



# Credit EDA Case Study: Doubts Session

**Course :** Data Science

**Lecture On :** Credit EDA CS

**Instructor :** Sumit Shukla

# What we will cover in this session?

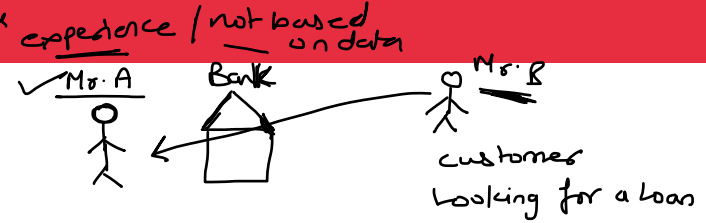
- 1 How to start with the “Credit EDA Case Study”
- 2 What are the important steps that should be included
- 3 Points to remember
- 4 Python Demo
- 5 QnA

# Problem Statement



You need to perform EPA and find out the various factors or variables that can help the bank to identify the person will default or not. **What to do?**

Data {   
 Loan approval   
 Loan rejection }



Two types of risks are associated with the bank's decision:

- If the applicant is likely to repay the loan, then not approving the loan results in a loss of business to the company.
- If the applicant is not likely to repay the loan, i.e. he/she is likely to default, then approving the loan may lead to a financial loss for the company.

**The company wants to understand the driving factors (or driver variables) behind loan default, i.e. the variables which are strong indicators of default**

- Reject Loan
- Approve Loan

Reject the application [ Mr. B is capable of paying back the loan. But Mr. A thinks that, he is not capable of paying back the loan ]

Approved the loan [ Mr. B is not capable of paying back the loan. But Mr. A thinks that he may repay the loan. ]

Interest loss →   
 credit loss

Data.   
 i) Application.csv   
 ii) prev.application.csv

# How to Solve?

[memory issue  
long time to execute]

Random sample from the complete data  
code

upGrad



## What to do?

Head  
Info  
Describe

1. Start by importing the 'application\_train.csv'.
2. Check the structure of the data (Normal routine check).
3. Data Quality Check and Missing values

- A. Find the percentage of missing values for all the columns.
- B. Remove columns with high missing percentage. ( $> 50\%$ )
- C. For columns which has less percentage (around 13% or so), you need to check what will be the best metric to impute the missing values? Like if the column you are checking is a categorical column check, which category you can use to fill the nulls. For others check does mean or median can be imputed or not. Others cases may be imputing with 0. You need to do this task for some variables and not for all, say 5.

- D. Check the datatypes of all the columns and change the datatype like negative age and date.
- E. For numerical columns check for outliers and report them for at-least 5 variables. Add observations and reasoning.

- F. Binning of continuous variables. Check if you need to bin any variable in different categories. Do this for at least 2 variables.

C: You need to only suggest/comment the best metric that can be used for imputation but you don't need to actually impute it. (5-columns)

Cx

C1  
C2  
C3  
NA

mode  
technique  
comment:

Cx-c

Have  
outliers

median  
else - mean

Demo  
Cx  
0  
2  
int column  
object

Demo  
Cx  
0  
2  
int column  
object

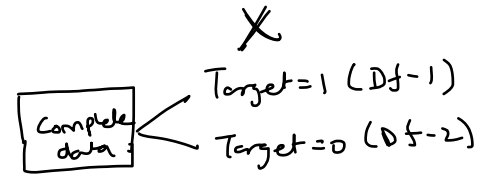
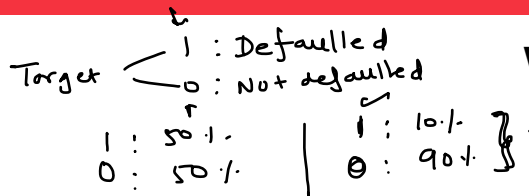
Conti → Ccateg.  
Income → Income-cat  
100 → 1000  
100-250 → Low income

150-250 : Low  
251-500 : Media  
501-750 : Above m  
751-1000 : High inc

# How to solve..Continue.

## What to do?

### 4. Analysis



Check the Imbalance percentage. No balancing technique required.

Divide the data into two sets, i.e. Target=1 and Target=0.

Perform univariate analysis for categorical variables for both 0 and 1. Compare the target variable across categories of categorical variables.

Find correlation for numerical columns for both the cases, i.e. 0 and 1.

Check the variables with highest correlation are the same in both the files or not?

Perform univariate for numerical variables for both 0 and 1. Compared the target variable across categories of continuous variables.

Perform bivariate analysis for numerical variables for both 0 and 1.

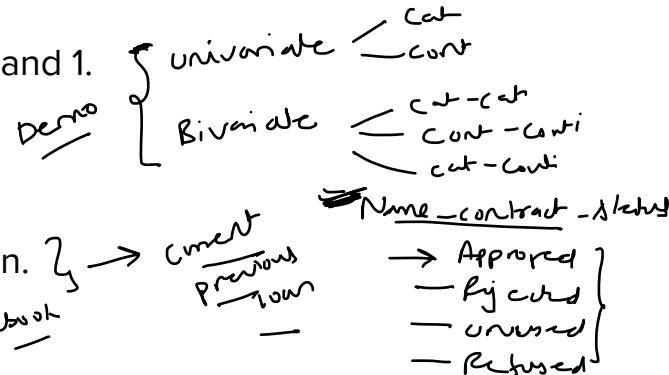
### 5. Read "Previous Application" data.

You can merge the files, but it's a challenge.

Perform univariate and bivariate analysis to find some pattern.

### 6. Final words

Based on your analysis, define the results and conclusion.



- Keep in mind “There is no correct or incorrect solution”.
- Every approach is correct if you are able to answer all the questions asked.
- The main objective of this case study is to learn and implement EDA techniques. So don't focus so much on columns and their descriptions.
- Remember you need to use plots and then understand the pattern. Then report your analysis in your notebook or PPT. No marks will be awarded if you have just plotted so many variables and have not explained the pattern.
- It's not possible to cover all the columns, so try to cover some of them. Based on the plots you get, try to identify the important variables.

PPT → Tech Business ] 10-15  
→ Explain some important plots  
→ find recommendations from this case study.

My data contains two classes 0 and 1. The total number of rows in my data is 150. The count of rows for class 1 is 40. What is my Imbalance percentage for class 1?

- A. 47%
- B. 38%
- C. 27%
- D. 26%

I have two files df1 and df2, both files have a common primary column as "ID". I need all rows from my df1 data and only the common rows from the df2 data. Which join should I use?(df1.merge(df2))

- A. Right
- B. Inner
- C. Left





# Python Demo. Let's Code



Thank You!