

Exploring Credit One Data

Project Overview

Credit One has recently seen an increase in the number of customers defaulting on their loans. The company is interested in understanding if it is possible to ensure customers will pay their loans and if possible, how to do so.

Dataset Details

Credit One provided a data set with 30,202 transactions and 25 features. After some cleaning efforts, the dataset was cut down to 30,000 transactions (202 duplicates or unintentional rows dropped) and 24 features (the ID feature was discarded). Loosely speaking, the features broke down into two types, customer demographics or loan information.

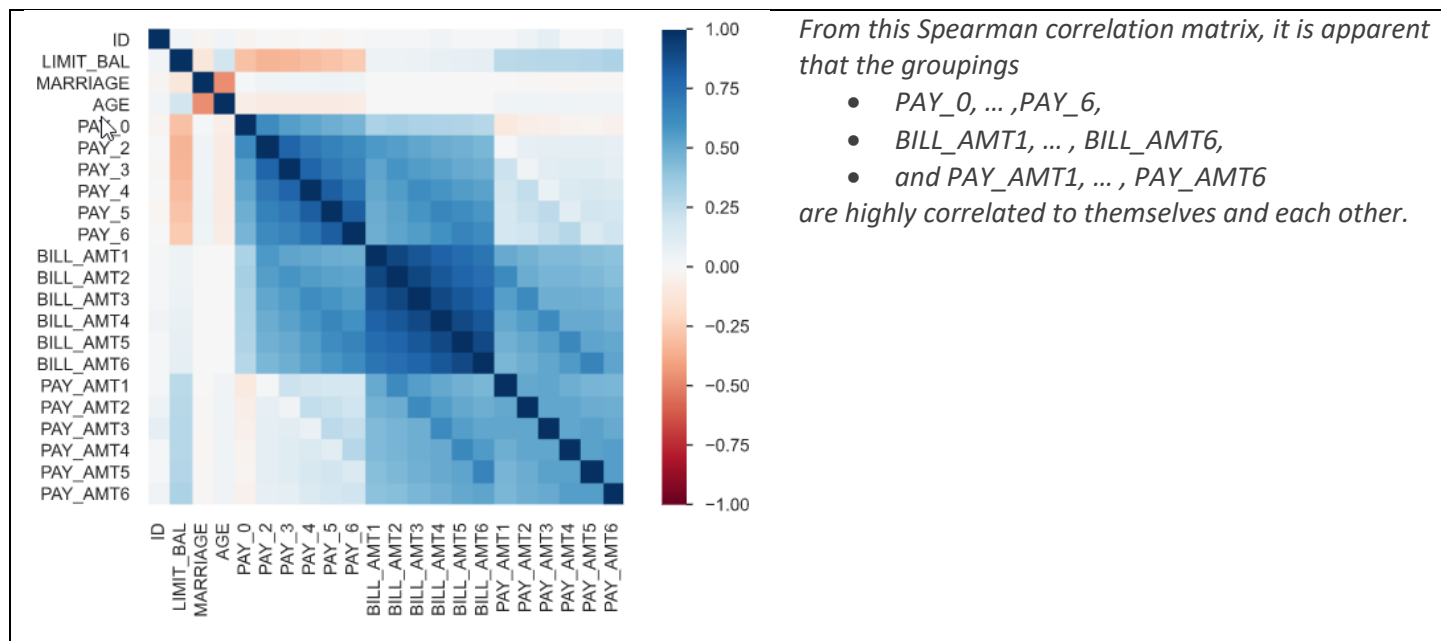
| | |
|---------------------------|--|
| Demographic Features | SEX, EDUCATION, MARRIAGE, AGE |
| Loan Information Features | PAY_0, PAY_2, PAY_3, ... , PAY_6, BILL_AMT1, BILL_AMT2, ... , BILL_AMT6, PAY_AMT1, PAY_AMT2, ... , PAY_AMT6, LIMIT_BAL, DEFAULT PAYMENT NEXT MONTH |

A full description of these features can be found at the end of the report

Analysis

Correlations

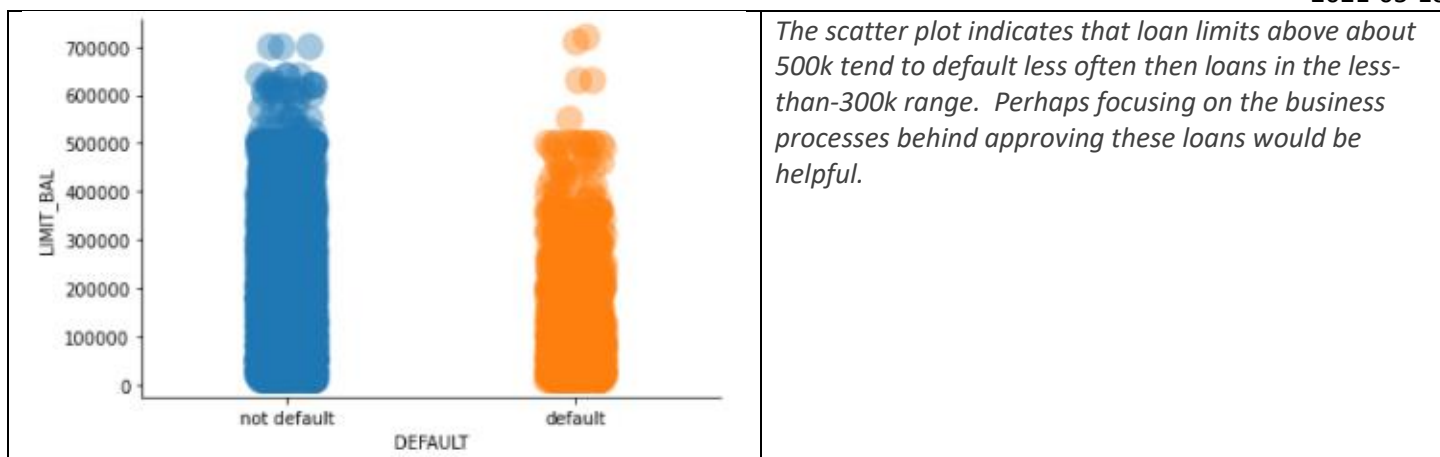
After using Pandas profiling to find any obvious patterns in the data, it became apparent that many of the loan information features were highly correlated to one another.



For future analysis, it could be helpful to remove or aggregate those strongly correlated features.

Larger loan limits

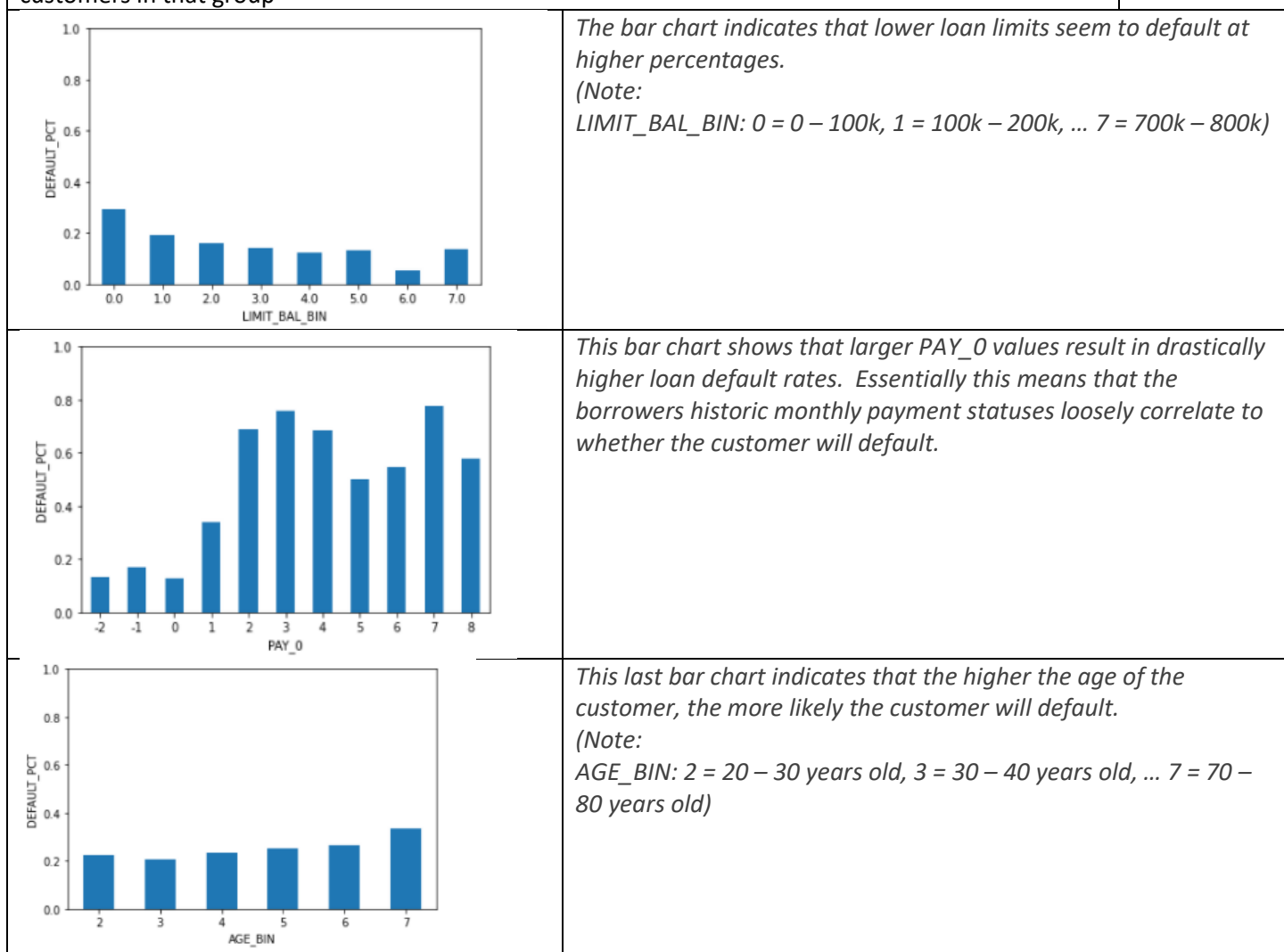
When comparing default status to the loan limit balance, the dataset seems to indicate that higher loan balances do not present a high risk of default.



Default percentages by group

Finally, to help understand which groups are defaulting the most, a new variable was added to the dataset:

DEFAULT_PCT = sum of customers who defaulted in a particular group divided by total number of customers in that group



Conclusions

Through the exploratory data analysis, it became clear that many of the loan information variables were correlated and will need to be reduced before further modeling can occur. Also, it seems clear that loans over \$500K default less frequently than loans under \$300k. Finally, age, payment status, and loan limit seem to be linked to how likely the customer is to pay back their loan.

Feature Descriptions

LIMIT_BAL: Amount of the given credit (NT dollar): it includes both the individual consumer credit and his/her family (supplementary) credit.

SEX: (1 = male; 2 = female).

EDUCATION: (1 = graduate school; 2 = university; 3 = high school; 0, 4, 5, 6 = others).

MARRIAGE: (1 = married; 2 = single; 3 = divorce; 0=others).

AGE: (year).

PAY_0 – PAY_6: History of past payments. (-2: No consumption; -1: Paid in full; 0: The use of revolving credit; 1 = payment delay for one month; 2 = payment delay for two months; . . .; 8 = payment delay for eight months; 9 = payment delay for nine months and above.)

BILL_AMT1 – BILL_AMT6: Amount of bill statement (NT dollar).

PAY_AMT1- -PAY_AMT6: Amount of previous payments (NT dollar).

DEFAULT PAYMENT NEXT MONTH: (default, not default)