# Agenda

- **The company**

- **Problem Statement**

- **Data Sources & Tools**

- **Hypothesis & Goals**

- **Exploratory Data Analysis**

- **Feature Engineering**

- **Model Building**

# 1.
# The company

**$5.2 billion**
Revenues!

**132,000 employees**
And a lot of customers!

**1997**
Founded in Czech Republic

# Maps

our office

10 Countries

**Broaden financial inclusion to provide comfortable and safe borrowing experience**

**Focuses on the clients with little to no credit history**

**Transactional information, annual income, family status, housing type, etc. in order to predict their clients' repayment abilities**

# 2.
# Problem Statement

Hypothesis: Clients in careers with historically worse job security are most likely to default on their loan payments

Hypothesis: Clients with many previous credits are more likely to default on loans

Goal: Establish a trustworthy algorithm to validate/invalidate these claims and reveal other trends among the clientbase

Goal: Communicate results of said algorithm in a comprehensible manner

# 3.
# Data Source and Tools

Data Source - **Kaggle**

Data Processing and modelling - **Pyspark and Python on the Databricks platform**

Data visualization - **Tableau , Draw.io**

## Bureau

All clients previous credits provided by other financial institutions

## Bureau Balance

Monthly balances of previous credits in the credit bureau

## POS_Cash

Monthly balance snapshots of previous POS and loans

## Credit Card Balance

Monthly balance snapshots of previous credit cards owned by the applicant

## Previous Applications

All previous applications for loans by the client

## Installment Payments

Repayment history for previously disbursed credits

## Applications

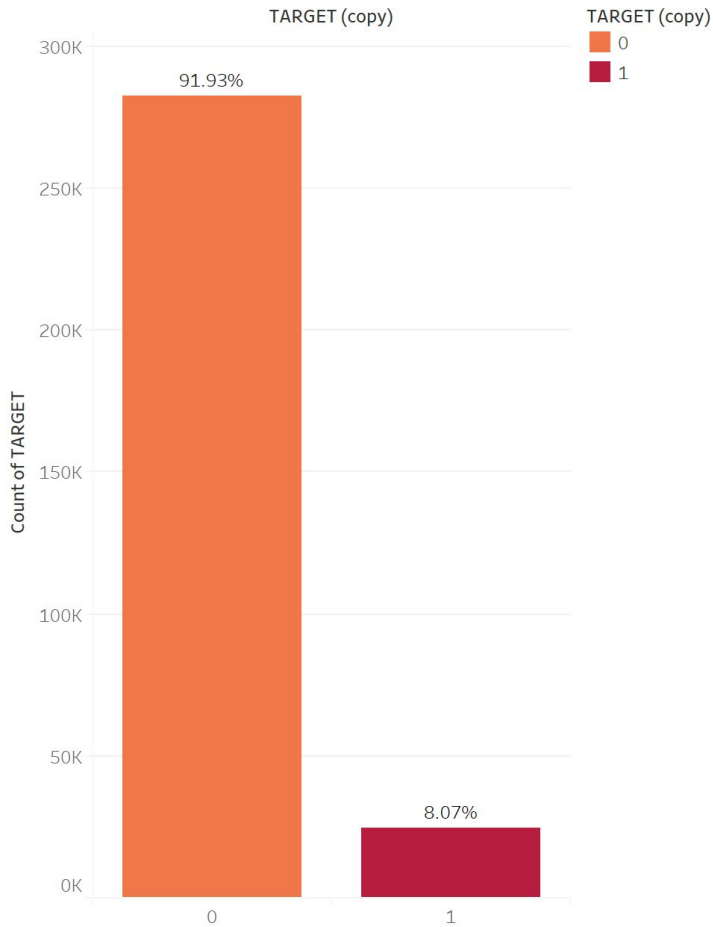Main table depicting current loan applications for each applicant
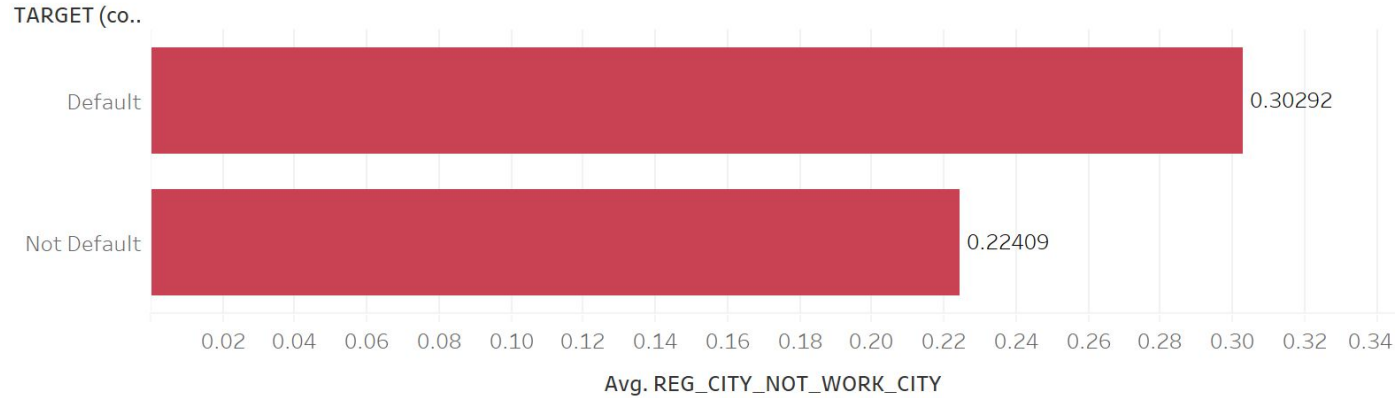
# 4.
# Exploratory Data Analysis

# Distribution of the Default Data



Count of TARGET for each TARGET (copy). Color shows details about TARGET (copy). The marks are labeled by % of Total Count of TARGET.

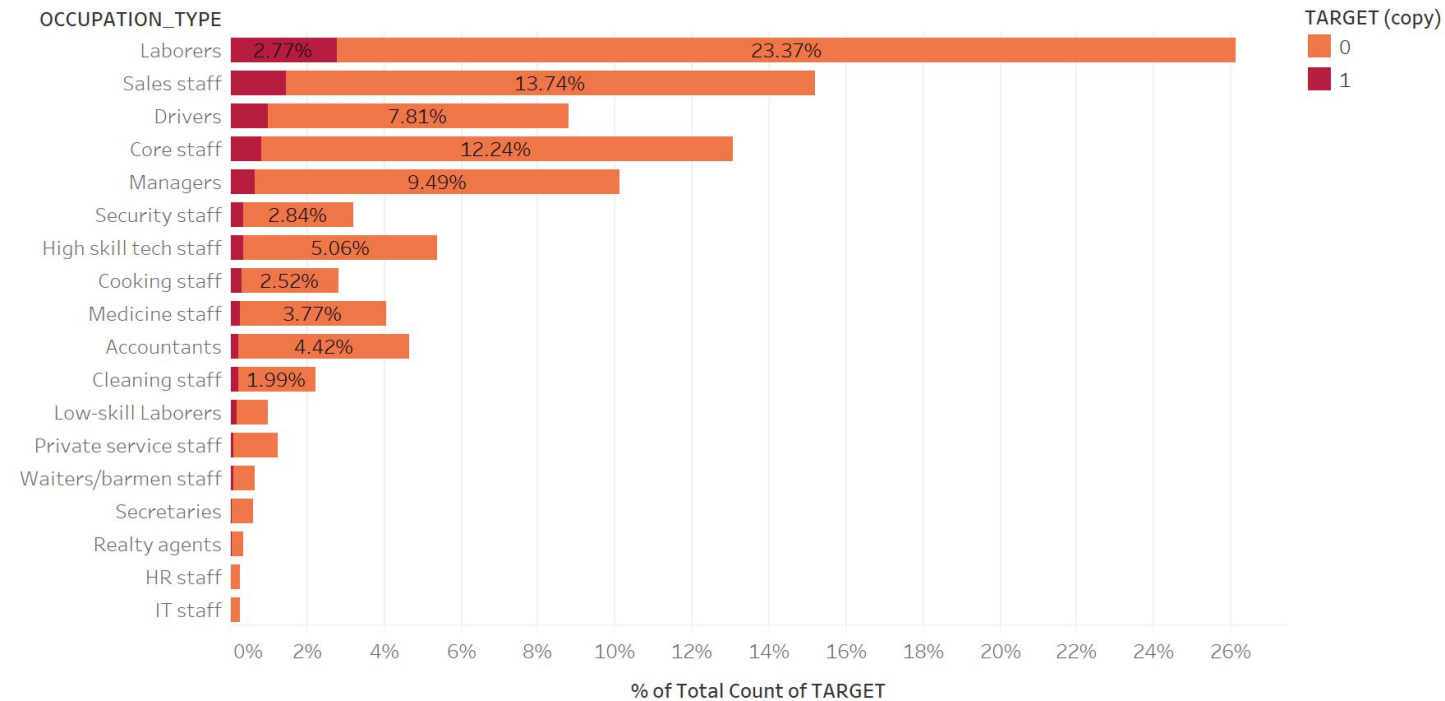## Proportion of discrepancy in residence and work location

**TARGET (co..**



Average of REG_CITY_NOT_WORK_CITY for each TARGET (copy). The marks are labeled by average of REG_CITY_NOT_WORK_CITY. The view is filtered on TARGET (copy), which keeps Not Default and Default.

15

Hypothesis 1:

Are clients with historically worse job security more likely to default on their loan payments?
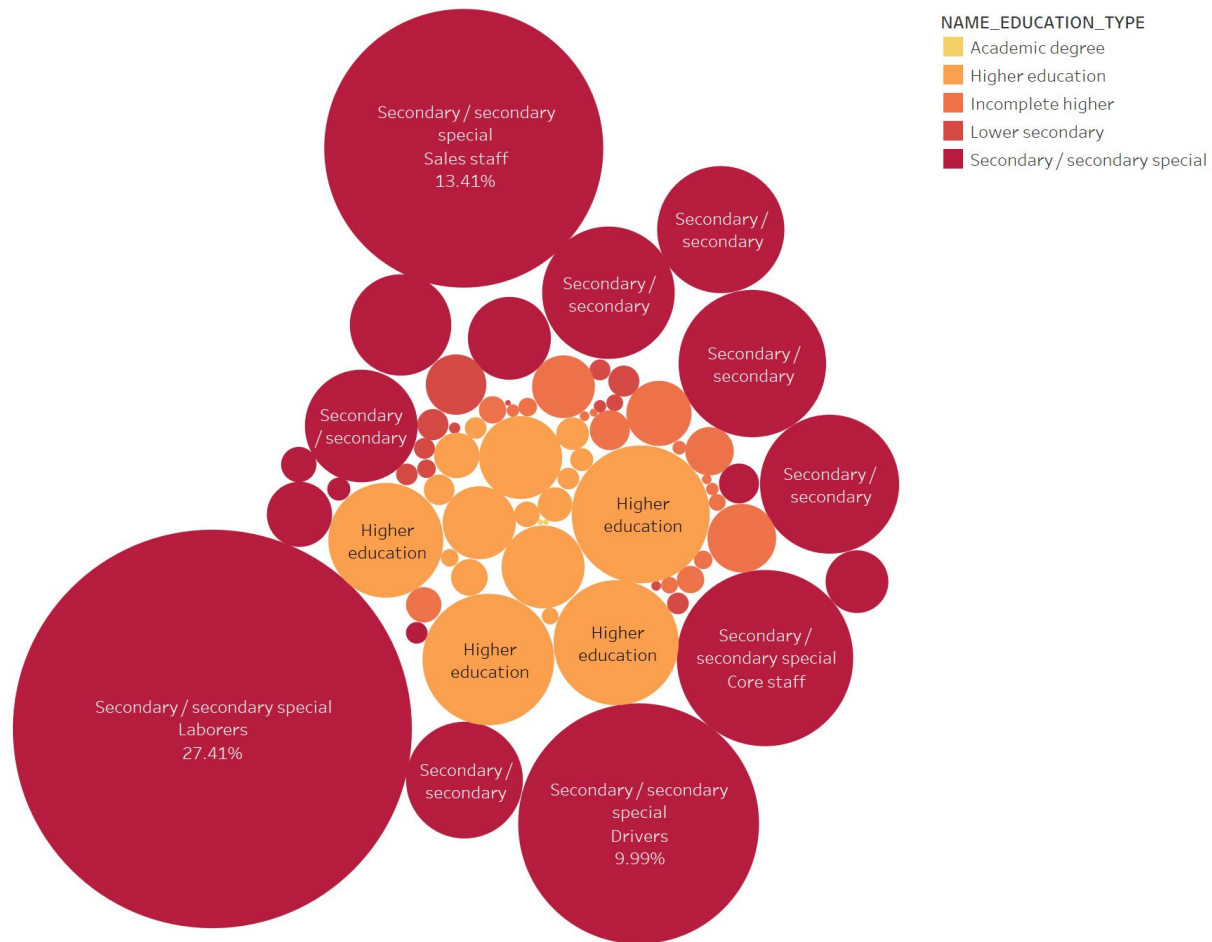
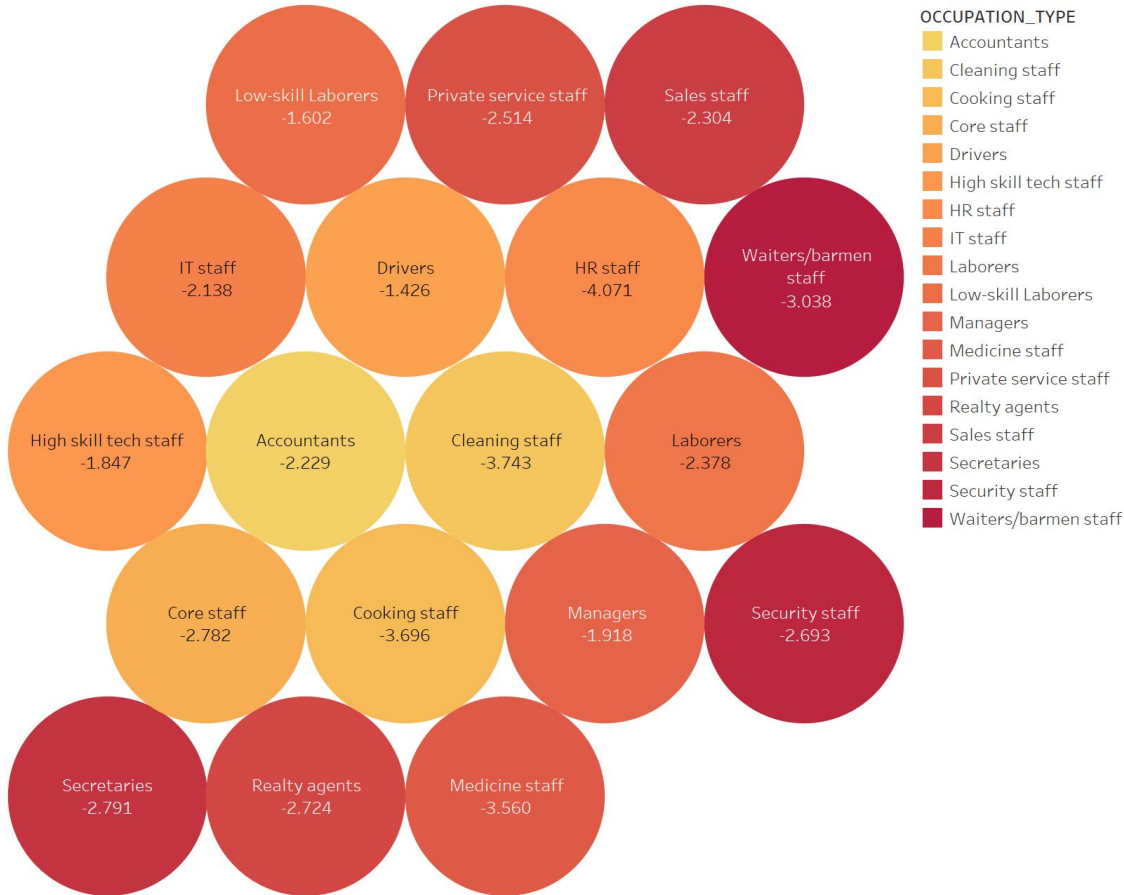# Number of people that default on a loan based on Occupation Types



% of Total Count of TARGET for each OCCUPATION_TYPE. Color shows details about TARGET (copy). The marks are labeled by % of Total Count of TARGET. The view is filtered on OCCUPATION_TYPE, which excludes Null.

# Number of Defaulters by Education Level and Occupation Type



**NAME_EDUCATION_TYPE**
- Academic degree
- Higher education
- Incomplete higher
- Lower secondary
- Secondary / secondary special

NAME_EDUCATION_TYPE, OCCUPATION_TYPE and % of Total TARGET.  Color shows details about
NAME_EDUCATION_TYPE.  Size shows % of Total TARGET.  The marks are labeled by NAME_EDUCATION_TYPE,
OCCUPATION_TYPE and % of Total TARGET. The view is filtered on OCCUPATION_TYPE, which excludes Null.

# Difference in Average Age between the Defaulters and Non Defaulters by Occupation Type



**OCCUPATION_TYPE**
- Accountants
- Cleaning staff
- Cooking staff
- Core staff
- Drivers
- High skill tech staff
- HR staff
- IT staff
- Laborers
- Low-skill Laborers
- Managers
- Medicine staff
- Private service staff
- Realty agents
- Sales staff
- Secretaries
- Security staff
- Waiters/barmen staff

Low-skill Laborers -1.602
Private service staff -2.514
Sales staff -2.304
IT staff -2.138
Drivers -1.426
HR staff -4.071
Waiters/barmen staff -3.038
High skill tech staff -1.847
Accountants -2.229
Cleaning staff -3.743
Laborers -2.378
Core staff -2.782
Cooking staff -3.696
Managers -1.918
Security staff -2.693
Secretaries -2.791
Realty agents -2.724
Medicine staff -3.560

OCCUPATION_TYPE and Avg Age Diff.  Color shows details about OCCUPATION_TYPE.  Size shows Avg Age Diff.
The marks are labeled by OCCUPATION_TYPE and Avg Age Diff. The view is filtered on OCCUPATION_TYPE, which
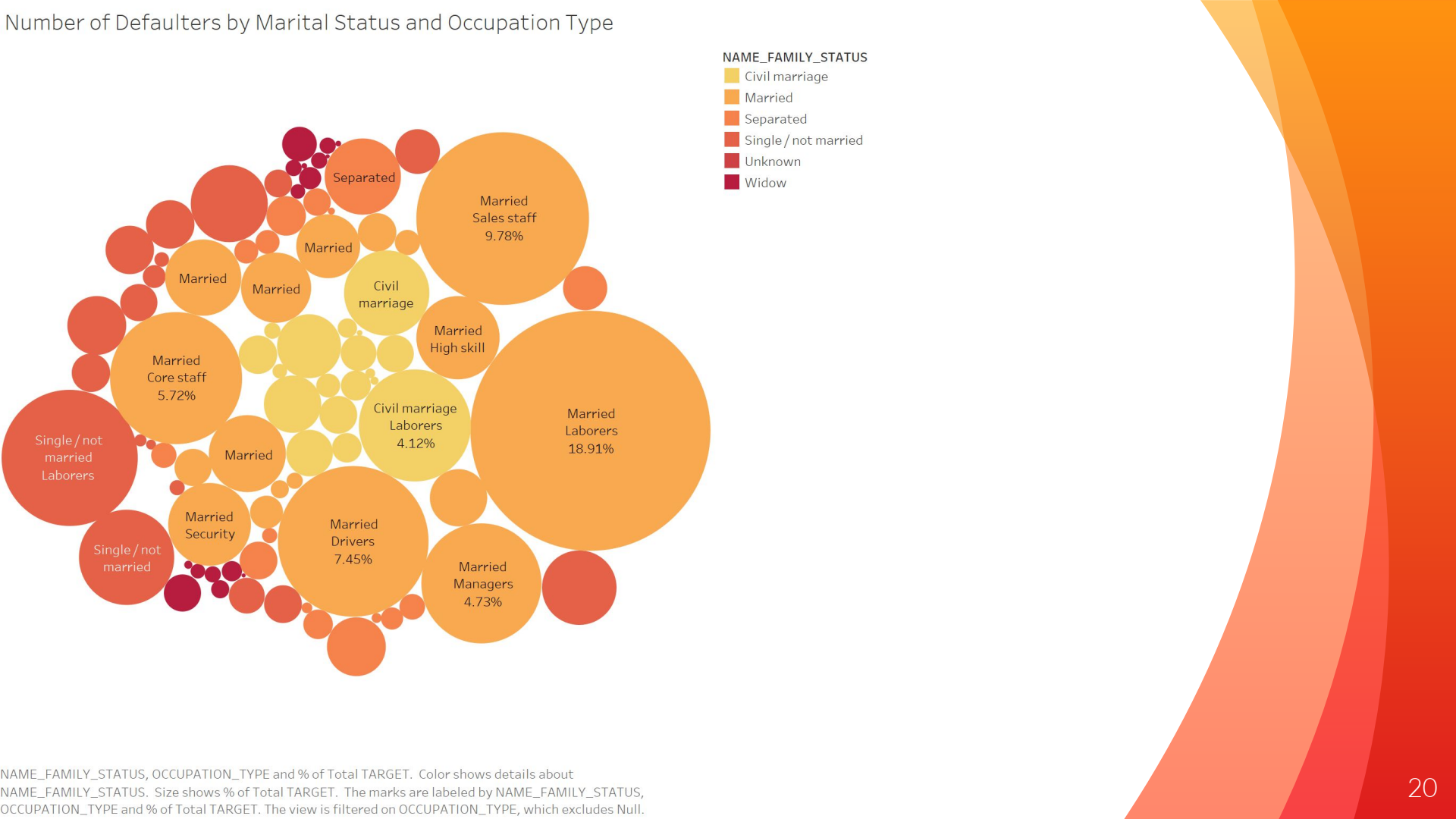excludes Null.

# Number of Defaulters by Marital Status and Occupation Type



**NAME_FAMILY_STATUS**
- Civil marriage
- Married
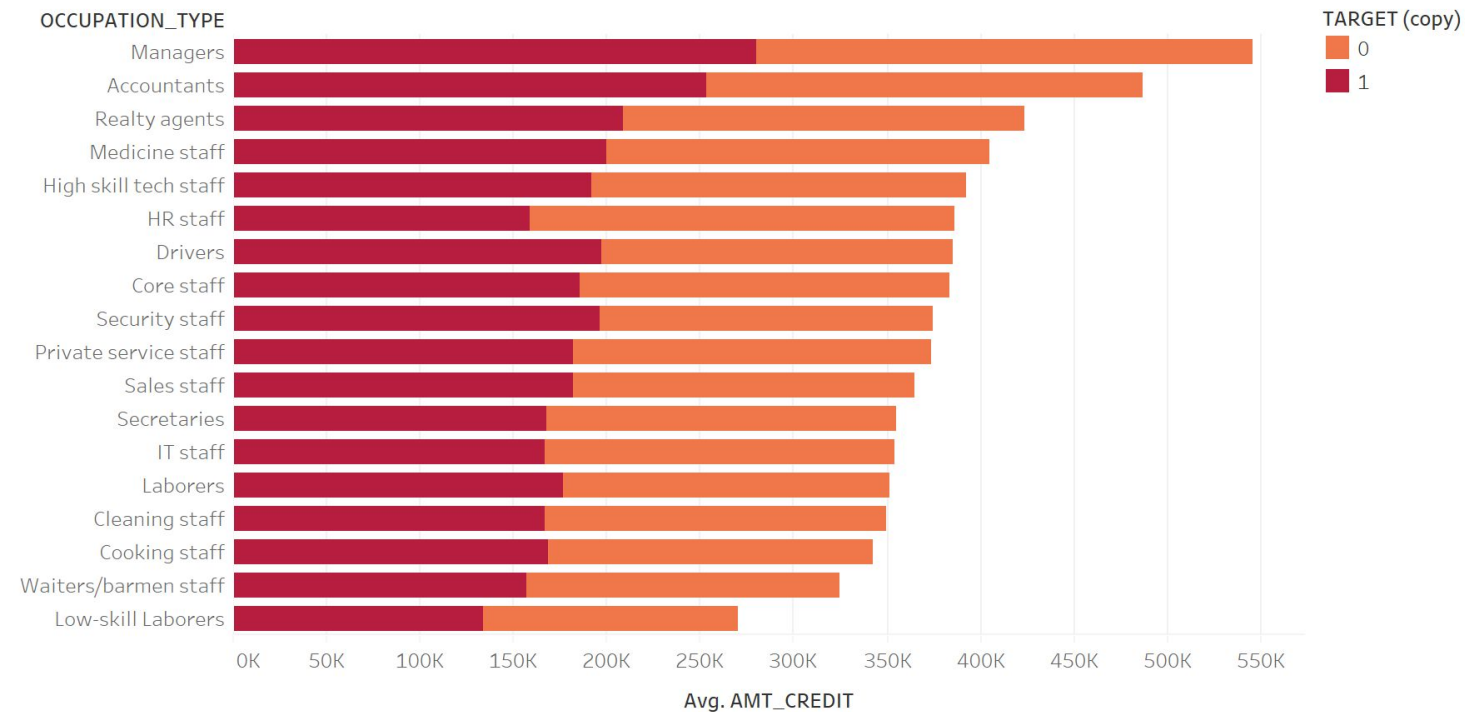- Separated
- Single / not married
- Unknown
- Widow

Separated

Married
Sales staff
9.78%

Married

Married

Married

Civil
marriage

Married
High skill

Married
Core staff
5.72%

Civil marriage
Laborers
4.12%

Married
Laborers
18.91%

Single / not
married
Laborers

Married

Married
Security

Married
Drivers
7.45%

Single / not
married

Married
Managers
4.73%

NAME_FAMILY_STATUS, OCCUPATION_TYPE and % of Total TARGET.  Color shows details about
NAME_FAMILY_STATUS.  Size shows % of Total TARGET.  The marks are labeled by NAME_FAMILY_STATUS,
OCCUPATION_TYPE and % of Total TARGET. The view is filtered on OCCUPATION_TYPE, which excludes Null.
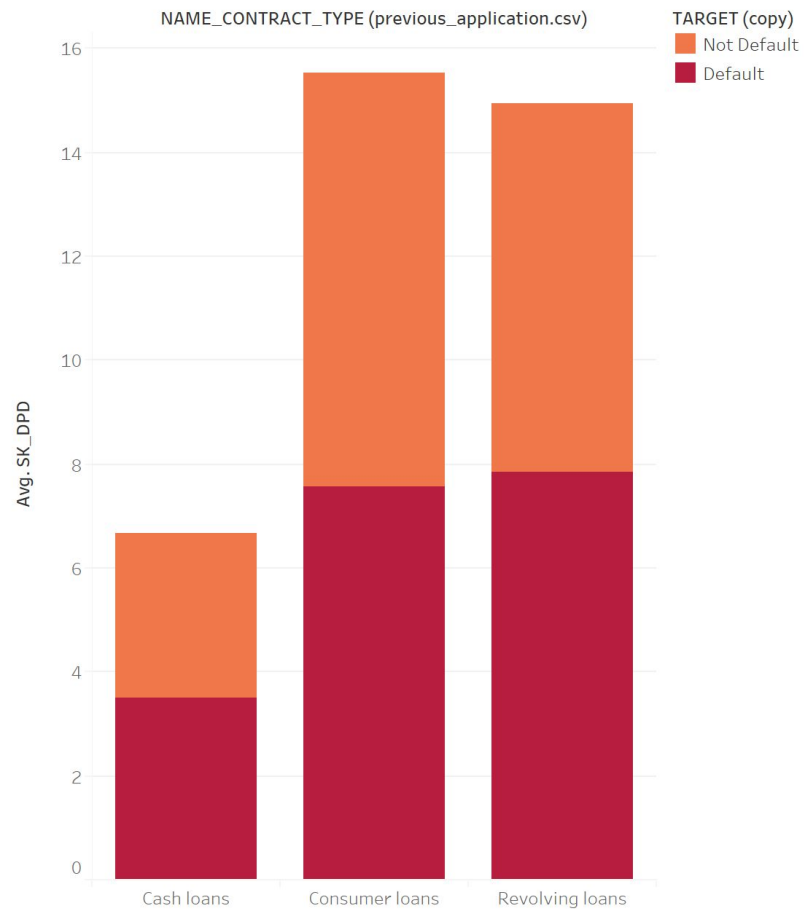
Hypothesis 2:

Are clients with many previous credits more likely to default on loans?

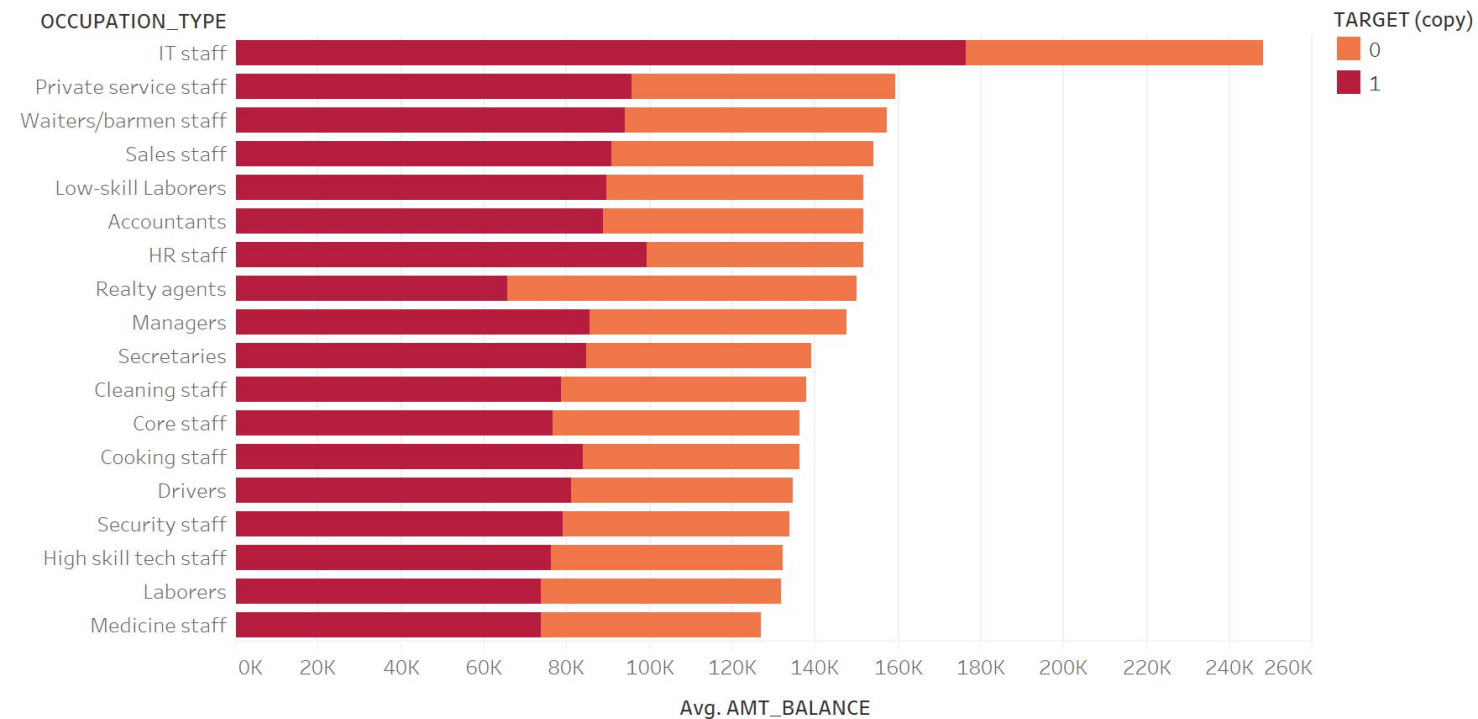# Average Amount of Previous Credit based on Occupation Type



Average of AMT_CREDIT for each OCCUPATION_TYPE. Color shows details about TARGET (copy). The view is filtered on OCCUPATION_TYPE and TARGET (copy). The OCCUPATION_TYPE filter excludes Null. The TARGET (copy) filter keeps 0 and 1.

# Days past due for different loan types



NAME_CONTRACT_TYPE (previous_application.csv)

TARGET (copy)
- Not Default
- Default

Avg. SK_DPD
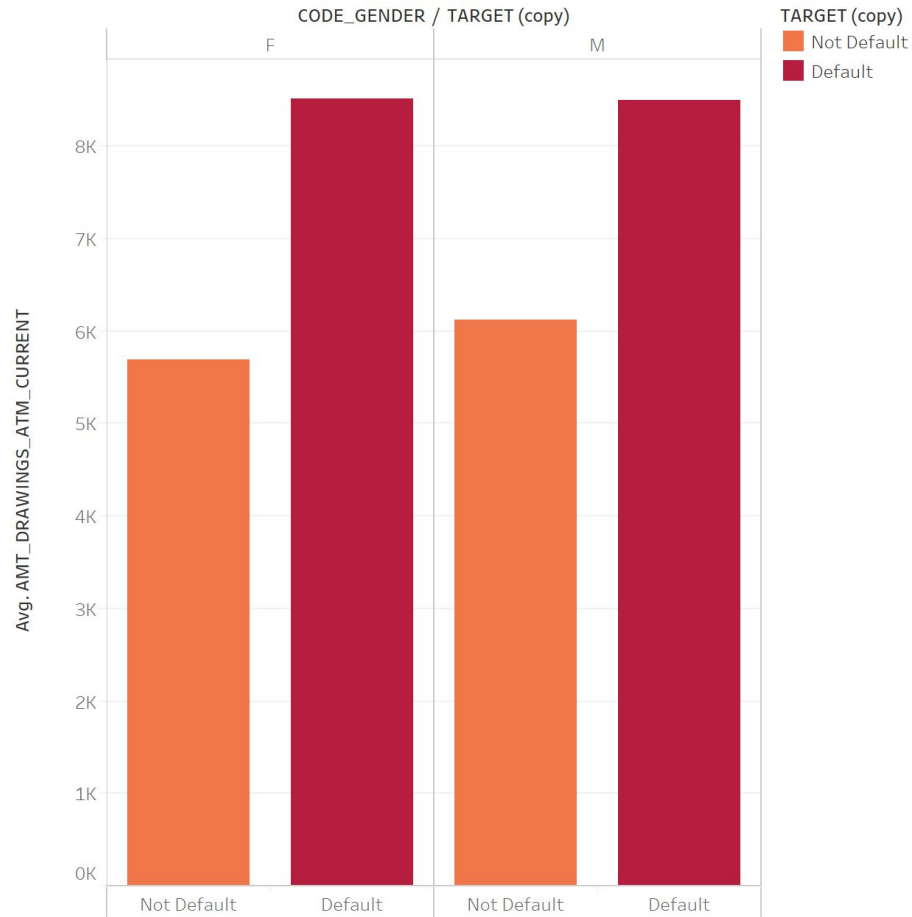
Cash loans   Consumer loans   Revolving loans

Average of SK_DPD for each NAME_CONTRACT_TYPE
(previous_application.csv).  Color shows details about TARGET (copy). The view
is filtered on NAME_CONTRACT_TYPE (previous_application.csv), which keeps
Cash loans, Consumer loans and Revolving loans.

# Average Credit Card Balance by Occupation



Average of AMT_BALANCE for each OCCUPATION_TYPE. Color shows details about TARGET (copy). The view is filtered on OCCUPATION_TYPE, which excludes Null.

24

# Average Amount Withdrawn from ATM by Gender and Default Status



Average of AMT_DRAWINGS_ATM_CURRENT for each TARGET (copy) broken down by CODE_GENDER. Color shows details about TARGET (copy).

# 5.
# Feature Engineering

- ‘Days’ variable made positive and in terms of years
- ‘Family size’ converted to binned categorical variable
- Technical information added to improve model performance
  - Credit term
  - % of days employed

- Missing values imputed with the median

- Data split into 80% train, 20% test

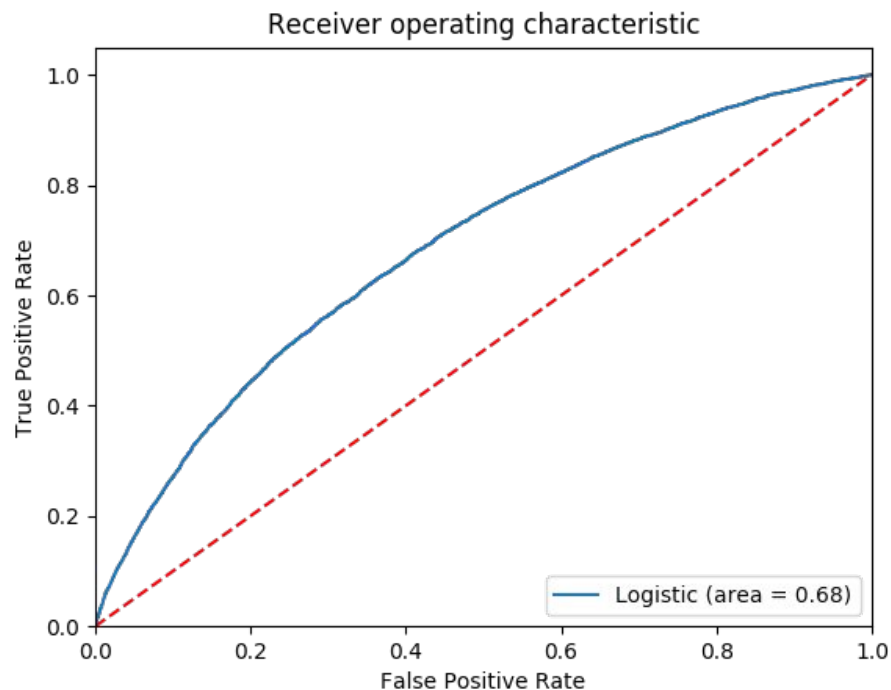- SMOTE package used for resampling to deal with imbalanced classes

# 6.
# Model Building

Insights and Recommendations

# Models Built: A Comparison

01 **Random Forests** — AUROC Score: 0.65

02 **Random Forests: Resampling** — AUROC Score: 0.63

03 **Logistic Regression** — AUROC Score: 0.68

04 **Cat Boosting** — AUROC Score: 0.69
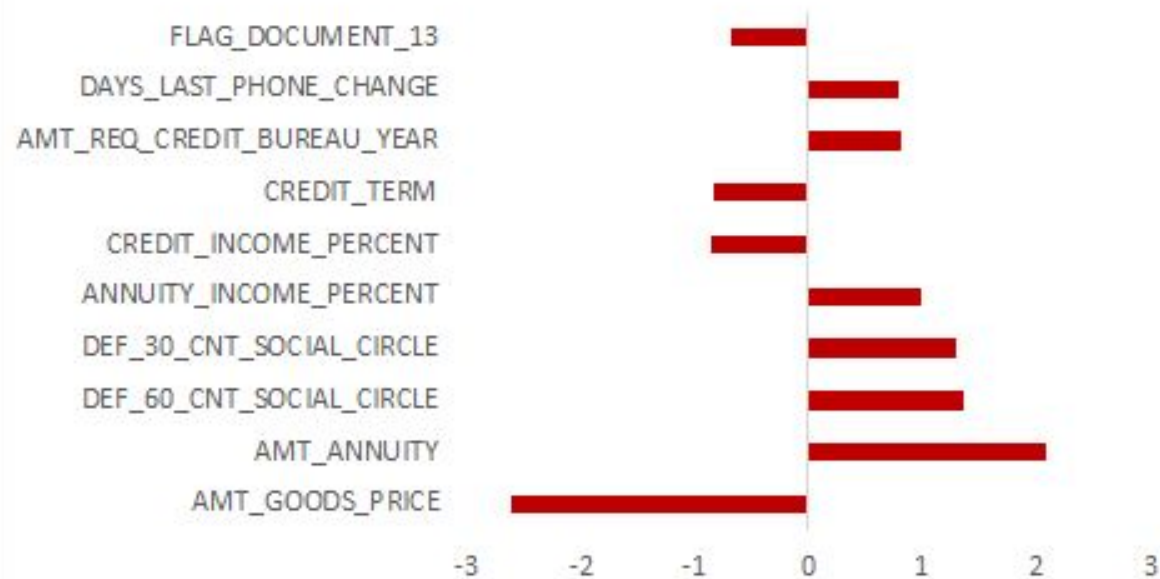
# Evaluation of Logistic Regression



Receiver operating characteristic

(True Positive Rate vs False Positive Rate)

Logistic (area = 0.68)

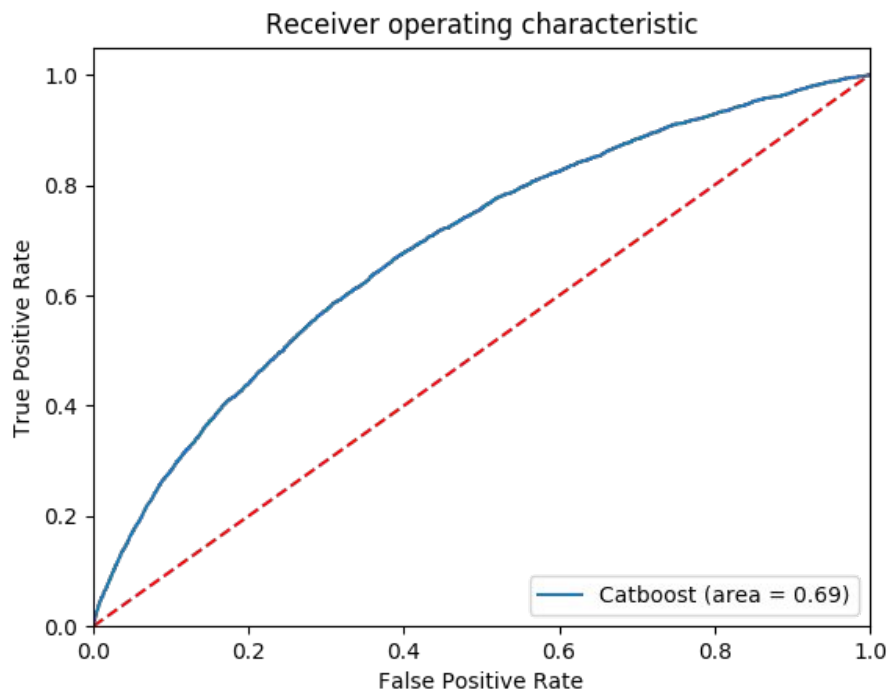|          | 0-Predict | 1-Predict |
|----------|-----------|-----------|
| 0-Actual | 55880     | 768       |
| 1-Actual | 4574      | 281       |

**Reduce False Negatives!**
People who default but the model predicts the client won't!

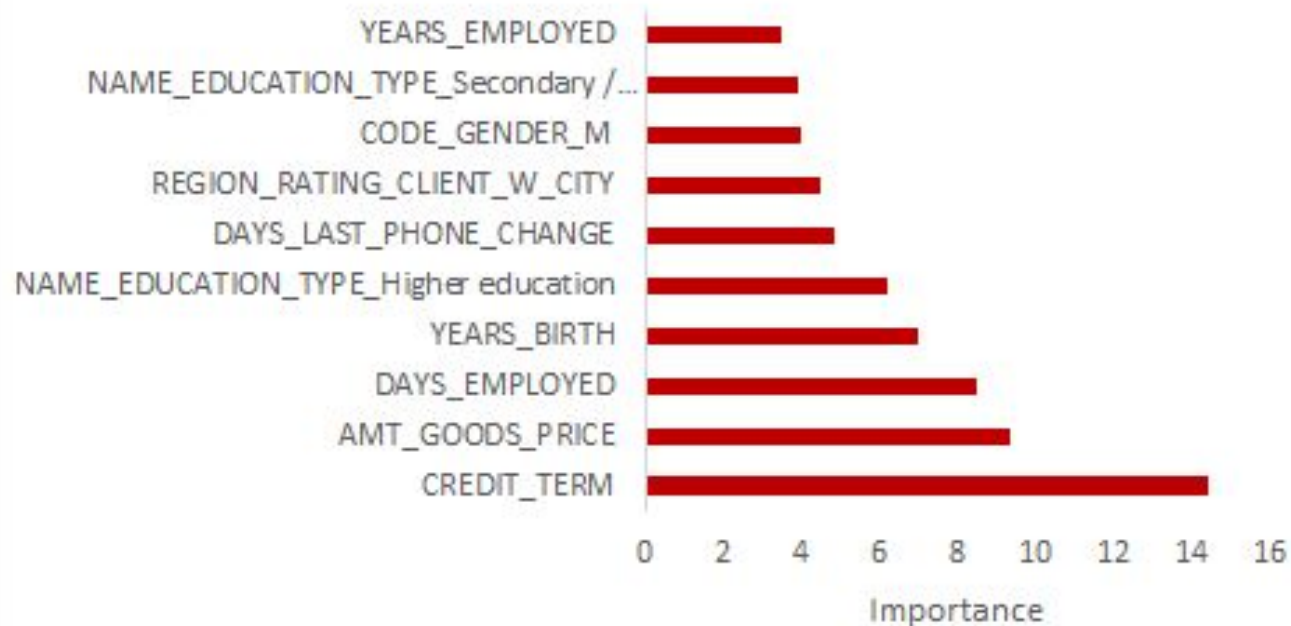Logistic Regression: Coefficients

# Evaluation of CatBoost


Receiver operating characteristic

|  | 0-Predict | 1-Predict |
|---|---|---|
| **0-Actual** | 54233 | 2415 |
| **1-Actual** | 4136 | 719 |

**Reduce False Negatives!**
People who default but the model predicts the client won't!

Feature Importance for CatBoost

# Insights and recommendations

- Be more cautious when you are lending to labourers and not highly educated clients
- The recent withdrawals from the ATM has an impact on the default risk
- Defaulting is not instant - If the credit balance increases over time, then the client is highly likely to default
- Region rating from the model as well as the discrepancy in the work and residence location
- Amts_goods_price the proposed loan purpose higher - more likely to default
- If a person has recently changed the phone number, then the propensity to default increases

# Questions?