

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

**Lecture – 13**  
**Dimension Reduction Techniques- Part I**

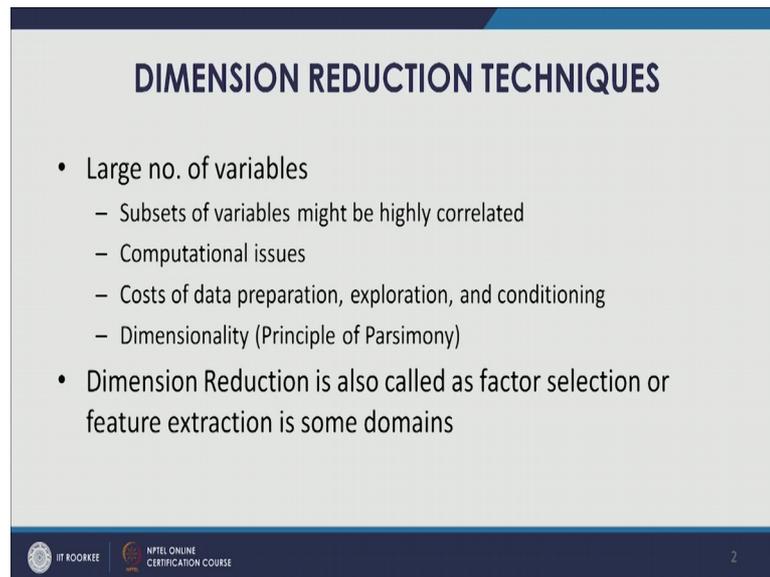
Welcome to the course business analytics and data mining modelling using R. So, we are going to start our discussion on next topic that is dimension reduction techniques. So, let us start. So, first, let us understand the context, so generally whatever discussion we have done till now, we have been talking about the different variables predictors that we generally require to build a model. So, sometimes if there are too many variables because we are also dealing with a large data set, when we talk about data mining modelling, we are generally dealing with large data sets. So, sometimes there are could be too many variables and subset of some of these variables might be highly correlated.

So, we are of course, looking for some amount of correlations, so that some link or association can be established, which could be used for prediction, but if the correlation is on the higher side then the information overlap could be problematic because few variables if they are highly correlated they can deep to is spurious relationships because they are going to dominate the model.

So, their coefficient for example, if we are running regression, if we are doing regression analysis then they will dominate; their coefficient will dominate and the others, other variables coefficient might come down. So, therefore, they will dominate the influence on the outcome variable that is y. So, we do not want that, we do not want our model to be spurious or results to be relationships, results to be meaningless. So, we want to get rid of this highly correlated variable problem.

So, dimension reduction techniques are actually used for this purpose, then there are computational issues as well, for example, if we are dealing with too many variables, sometimes computation might take more time, especially if you are running different variants of your techniques and so many candidate models are there. So, it might be time consuming too much of time would be required to do your modelling exercise. So, for that also we would like to reduce the number of variables.

(Refer Slide Time: 02:35)



**DIMENSION REDUCTION TECHNIQUES**

- Large no. of variables
  - Subsets of variables might be highly correlated
  - Computational issues
  - Costs of data preparation, exploration, and conditioning
  - Dimensionality (Principle of Parsimony)
- Dimension Reduction is also called as factor selection or feature extraction in some domains

IIT ROORKEE | NPTEL ONLINE CERTIFICATION COURSE

Now, another reason could be cost of data preparation exploration and conditioning. So, we talked about this space and the kind of things that we are required to do in data preparation exploration and conditioning. So, that might take more time, if there are more number of variables for example, visual techniques that we covered, that might be, that might require more time and so the cost of modelling would go up, so that could be another reason. So, dimensionality, so we talked about in previous lectures, we talked about principle of parsimony, where we said that we want a variable to have as few variables as possible and still we are able to explain most of the variation in the outcome variable, so that is the desirable property.

So, therefore, we would hardly look to reduce the dimensionality of a model. So, therefore, large number of variables we always look to reduce the number of variables, but we want to retain the important variables in the model. So, we need to identify them and some of the problems for example, some variables been highly correlated, we want to get rid of that also. Dimension reduction is also called factor selection or feature selection in feature extraction in some domains, specially machine learning, artificial intelligence, so they are some other names are also being used.

(Refer Slide Time: 04:02)

**DIMENSION REDUCTION TECHNIQUES**

- Dimension Reduction Techniques
  - Domain Knowledge
  - Data Exploration Techniques
  - Data Conversion Techniques
  - Automated reduction Techniques
  - Data Mining Techniques

IIT ROORKEE | NPTEL ONLINE CERTIFICATION COURSE 3

So, what are the different dimension reduction techniques that can be used? So, these are the few that we are going to cover domain knowledge, data exploration techniques, data conversion techniques, automated reduction techniques and data mining techniques. So, domain knowledge is the one, that we generally that can be very handy, if you understand the phenomena, if you understand the underlying theories that are explain that particular phenomena, you understand the key variables construct and keeping in mind the kind of task that you have to perform, you can identify; easily identify key variables because you have been working on that area. So, because of the knowledge that you have throughout the phenomena about the variables construct and everything else, that can help you, that can put you in a better position in a identifying key variables.

So, in that fashion you can reduce the dimension, you might not use other variable state and your modelling and that can is that is going to help. Removing redundant variable, some variables through that same process if you understand your area, you will immediately understand that these variables are not useful, so you can get rid of those variables.

Sometimes, if you have been involved in data preparation exercise and the collection exercise you would also know some of the variables which might be having a redundant values because of the way the data might have been collected, it might not you know, those variables might not be suitable because some errors will come into your model and

the results that you actually might expect might not come and your own arguments and theories might not be justified. So, therefore, we can also identify these redundancy variables and if we can handle them, we should go for that, otherwise we can get rid of these variables as well.

Measurement issues for variables, so sometimes we also have to see, if a variable if there is an we need to repeat the exercise can the variable be measured again, is it easy to measure that particular variable, so those kind of things also. So, that is from the project point of view, if you are running a project for your business organisation and there is some business roles and the relevant business analytics exercise that you have to perform.

So, you will also have to see, whether a particular variable can be again it can be data on that variable can be again measured and collected. So, that is also you have to see whether the exercise can be repeated or not, so those kind of issues could also be there. Next is data exploration techniques, so some of them, some of these things we have already discussed in previous lectures for example, descriptive statistics, that can also help us in reducing the dimension for example, summary statistics.

Sometimes it might be easily able to see, the you know, from the summary itself, the we can identify important variables and we can also identify the variables which we can get rid of, pivot tables in excel that can also be used to, again look at the summary label data and make further decisions about your modelling exercise. Correlation analysis can also help you identifying the variable which are correlated and should be included in the model, highly correlated variable can also be identified there that can again help you in determining which variable you need to, you have to perform closer inspection and so that is also, so correlation analysis is also going to be quite helpful in reducing the dimension.

Then we have visualization techniques that we have covered in the in previous lectures. So, let us go through some exercises for these techniques. So, let us open R studio.

(Refer Slide Time: 08:16)

```
1 library(xlsx)
2
3 # UsedCars.xlsx
4 df=read.xlsx(file.choose(), 1, header = T)
5 df=df[, !apply(is.na(df), 2, all)]
6
7 # First 10 observations
8 df[1:10,]
9
10 Age=2017-df$Mfg_Year
11 df=cbind(df, Age)
12 dfb=df
13 df=df[,-c(1,2,3)]
14 str(df)
15
16 #summary statistics
17
```

	Brand	Model	Mfg_Year	Fuel_type	SR_Price	KM Price	Transmission	
1	Hyundai	verna	2013	Petrol	8.88	75.000	5.60	
2	Mahindra	Quanto	2012	Diesel	6.99	49.292	3.95	
3	Maruti	Suzuki	SX4	2011	Petrol	7.18	48.000	2.99
4	Chevrolet	Beat	2013	Petrol	4.92	41.000	2.35	
5	Honda	Civic	2008	Petrol	13.50	110.000	3.65	
6	Honda	Brio	2012	Petrol	5.74	60.000	2.99	
7	Hyundai	i20	2011	Diesel	8.42	56.000	3.87	
8	Skoda	Rapid	2014	Diesel	10.49	27.500	5.80	
9	Mitsubishi	Pajero	2007	Diesel	19.50	101.000	5.50	
10	Hyundai	i10	2013	Diesel	6.25	94.770	3.10	

Let us, load this particular library that we have been using in every exercise. So, used cars data set this is the 1 that we are going to use this for this exercise. Let us import this data set again, now first an observation we have already seen this data in previous lecture, but let us have a look again. So, these are the variables, as you can see you are by now you should be familiar with this particular data set, brand model, manufacturing year, fuel type, showroom price, kilometre price, transmission, owner's airbag and C price. So, these are the variables.

(Refer Slide Time: 09:07)

```
10 Hyundai i10 2013 Diesel 6.25 94.770 3.10 1
Owners Airbag C_Price
1 1 0 1
2 1 0 0
3 1 0 0
4 1 0 0
5 2 0 0
6 1 0 0
7 1 1 0
8 1 1 1
9 1 1 1
10 1 1 0
```

Now, another variable of interest that we have been using in other lectures as well is age. So, we are going to compute age like we did in previous lectures from manufacturing year, let us append this to our data frame, let us take a back up and first 3 variables are we were not interested.

So, let us drop them and now let us have a look at the structure. So, these are the variables of interest, 9 variables, you can see 1 being the categorical, others being numerical. Now, let us look at the summary statistics.

(Refer Slide Time: 09:44)

```

10 Age=2017-df$Mfg_Year
11 df=cbind(df, Age)
12 dfb=df
13 df=df[,c(1,2,3)]
14 str(df)
15
16 #Summary statistics
17 countblank=function(x) sum(x=="")
18
19 dfsum=data.frame(Average=sapply(df[, -1], mean), Median=sapply(df[, -1], median),
20                 Min=sapply(df[, -1], min), Max=sapply(df[, -1], max),
21                 Std=sapply(df[, -1], sd), Count=sapply(df[, -1], length),
22                 countblank=sapply(df[, -1], countblank))
23 round(dfsum, digits = 2)
24

```

```

'data frame': 79 obs. of 9 variables:
 $ Fuel_type : Factor w/ 3 levels "CNG","Diesel",...: 3 2 3 3 3 3 2 2 2 2 ...
 $ SR_Price  : num  8.88 6.99 7.18 4.92 13.5 ...
 $ KM        : num  75 49.3 48 41 110 ...
 $ Price     : num  5.6 3.95 2.99 2.35 3.65 2.99 3.87 5.8 5.5 3.1 ...
 $ Transmission: num  0 0 0 0 1 0 1 0 0 1 ...
 $ Owners    : num  1 1 1 1 2 1 1 1 1 1 ...
 $ Airbag    : num  0 0 0 0 0 0 1 1 1 1 ...
 $ C_Price   : num  1 0 0 0 0 0 0 1 1 0 ...
 $ Age       : num  4 5 6 4 9 5 6 3 10 4 ...

```

Now, this is 1 function that is also part of summary statistics that we are going to create a count blank. So, this is going to count the number of a blank cells that are there in a particular variable, in a particular column. So, that we can identify there are any missing values that are represented through blanks. So, if they are represented through NA or some other format, so we can use that and create our; write our own function and find out, if there are any such values.

So, let us create this function, so we can see a function has been created; count blank function has been created. Now, this is the data frame that we are going to create that is based on summary statistics that we are going to compute. So, mean is 1, then we are going to compute median, you would see that first column we have eliminated because this being fuel type being categorical. So, there is summary statistics are mainly applicable to numerical variables, mean, median then a minimum value and the

maximum value. So, that will tell us the range, then the standard deviation, then the length for that particular column, how many observations are there, then the count blank also. So, that we can easily spot, which variables are having all the values and if there are any missing values at all. So, let us execute this line and find out more about this variables, let us look at the results.

(Refer Slide Time: 11:18)

```

14 str(df)
15
16 #Summary statistics
17 countblank=Function(x) sum(x=="")
18
19 dfsum=data.frame(Average=sapply(df[, -1], mean), Median=sapply(df[, -1], median),
20                 Min=sapply(df[, -1], min), Max=sapply(df[, -1], max),
21                 Std=sapply(df[, -1], sd), count=sapply(df[, -1], length),
22                 countblank=sapply(df[, -1], countblank))
23 round(dfsum, digits = 2)
24
25 # Correlation table
26 M=cor(df[, -c(1,5,8)]); M
27 M[upper.tri(M)]=NA; M
28 print(round(M, digits = 2), na.print = "")
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100

```

```

> round(dfsum, digits = 2)
  Average Median  Min Max  Std Count Countblank
SR_Price  10.80  7.76  3.11 116 13.17  79      0
KM         77.71 78.90 19.00 167 27.96  79      0
Price       4.94  3.85  1.15  72  7.89  79      0
Transmission 0.20  0.00  0.00  1  0.40  79      0
Owners       1.14  1.00  1.00  2  0.35  79      0
Airbag       0.22  0.00  0.00  1  0.41  79      0
C_Price      0.39  0.00  0.00  1  0.49  79      0
Age          5.67  5.00  2.00  10  2.07  79      0

```

Now, you can see all these variables are having 79 as count. So, therefore, 79 observations are there and no variable is having any missing values, count blank numbers or 0 for each variables. So, we are on the safer side on these 2 counts.

Now, let us look at the range minimum and maximum value, you can see the range for different variables SR price and price, kilometres, so different range you can see, you can also look at the median values and the mean values. So, these values can also help you understand, average values are going to give you the central sense of data centralise, a numbers of the data, a minimum and maximum value will give you the range of the data.

If the, you can also compare average and median value, so if the average is higher than the median value then probably the data is rightly right skewed, if the it is average mean value is less than the median value then the probably the data is left skewed. So, that kind of thing can also be understood from these statistics. Now, this brings us to our next part that is pivot tables. So, pivot tables can also be used, so we can if there is some label of interactivity that can be achieved using excel pivot tables and some kind of combining



If you look at the first pivot table, so once you create the pivot table, you would have this kind of a view. So, where all the variables would be listed on this particular right top of your excel file and all the variables you can see.

Now, you would have report filter column, label, row labels and a values. So, in this particular example we have taken transmission variable in the row side and the in the values section we have taken the price and for the price also it is the counting of those values that has been done. So, if we look at this particular table, you would see transmission 0, 63 such cars are there and for transmission 1 that is for automatic 16 cars are there.

If you look at the another pivot table that is in more detail you would see we have now 2 labelling, 1 for kilometre is on the row side and the transmission is on the column side. So, transmission 0 and 1 and then different ranges have been created. Once you create a pivot table you will have to group them. So, you will have group them to have this kind of result that we have, if you are in your own time you can learn more about the pivot tables.

So, different billing of a kilometres have been billed into these groups and you can see the values for these bins, across these 2 categories of transmission and if you look at again the price, average of price values have been taken, so the values that you see in these cells, these are actually the average price you can see average of price value 5.8 and row is this 1929. So, these kind of values they can easily have here. From there, from this, these kind of values you can easily spot, if some of the categories can be combined, for example, if KM was a variable which we wanted to convert into a categorical variable, therefore, how billing can be done and the kind of sizes of those bins and number of such bins that those kind of decision can be performed. For example, if you look at the 79 to 89 range and 89 to 99 range, if you look at the average pricing for these 2 range this is quite close. So, probably we could have clubbed these 2 bins.

(Refer Slide Time: 15:59)

The screenshot shows an Excel pivot table with the following data:

Transmission	Count of Price
0	63
1	16
<b>Grand Total</b>	<b>79</b>

Average of Price	Transmission	Grand Total
19-29	0	38.9
20-29	0	2.95
30-39	0	8.65
40-49	0	3.87
50-59	0	13.55
60-69	0	5.55
70-79	0	3.82
80-89	0	4.32
90-99	0	3.85
100-109	0	6.2
110-119	0	2.5
120-129	0	1.1
<b>Grand Total</b>	<b>3.747619048</b>	<b>9.644375</b>

Similarly, if you look at 29 to 39 and we look at 39 to 49 the average pricing is quite close. So, therefore, we could have combined these 2 bins. So, this kind of combining can be done because as you know, when you have categorical variable and when you do your modelling, you will have to create dummy variables depending on the number of categories that are there in the variable. So, that would actually increase the dimensionality of your model and you are always looking to reduce the dimensionality, if there are few categories that is better for the model.

So, we are always looking to reduce the number of categories for a variable. So, we are always looking to combine some of the categories or drop some of the categories if required, so that kind of thing can be easily done using pivot tables. Now, next thing that we can do is correlation analysis. So, as we talked about correlation analysis can be helpful in identifying the highly correlated variables. So, let us open R studio and will understand through an exercise. So, again you would see the same data frame, the used cars data set we are going to use, you would see the categorical variables we are not going to consider for correlation tables because the variables have to be numerical. So, let us, so core car is the function. So, this is something that basic stat is something that we have covered in our supplementary lectures. So, you are advised to actually go through those videos, those lectures. So, correlation let us compute between the variables.

(Refer Slide Time: 17:48)

```

20 # Sum=us.c(1:1) # sum(average=sapply(df[,1:2], mean), median=sapply(df[,1:2], median),
21       Min=sapply(df[,1], min), Max=sapply(df[,1], max),
22       Std=sapply(df[,1], sd), Count=sapply(df[,1], length),
23       countblank=sapply(df[,1], countblank))
24 round(dfsun, digits = 2)
25
26 # correlation table
27 M=cor(df[, -c(1,5,8)]); M
28 M[upper.tri(M)]=NA; M
29 print(round(M, digits = 2), na.print = "")
30 symnum(M)
31
32 # Reducing Categories
33 Age_groups=levels(as.factor(df$Age))
34 Age_groups2=as.numeric(Age_groups)
35 C_PricebyAge1=NULL

```

```

Console C:/Users/user/Desktop/MOOC January 2018/Dr. Gaurav Dixit/Session 4/
Airbag  0.42830900 -0.05869390  0.35284463 -0.21060588  1.00000000 -0.12563911
Age     -0.02907086  0.53381082 -0.17338213  0.29491101 -0.1256391  1.00000000
> M[upper.tri(M)]=NA; M
      SR_Price  KM Price Owners Airbag Age
SR_Price  1.00000000 NA      NA      NA      NA      NA
KM        -0.06606499 1.00000000 NA      NA      NA      NA
Price     0.95517032 -0.2193900  1.00000000 NA      NA      NA
Owners   -0.03125270  0.3248002  -0.08488456  1.00000000 NA      NA
Airbag    0.42830900 -0.0586939  0.35284463 -0.2106059  1.00000000 NA
Age      -0.02907086  0.5338108  -0.17338213  0.2949110  -0.1256391  1

```

You would see SR price, so this matrix you would see in 79 rows and columns we have the same names, all the numerical variables have been selected SR price, KM price, owners airbag and age and you would see. Now, this particular metrics is symmetrical from the diagonal all the diagonal values are 1 because the variable is going to be 100 percent correlated with itself, therefore, all these values are gone.

(Refer Slide Time: 18:28)

```

20 # Sum=us.c(1:1) # sum(average=sapply(df[,1:2], mean), median=sapply(df[,1:2], median),
21       Min=sapply(df[,1], min), Max=sapply(df[,1], max),
22       Std=sapply(df[,1], sd), Count=sapply(df[,1], length),
23       countblank=sapply(df[,1], countblank))
24 round(dfsun, digits = 2)
25
26 # correlation table
27 M=cor(df[, -c(1,5,8)]); M
28 M[upper.tri(M)]=NA; M
29 print(round(M, digits = 2), na.print = "")
30 symnum(M)
31
32 # Reducing Categories
33 Age_groups=levels(as.factor(df$Age))
34 Age_groups2=as.numeric(Age_groups)
35 C_PricebyAge1=NULL

```

```

Console C:/Users/user/Desktop/MOOC January 2018/Dr. Gaurav Dixit/Session 4/
Airbag  0.42830900 -0.0586939  0.35284463 -0.2106059  1.0000000 NA
Age     -0.02907086  0.5338108 -0.17338213  0.2949110 -0.1256391  1
> print(round(M, digits = 2), na.print = "")
      SR_Price  KM Price Owners Airbag Age
SR_Price  1.00
KM        -0.07  1.00
Price     0.96 -0.22  1.00
Owners   -0.03  0.32 -0.08  1.00
Airbag    0.43 -0.06  0.35 -0.21  1.00
Age      -0.03  0.53 -0.17  0.29 -0.13  1

```

Now, we want to get rid of the upper triangular, so this we can do through this particular function upper tri, we will set it to NA and then later on will get rid of this, print is the

function that can be used, so we can live NA as blank and will have this particular correlation table, which is now showing just a 1 half lower triangular matrix and we can see the numbers. If we look at the numbers, we can see price and showroom price they are very highly correlated 96 percent correlation 596 is the correlation coefficient here.

Similarly, we can look at some other numbers for example, airbag and SR price 43 and then we can have look at other numbers like age and KM a 53, similarly airbag and price they are 35. So, there are some a highly correlated numbers, price and showroom being very highly correlated. So, these highly correlated numbers we can easily spot and then we can try to understand if some of these variables can be, some of the dimensionality can be reduced, but in this case SR price is an important variable, it being correlated with the price is because of the, that SR price is mainly determining the used price of a used car. So, a because of that this high correlation is there. So, we do not need to get rid of because this being a key variable and important variable and it is relationship is again with the outcome variable of interest that is price.

Now, there is another function that can be used to perform the same thing that is sym num, sym num is 1 function that again displays the, a correlation table, you can see here.

(Refer Slide Time: 19:56)

```

20  Min=sapply(df[, -1], min), Max=sapply(df[, -1], max),
21  Std=sapply(df[, -1], sd), Count=sapply(df[, -1], length),
22  countblank=sapply(df[, -1], countblank))
23  round(dfsun, digits = 2)
24
25  # correlation table
26  M=cor(df[, -c(1,5,8)]); M
27  M[upper.tri(M)]=NA; M
28  print(round(M, digits = 2), na.print = "")
29  symnum(M)
30
31  # Reducing Categories
32  Age_groups=levels(as.factor(df$Age))
33  Age_groups2=as.numeric(Age_groups)
34  C_PricebyAge1=NULL
311 Top Level

```

```

> symnum(M)
  SR_Price 1
  KM        1
  Price     B 1
  Owners    . 1
  Airbag    . . 1
  Age       . 1
attr(,"Legend")
[1] 0 . 0.3 . 0.6 . 0.8 + 0.9 +* 0.95 +B 1 \t ## NA: '?'
>

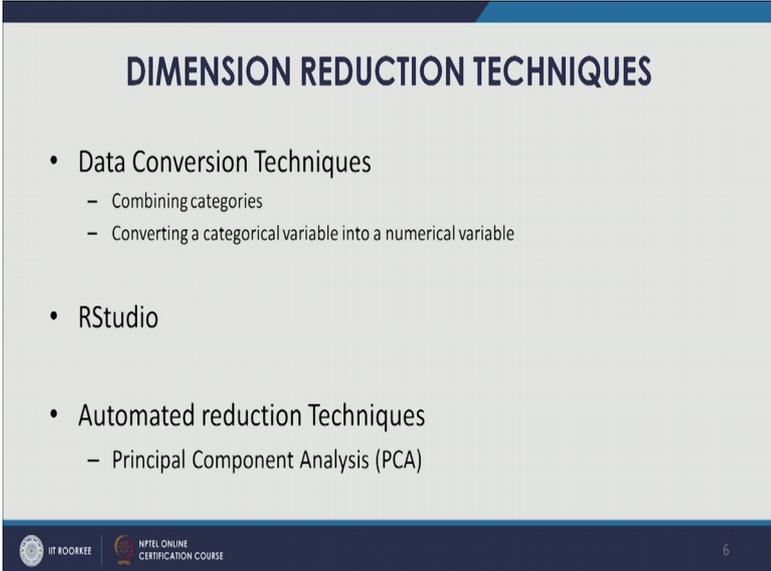
```

But the instead of the actual numbers they display the different symbols for different range. For example, any correlation value that is between 0 and 0.3 that is left blank, then any value that is between 0.3 and 0.6 they would displayed using dots, you would

see 4 dots there. So, 4 correlation values are dots. So, this is easier for because we are generally looking to identify the variables which are slightly on the higher side, so more than 0.3.

So, therefore, you would see any values, any correlation value which is more than 0.3, we are interested in knowing those 2 variables. So, 0.3 and 0.6 are being depicted by dot, there are 4 such values, then 0.6 to 0.8 they can be depicted by comma, we do not have any such value and 0.8 to 0.9 y plus then 0.9 to 0.95, that is why as trick. So, we do not have any such value, but 0.9521 that is depicted by B. So, that we have 1 value in this category that is price and showroom price. So, immediately we can spot through these symbols because we are interested in identifying the highly correlated values. So, this can actually help. So, this is another function. So, let us go back to slides.

(Refer Slide Time: 21:15)



**DIMENSION REDUCTION TECHNIQUES**

- Data Conversion Techniques
  - Combining categories
  - Converting a categorical variable into a numerical variable
- RStudio
- Automated reduction Techniques
  - Principal Component Analysis (PCA)

IIT ROORKEE | NPTEL ONLINE CERTIFICATION COURSE | 6

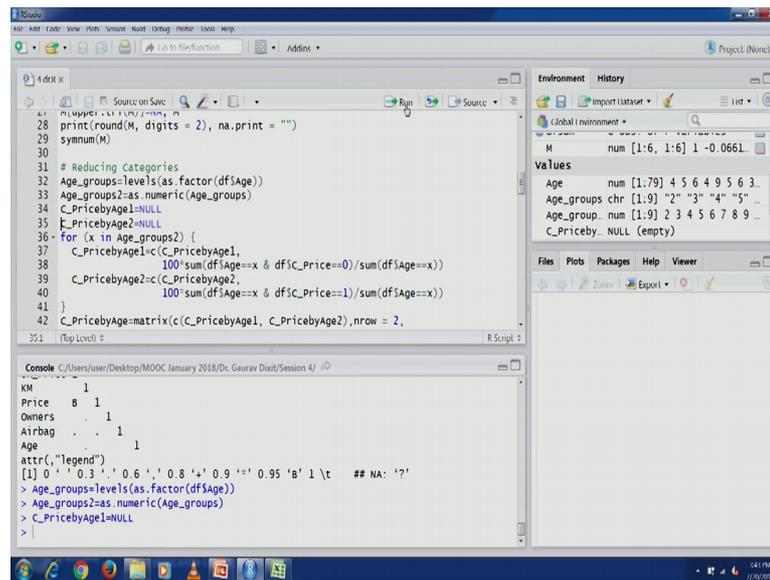
Now, another method for reducing dimension could be data convergent techniques. So, if we have a categorical variable and many categories are there and it will we can have a look at the data and identify the categories which can be combined. So, combining categories as we discussed that can reduce dimensions. Similarly, if we have a categorical variable and it is possible to convert it in to a numerical variable, so that can also help us in reducing the dimension. For example, suppose we have a collected data so on age, age groups.

So, we have not been able to collect the data on actual age of a individual, rather we have age group information. So, particular record belonging to a particular age group, so in that case depending on the number of age group that we have, so those many dummy variables have will have to be created for our modelling exercise and that will increase the dimension. So, we can reduce this scenario because age is essentially is a continuous variable. So, we can, if there are many more categories for age groups we can treat it as a numerical variable. So, the middle value of a particular range can be assigned, that mid value of particular range can be assigned as the value for that particular, individual age value for that particular individual and we can convert this categorical variable into a numerical variable and thereby reducing the dimensions, so this can be done.

Similarly, for categorical variable as we talked about in the previous exercises as well, a for 2 groups if there are similar kind of numbers are there with respect to our outcome variable, then probably we can combine them. If there is 1 group which is having very few number of observation for that group, we can also drop it or we can combine it with the major category. So, that kind of combining can be performed using data convergent techniques. So, let us open R studio and let us do the same thing through an exercise.

So, in this particular case again same data set we are going to use. So, age is the variable that we have in our data set; age of used cars. So, let us convert it into age groups. So, if age this was numerical variable we need to use the hash dot factor function to convert it into factor variable or categorical variable and now we are trying to extract the labels, so of this particular categorical variable once created.

(Refer Slide Time: 23:59)



```
28 print(round(M, digits = 2), na.print = "")
29 symnum(M)
30
31 # Reducing Categories
32 Age_groups=levels(as.factor(df$Age))
33 Age_groups2=as.numeric(Age_groups)
34 C_PricebyAge1=NULL
35 C_PricebyAge2=NULL
36 for (x in Age_groups2) {
37   C_PricebyAge1=c(C_PricebyAge1,
38                 100*sum(df$Age==x & df$C_Price==0)/sum(df$Age==x))
39   C_PricebyAge2=c(C_PricebyAge2,
40                 100*sum(df$Age==x & df$C_Price==1)/sum(df$Age==x))
41 }
42 C_PricebyAge=matrix(c(C_PricebyAge1, C_PricebyAge2),nrow = 2,
```

Environment History

Global Environment

M num [1:6, 1:6] 1 -0.0661

Values

Age	num	[1:79]	4 5 6 4 9 5 6 3...
Age_groups	chr	[1:9]	"2" "3" "4" "5" ...
Age_group...	num	[1:9]	2 3 4 5 6 7 8 9 ...
C_Priceby...	NULL	(empty)	

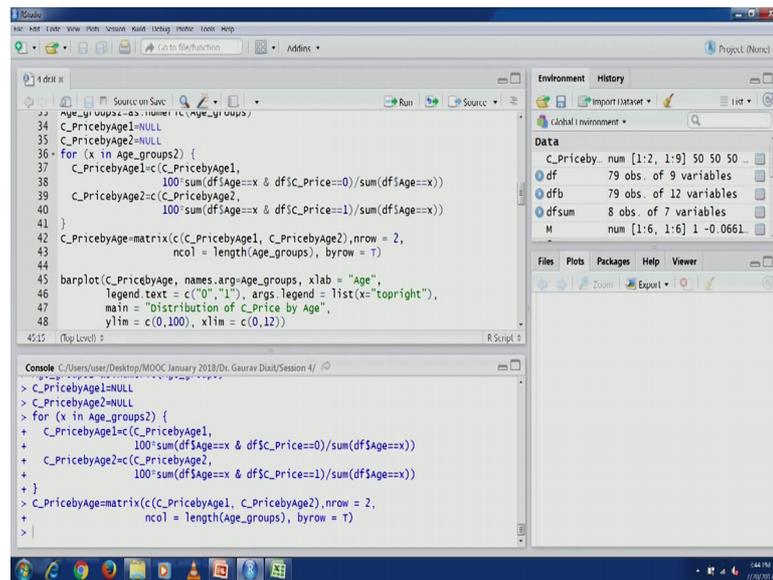
Console

```
C:/Users/user/Desktop/MOOC January 2018/Dt. Gaurav Dast/Session 4/
KM      1
Price   B  1
Owners  .  1
Airbag  .  1
Age     .  1
attr(,"Legend")
[1] 0 "0.3" "0.6" "0.8" "0.9" "0.95" "B" 1 \t ## NA: '?'
> Age_groups=levels(as.factor(df$Age))
> Age_groups2=as.numeric(Age_groups)
> C_PricebyAge1=NULL
>
```

So, you would see age groups has been created and different labels can be seen here in the environment or data section 2, 3, 4, 5 and many others. Now, for further coding we would also require these groups in a numeric format. So, we are trying to create another variable to store the same information. Then, we want to create a bar plot, which will see as later on. So, for that we need to create this C price that categorical pricing that we had created for this particular data set by age. So, let us do that.

So, there are 2 groups for categorical price 1 is 0 and 1 is for 1. So, used cars with less than 4 lakh value and then used car is greater than or equal to 4 lakh value. So, these 2 groups we have, one is, so 2 variables we need to create, for each of these 2 variables and for different age categories, we are trying to find out the these values, we are trying to sum them up. So, you would see first C price age 1 is actually, we are trying to find out the percentage values, where age is a particular belong, belongs to a particular age is there and it belongs to particular age group and the C price is also 0. So, we want to find out the percentage of that, similarly for C price 1, we are going to perform the same thing.

(Refer Slide Time: 25:50)



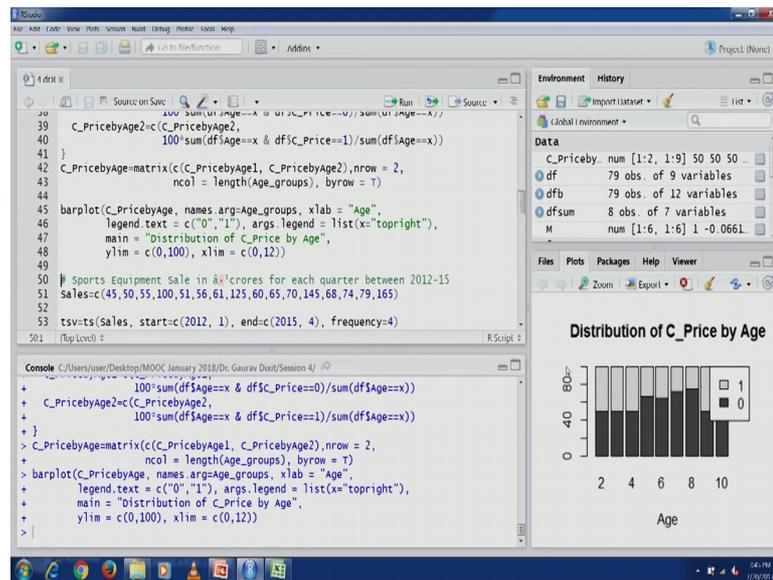
```
34 C_PricebyAge1=NULL
35 C_PricebyAge2=NULL
36 for (x in Age_groups2) {
37   C_PricebyAge1=c(C_PricebyAge1,
38     100*sum(df$Age==x & df$c_Price==0)/sum(df$Age==x))
39   C_PricebyAge2=c(C_PricebyAge2,
40     100*sum(df$Age==x & df$c_Price==1)/sum(df$Age==x))
41 }
42 C_PricebyAge=matrix(c(C_PricebyAge1, C_PricebyAge2),nrow = 2,
43   ncol = length(Age_groups), byrow = T)
44
45 barplot(C_PricebyAge, names.arg=Age_groups, xlab = "Age",
46   legend.text = c("0","1"), args.legend = list(x="topright"),
47   main = "Distribution of C_Price by Age",
48   ylim = c(0,100), xlim = c(0,12))
```

```
> C_PricebyAge1=NULL
> C_PricebyAge2=NULL
> for (x in Age_groups2) {
+   C_PricebyAge1=c(C_PricebyAge1,
+     100*sum(df$Age==x & df$c_Price==0)/sum(df$Age==x))
+   C_PricebyAge2=c(C_PricebyAge2,
+     100*sum(df$Age==x & df$c_Price==1)/sum(df$Age==x))
+ }
> C_PricebyAge=matrix(c(C_PricebyAge1, C_PricebyAge2),nrow = 2,
+   ncol = length(Age_groups), byrow = T)
>
```

So, let us run this, you would see that different 2 different variables C price by age 1 and C price by age 2 have been created. So, 9 groups were there, you would see 9 groups were there and 9 values have been computed.

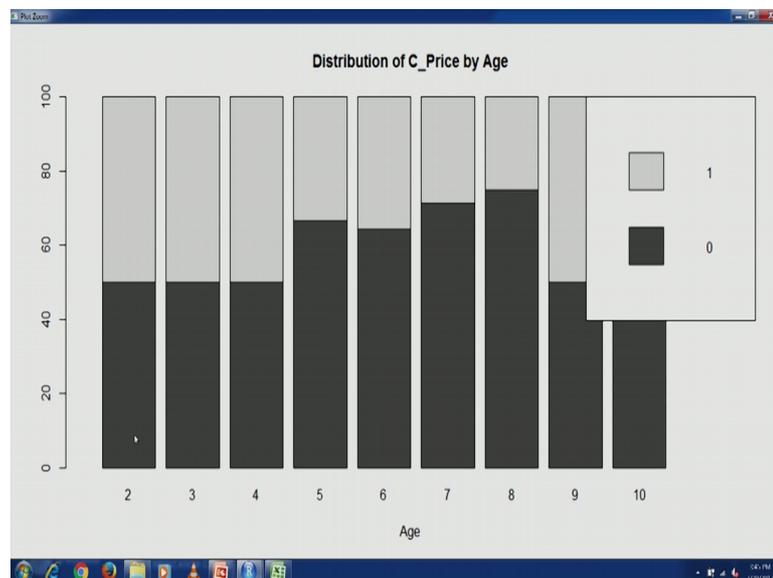
Now, let us create a matrix of these 2 variables, once this is done. Now, this particular information, this particular variable can be used to draw a bar plot, the x axis is again going to be categorical as we have been doing before. So, age groups; the labels from the age groups would be taken for this labelling, limits have been is specified every appropriately because this is going to be percentage number. So, 0 to 100 is the range for y limit and then because there are 9 groups, so this, they are going to be captured in this range for x axis.

(Refer Slide Time: 26:57)



So, let us execute this.

(Refer Slide Time: 27:04)

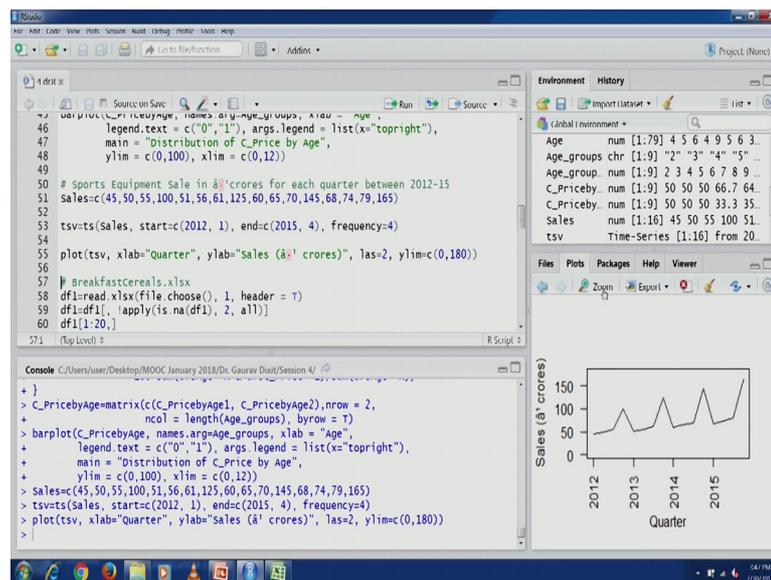


You would see that a bar plot has been generated, now you can see that 0 is represented by this higher intensity shade of grey colour and then 1 is represented by light grey colour. So, you can see for age group 2, 3 and 4 you would see that proportion of these values, this is number of cars for these is same and count of these cars for other groups this is slightly different, if you look at these 2, 3 and 4 from here, if we can see that probably age group 2, 3 and 4 they can be combined together and similarly age group 9

and 10 they can also be combined, if you look at age groups 5 6, they can be combined similarly, age group 7 and 8 can also be combined, that is mainly because the proportion does not change, so therefore, we can combine these categories.

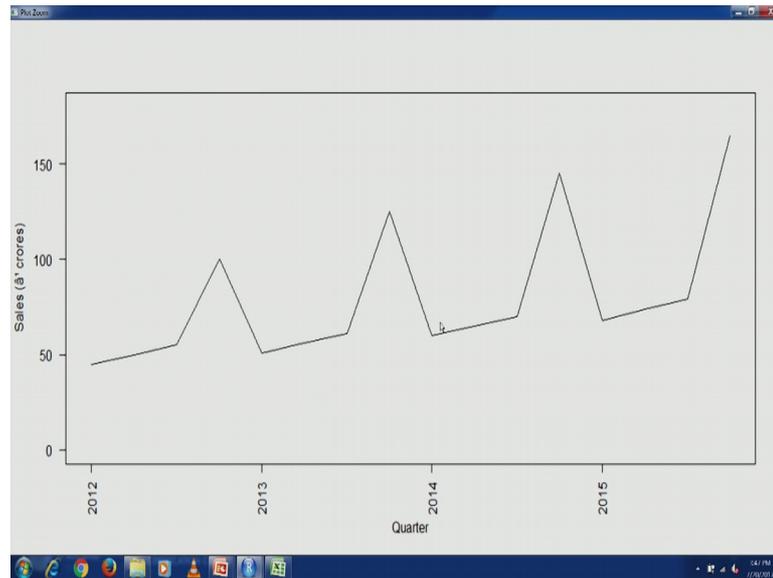
So effectively, we can reduce the number of categories from, right now they are 9 to 1, 2, 3 and 4. So, we can reduce from 9 categories to 4 categories, so thereby reducing a dimension by 5. So, this kind of combining categories can always be done as we talked about. Now, let us do 1 more exercise for time series data, so this hypothetical data we are going to create this is a time series data, so this is sales, a for force equipment sales in rupees corrodes for each quarter between 2012 to 2015 and for each quarters these numbers are given, let us create this variable sales and a time series is being created for the same.

(Refer Slide Time: 28:53)



You can see the frequency is 4 because this is quarterly data and the sales data is being used from 2012 onwards to 2015 and this time series has been created. Now, let us plot this.

(Refer Slide Time: 29:15)



Now, let us have a look at this particular plot. Now, you would see that let us a look at the 2012 data. So, there are going to be 3 points as we saw in the data itself, 1, 2nd point somewhere here and 3rd point here and then 4th point. So, it is always at the 4th point, where some peak is happening. So, it is at the in the 4th quarter were then there is huge spike in sales, but if we look at the a 1st and 2nd and 3rd quarter the numbers are quite close to each other. So, we are looking to combine categories, so we are looking at we have sales data for different quarters and we are interested in analysing this the same quarter wise performance.

So, probably we can combine the data for quarter 1, quarter 2 and quarter 3 because the kind of numbers that we have for this time series you would see that, their performance is quite similar. So, we can combine group them into 1 and the peak 1 could be the another group. So, this is how, combining can be done even for time series data. So, will stop here and in the next lecture will start our discussion on automated reduction techniques and the 1 that we are going to cover in this particular course is principle component analysis, so will start from there.

Thank you.