

# **Supervised Learning Capstone**

## **Loan Risk Assessment**

By Sebastian Rosado



---

1

# Introduction

Motivation, Data Details



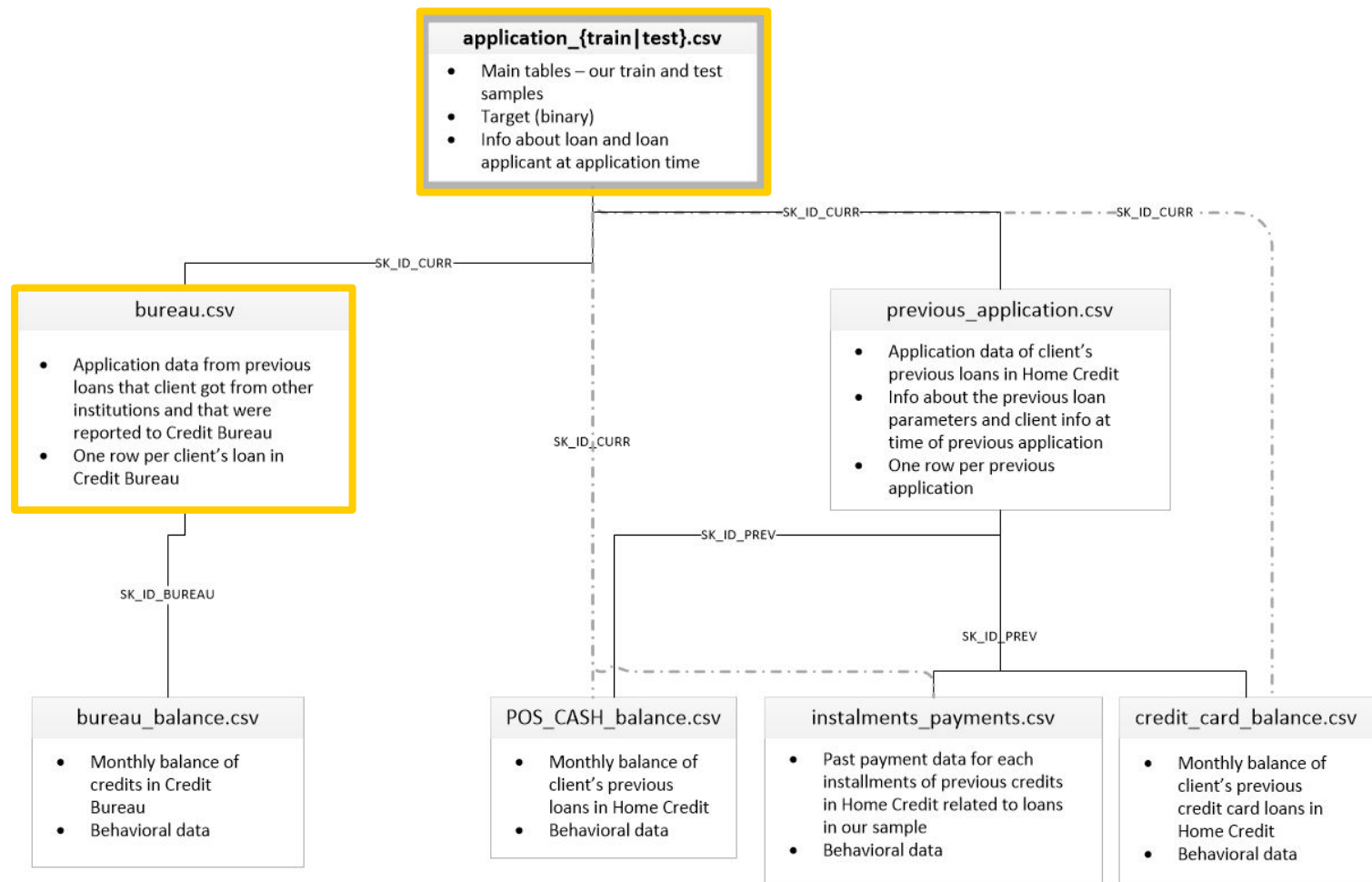
## Motivation

- ◉ What variables are the most predictive of a client's difficulty to pay back their loans?
- ◉ Banking the unbanked
  - Machine learning + finance
- ◉ Real data
  - Direct application
  - Challenging
- ◉ Interesting, unconventional variables
  - Apartment sizes, social circle creditworthiness



## Description of the Data

- Home Credit → non-bank financial institution
  - Applicants with little / no credit history
- Two merged DataFrames (one-to-one)
  1. Clients with active Home Credit loans
  2. All previous loans for those clients
- 236,630 rows, 194 columns
- Target: Payment Difficulties (1/0)



---

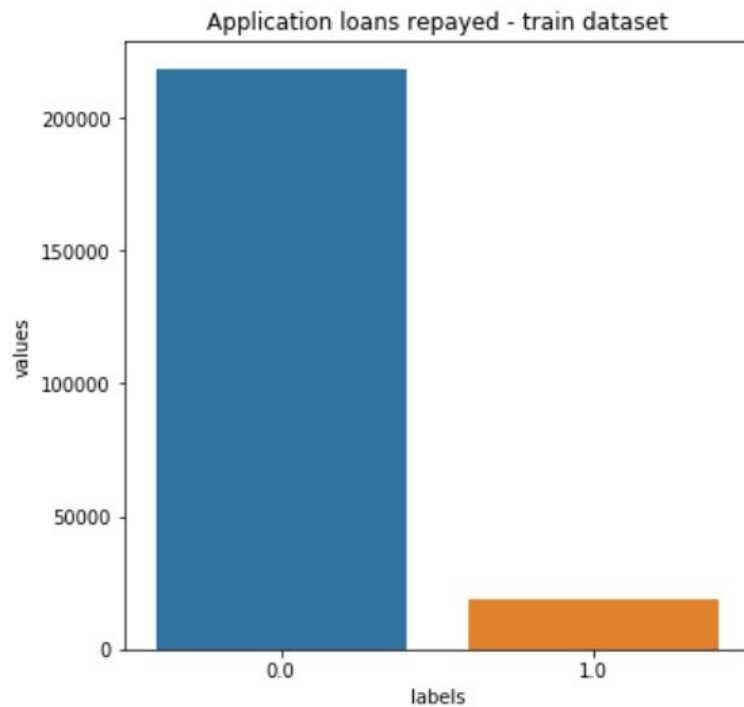
2

## **Data Exploration & Cleaning**

---



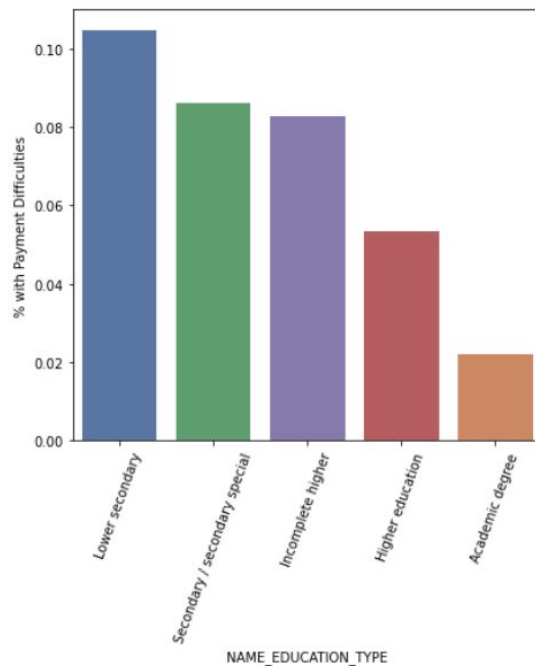
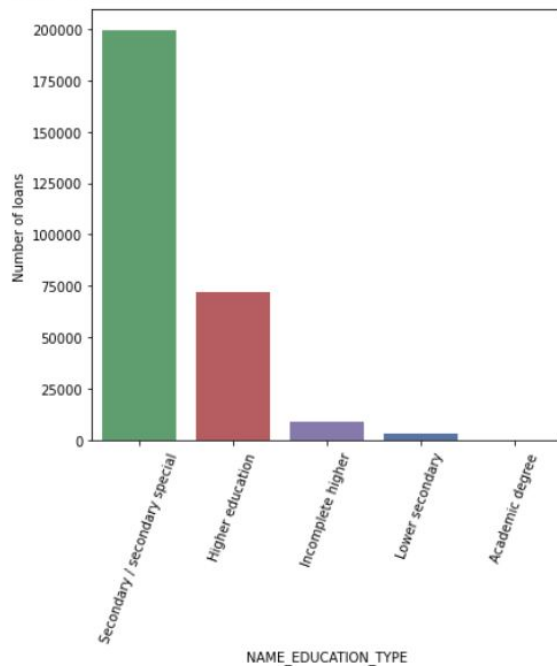
## Target Variable





# Educational Achievement

Description of NAME\_EDUCATION\_TYPE: 13      Level of highest education the client achieved  
Name: Description, dtype: object

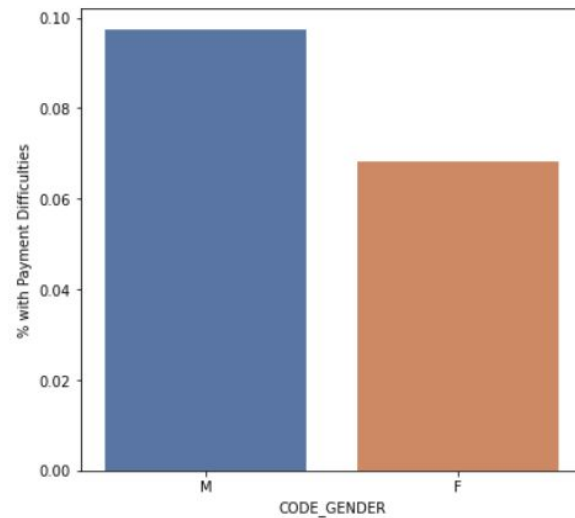
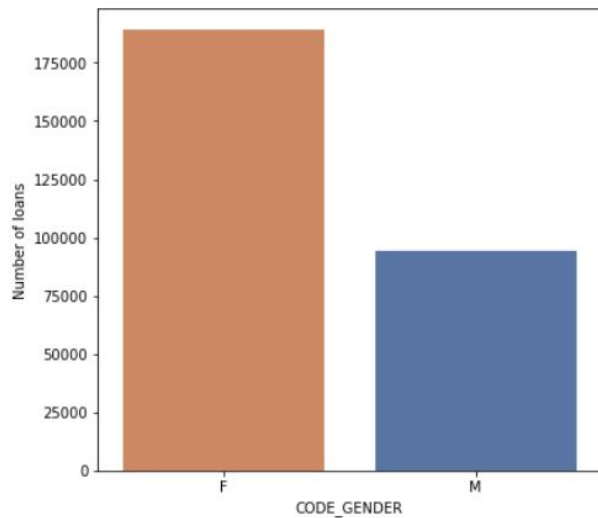






# Gender

Description of CODE\_GENDER: 3 Gender of the client  
Name: Description, dtype: object

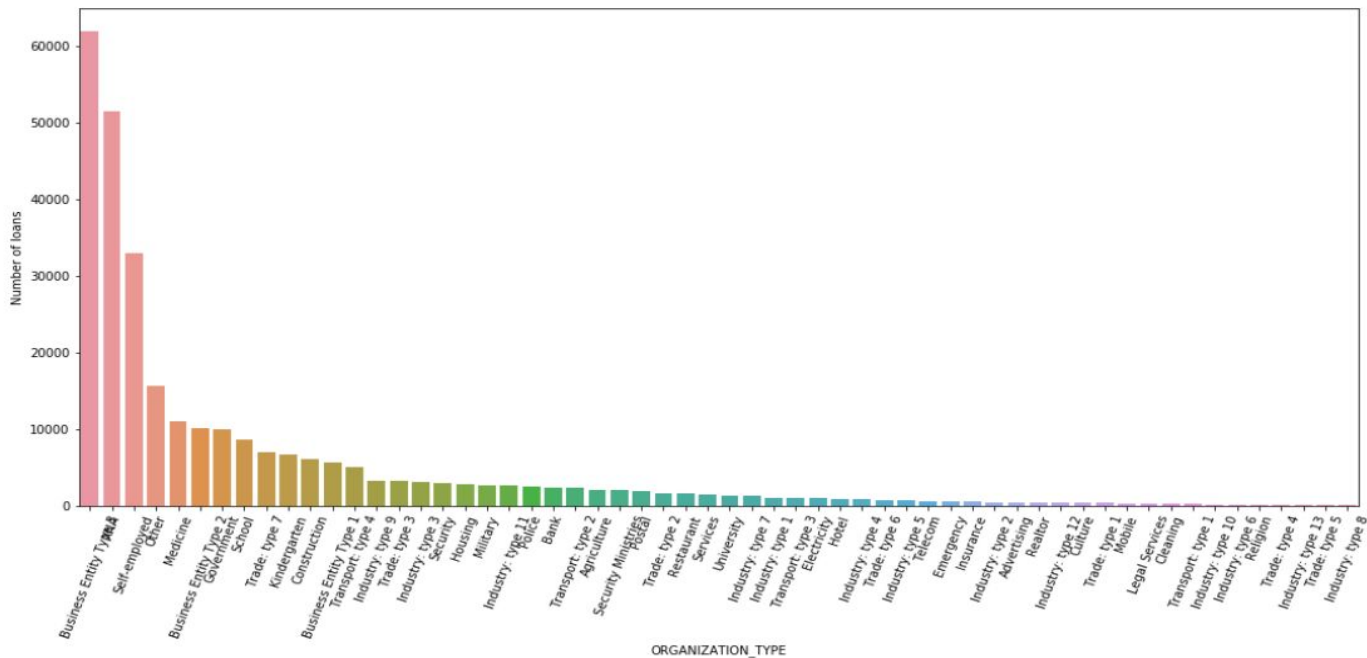




# Organization Type

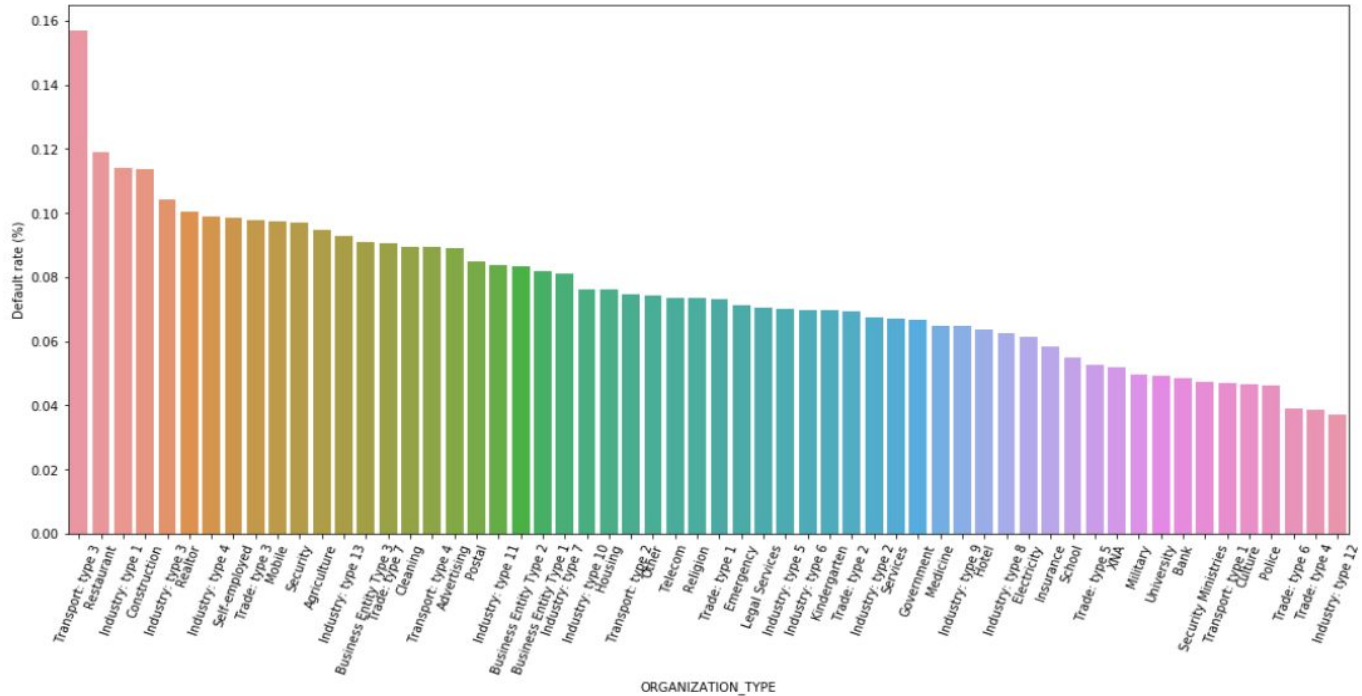
Description of ORGANIZATION\_TYPE: 40  
Name: Description, dtype: object

Type of organization where client works



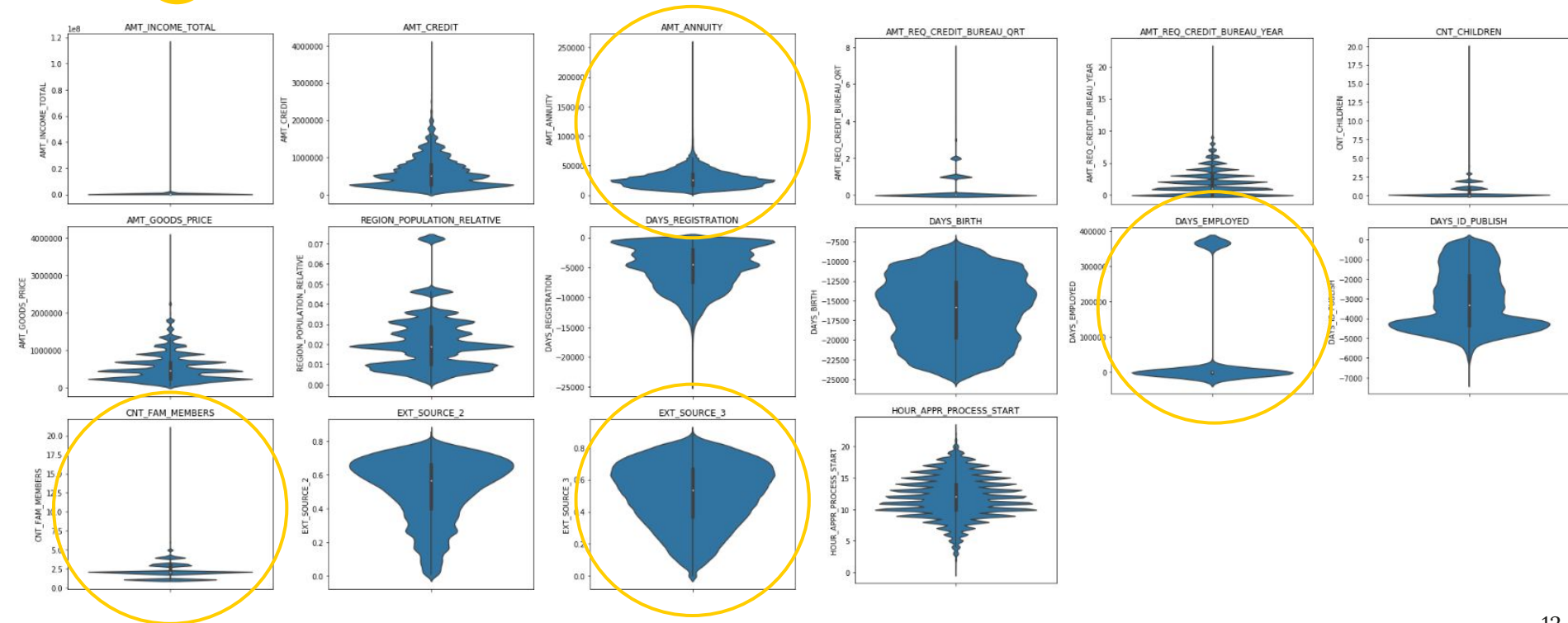


# Organization Type





# Continuous & Discrete



3

## **Feature Engineering**



## Created Features

- Third Degree Polynomial Transformations:
  - 'EXT\_SOURCE\_2', 'EXT\_SOURCE\_3', 'DAYS\_BIRTH', 'bureau\_DAYS\_CREDIT\_mean'
- Financial Variables
  - Annuity payment as % of income (kept)
  - Credit payments as % of income
  - Credit overdue as % of total credit
- PCA
  - Tried four, kept one (['AMT\_INCOME\_TOTAL', 'AMT\_ANNUITY', 'AMT\_CREDIT', 'AMT\_GOODS\_PRICE'])

---

4

## **Modeling & Tuning Process**

---



## Modeling Process

	SMOTE-Balancing	Imbalanced Original
Random Forest	Y	Y
Logistic Regression	Y	Y
XGBoost	Y	Y
Dimen.- Reduced XGBoost	Y	N
Dimen.- Reduced RF	Y	N



---

5

# Results

---



## Metric Comparison - Recall

	SMOTE-Balancing	Imbalanced Original
Random Forest	87.75%	0.0%
Logistic Regression	54.41%	0.01%
XGBoost	92.86%	0.17%
Dimen.- Reduced RF	84.62%	N
Dimen.- Reduced XGBoost	90.2%	N

---

6

# Conclusions

Key Takeaways and Lessons Learnt



## Key Takeaways

- Most Important Variables:
  - Completed Higher Ed. (26.37%)
  - Owns phone (7.95%)
  - Applied on Sunday (5.88%)
  - Active external loans (5.64%)
  - Owns a car (5.51%)
- Best Model (judging efficiency vs. score)
  - Dimen. Reduced XGBoost



## Lessons Learned

- SMOTE increases recall, though sometimes at the expense of accuracy
- Averaging feature importances reduces blind spots (don't trust a single algorithm)
  - This can be used for efficient dimensionality reduction
- Visualization functions have a high ROI
- Good data science is a marathon, not a race



## Room for Improvement

- Precision-recall curves
- More feature engineering
- Add other models (e.g. SVM & KNN)

---

7

**Thank You**

---



## Data Source

- ◉ <https://www.kaggle.com/c/home-credit-default-risk>