other variables are standardized. Like if beta 1 is not just with it is X i rather X i is taken from their mean value divided by standard deviation. That means the X variable X i variable has been standardized. X 1 variable has been standardized here also X 2 variable has been standardized.

The coefficient that we derived that is beta 2 star and beta 1 star are actually called standardized coefficients. And there is no beta naught because if you standardize that since it is already a constant. And if you take its mean value from its constant value, that we mean and the value itself is the same. So, the difference is going to be 0, so there will be no question of beta naught.

Therefore, when we have a standardized coefficient. So, beta 1 star is equal to beta 1 times the standard deviation of X1 divided by the standard deviation of the dependent variable. So, this is very important so far as any quiz question is concerned you can just note and you will understand further.

(Refer Slide Time: 56:59)



To get the standardized coefficient include the beta option at the end of the command. So, the beta if you are just by adding comma then followed by beta, the beta coefficient is going to give you standardized coefficient value. The rest of the commands are the same.

(Refer Slide Time: 57:16)

Standa	rdized and Non standa	ardized Coefficient
	Unstandardized β	Standardized β
Definition	Obtained after running a regression model on variables measured in their original scales	Obtained after running a regression model on standardized variables (Z- score) (i.e. – mean =0, SD =1) Or all variables in same scale.
Interpretation	A change of 1 unit in the independent variable X is associated with a change of β units in the outcome of Y	A change of 1 SD in X is associated with a change of β SD in Y
Used for	Interpreting the individual effect of X on Y	Comparing the effects of different predictors X_i on the outcome of Y
Misleading when	We want to compare the importance of a variable X_i with other variables in the model (since these variable are in different scales)	Variables in the model have different standard deviations or follow different distributions
For binary variables	The coefficients has an intuitive interpretation	The coefficients is not interpretable
swayam 🧕		

I am not calculating step by step; we will certainly come up with these in the next episode at this moment we are just clarifying it. Standardize and non-standardized coefficient by definition by interpretation, how it is misleaded, how these are regarded in case of binary variables. All those things I am sure if you read you can able to enjoy and understand the difference between non-standardized beta and standardized beta.

(Refer Slide Time: 57:47)

	Some points to consider while fitting or estimating a model:
	Fit the maximal model-
	The maximal model is the large model that contains all the explanatory variables of interest as predictors.
	This model also contains all the covariates that might be of interest.
	It also contains all the interactions that might be of interest. Note the amount of variation explained.
	Begin simplifying the model-
	Inspect each of the terms in the "maximal" model with the goal of removing the predictor that is the least significant. Drop from the model the predictors that are the least significant.
0	swayam 🔮 42

I am not going going through them in detail these are all explaining the basic meaning I have already explained I am sure you will understand this. Some points to consider while fitting or estimating a model like feeding the maximal model beginning simplifying the model etcetera can be done like.

The maximal model is the large model that contains all the explanatory variables of interest as predictors. That is called a maximal model. This model also contains all the covariates that might be of interest. It also explains or it also contains all the interactions that might be of interest note the amount of variation explained is also captured correctly.

Then begin simplifying the model inspect each of the terms in maximal model with the goal of removing the predictor, that is this least significant one. Then drop those from the model and those are least significant and then further you run the regression.

(Refer Slide Time: 58:51)



Now, feed the reduce model compare the amount of variation explained by the reduced model with the amount of variation explained by the maximal model. So, the strategy to keep or draw variables is like predictor not significant and has the highest expected sign, in that case you can keep the command we have already set keep or drop command that you can show up or keep. And then further review and find out which one is in fact your better model.

(Refer Slide Time: 59:20)



Then you can actually also predict by predicting you have to give the variable command is predict then that variable with a hat is there. Now, how good the model is dependent on how

well it predicts the linearity of the model and the behavior of the residuals etcetera generating predicted values of log medical expenditure.

After running regression simply, you give predict at the command. So, basically whatever the result we have got that is in terms of the hat. By hat it should save it as by adding with a hat in addition to that name. Then predict r and resid that will give you the distribution of his residual terms.

(Refer Slide Time: 60:13)



Checking model assumption and in fitting how clearly these are going to give it.

(Refer Slide Time: 60:18)



(Refer Slide Time: 60:19)

Check for Normality inResiduals: Three graphs will help us check	
For normality in residuals:	Kernel density: kdensity r, normal
Kdensity, pnorm and qnorm pnorm r, title("normal check") qnorm r, title("normal check")	Kernel density estimate
) swayan 👰	47

(Refer Slide Time: 60:20)



(Refer Slide Time: 60:20)

	Check for mult The Stata comi inflation factor	icollineari nand to cl). Right aft	ty: heck for i er runnir	multi Ig the	collinearity is vif (variance e regression type:	
	. vif Variable	VIF	1/VIF)	A VIF > 10 or a 1/VIF < 0.10	
	new_edu square_age Sector SocialGroup nature_ail~t HouseholdS~e reimbursed~r kamal_trad	1.11 1.09 1.09 1.08 1.02 1.00 1.00	0.899521 0.914909 0.917311 0.919220 0.928342 0.981031 0.996283 0.997673	$\left\langle \right\rangle$		
0	Mean VIF	1.06)	/	49	

(Refer Slide Time: 60:21)



I am just explaining the overview of it. So, checking whether the model assumptions are fit or not. We have already started with some model called graph, if you go by graph two-way scatter plot of all those things will give you a graph. So, graph then the fitted line can also be derived.

That fitted line if it is perfectly going now with the distribution then your model is fine, here it seems your model is not perfectly fit. Similarly normality in residuals I have already told you one command called kdensity, but in this case pnorm and or and qnorm to be to be given to check whether they are following normal distribution or not.

So, then kdensity and normal if you do it will be plotting like this and accordingly, we can understand. Similarly, standardized normal probability plot that is pnorm checks for non-normality in the middle range of residuals. In that case pnorm r then you can give title of that particular a chart with an in with the bracket about the normality normal check, it will give you information.

Some other assumptions checked are through multicollinearity I know we will explain in our other models as well, but at this moment if you go by the statistics books called VIF Variance Inflation Factor, the mean VIF value should be no less than 10.

If it is greater than 10 then there will be multicollinearity. So, the mean VIF value is here if it is very less; that means, your model is perfectly fine. So, the VIF command is an important

VIF that will give you the result like this. I am not drawing it, we will draw it sometimes in our other modules. Similarly testing for homoscedasticity usually we go for the Breusch Pagan test.

So, that basically detects heteroscedasticity in the model the null hypothesis is that the results are homoscedastic whether; which means the having constant variance. So, that command is called estat hettes. estat hettest my previous module on large scale data handling we have also included this, if you want further reading you can also refer to that as well.

But otherwise, if you simply run with this command you will get this kind of information. Our assumption is that our residuals are homoscedastic there is no variation in the residuals the standard deviation is constant. But here it seems that your result is significant the assumption is constant variance; that means, it is deviating it is actually rejecting your assumption. So that means, there is heteroscedasticity.

So, this is what is interpreted, and the rest you can just read between the line and I am sure you will enjoy reading.

And this is a very big lecture and we have skipped some of this section deliberately and I am I know that, if you go through you will have lots of queries and we will be happy to deal with them in our live session class. So, with this, I think it is time to close and I look forward to your queries.

Thank you.