# UNIVERSIDAD CENFOTEC

**Data Analytics / Big Data 2018c-3**

Program 5 Task 3

## Data Science with Python

**Build and Evaluate Models**

PRESENTA:

**Alberto Madrigal Jiménez**

ASESOR:

**Gustavo Rojas**

San José, CRC.                                                    2020

# Credit One Report

## Main Problem:

An increase in customer default rates is bad for Credit One since its business is approving customers for loans in the first place. This is likely to result in the loss of Credit One's business customers.

## Proposed Solution:

Make predictions with the Default Payment Field, in order to understand the behavior of the customers answering the following questions.

1) How do you ensure that customers can/will pay their loans?

   With the information provided in the dataset as: sex, education, marriage, age, the last 6 month pays, and the last 6 billing amounts, we can predict the customer behavior, if they will default the loan or not.

2) Can we approve customers with high certainty?

   Based on the story provided by the data, we can understand the behavior of the current customer with at least six month with historical information. We can predict and say with 82% of accuracy the answer about the default loan.

# Cleaning and Pre-processing the Dataset

The following processes were applying over the data to clean and show the final dataset, we did not add a lot of technical information in this section, because, it information can be checked in the python doc.

### 1) Fit

In the following image we can see how using fit function, we can fit an estimator to be able to predict.

Example:

```
In [98]: #SVR
         modelSVR.fit(X_train,y_train)
Out[98]: SVC(C=1.0, cache_size=200, class_weight=None, coef0=0.0,
             decision_function_shape='ovr', degree=3, gamma='auto', kernel='sigmoid',
             max_iter=-1, probability=False, random_state=0, shrinking=True, tol=0.001,
             verbose=False)
```

### 2) Mapping to a Uniform distribution

Function: train_test_split is for splitting a single dataset for two different purposes: training and testing. The testing subset is for building your model. The testing subset is for using the model on unknown data to evaluate the performance of the model.

Example:

```
X_train, X_test, y_train, y_test = train_test_split(features, depVar, train_size=0.6, random_state=0)
print(X_train.shape)
print(X_test.shape)

(17999, 24)
(12000, 24)
```
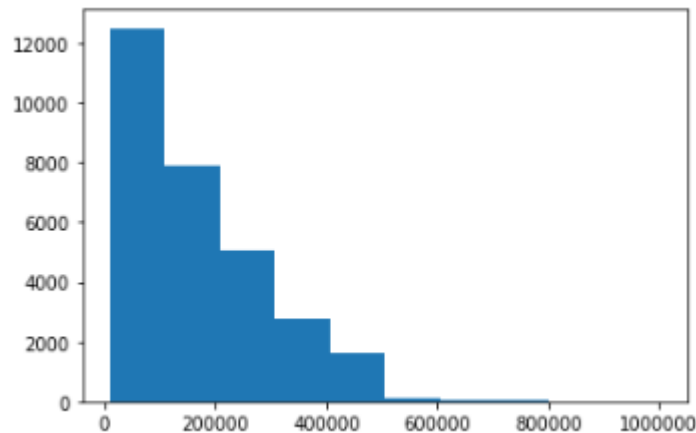
### 3) Encoding categorical features

Categorical Data: Often in real-time, data includes the text columns, which are repetitive. Features like gender, country, and codes are always repetitive. These are the examples for categorical data.

```
#Setting variables to categorical
rawData['SEX'] = rawData['SEX'].astype('category')
rawData['EDUCATION'] = rawData['EDUCATION'].astype('category')
rawData['MARRIAGE'] = rawData['MARRIAGE'].astype('category')
rawData['DPNM'] = rawData['DPNM'].astype('category')
rawData['AGE'] = rawData['AGE'].astype('category')
```
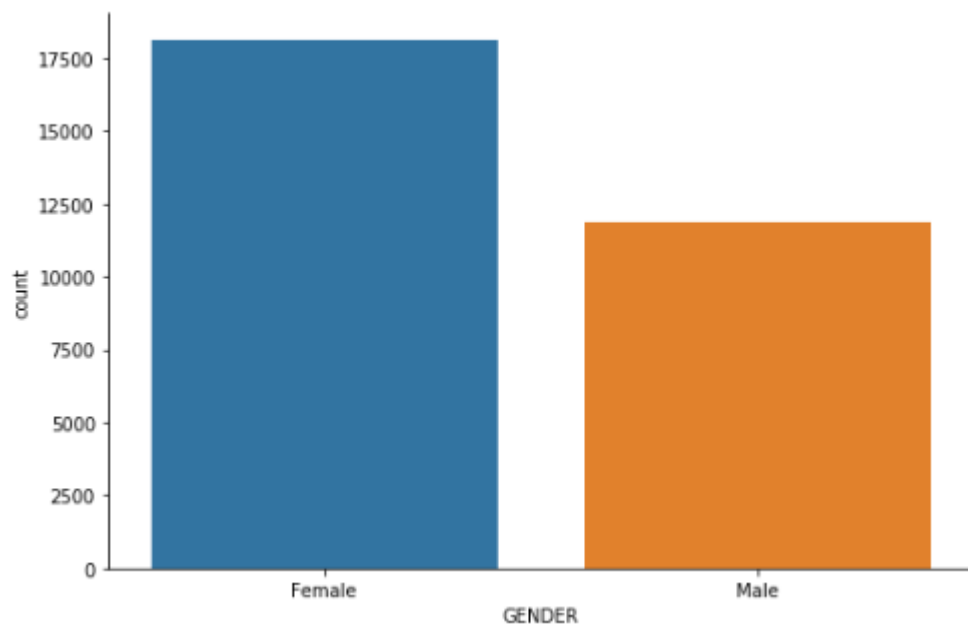
# EDA

Exploratory data analysis (EDA) is an approach to analyzing data sets to summarize their main features, often with visual methods. A statistical model can be used or not, but primarily EDA is for identifying what the data can tell us beyond the formal modeling or hypothesis testing task.
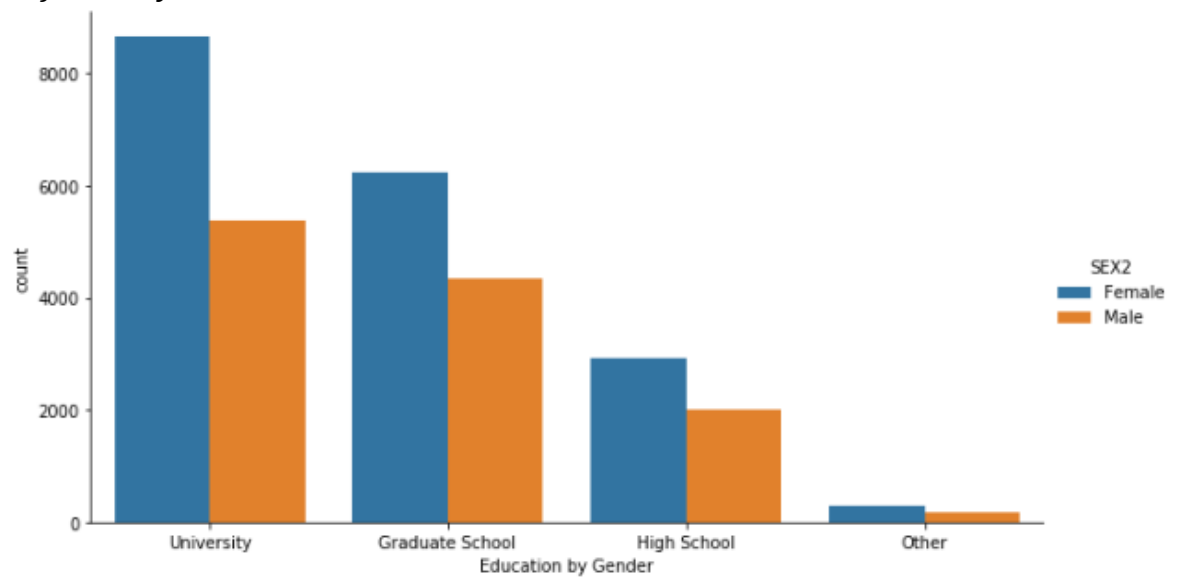
*Balance of the limits:*



- o   As we can see in the previous chart, the highest count of amount are between zero and 200000.
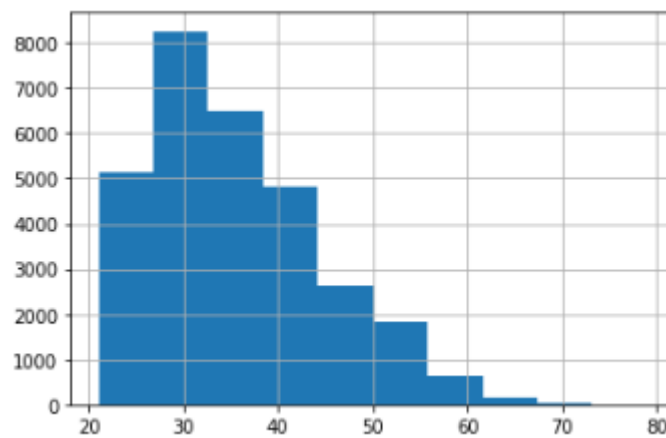
*Count of rows by Gender:*



- o   The Female count is higher than Male count

## *Analysis of Count by Education*



- o The previous image shows how the university grade school has the highest count of people, also, the female domains the volume.

## *Distribution Clients by Ages*



- o The highest numbers of people are around of 30 years old.

# Classification Methods Applied

**SVM (SVC):**

Vector support machines, support vector machines, or support vector machines are a set of supervised learning algorithms. These methods are properly related to classification problems.



**Outcome:** But the results obtained are not the best, because this method is not properly due to, all the values obtained are zero.

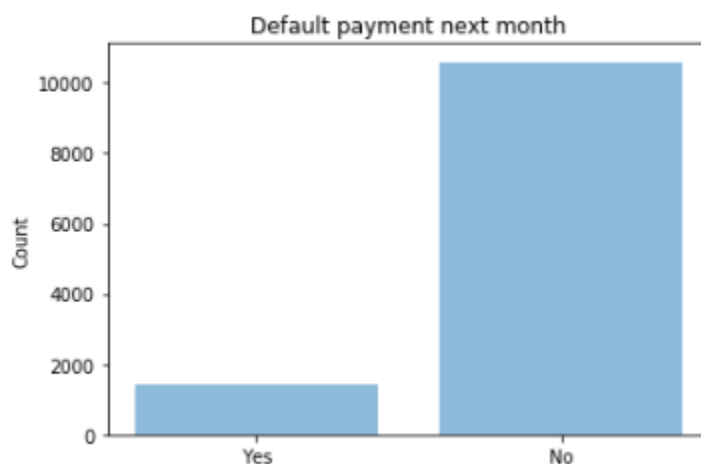**Random Forest Classifier (Classification Tree):**

Random forest also known in Spanish as '"Random Forest"' is a combination of predictor trees such that each tree depends on the values of a random vector tested independently and with the same distribution for each of these.



**Outcome:** the results obtained look pretty nice, with high percent of accuracy on the results.

**Logistic Regression:**

Logistic Regression is a Machine Learning classification algorithm that is used to predict the probability of a categorical dependent variable. In logistic regression, the dependent variable is a binary variable that contains data coded as 1 (yes, success, etc.) or 0 (no, failure, etc.).



**Outcome:** But the results obtained are not the best, because this method is not properly due to, all the values obtained are zero.

## Method Selected:

After applied all the methods, the best method selected was the Random Forest, we achieve with **82% of accuracy** the results, and also related with the higher indicator we got **81% of the better performance model**.

*Model used through thru the process:*

| | Higher is an indicator of a better performing model | Score | Accuracy |
|---|---|---|---|
| RANDOM FOREST | 81% | 100% | 82% |
| LINEAR REGRESSION | 78% | 78% | 78% |
| SVR | 77% | 78% | 78% |

The results obtained with the test data shows the following results.

Default payment next month

## Observations:

1) On this project we did not delete any column, based on the Correlation Matrix, we saw some weird behaviors, with high level of correlations between the independent variables, but, we did not consider drop them off due to the dataset lost the sense, basically, we need to understand the customer and the payment, when we deleted some columns, we lost the north of the investigation and we obtained abnormal results.

Correlation Matrix Image:

2) In the data source document where this explains every single field, we found some duplicated information in the Education field.  To solve this situation, and improve our analysis, we replace all the values in 4, 5 and 6 by 0. With this change we could get better results and minimize the errors in the results.

X3: Education (1 = graduate school; 2 = university; 3 = high school; 0, 4, 5, 6 = others).

## Recommendations:

- If the company make the decision and put in action the new model, we can improve the current process to define the default and no default customers.
- The machine learning is a continuous improvement process where it very important understands every single day the new features and improves our models.
- The method selected it has 82% of accuracy, and the final results are strong and consistent, consider that a new tool in the way to answer the current and the future questions, but remember to combine this results with the feedback of the Subject Matter Experts.

## End.