**Lecture – 43**
**Case study on regression**

Welcome to the Case study on regression. In this case study we are going to take a problem that is on predicting price of pre owned cars.

(Refer Slide Time: 00:19)



So, let us go ahead and take a look at the problem statement. Now storm motors they are any commerce company and they act as mediators between parties who are interested in selling or buying a pre owned cars. Now these are second hand cars and storm motors act as mediators. Now in specific for the year 2015 to 2016, storm motors they have recorded the data about the seller and car details.

Now, these set of details includes specification details, condition of the car, seller details, registration details, web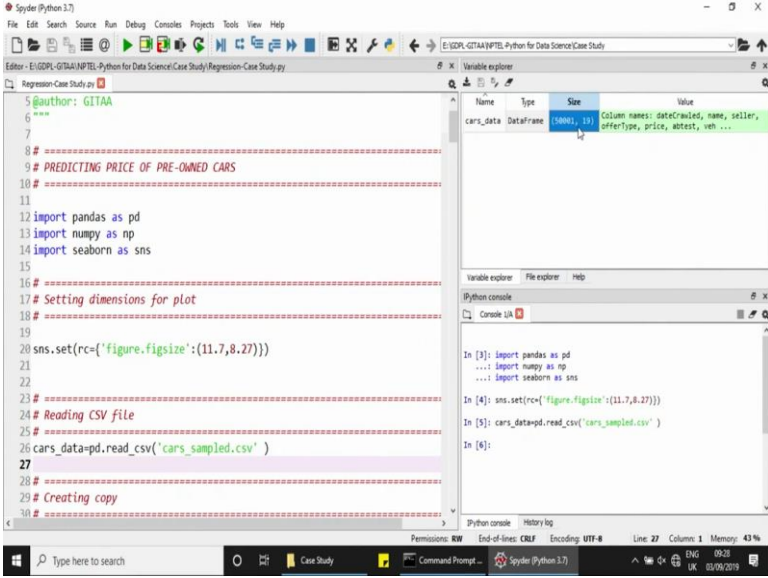 advertisement details, make and model information and price. Now these are several buckets and there are specific parameters or variables that storm motors as collected and each of these variables will definitely fall into one of these buckets. Now storm motors what they want to do is that, they want to develop an algorithm that will help them predict the price of the pre-own cars and this can be based on the various attributes that are associated with the car.

(Refer Slide Time: 01:23)



So, let us see how to solve this case study in python. So, let us start by importing a few required packages, first we have going to import pandas for reading all the files from the CSV format, then I am going to do a few numerical operation and hence I am going to import numpy as well. And I am going to visualize the data and derive some insides and going to import seaborn for it. So, I am importing pandas as pd numpy as np and see born as sns. So, let us just import these packages.

So, these packages have been imported now. Now I am setting a dimensions for all the plots that I am going to be generating, I am using the function sns touch set is a function from the seaborn package and within the parenthesis, I am giving the figure size and this sets the dimensions for rather. And this will set the dimensions for all our plots. So, let us begin by importing the CSV file. So, you have been given a data called cars_sample.csv. Now I have already set my working directory and the data is also residing in my working directory.

So, I am just going to right away read it using the pd.read_csv function. So, the data has been read and we have 50001 records and 19 columns let us just open and see what does the data have to say.
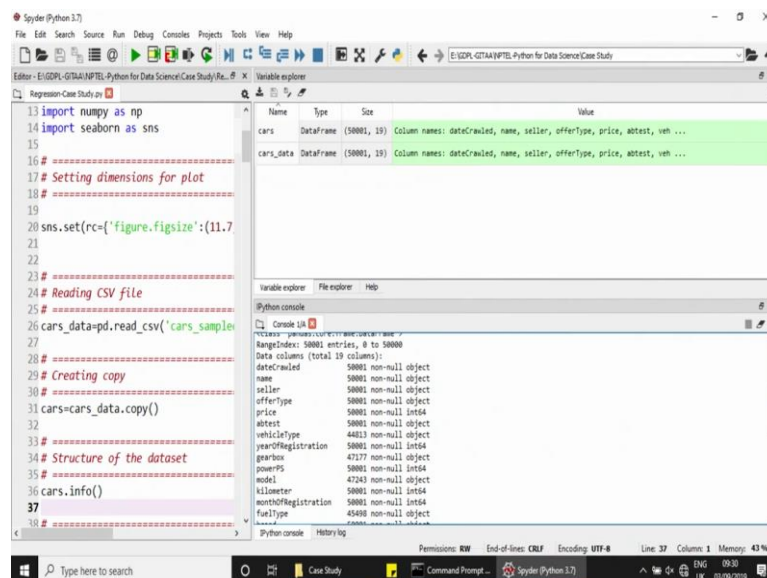
So, you have index from 0 to 50,001 you have the date crawled you have the name of the car; the seller type, the offerType the price you have a b test then vehicle type and year of registration and there are couple of other variables here, all these have been explain to you before the description of the variables. So, now, that we have read let us just create a copy of it. So, you have already seen how to create a copy of your data. So, you can do a deep copy and any change that you made to cars will not be reflected back in cars_data.

So, I made a deep copy and I saved it as cars and we are going to be using cars to work here after. So, let us start by looking at the structure of the data, I am using the cars.info. So, this gives you the structure of the data. Now this tells me that is the data frame it has entries ranging from 0 to 50,000 and the total number of columns are 19 here what you also have is a number of non null objects and what is the data type of each of the column.

So, date crawled has 50,001non null object it is an object data type similarly you have name for seller type, offerType, price and abtest all these are full and if you take vehicle type that is 44,813 non null entries. So, you have some missing values there and similarly even in gearbox model and fuel type this will give you just of all the column names and under each column the number of non null entries and what is the data type of each of the columns.

(Refer Slide Time: 04:49)



So, far as we have 6 columns which are of in 64 data type and there are 13 column which are of object data type. So, now, let us summarize the data I am going to be using the function cars.describe for this.

So, when I used cars.describe you can see that there is a. So, when I use a we are do a.describe function, you can see the summary is given only for a few variables and it is not given for every variable there are some dots that you see in between; now to get rid of this you can use the.

So, when you do a summary you can see the summary is given when you do a.describe, you can see the description gives you the output in a scientific notation what you also see is that it displays the summary of only a few variables and the remaining variables are displayed as '…'. So, let us see how to get rid of the scientific notation first, I am going to use the function.set_option.

Within parenthesis what you have to give is displayed.float format because these are all float values and it is for these that you want to change the output. Now there is a lambda function that we are declaring inside. So, what I am saying here is converted to a three decimal place float value. So, let us just run this function again and then let us rerun cars.describe.

So, here you can see all the values have been rounded off to three decimal places, but we still and not able to get the description for all the columns. So, for that you need to use again the same function pd.set_option. So, within parenthesis you can give the parameter display.max_columns and it will display the maximum number of columns and the next parameter is basically how many of a columns you wanted to be display.

So, I have given a big number here 500, but usually you will not be seen 500 columns in a data, but this is just to set the maximum number of columns that can be displayed.

So, let us do that and let us read on cars.describe. So, now, you can see that all your columns are displayed here and you have a summary for all of it. Now price tells you that there are 50,000. So, under price if you check the count there are 50,001 records.

So, all the records having field and the mean is about 6559 and there is a standard there is a huge standard deviation, it is about 85,818 the minimum value is 0 the first quartile value is 1150, the second quartile which is also the median is 2950, the third quartile is 7190, but the maximum is very very huge and if you look at the difference between the mean and the median which is the second quartile, you can see there is a huge difference between the mean and the median there is itself shows you that price is very very secured.

And if you take care of registration you have 50,001 non null values the mean of year of registration is 2005, but this does not make sense because year of registration is an integer and it cannot be rounded off to a decimal.

But if you take the minimum of it you can see there are years ranging from 1000. Similarly, if you look at powerPS it again has 50,001 non null values, the mean being around 117 and if you look at the median, the mean being around 117, but if you look at the minimum value the minimum is 0 which does not make sense again and the first quartile is at around 70, the second quarter which is the median is around 105 and maximum is very very high, the maximum is around 19312 horse power. Now then you have kilometer, kilometer has again 50,001 non null entries, it is mean is about 1,25,613 kilometers, the minimum is again 5,000, the maximum being 1,50,000.

(Refer Slide Time: 08:57)



Next you have month of registration. Month of registration does not have any missing values the minimum again is 0. But 0 does not make sense and postal code is something that we will be getting rid off because we are not going to be using for this case study, but you can you are free to use this variable if you want to do a more spatial based analysis right. So, again even for postal code, they are all integers and you cannot have decimal values.

(Refer Slide Time: 09:23)

So, let us just start by dropping unwanted columns up front. Now you have variables like name date crawled date created postal code and last seen, now we are not going to using these variables for our analysis and I am going to just drop them. So, I am creating a list of all the variable names here and once I run it all of them are stored under col and I am just going to use the.drop function and to drop them.

So, I do a cars.drop and I say columns is equal to col. So, these are the columns I want to drop, since these are columns I am giving the axis as 1. So, from 19 you can see the number of columns have come to 14, I have dropped 5 columns here. Now, we are going to see if there are any duplicate records in our data right and if there are any duplicate records, we are just going to keep the first occurrence of that such records. So, we started out with 50,001 records let us see when you do a.drop_duplicates how many are we will using. So, one 50,001 it is come down to 49531. Now that we have removed it let us get into data cleaning now.

So, my first task is to compute the number of missing values under each column, now from the.info function we saw not all column were filled right. So, now, we are going to find the number of missing values and each column and later on we will come up the method of how to logically fill such missing values. So, the first step is to basically calculate the number of missing values under each column. So, I am using the isnull function. So, isnull will basically return a data which is true or false.

It will mark the cells which are missing with true, but I am going to sum of the occurrences of the trues. So, I am just going to do it.sum for it. So, if you run the function you can see on my right that seller of a type price, ab test do not have any missing values now vehicle type has 5152 missing values, year of registration is again full, gearbox; however, has 2765 records missing.

And under model you can see you have about 2730 records missing. And under fuel you can see you have about 4467 records missing and about 9640 records are missing and not repaired or damaged. So, this is a column wise count of the number of cells that are missing.

(Refer Slide Time: 12:01)



So, now, let us just take year of registration first. So, in 975 what I am trying to do is I am trying to find year wise count based on this variable. Now, you can use the.value_counts here and I am sorting the index to make sure that, here year of registration is sorted based on the years and not on the counts. So, by default if you do a.value_counts, you would have seen that.

So, it gives you the category with the highest frequency on top, but I do not want that I basically wanted sorted placed on the years. So, let us just run this, this is too big output. So, I am saving it onto year wise under count variable.

(Refer Slide Time: 12:39)



Now, if you open and see. So, let us just see what this data frame has. It has index and it has year of registration index is the set of all years and year of registration on the that column you basically have the frequencies of each of these years. So, you have years that are ranging from 1000 and then you have year from 1910 and what you can also see if you scroll down is that, you have years which are in the future that is which are after 2019.

So, most of these do not make sense so, that is considerable amount of cleaning that we will have to do for this column. So, let us see how to come up with the strategy to clean them.

(Refer Slide Time: 13:19)



So, if you look at the sum of records which are greater than 2018 then that is only 26 right. This is 2 lesser number and we cannot predict in the future. We are in 2019, but I am not considering that because it is occurred only twice and it would not make sense to consider it. So, if you look at the number of cars where registration is greater than 2018.

Now, you have about 26 of them 26 records have a year that is greater than 2018. Now similarly if you do for a lower and check what you can also see is, I am taking the limit 1950. So, if you look at the cars that are made before or that are registered before in 1950, then you have about 38 records of that kind right. Now these are too less correct and they are going to smear or effect or you know sway the effect of the model. So, we are going to get rid of them. So, the working range that I am going to set is between 1950 and 2018. So, just to reconfirm this if you do a scatter plot I am using the sns.regplot.

(Refer Slide Time: 14:31)



If you look at this scatter plot, I have the year of registration on x axis and I have the price on the y axis, let me just drag it to the size just to show you a better plot. So, you can see that the plot is not explaining anything you have all only you have only dots here and there and that is because there are values which are very very high and they are just smearing out the effect on the other points. So, we have to clean this column to really understand what is the effect of year of registration on price.

(Refer Slide Time: 15:09)

So, that is as far as year of registration goes, the next is the variable price. Now if you again do a.value counts and then you sort based on the ascending order of the price and I am just saving it on to the variable price_count.

(Refer Slide Time: 15:25)



So, the left is basically the index and that is the value of the price of the car and the right is basically the frequencies of each occurrence. Now that is a too huge range. So, we have to do tweak it down and even price equal to 0 we will not make any sense, I mean it can be because if you want to just sell it for a few dollars or a few for a very low amount you can.

But then we are looking to generalize this model for a workable range of data and these values are these values can be very very extreme right. So, there is nothing wrong in selling a car for a lower price, but then we would like to generalize a model that would come up with the better prediction, but however, these if you have these values when they are going to smear the effect out. So, we are going to do check on this as well and we are going to arrive at a working range that is easy for us to understand right.

(Refer Slide Time: 16:29)



So, now if you do histogram which is using the distplot right.

(Refer Slide Time: 16:31)



You can see that because there are a lot of entries under the zero prize category. So, this is something that we have to work on.

(Refer Slide Time: 16:47)



So, let us just quickly reiterated by even looking at the describe. So, describe also tells you that the mean is about 6567, the median is way of the medium is around 2950. Now there is a huge different than this itself accounts for the skewness in the data and the minimum here you can see is again 0 and the maximum value is also given. So, the range is really really wide and we have to narrow it down to come up with the generalized model.

(Refer Slide Time: 17:15)

So, let us just quickly do a box plot. Now a box plot also tells you that there are some outliers here right if I do a box plot for the price, you can see you cannot even see the box. In fact, what you just see is a line. Now this itself tells you that there are considerable outliers in your data which are very very extreme in nature and hence you not able to really see the behavior of the variable. So, let us just do a quick check of the range.

So, I am setting the range between 100 to lakh and 50,000 dollars. So, the first is to see how many cars are price above 1,50,000 dollars and that is about 34 and if you do a price check on the lower end you will see that 1748 cars a price lower than 100, but this is because 100 dollars to lakh and 50s decent range to work with. So, the next variable that we are going to look into is power ps. Now again if you do a.values count and then you sort it and now I have saved it under power count.

(Refer Slide Time: 18:27)



So, here again you will see that the power values are on the extreme left and the right hand side is the frequencies of each. So, 5533 records have a power 0; again we have the same problem because the range is too diverse and we cannot really inform much. So, we have to again narrow it down. Somewhere in between you have a lot of occurrences of once some of them are two extreme for instance you have 19,312 horsepower that is very very extreme.

And which is occurred only once and so, we have to narrow down the range of power PS as well. So, let us just do a distplot just to see if there is any skewness yes there is you can see that from here itself because of entries that are of 0. So, this is something that we have to take into account we have to clean it up.

(Refer Slide Time: 19:23)



And if you do a.describe, this will also show you this skewness we have already done this before just reiterating it here. So, the mean is about 116 and if you look at the minimum it is 0 and if you look at the first quartile it is around 76 it is around 69 and if you look at the median, that is around 105. Though the mean and median are not far away the standard deviation is quite huge.

Now for a mean value of 116 if you are going to deviate by around 200 units then that is a lot right. So, though the mean and the median are not very far, but they are still far there is still some skewness. But if you look at the range of powerPS it is too diverse and you also have values that are around 0.

So, you just even if your mean and median are very very close you also have to take into account the standard deviation in the data and in this case it is very very huge right. So, we are going to take up this we are going to tweak the range down, but before we do that let us just quickly see if this is also getting reflected in the box plot.

(Refer Slide Time: 20:27)



This plot is a lot better compared to the price plot that we are earlier did there we were just able to see a line at least here you are able to see a small box. So, yes of course, you have some extreme values here that are actually compressing the box. So, we will have to deal with this, we have to come with a workable range of data. So, let us just plot and see if powerPS has an effect on price before we clean it up.

(Refer Slide Time: 21:03)



Again it is all bundled up together on the lower end and I think that is because of the values that are 0, we have to clean it up before we further do anything with this variable

because otherwise we will not be able to see the effect of this variable on price. So, let us just fixe a range. So, I am fixing a range between 10 and 500, now I have arrived at all these ranges by trial and error also the idea is that you do not want to let go of too many records correct.

So, I am going to check the number of cars that have a powerPS of greater than 500 and that is about 115 and if you take the lower range and if you take the lower value the number of cars that are less than 10 it is about 5565. These are from trial and error and that is considerable reading that has been done to even check you know what is the minimum power that you need to start any vehicle.

So, the working range year of registration is between 19 to 50, the working range for year of registration is between 19 50 to 2018 and for price we are stick into 100 to 1, 50,000 dollars and the working range for powerPS variable is from 10 to 500. So, now, that we have checked the working range for three variables we are now going to clean the data by giving these variables and any further modifications we are going to do in the clean data.