## Business Analytics & Data Mining Modeling Using R Dr. Gaurav Dixit Department of Management Studies Indian Institute of Technology, Roorkee

## Lecture – 23 Multiple Linear Regression-Parts II

Welcome to the course business analytics and data mining modeling using r. So, in the previous lecture we started our discussion on this particular technique multiple linear regression. So, we discussed the linear regression the equation, the different coefficient that we are required to estimates. We talked about the exploratory modeling and the predictive modeling a few differences and we also try to understand the application of this particular technique multiple linear regression, in a statistical technique and how it is different in a data mining environment. So, we also talked about some of the assumptions you know that when we apply OLS to estimate those coefficient those beta zeros and sigma.

Then what are the underlying assumption that we have to follow, and how those assumptions are different you know how in first assumption especially noise following normal distribution, how we get some relaction relaxation for from that assumption in a data mining setting. We talked about all those things then we will again go through an exercise to understand how linear regression modeling is done and how different concept can be put into practice.

(Refer Slide Time: 01:54)

8 Nudio		
Hie Hilt Code View Plots Sessen Kuid Debug Profile Loois Help		
💽 🔹 😴 🔹 📄 🔚 🔚 🥻 to to file/function 🔤 👘 Addins 🔹		Project: (None) •
● 6 regress.R ×	-0	Environment History
💠 🗇 🙍 🚍 🗖 Source on Save 🔍 🖉 📲 🔹	📑 Bun 💽 🕞 Source 🔹 🛎	🞯 🔒 🖃 Import Dataset 🔹 🥑 📃 List 🔹 🎯
<pre>1 library(x1sx) 2 # usedCars.x1sx 4 df=read.x1sx(file.choose(), 1, header = T) 5 df=df(, lapp)(is.ma(df), 2, all)) 6 head(df) 7 8 Age=2017-dfMfg_year 9 df=bind(df, Age) 10 10 dfb=df 12 df=df[,-c(1,2,3,11)] 13 14 str(df) 15 df5Transmission-as.factor(df5Transmission) 16 df5Transmission-as.factor(df5Transmission) 17 # meaningment (60% 40%)</pre>		Clobal I nuironment   Fries Plots Packages Help Viewer  Compared to Separate Compared to Sep
.k1 (lop i evel) 0	R Script \$	
Console Cytoursynaer/Dektog/MOOC Innuary 2014/DF Gaurae Obhr/Seadon 6/ PO R is a collaborative project with many contributors. Type 'contributors()' for more information and 'citation()' on how to cite R or R packages in publications. Type 'demo()' for some demos, 'help()' for on-lne help, or 'help, start()' for an HTML browser interface to help. Type 'q()' to quit R. > library(xlsx) Loading required package: rlava Loading required package: xlsxjars >		

- So, let us open r studio. So, as usual we will load this particular library x l s x because the data that the this particular data used cars data set, that we have it is in the excel file.
- So, let us import this data set. So, the function that we are going to use is read dot x l s first argument is as usual we are going to browse for this particular file and then the first worksheet will import the data of the first worksheet and the header is true because we have name of all the variables there in the data set. So, let us execute this line you can see in the environment section this particular data set has been imported and you can see 79 observation of eleven variables.

(Refer Slide Time: 02:54)

3	Greatess R #	dfx									_	Environment History -
5	() (n) V	Lilter							Q			🧭 🔲 🖉 Import Datacet y 🧃 📃 list y
	Brand *	Model	Mfg_Year	Fuel_type	SR_Price	KM ÷	Price	Transmission	Owners	Airbag	C_Price	Global I nvironment • Q
	Hyundai	Verna	2013	Petrol	8.88	/5.000	5.60	0	1	0	:	Data
	Mahindra	Quanto	2012	Diesel	6.99	49.292	3.95	0	1	0	-	• odf 79 obs. of 11 variables
	Maruti Suzuki	SX4	2011	Petrol	7.18	48.000	2.99	0	1	0		
	Chevrolet	Beat	2013	Petrol	4.92	41.000	2.35	₽ o	1	0		
	Honda	Civic	2008	Petrol	13.50	110.000	3.65	1	2	0		
5	Honda	Brio	2012	Petrol	5.74	60.000	2.99	0	1	0		
,	Hyundai	i20	2011	Diesel	8.42	\$6.000	3.8/	I	1	1		Files Plots Packages Help Viewer
	Skoda	Rapid	2014	Diesel	10.49	27.500	5.80	0	1	I		🕼 📣 🖉 Zoom 🍱 Export • 🍳 🎸
,	Mitsubishi	Pajero	2007	Diesel	19.50	101.000	5.50	0	1	1		
10	wing I to I0 of	/9 entries	00 January 20	118/Dr Caur	w Divit/Car	tion 61						
	<pre>itation()' pe 'demo()' alp.start() pe 'q()' to library(xls ading requi ading requi df=read.xls view(df)</pre>	on how to c for some c for an HT quit R. x) red package x(file.choc	cite R or demos, 'he TML browse e: rJava e: xlsxjar ose(), l,	R packag lp()' fo r interf s header =	es in p or on-li ace to : T)	ublicat ne help help.	ions.					

- Let us have a look at this data in the r environment. So, this is small icon that you see in the environment section once the data has been imported. So, once you click on this particular icon you would be able to see another tab that would open and you would be able to see the data just like you see it in the excel files. So, you can see the variable names brand models, manufacturing, your fuel types. So, there are 3 fuel types for these used cars petrol diesel CNG and then we have showroom price for each of these used cars. So, when these cars were what first time what; these right as the new cars. So, what was the price and then the kilometers since these cars have been running on road, the kilometers that have been accumulated from the starting years from the purchase year
- And then the price this particular price is the used the offered price for the these used cars. The cars whether the transmission is manual or automatic that is also we have information on. So, 0 representing the manual, and one representing the automatic transmission.
- Now, next variable is on owners where each number is representing the number of owners that actually owned the number of people number of individuals who actually owned this particular car. So, that is also there. You also have information on some of the security features for example airbags.

(Refer Slide Time: 04:31)

a										-	_	I subsection attenue
6 regress.k	K dt x							(0)		=		Environment History
ଡ଼ା <u>ଯା</u> ଜଣା ଶା	Y Lilter	146- V	Fuel hard			-		Current	Alabañ	c mini	4	🐨 🖬 🔄 Import Dataset 🔹 🦉 📃 List 🔹
srand	Model	MIG_Tear	Fuer_type	SK_Price	KM	Price	Transmission	Owners	Airbag	C_Price		Gobal I mironment •
iyundal	verna	2013	Petrol	8.88	75.000	5.60	0	1	0			odf 79 obs of 11 variables
Aahindra	Quanto	2012	Diesel	6.99	49.292	3.95	0	1	0	0		
Maruti Suzuki	SX4	2011	Petrol	7.18	48.000	2.99	0	1	0	0		
hevrolet	Beat	2013	Petrol	4.92	41.000	2.35	0	1	0	0		
Ionda	Civic	2008	Petrol	13.50	110.000	3.65	1	2	0	0		
Ionda	Brio	2012	Petrol	5.74	60.000	2.99	0	1	0	0		
Iyundai	i20	2011	Diesel	8.42	56.000	3.8/	1	1	1	0		Files Plots Packages Help Viewer
skoda	Rapid	2014	Diesel	10.49	27.500	5.80	0	1	1	1		🖕 🧅 🔎 Zoom 🛛 🗷 Export 🔹 🍳 🏒
Mitsubishi	Pajero	2007	Diesel	19.50	101.000	5.50	0	1	1	1		
howing I to II	0 of /9 entries			-	_			_		Q 17		
citation() ype 'demo help.star ype 'q()' library() oading rec oading rec oading rec	)' on how to ()' for some t()' for an to quit R. xlsx) quired packa quired packa xlsx(file.ch	e cite R e demos, HTML bron ge: rJav. ge: xlsx, boose(),	or R pack 'help()' wser inte a jars 1, heade	for on- erface f	n public -line he to help	ation	ns. or				•	

So, number of airbags that are there in the car. So, that that information is also available.

- We have another variable in the data set that is c prices, that is this is this variable what is generally present for the classification task where any car having a less than offered value of a 4 lakhs is represented by 0 and the cars having value of equal to or more than 4 lakhs are represented by 1 So, this is the data set. So, let us close this particular tab and let us. So, first thing would be as usual we would like to remove the n a columns.
- So, we used these within the brackets and the for the column value. So, we have applied this first particular function. So, apply what is going it is going to do is it will find out the using is dot n a function will find out which of the columns in this particular data frame d m d f they are having any values right. So, those particular columns would be selected. So, 2 is indicating that this particular function is being applied or you know column wise on this particular data frame.
- So, therefore, those particular columns n a columns would be selected and the all function would then be applied on those columns. So, so other columns which have which do not have any values, they would be true and then the columns which I have n a values which will have false values. A logical operator logical vector would be written from this function and then we have this another operator not that would apply on this logical vector logical result, and then all the true values would be converted to false and all the false values would be converted into true. So, therefore, all the n a columns which were identified using the apply function and they were indicated as false now they would become.

Now,. So, the reverse will happen. So, the all the columns which do not have any values they would be returned as false using apply function, and when not is applied then they would become true, and the other columns they would be returned as true because they were any columns and when not is applied they will become false to those columns would actually be dropped. So, this is how this particular line we have been using quite often. So, this is how it will operate. So, in this particular data set we did not have any says column. So, the result would remain same.

(Refer Slide Time: 07:40)

Studio		Contract of the local division of the local			- 6 - x
Hie Hat Lode View Plots Sesson Kuld Libebug Profile Loois Help	Adding *			Ť.	Duries I: (Name)
€] 6 regress.R ≈				Environment History	
🗇 🔿 🙍 🔒 🗖 Source on Save 🛛 🧕 🖉 📲	•	📑 Run 📑 📑 Source	- 2	🞯 🔒 📑 Import Dataset 🔹 🍕	≣ list • 🞯
1 library(xlsx)				🚳 Global I nvironment 🔹 🔍	
2 3 # UsedCars vlsv				Data	
<pre>4 df=read.xlsx(file.choose(), 1, hea 5 df=df[, !apply(is.na(df), 2, all)] 6 head(df) 7</pre>	der = T)			• df 79 obs. of 11 varia	bles 📃
<pre>% Age=2017-df\$Mfg_Year 9 df=cbind(df, Age) 10</pre>					
11 dfb=df					-
12 df=df[,-c(1,2,3,11)]				Files Plots Packages Help Viewer	-
14 str(df)				🔷 🧅 🔎 Zoom 🛛 🗷 Export 🔹 🍳	
15 df\$Transmission=as.factor(df\$Trans	mission)				
10 17 # Partitioning (60%:40%)					
81 (lop i evel) \$			R Script \$		
Canada Cultura (Inclusion (Inclusion Current) 2018/00-Current	- Diductation (1. (2)		- 1		
control courter package. Alaxing 2016/01. data	V Dixity session of the				
<pre>&gt; df=read.xlsx(file.choose(), 1, header =</pre>	т)				
<pre>&gt; df=df[, !app]y(is.na(df), 2, all)]</pre>					
> head(df)					
Brand Model Mfg_Year Fuel_type	SR_Price KM Price	Transmission Owner	rs		
1 Hyundai Verna 2013 Petro 2 Mahindra Quanto 2012 Diece	6.88 /5.000 5.60	0	1		
3 Maruti Suzuki SX4 2011 Petrol	7.18 48.000 2.99	0	1 =		
4 Chevrolet Beat 2013 Petrol	4.92 41.000 2.35	0	1		
5 Honda Civic 2008 Petro	13.50 110.000 3.65	1	2		
6 Honda Brio 2012 Petro	5.74 60.000 2.99	0	1		
All bag C_FI ICC			•	L	

- Now, let us look at the first 6 observation of this particular dataset, you can see the name these variables as we saw through other options in r studio by clicking this particular icon, and looking at the full file in one go. So, using head function we can look at the 6 observation we need not look at all the observation, that would require the whole file to be loaded into memory and therefore, if your device where you are running your r studio and if it does not have sufficient ram built into it then probably you would are you are better of running head function, and you would just be loading 6 observations into memory.
- So, these are the observations the variables we have already discussed. So, one particular third column that we can see is manufacturing here, when the car was actually manufactured. So, the age of car can actually be computed using this particular vector. So, this is what we are going to do next. So, you can see age variable and if the if we if these offered prices if all the information is in the context of year 2017. So, therefore, the current here is 2017. So, we can subtract all this particular vector from 2017 and all the values all the observations for all the observations will get the difference and therefore, the age.

(Refer Slide Time: 09:19)



- So, let us execute this line you can see in the environment section is variable has been created this numeric vector, having 79 values right. So, age of all the used cars have been computed. Now age could be a relevant variable because as we discussed the task here is the prediction task we are trying to we would be trying to build a model, trying to we will build a model to predict the price of a used car right offered the price of a used car. So, age could be an important variable in terms of explaining in terms of predicting that particular price. So, therefore, you would like to have age also in our model.
- So, let us append this particular variable this particular vector and this data frame. So, c bind is the function as we have talked about in previous lectures also, this can be used to an already existing you know a data frame to append this particular variable. By default it would append this variable at the end of all the all other columns instead data frame. So, let us execute this line. So, now, this particular variable has been appended. If you want to check again you can run the head you can call the head function again and you would see the last column age has been created and you can also see the values for first 6 observations.
- Now, if we look at this particular data set, then brand name first 2 columns brand name model and also third one also manufacturing here they do not seem to be relevant for our analysis for our manufacturing here we have already we transformed this particular variable into an age variable, age of the used car. So, therefore, we would not be requiring this particular variable right. So, we will get rid of this variable also since we would be building a prediction model. So, therefore, this c underscore price which was the outcome variable mainly designed for the classification task, we also would not be requiring this variable.

So, we can get rid of these four variables; first one brand then model than manufacturing year and then this one c price, and the remaining variables would be the outcome variable that is price or the relevant predictors, that we want to include in our model. So, first let us take a backup of the existing data frame.

(Refer Slide Time: 12:01)

3 Nourio	- 6 <del>- </del> 8
Hie Hitt Code Wew Plots Session Ruid Debug Photoe Tools Help	
💽 • 🚰 • 🔒 🔒 🥻 Addins •	Project: (None) •
0 6 regressik x	Environment History
	n • E 🔿 🕞 🕞 Import Dataset • 🖌 📃 List • 🚱
1 library(xlsx)	Gohal Invironment *
2	
3 # UsedCars.xlsx	odf 70 she of 12 uppishlas
5 df=df[, !apply(is.na(df), 2, all)]	Odfh 79 obs. of 12 variables
6 head(df)	Values
7	Aco num [1:70] 4 5 6 4 0 5 6 2
8 Age=201/-dT3MTg_Year 9 df=chind(df_4ge)	Age 110111 [1.79] 4 5 0 4 5 5 0 5
10	
11 dfb=df	
12 df=df[,-c(1,2,3,11)]	Files Plots Packages Help Viewer
13 14 str(df)	🧼 🧅 🔎 Zoom 🛛 🗷 Export 🔹 🍳 🥑
<pre>15 df\$Transmission=as.factor(df\$Transmission)</pre>	
16	
17 # Partitioning (60%:40%) 121 (lop(evel) =	R Script \$
Console C:/Users/user/Desktop/MOOC January 2018/Dr. Gaurav Dixit/Session 6/	
4 Chevrolet Beat 2013 Petrol 4.92 41.000 2.35 0	1 '
5 Honda Civic 2008 Petrol 13.50 110.000 3.65 1	2
Airbag C_Price Age	1
1 0 1 4	
2 0 0 5	
5 0 0 9	
6 0 0 5	
> dfb=df	
>	-

- So let us take a backup you can see this particular data frame has been created now let us eliminate these columns. So, combined function can be used and minus before combined function indicating that we do not want to subset, that we want is other columns that are mentioned here in the combined function.
- So, if you want to have a look at now if we are interested, we can have a look at the structure of the data frame, now in the present data frame this data frame that we have is we have all the variables of interest.

(Refer Slide Time: 12:44)



- You can see now we have 79 observation of 8 variables. So, all the variables are now of interest to us whatever modeling exercise. So, fuel type right that would also help us determine the price of offered price of a used car because the CNG whether the car is running on a CNG fuel or diesel or petrol.
- So, these cars are when they are purchased for the first time, they are priced differently and therefore, depreciation and other factors they work differently. Therefore, prices of these cars based on the fuel type that they have are going to be different therefore, fuel type being an important predictor for us in this prediction modeling exercise. Now as s r price is actually the show room price. So, this is equivalent of similar to the first time this particular and that particular car was purchased what was the price that was paid right this is the price over the years, the depreciation would be applied and the depending on the condition of the car and other variables some of them some of those variables, we have in this data set as well the offered price would be adjusted.
- Now, how this offered price is being determined by the individual is part of this exercise of this model. Now the another important variable that we have is k m. So, the number of kilometers that a car has accumulated that also tells about the wear and tear that the car might have gone through, because of the those number of kilometers covered right. So, therefore, kilometer might also indicate the value of a car that is you know that should be depreciated that should become part of the depreciation. If the car has been raven less then probably it has not gone through that much of wear and tear.

- But if it has traveled more than probably you know more wear and tear might have happened and therefore, the price might be on the lower side therefore, kilometer is also an important variable in this exercise. Price is the variable the offered price that is the variable that we are trying to predict this is the outcome variable of interest to us, next variable that we have is transmission. So, transmission right now as indicated in this output this is right now numerical vector as shown here, but this variable can have only 2 values or let us say 2 labels, because this is a categorical variable because we have just 2 labels whether the car is automatic or the car is manual.
- But here in this data set it is being shown in it is a numeric vector here zeros and ones. So, we would like to convert this vector from numeric to categorical variable to factor variable. So, therefore, you can see in the next line we have used as dot factor function to coerce this particular numeric variable into a factor variable in r environment. Because this variable has 2 labels 2 categories that is the car is automatic or manual.
- So, let us execute this particular code and once it is done you can again run this structure function and you would see a change there the transmission variable, now you can see it has been converted into a factor variable with 2 levels 0 and 1.
- So, now the values 0 and 1 they are being treated appropriately because this being a categorical variable. So, now, the these values have become numeric codes indicating 2 different labels now in the structure function you would also see that the some few initial observation that are being shown here they are in this format one category 1, category 2 1 and 1 and then 2 and then 1, but the actual values at we had saw before in the using the head function they are either 0 and 1.
- So, this is just the way output of a structure function is presented the values have not changed because of that conversion that we just did if you want if you are interested.

(Refer Slide Time: 17:36)

Man Lande Ware 1995 Second Reid Debug 1995 August Helle	-0
• Control field states and states	👔 Project: (Ni
	() Hyper (in
6 regressit x	Environment History
Solar G Source on Save Q Z + E +	🔹 🧟 😭 🕞 import Dataset 🔹 🖌 📃 List 🔹
<pre>4 dT=read.xisx(T1ie.cnoose(), 1, neader = i)</pre>	· Gobal Invironment • Q
5 df=df[, !apply(is.na(df), 2, all)]	
b head(dT)	E Data
8 Ace-2017-dfSMfo year	Odf 79 obs. of 8 variables
9 df=cbind(df, Ane)	Odfb 79 obs. of 12 variables
10	Values
11 dfb=df	Age num [1:79] 4 5 6 4 9 5 6 3
12 df=df[,-c(1,2,3,11)]	
13	
14 str(df)	
1) df%Transmission=as.factor(df%Transmission)	Files Plots Packages Help Viewer
17 # Partitioning (60%-40%)	A A D Zroup J Evoud + O
18 partidy=sample(1:prow(df) 0.6*prow(df) replace = E)	
19 dftrain=df[partidx.]	
20 dftest=df[-partidx,]	
15:1 (lop Level) \$	R Script ©
	-0
Concerner in the second s	
S Airbag : num 0000001111	
S Age : num 4 5 6 4 9 5 6 3 10 4	
head(df)	
Fuel_type SR_Price KM Price Transmission Owners Airbag Age	
Petrol 8.88 75.000 5.60 0 1 0 4	
Diesel 6.99 49.292 3.95 0 1 0 5	
Petrol 7.18 48.000 2.99 0 1 0 6	
Petrol 4.92 41.000 2.35 0 1 0 4	
Petrol 4.92 41.000 2.35 0 1 0 4 Petrol 13.50 110.000 3.65 1 2 0 9	
Petrol         4.92         41.000         2.35         0         1         0         4           Petrol         13.50         11.000         3.65         1         2         0         9           Petrol         5.74         60.000         2.99         0         1         0         5	

- You can again have a look at the actual values you can see the transmission variable is still having the same 0 and values 0 and 1 value it is only the structural function that is presenting the output in that fashion, that for factor variables if for fact for factor variables it generally shows 1 2 and 3 for different categories, different labels that we might have. You can see even for the fuel type variable CNG diesel and petrol these are the 3 labels that we had the in the structure output we see 3 2 3
- So, these representing the 3 classes 1 2 and 3, but actual values are same they are not disturbed. So, you can see fuel type you can see it the these strings these characters patrol diesel and therefore, those values are not changed it is just the representation in the output of the structure function.

(Refer Slide Time: 18:40)

Kludio	the second se				5	. 0 X
Hie Hait Ce	de View Plets Session Ruild Debug Poehle Tonis Help					
0.6	🔭 🖶 🗐   🚔   🍌 (a to file/function 🔤   🔯 📲 Addins 🔹				🕓 Proje	t (None)
0 6 rec	peccil at	-0	Environment	History		-0
A		-	ania		4 =	- R
4	dT=read, XISX(T) IE, Choose(), 1, header = 1)	Source • -		Import Maser *	2 = 0	1.0
5	df=df[, !apply(is.na(df), 2, all)]		Global Lrw	ironment *	(u	
6	head(df)	I	Data			
7			0 df	79 obs. of	8 variables	
8	Age=201/-dt3Mtg_Year		0 dfb	79 obs. of	12 variables	
10	ur-connu(ur, Age)		values			
11	dfb=df		Ane	num [1:79]	456495	6.3
12	df=df[,-c(1,2,3,11)]		Age	from (2175)		• • •
13						
14	str(df)		The second se			
15	df%Transmission=as.factor(df%Transmission)		Files Plots	Packages Help	Viewer	-
17	# Partitioning (60%-40%)			Zourn Every	. 0 4	
18	partidx=sample(1:nrow(df), 0.6*nrow(df), replace = F)		W WIC	200mm Paper		
19	dftrain=df[partidx.]					
20	dftest=df[-partidx,]					
15:1	(inplevel) 0	R Script \$				
Console	C:/Users/user/Desktop/MOOC January 2018/Dr. Gaurav Dixit/Session 6/					
'data.	frame': 79 obs. of 8 variables:	•				
\$ Fue	el_type : Factor w/ 3 levels "CNG", "Diesel",: 3 2 3 3 3 2 2 2 2					
S SK	Price : num 8.88 0.99 /.18 4.92 13.5					
S RM	num 5 6 2 05 2 00 2 25 2 65 2 00 2 87 5 8 5 5 2 1					
\$ Tra	insmission: Factor w/ 2 levels "0" "1": 1 1 1 1 2 1 2 1 1 2					
\$ Own	vers : num 1111211111					
\$ Air	bag : num 00000011111					
\$ Age	: num 4 5 6 4 9 5 6 3 10 4					
> head	l(df)					
Fuel	_type SR_Price KM Price Transmission Owners Airbag Age					
1 F	Petrol 8.88 75.000 5.60 0 1 0 4					

- Now, another important thing that we need to understand here is, that the factor variables that we see the are nominal variables the way they have been created these are nominal variables, they are not ordinal and the way labeling is done here.
- For example if we check the be run d f of function for the first variable that is fuel type, you would see these values first 6 values petrol diesel, petrol petrol and petrol all and then labels CNG diesel and petrol these labeling in r environment is done alphabetically is therefore, CNG because it start with c which is and then diesel it is start with d and petrol p. So, the alphabetical in alphabetical fashion the ordering of labels is in that fashion, but when we do when we run a classification task, we will have to decide our reference category. So, in that case if we happen to select our reference category as petrol that would not be by default made as the reference category here in this case.
- We will have to relabeled this variable. So, some this kind of exercise we will do when we you know discuss a particular technique that is suitable for classification or used for classification. So, next thing that would be required, once we have done on all very you know variable transformation, we have checked the variable type appropriately transform them, now all the variables are ready and we are ready for the modeling exercise. So, before we go ahead we need to partition on this particular sample. So, in this particular exercise will partition this particular dataset into 60 percent for the training set and the 40 percent for the test set.
- So, we would not be having validation set we would be building our model on the training set and then we would be testing our model on this test partition. So, sample is the function that could be used to perform this partitioning. So, in the first argument as we have discussed before, we need to specify in a numeric vector form number of observations. So, in this case 1 2 number of rows that are there in this particular data frame representing the number of observations that are there in the data frame. So, that being indicated then the size of the sample, that is being indicated by this 0.6 into the length of the data frame.
- So, because we want 60 percent of the observation 2 randomly selected observations to go into our training partition. So, therefore, 0.6 into this the full size, and the replacement is false because we want to do our sampling without replacement which is the typical way of sampling in a modeling exercise.
- So, this particular sample function would return as the indexes indices of those observation which have been randomly drawn; now to further partition once. So, let us execute this line. So, you would see that part i d x this index variable has been created this is integer because these are the indices.

(Refer Slide Time: 22:09)



- That have been returned by the sample function and you would see that 60 percent of the other observation, the indices of 60 percent of the observation randomly drawn from that particular data set have been written.
- Now, we need to partition the data set. So, we can do using these brackets functions, we can subset these particular observations from the full data set. So, in the rows value in the for the row value within the brackets, we can mention this particular variable and all those indices would then be selected subsetted for training partitions. So, let us execute this code. So, we have been able to create d f train, you can see forty seven observations of 8 variables, which is same as the part i d x which was also 47 have which had 47 indices in the first place. Now since we are getting just 2 partitions therefore, all the remaining observation can go through can actually be left for the test partition.
- So, in within the brackets in the for the row value we can mention minus part i d x. So, the remaining indices they would be subsetted for the test partition d f test. So, let us execute this code.
- Now, once the partitioning exercise is over, now you can see in the environment section d f test another partition has been created having 32 observations of 8 variables. Now once this partitioning exercises has been done the function we can move to our linear modeling linear regression modeling. So, the function that is available in r for the to perform this modeling is 1 m. If you are interested in finding more details about this further function you can go to the help section and type 1 m and enter and you would see that in the help section it talks about this particular function, 1 m is about fitting linear models.
- So, you can see in the description that it talks about that l m is used to fit linear models; it can be used to carry out regression, single stratum analysis of variance analysis of covariance right. So, all those statistical tests can be performed using this function.

(Refer Slide Time: 24:33)



- Now, if you look at the usage, you can look at the function and the arguments that can be fast first one is a formula, formula that is going to represent the linear regression model right. So, we need to pass this particular formula and then the second important variable is data. So, there are many other arguments also. So, you can on your own time you can go through some of these variables which are not typically used.
- So, the formula is written in this particular format. So, the output outcome variable of interest it is written first and then we use the tilde operator and then we can write the all the names of the predictors that we have in our data set, and that and that we want the variables predictors that know which we want to include in our model. Or we can simply type dot if we want to include all the variables that are that are available in the data set. So, if you remember then the data frame that we are using now the that we partitioned, we had already excluded the variables that it we did not want in the first place and therefore, all the remaining variables are the come into our set of predictors and we would like to have all of them in our model.
- So, our formula is going to be priced tilde dot; dot indicating to this function that all the all other variables should be part of the set of predictors. Now the data set that we are using is d f train that is the training partition. So, we would be building this model on training partition.

(Refer Slide Time: 26:14)



- So, let us execute this line and you would see that in the environment section a mod variable has been created if you are not so, it has it is this this particular variable is actually a list of 3. So, it has information on thirteen elements, if you are interested in finding out the all the names of these this particular list, you can the names come on.
- So, these are the thirteen elements that are there you can see. The first one is about coefficients, then second one is about residuals then effects, rank, fitted values, assignment you know similarly so many other details have been computed by this particular function.
- If you are understand finding out all these values, you can again go to the help section and find out and it is called down in this particular section, and you would see that there is going to be discussion under the value sub section. The kind of values that are returned by 1 m function and within this you will have details what are coefficients a named vector of coefficients of the coefficient that we have residuals fitted values.

(Refer Slide Time: 27:21)

We be to the two must long multi lo	Ktudio		
Imagesci #       Imagesci #         Imagesci #       I	Hdit Co	de View Plots Session Kuild Debug Plothie Loois Help	
Importing intervelop	•	🕐 🔒 🔒 🎽 🖍 Co to file/function	Project: (None) •
Image: Source wise in the source wise in the source wise intervention interventintervention intervention	6 reg	ressR x	Environment History
10       0		AT D T Source on Save Q Z + E +	🔿 🔲 🕐 Import Dataset 🔹 🖌 📃 List 🔹 🚱
11       dfb-df         2       dfb-df         3       dfb-df         13       dff-df         14       str(df)         15       dffTransmission-as.factor(dfSTransmission)         16       dffadf         17       # partitioning (GM: 40%)         18       partide:sample(linrow(df), 0.6*nrow(df), replace = F)         10       dftsst-df[-partix,]         21       odftsst-df[-partix,]         22       odftsst-df[-partix,]         23       sumary (mod)         24       anova(mod)         25       # List of 13         26       # Measures of Goodness of fit         27.11       (pre/rew) ±         26       # Measures of Goodness of fit         27.11       (pre/rew) ±         27.12       (pre/rew) ±         26       # Measures of Goodness of fit         27.11       (pre/rew) ±         27.12       (pre/rew) ±         28       (pre/rew) ±         29       (pre/rew) ±         20       (pre/rew) ±         21       (pre/rew) ±         21       (pre/rew) ±         22       (pre/rew) ±         23	TU		Clabal Laurenmant =
12       draft ,-(1,2,3,11)         13       str(df)         14       str(df)         15       dfirst 32 obs. of 8 variables         16       dfirst 32 obs. of 8 variables         17       # partitioning (60%:40%)         18       partitioning (60%:40%)         19       dfirst 32 obs. of 8 variables         10       dfirst 32 obs. of 8 variables         11       dfirst 32 obs. of 8 variables         12       dfirst 32 obs. of 8 variables         13       gfirst 32 obs. of 8 variables         14       strant, 47 obs. of 8 variables         15       gfirst 32 obs. of 8 variables         16       gfirst 32 obs. of 8 variables         11       gfirst 32 obs. of 8 variables         12       gfirst 32 obs. of 8 variables         13       gfirst 32 obs. of 8 variables         14       title 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,	11	dfb=df	dth /9 obs of 12 variables
13       str(df)         14       str(df)         15       str(df)         16       dfstransmission-as.factor(dfstransmission)         16       # Partitioning (60%:40%)         18       partidax sample(1:nrow(df), 0.6*nrow(df), replace = F)         19       dftrain df (partidx, ]         20       odd=lm(Price ~ ., dftrain)         21       mod-lm(Price ~ ., dftrain)         22       anova(mod)         23       summary(mod)         24       anova(mod)         25       # Measures of Coodness of fit         711       (por iew) :         721       (por iew) :         722       (por iew) :         723       (por iew) :         724	12	dT=dT[,-C(1,2,3,11)]	Odftert 22 abs of 8 uppishles
15       difframenission-as.factor(dfiTransmission)         16       difframenission-as.factor(dfiTransmission)         16       difframenission-as.factor(dfiTransmission)         16       # Partitioning (60%:40%)         18       partidx-sample(inrow(df), 0.6*nrow(df), replace = F)         19       dftrain/dfipartidx,]         20       dftst-df[partidx]         21       od=ln(Price, dftrain)         23       sumary (mod)         24       anova(mod)         25       # Measures of Goodness of fit         26       # Measures of Goodness of fit         27.11       (Mprimel) ±         26       # Measures of Goodness of fit         27.11       (Mprimel) ±         27.12       (Mprimel) ±         26       # Measures of Goodness of fit         27.11       (Mprimel) ±         27.12       (Mprimel) ±         26       # Measures of Goodness of fit         27.11       (Mprimel) ±         27.12       (Mprimel) ±         28.12       (Mprimel) ±         29.12       (Mprimel) ±         20.12       (Mprimel) ±         21.12       (Mprimel) ±         21.12       (Mprimel) ±	14	str(df)	ditest 52 obs. of 6 variables
16       # Partitioning (60%:40%)         17       # Partitioning (60%:40%)         18       partidx-sample(linrow(df), 0.6*nrow(df), replace = F)         20       oftast=df[-partidx,]         21       nod-in(Price, dftrain)         23       summary(mod)         24       heasures of Goodnass of fit         21       (rep:#en) :         22       rep:#end         23       summary(mod)         24       heasures of Goodnass of fit         211       (rep:#en) :         221       rep:#end         23       summary(mod)         24       heasures of Goodnass of fit         211       (rep:#en) :         221       rep:#end         231       (rep:#end) :         241       heasures of Goodnass of fit         211       (rep:#end) :         212       (rep:#end) :         213       (rep:#end) :         214       (rep:#end) :         215       heasures of Goodnass of fit         216       (rep:#end) :         217       (rep:#end) :         218       (rep:#end) :         219       (rep:#end) :         2111       (rep:#end)	15	df\$Transmission=as.factor(df\$Transmission)	offrain 47 obs. of 8 variables
17       # Partitioning (60%:40%)         18       partidx:sample(1:nrow(df), 0.6*nrow(df), replace = F)         19       dftrain-df(partidx,]         20       mod-ln(Price ~., dftrain)         21       mod-ln(Price ~., dftrain)         25       anova(mod)         26       # Measures of Goodness of fit         211       (prime)         25       mod-ln(Price ~., dftrain)         26       # Measures of Goodness of fit         211       (prime)         25       mod List of 13         26       # Measures of Goodness of fit         211       (prime)         25       # Norgets         Consel (Cistra-Uppe)       # Norgets         (1) Packages Hulp Verwer       # Out in Out in Output Verseor         26       # Measures of Goodness of fit         211       (prime)       # Out in Output Verseor         212       test of 13       mod List of 13         213       test of 14       test of 16         214       test of 12       # Output Verseor         215       test of 12       # Output Verseor         226       test of 12       # Output Verseor         237       test of 12       # Output Verseor <td>16</td> <td></td> <td>Values</td>	16		Values
18       partidx:sample(l:nrow(df), 0.6*nrow(df), replace = F)         20       dftrain-df[partidx,]         21       mod=lm(frice ~, , dftrain)         25       summary(mod)         26       # Weasures of Goodness of fit         27       # Weasures of Goodness of fit         28       winey(mod)         29       # Weasures of Goodness of fit         21       (moj=wei)         26       # Weasures of Goodness of fit         21       (moj=wei)         26       # Weasures of Goodness of fit         271       (moj=wei)         28       winey(mod)         29       # weasures of Goodness of fit         20       model         210       (moj=wei)         211       (moj=wei)         212       (moj=wei)         213       (moj=wei)         214       (moj=wei)         215       (moj=wei)         216       (moj=wei)         217       (moj=wei)         218       (moj=wei)         219       (moj=wei)         211       (moj=wei)         212       (moj=wei)         213       (moj=wei)         214	17	# Partitioning (60%:40%)	Age num [1:79] 4 5 6 4 9 5 6 3
19       oftrain_of[partidx.]         20       offrain_of[partidx.]         21       nod=lm(Price ~ , dftrain)         23       summary(mod)         24       anova (mod)         25       # Measures of Coodness of fit         211       (mod=lm(row det, model and the residues, their sequences minutes)         26       # Measures of Coodness of fit         211       (mod=lm(row det, model and the residues, their sequences minutes)         212       Income Crower det, model and the residues, their sequences minutes minutes)         211       Long Level 3: (who first) MOOC lanuary 2010/07. Gauver Obtrivisedon Sy PO         212       Long Level 3: (who first) Petrol Petrol Petrol Petrol Level 3: (who first) Petrol Discal Petrol Discal Petrol Discal Petrol Petr	18	<pre>partidx=sample(1:nrow(df), 0.6=nrow(df), replace = F)</pre>	0 mod List of 13
21       mod-lm(Price ~ ., dftrain)         22       mod-lm(Price ~ ., dftrain)         23       summary(mod)         24       anova(mod)         25       # Measures of Goodness of fit         21       (privel) ≥         25       # Measures of Goodness of fit         21       (privel) ≥         25       # Measures of Goodness of fit         21       (privel) ≥         25       # Measures of Goodness of fit         21       (privel) ≥         25       # Measures of Goodness of fit         21       (privel) ≥         26       # Measures of Goodness of fit         21       (privel) ≥         21       (privel) ≥         22       # Measures of Goodness of fit         23       # Measures of Goodness of fit         24       # Measures of Goodness of fit         25       # Measures of Goodness of fit         26       # Measures of Goodness of Fit         21       # Measures of Goodness of Errol         21       period [Gittrack and [Detrol Petrol Petrol         21       # Goodness of Errol         29       # Goodness of Errol         20       # Goodness of Errol	19	dftrain=df[partidx,]	partidx int [1:47] 60 21 42 15 76
22       and-ln(price - , dfrain)         23       summary(mod)         24       anova(mod)         25       # Measures of Goodness of fit         26       # Measures of Goodness of fit         27       (repress)         28       > Measures of Goodness of fit         211       (repress)         22       > Measures of Goodness of fit         23       > Measures of Goodness of fit         24       In repression         25       # Measures of Goodness of fit         211       (repress)         220       File         23       Summary(mod)         24       In repression of fit         211       (repress)         220       File         23       Summary(mod)         24       In repression of fit         25       # Measury 2018/07 Gausery Duit/Session S/ ©         26       File         211       repression of file         221       (repress)         23       Sumary (repression of file         24       (repression of file         25       File         26       File         26       File <t< td=""><td>21</td><td>urcesc-ur[-parcrox,]</td><td></td></t<>	21	urcesc-ur[-parcrox,]	
23       summary(mod)         24       anova(mod)         25       # Measures of Goodness of fit         26       # Measures of Goodness of fit         21       (perweit)=         20       # Weasures of Goodness of fit         21       (perweit)=         25       # Measures of Goodness of fit         26       # Measures of Goodness of fit         27       # Weasures of Goodness of fit         28       # Weasures of Goodness of fit         29       # Measures of Goodness of fit         20       # Weasures of Goodness of fit         21       (perweit)=         20       # Weasures of Goodness of fit         20       # Weasures of Goodness of fit         21       (perweit)=         21       pervectores         21       pervectores         21       pervectores         21       pervectores         21       pervectores         22       # State of Pervectores         23       # State of Pervectores         24       # State of Pervectores         25       # State of Pervectores         26       # State of Pervectores         27       State of Pervectore	22	mod=lm(Price ~ dftrain)	Files Plots Packages Help Viewer
24     anova (mod)       25     # Weasures of Coodness of fit       26     # Weasures of Coodness of fit       211     (ppixel) 2       25     # Verget 2       26     # Weasures of Coodness of fit       27.11     (ppixel) 2       27.12     (ppixel) 2       28.12     # Verget 2       29.12     # Verget 2       20.12	23	summary(mod)	🖕 🧅 🏠 🚍 🙇 🔍 Im 🛛 🚱
25       # Measures of Goodness of fit       Tole works of the transmission	24	anova (mod)	It: Litting Linear Models + Find in Topic
20 m Medsures of coordinates of rite     interpretation       21 (pre-Med)     (N-TM)       2 medsures of coordinates of rite     (N-TM)       2 medsures of rite     (N-TM)       2 model     (N-TM)	25	A Maximum of conducts of file	
Consel C/User/User/User/Desktop/MOOC January 2018/OF Gaurary ObstySesion 6/ ↔ bead(dfStruel_type) [1] petrol Dises! Petrol Petrol Petrol Petrol Levels: CK Dises! Petrol > partid/scsample(L:nrcw(df), 0.5*nrow(df), replace = F) > dfrain-df(partidx.] > offseti-df(partidx.] > offseti-df(partidx.]	26	# Measures of Goodness of fit	following components:
Console C/User/User/User/Desktop/MOOC January 2011/DF Gaurary Disht/Deskdon 6/ 00 minuterine and a feature and a f	23.1	(lop rever) + K script +	www.tilia.iumitus_a named vector of coefficients
<pre>&gt; head(dfSruel_type) restdual a he resduals, that is response minus [1] Petrol Dissel Petrol Petrol Petrol Petrol Levels: (xx Dissel Petrol) &gt; partidx=sample(1=nrw(df), 0.6=nrw(df), replace = F) &gt; dfrain=df(partidx,] &gt; offsetadf(partidx,] &gt; offsetadf(partidx,] &gt; offsetadf(partidx,]</pre>	Console	C:/Users/user/Desktop/MOOC January 2018/Dr. Gaurav Dixit/Session 6/	coefficiences a named vector or coemclents
[1] Petrol Diesel Petrol Petrol Petrol Petrol Petrol     fifted values       Levels: CNG Diesel Petrol     fitzed.ixibex the fifted mean values       > partidix-sample (Linzow(df), 0.6°nrow(df), replace = F)     fitzed.ixibex the fifted filted linear model       > dftrain=df[partidx,]     model       > odfuls[crincedfrain]     model	> head	(dfSFuel type)	residuals the residuals, that is response minus
Levels: ONG Diesel Perol > partidx=sample(l:nrow(df), 0.6*nrow(df), replace = F) > dfrasin=df[partidx,] > dfstzetdf[-partidx,] = odd[partidx_c = dfrasin]	[1] Pe	trol Diesel Petrol Petrol Petrol	fitted values.
<pre>&gt; partidessample(1:nrow(df), 0.6*nrow(df), replace = F) &gt; dfrain-df(partidx,] &gt; dfrain-df(partidx,] &gt; model </pre>	Levels	: CNG Diesel Petrol	fitted.values the fitted mean values.
> dftrained[partidx,] uneinument dank of une incomment > dftest_df[-partidx,] model. model.	> part	idx=sample(1:nrow(df), 0.6*nrow(df), replace = F)	the sumarie rank of the fitted linear
> dftest=d[-partix,]	> dftr	ain=df[partidx,]	model
	> arte	st=of[-partiox,]	model.
names (only for weighted ins) the specified weighted ins) the specified	> mod=	s(mod)	weights (only for weighted fits) the specified
[1] "coefficients" "residuals" "effects" "rank" "fitted.values" weights.	[1] "	coefficients" "residuals" "effects" "rank" "fitted.values"	weights.
[6] "assign" "qr" "df.residual" "contrasts" "xlevels" df.residual the residual degrees of freedom.	[6] "	assign" "qr" "df.residual" "contrasts" "xlevels"	di.residual the residual degrees of freedom.
[11] "call" "terms" "model"	[11] "	call" "terms" "model"	call the matched call
> Call ule matched call	>		the materied call.

So, details about all these written values would be there. So, we might not be interested in all these written values. So, somebody is one important function for the for us to get the relevant output from this exercise. So, let us execute this.

(Refer Slide Time: 27:57)

RStudio	an hard 🦯			-	6 x
9 • 🗃	de Wew Plets Sesson Ruid Lichug Photie Loois Hop			I Project:	(None)
0 6 reg		Environment	History		-0
001	a a Saura an Saura a 🖉 - E		Import Dataset *	/ ≡ liet	. 6
11			import Maser -		
12	df=df[,-c(1,2,3,11)]	Global I nvi	/9 obs of	12 yanabler	
13		O ULD	79 005. 01	12 variables	-
14	str(df)	ortest	32 ODS. OT	8 variables	-
15	dtstransmission=as.tactor(dtstransmission)	O dftrain	47 obs. of	8 variables	
17	# Partitioning (60%-40%)	values			
18	partidx=sample(1:nrow(df), 0.6*nrow(df), replace = F)	Age	num [1:79]	4 5 6 4 9 5 6	3
19	dftrain=df[partidx,]	0 mod	List of 13		
20	dftest=df[-partidx,]	partidx	int [1:47]	60 21 42 15 7	6
21					
22	<pre>mod=lm(Price ~ ., dftrain)</pre>	Files Plots	Packages Help	Viewer	-
23	summary (mod)		a .	O Im C	DIR
25	anova (mou)				010
26	# Measures of Goodness of fit	It: Litting Linear	Models • Find in Top		
27	gf=c(mod\$df.residual, summary(mod)\$r.squared, summary(mod)\$sigma,	residuals	the residuals.	that is response mi	nus
24:1	(lop Level) © R Script ©		fitted values.		
		1511 and and	the fitted mean	values	
Console	C:/Users/user/Desktop/MOOC January 2018/Dr. Gaurav Dixit/Session 6/ 🔗 📼 🗖	fittea.vai	ues the fitted filed	I values.	
coeffi	cients:	rank	the numeric ra	nk of the fitted lines	ar
	Estimate Std. Error t value Pr(> t )		model.		
(Inter	cept) 4.81532 1.56893 3.069 0.00395 **	weights	(only for weigh	ted fits) the specifi	ed
Fuel_t	ypeDiesel -0.46506 1.04170 -0.446 0.65781		weights.		
Fuel_t	ypePetrol -1.51559 1.02256 -1.482 0.14654				
SR_Pri	ce 0.27336 0.05278 5.179 7.58e-06 ***	df.residua	1 the residual de	grees of freedom.	
KM	-0.01219 0.009/6 -1.249 0.21946	call	the matched c	all.	
fransm	15510N1 0.43457 0.657/9 0.681 0.4997/		the sum shi	and used	
Aichao		Lerms	ule terms obj	eci useu.	
Ane	-0.28268 0.13213 -2.139 0.03888 *	contrasts	(only where re	levant) the contrast	ts
		-	used.		
	1				•

So, let us look at the summary output. So, you can see that the call explaining the formula that we had given there that we have the arguments that we are given to the 1 m function, then we have residuals some descriptive statistics about the residuals you can see that a minimum value median value residual, that this particular residual residual is of this models are having min max medium first quartile third quartile.

- Now, let us come to the important part that we would like to discuss first coefficients; the coefficients you can see for all the predictors or we can see here first one is the intercept that is the constant. So, this particular estimator is representing the beta 0 that we had in our slide.
- So, this particular intercept against intercept this particular value that we have this is beta 0, the corresponding standard error for this particular estimate is also given there. The t value and p value have the same meaning which we described in our supplementary lecture on introduction to basic statistics right. So, those values remain same.
- We have also discussed a few more details about these values there. So, you can watch those particular lectures. Now the next important variable that we can see is fuel type. Now you would see that instead of one we have 2 variables for fuel type fuel type diesel and fuel type petrol why this has happened is, because fuel type is a categorical variable and it had 3 categories. So, for all those categories because the way softwares are implemented, the way these techniques are implemented, they cannot handle the textual data and therefore, dummy codes dummy coding has to be performed on these variables categorical variables, and depending on the number of categories in a categorical variable we will have to create equal number of dummy variables.
- So, dummy variables are actually representing different categories. So, for example, fuel type we had 3 categories diesel, petrol and CNG.

(Refer Slide Time: 30:18)



- So, for each of these categories right diesel patrol and CNG. So, we had these 3 labels for fuel type for fuel type labels, for each of these labels will have to create dummy variables. So, dummy variables would indicate presence or absence of that particular value that particular label. So, if the value is one; that means, that car is having fuel type of diesel if the value is 0; that means, that car is not having fuel type of diesel.
- Similarly, the variable the dummy variable corresponding dummy variable for petrol having value 1 will mean that the car is running on the fuel type of petrol, and if it is 0 it is not running on fuel type petrol similarly for this one. So, therefore, for each label will have to create equal number of dummy variable. So, we will have 3 variables. So, for the fuel type instead of having one variable in our model we will end up with 3 dummy variables to represent this particular categorical variable. But if we look at the variables presence of if you know if we have information on 2 variables, any 2 variables out of 3; if we have these any 2 variables out of these 3 then the other variables information is automatically known.
- So, therefore, in our modeling exercise right if a particular value is if a particular car is not having diesel or petrol then of course, it will have CNG. So, because of that if we have information on 2 dummy variables the third one is automatically known. So, therefore, in our model we just have to include 2 dummy variables. So, if there are n classes then we will have to include n minus one dummy variables in the model. Now what happens to do the remaining label right.
- So, for example, petrol and diesel are selected here, and we had petrol diesel and CNG. So, the remaining label it becomes the reference category that we have been talking about. So, any results that we get for these 2 variables now these 2 dummy variables p and d, those results would have to be interpreted with respect to this res reference category. More on this will stop here more on this particular dummy coding we will discuss in our next lecture.

Thank you.