



Credit EDA Case Study: Doubts Session

Course : Data Science

Lecture On : Credit EDA CS

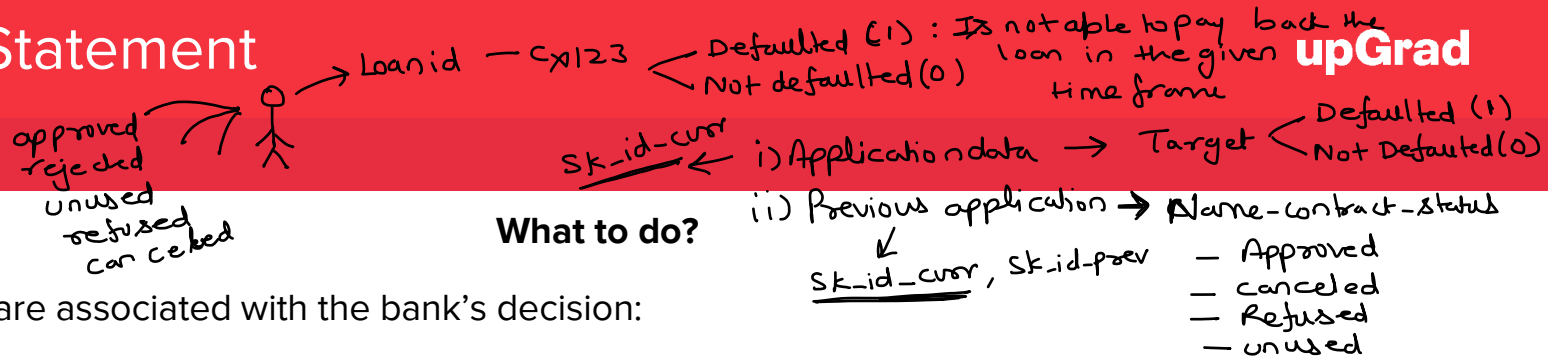
Instructor : Sumit Shukla

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

→ Steps solve this CS
+
→ Demo Python

Problem Statement

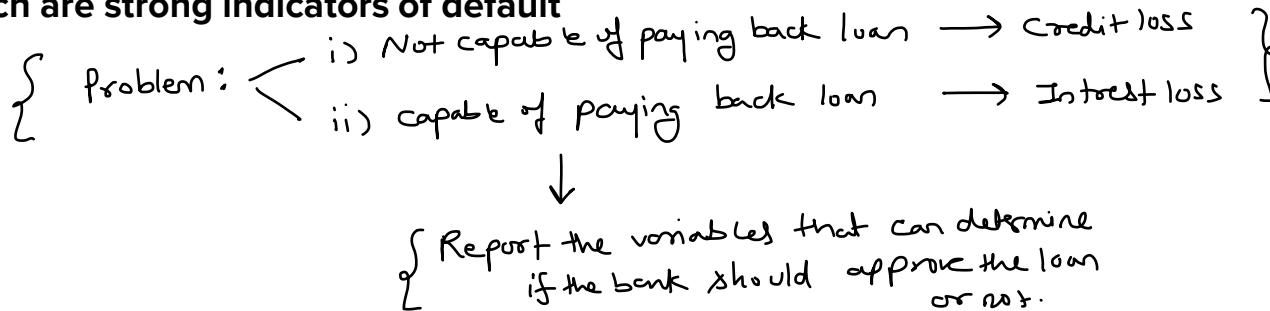


What to do?

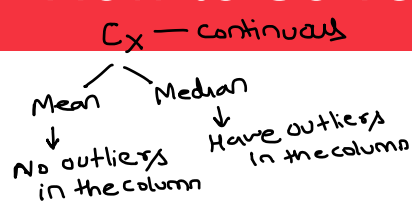
Two types of risks are associated with the bank's decision:

- 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



How to Solve?



Cx — categorical
mode

memory issue

What to do?

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

data

- df.shape
- df.dtypes
- df.describe()
- df.info

- 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.

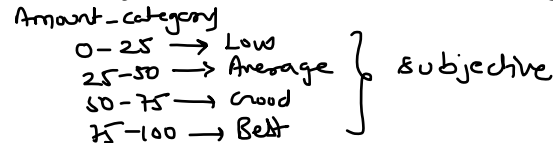
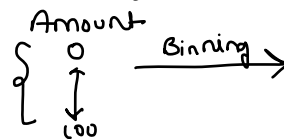
code demo

Categorical → object
Continuous → int, float

- 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.

make it as positive

load the date column as pandas datetime format



3.C: You don't need to actually impute the data, but report the best value that can be used to impute the missing data (At least 5 columns) + comments, markdown

df.isna()

What to do?

↓
- (age) = the age

4. Analysis Target < $\begin{matrix} 1 : 40\% \\ 0 : 60\% \end{matrix}$ $\left\{ \begin{matrix} 100 \text{ Rows} \\ 0 : 60 \\ 1 : 40 \end{matrix} \right\}$ $\begin{matrix} 1 : 50\% \\ 0 : 50\% \end{matrix}$

- Check the Imbalance percentage. No balancing technique required.
- Divide the data into two sets, i.e. Target=1 and Target=0. (Application data dataframe)
 - Target 1 Df ✓
 - Target 0 Df ✓
- Perform univariate analysis for categorical variables for both 0 and 1. Compare the target variable across categories of categorical variables.

- code demo
- 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. + Application data

$\left\{ \begin{matrix} 25-30 \\ \text{columns} \end{matrix} \right\}$
 $\left\{ \begin{matrix} \text{univariate} \\ \text{Bivariate} \end{matrix} \right\}$
 $\left\{ \begin{matrix} 10 \text{ continuous} \\ 10 \text{ categorical} \end{matrix} \right\}$
 $\left\{ \begin{matrix} 10 \text{ continuous} \\ 10 \text{ categorical} \end{matrix} \right\}$
 with respect to the target Column

combined-df

- You can merge the files, but it's a challenge.
- Perform univariate and bivariate analysis to find some pattern.

univariate } with respect to
 Bivariate } name contract status
 Approved Rejected unused
 loan application

6. Final words + Report the final set of variables that can be used by the bank for approving or rejecting any loan application


Jupyter notebook + conclusion

- Based on your analysis, define the results and conclusion.

What to keep in mind

(PPT) → PDF → Major Analysis + conclusion + Recommendation
to the Bank (up to 15)

ZIP ← Jupyter notebook
PPT → PDF

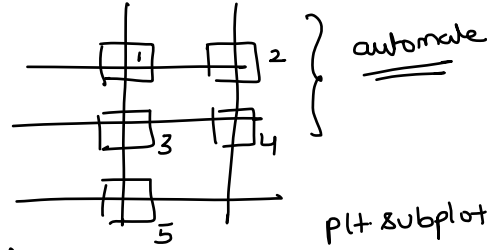
- Keep in mind “There is no correct or incorrect solution”.
- Every approach is correct if you are able to answer all the questions asked.
- The main objective of this case study is to learn and implement EDA techniques. So don't focus so much on columns and their descriptions. } 80+ columns (25-30)  } Comments explaining Notebook
- 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.

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

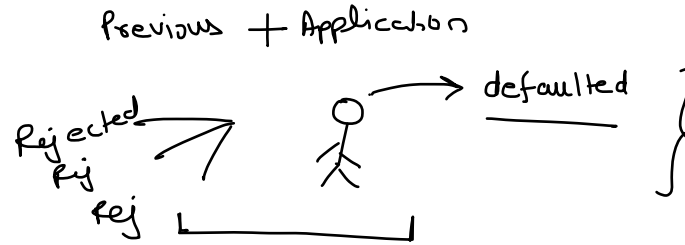
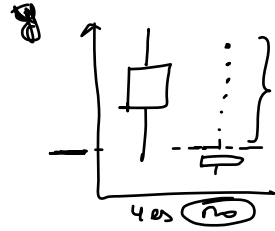


`plt.subplot(# of Rows, # of columns
/ figure number)`

`plt.subplot(3,2,1)`
Code → plot 1
`plt.subplot(3,2,2)`
Code → plot 2



Python Demo. Let's Code



`dropna()`

Thank You!