



In this module, we'll talk about how we execute data exploration and preparation tasks inside R. To demo the process of data exploration and preparation, we'll be looking at a credit card dataset. Let us, first of all, understand what are the variables present in credit card dataset by looking at the data dictionary for this particular dataset.

This is the data dictionary for the credit card dataset that we'll be using throughout this session. As you can see, since this is a credit card dataset, this is a financial dataset; we have some financial variables like Debt Ratio, monthly income. We also have some demographic variables like gender, region and occupation. Now the variable of interest for us would be NPA status. This variable talks about if a person has been delinquent in the past 180 days or not. Essentially, it gives an idea of who are my good customers and who are my bad customers.

So the person who has been delinquent in the past 180 days would be a bad customer for me and the person who has not been delinquent in the past 180 days would be a good customer for me. Let us now read the dataset inside R.

Since we will need to manipulate the data a lot, I have also loaded the library (dplyer). This command here `options(scipen=999)` switches off the scientific notation in terms of number representation. As far as the process of data exploration and preparation is concerned, we will be touching upon these particular topics. We will



## MY CLASS NOTES

The first thing that we will do as far as sanity checks are concerned is to look at the names of the columns that have been read in. So I will use the *names* command to display the names of the column in the dataset that I just read in. Now an interesting thing that I see here is that the 12<sup>th</sup> column has name called “MonthlyIncome.1” and the 6<sup>th</sup> column has the name “MonthlyIncome”. So probably these two columns have the same values and they are duplicate columns and I might need to verify that. Also if I take a look at the names of the columns, the first column, I see the first column as “NPA.Status” and last column is “Good Bad” column.

If I go back to my data dictionary, I see that I don't have a column called "Good\_Bad" in my dataset. So the original dataset that was given to me did not have the column called "Good\_Bad". The dataset that I have just read in might have been manipulated before it reached me and



## MY CLASS NOTES

[illegible]

Also if I take a look at the `MonthlyIncome.1` column, it has these values and if I take a look at the `MonthlyIncome` column, it has these values which seem to be similar. Let me go back to my Data dictionary and see if I have a column called “`MonthlyIncome.1`”. As you can see there is just one column which talks about Monthly Income and the name of that column is “`MonthlyIncome`”. So there is no column called “`MonthlyIncome.1`” as far as the data dictionary is concerned. So I can conclude that these two are duplicate columns so my data was manipulated in some way before it reached me. So I’ll need take this variable out from my dataset.



The way I would take this variable out is by using this syntax - `(-c(1,12))`, because `NPA.Status` is the first column and `MonthlyIncome.1` is the 12<sup>th</sup> column. So I am removing the duplicate variables out of my dataset. Now as you can see, the number of columns has reduced 18 to 16.

Let us look at the summary statistics of all the columns in my dataset by using the **Summary** command. Summary command produces 5 point summary statistics for numerical variable and frequency distribution for categorical or factor variables. Now one thing that I notice from the summary statistics is that in my good bad column, I have two missing values. Whenever we have credit card kind of dataset where we have information about people who have defaulted and information about people who have not defaulted, the aim is to build a model to identify good and bad customers. This `Good_Bad` variable will become a dependent variable if I have to build a model on it.

Now, in my dependent variable, there are two missing values. So I'll need to take these two missing values out of my dataset. The way to do that is first of all identify, for which rows, I have missing values in the `Good_Bad` column and for that I'll use this command - `(which(is.na))` of my dependent variable which has missing values. So the index object here is going to store the row number where I have missing values in the `Good_Bad` column. I will take these rows out of my data by using the minus symbol in the row index for the data frame `cr`.



## MY CLASS NOTES

Now as you can see I have two less observations now; I have now 1,50,000 observations. Now what we will do is we will run individual summary for all the variables that we have in our dataset.

I will do a simple check on this variable and figure out how many observations are zero and how many observations have really large magnitude. So I will be using this simple dplyr syntax here, so this syntax here tells me how many rows I have in my data for which this variable has a value of zero. I can see there are around 10,000 rows in my data, which have value of zero for RevolvingUtilizationOfUnsecuredLines. Also since this is ratio variable or percentage variable, ideally its value should not be more than 1 or 2. So Let us take a look at how many observations in my dataset corresponding to this dataset have a



magnitude greater than 0.99. Again 14,000 observations have a value greater than 0.99.

Another thing that we can do it, we can do a percentile wise break up for this variable and take a look at how this percentile break up for this particular variable looks like. I will use a quantile command here. What this command does is, it tells me what percentage of my data has a particular magnitude. Here I see 99% of my data has a magnitude less than equal to 1.09 and there are just 1% observations, which have magnitude greater than this number. What I will do is, I will go back to my stakeholder and have a discussion with my holder and tell him that this is what I found out when I explored your data. So he might give me an input that “Usually this variable never has a magnitude greater than 2”. So based on this business knowledge that has been provided by my stakeholder, I will cap the values of this particular variable.

So first of all I will see how many rows in my dataset I have which have a magnitude which less than or equal to 2. There are 1,49,000 rows and the next thing I’m going to do is to filter out those observations from my data whose value is more than 2 for this particular variable. As you can see, I have reduced number of rows in my dataset, because I have taken out the values in my dataset for this particular variable whose magnitude is more than 2.

This is one of the ways we take a look at the quality of a dataset by just observing the summary



statistics. Then we go back to the stakeholder and ask them if the numbers make sense. Basically, they provide us with the business context and tell us where we can cap this variable and where should not cap this variable.

Let us take another example. Let us talk about the age variable. Before we talk about the Age variable, let us clear the console using CTRL + L. Let us take a look at the 5 point summary for the age variable. If I look at the age variable, another odd thing I notice here is that the minimum value for age variable is 0. This would imply that in my dataset, I have some people whose age is 0 and they still have a credit card. Now this does not make sense. It would mean that someone who is not even adult and he still has a credit card, which should not be the case.

So I will first of all find out how many observations in my data have a magnitude of 0. So I will use this dplyr syntax. What it does is it tells me how many rows in my data have a magnitude of 0. I see there is just 1 row. This could have been a data entry error and I can neglect this single observation. But, before I proceed further, I can do quantile wise analysis and see how the distribution of the age variable looks like. As you can see the first 1% of the observations are less than or equal to 24, which makes sense. And the maximum is 109; again I don't see anything really odd here. The only odd thing was that I have one observation in my dataset whose magnitude is 0. So I will filter out those rows where I have a value of 0 for Age



## MY CLASS NOTES

[illegible]