

# HOME CREDIT



**Sijo VM**  
**Pooja Shah**  
**David Owen**  
**Saurabh Bodas**  
**Alisha Fernandes**

# Agenda

- The company
- Problem Statement
- Data Sources & Tools
- Hypothesis & Goals
- Exploratory Data Analysis
- Feature Engineering
- Model Building

# **1.**

# **The company**

The background of the slide features abstract, wavy shapes in various shades of orange and red, creating a modern and dynamic visual effect. The colors transition from lighter tones on the left to darker, more saturated tones on the right.

 **\$5.2 billion**

Revenues!

 **132,000 employees**

And a lot of customers!

**1997**

Founded in Czech Republic

# Maps



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**

The Kaggle logo, featuring the word "kaggle" in a lowercase, blue, sans-serif font.

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

The Databricks logo, consisting of a red icon of three stacked cubes above the word "databricks" in a lowercase, grey, sans-serif font.

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

The Tableau logo, featuring a cluster of small, multi-colored plus signs to the left of the word "tableau" in a lowercase, blue, sans-serif font.The PySpark logo, with "Py" in orange and "Spark" in black, both in a sans-serif font, followed by an orange star icon.

## **Bureau**

All clients previous credits provided by other financial institutions

## **Credit Card Balance**

Monthly balance snapshots of previous credit cards owned by the applicant

## **Bureau Balance**

Monthly balances of previous credits in the credit bureau

## **Previous Applications**

All previous applications for loans by the client

## **Applications**

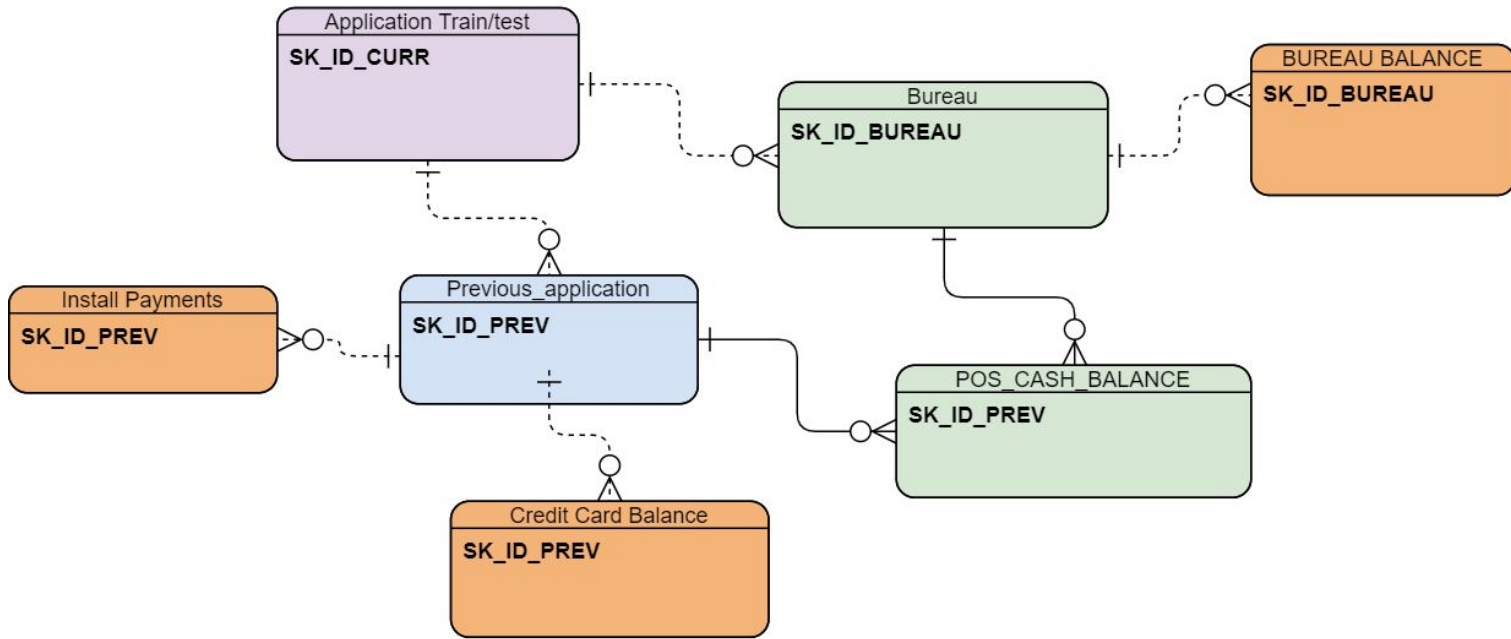
Main table depicting current loan applications for each applicant

## **POS\_Cash**

Monthly balance snapshots of previous POS and loans

## **Installment Payments**

Repayment history for previously disbursed credits

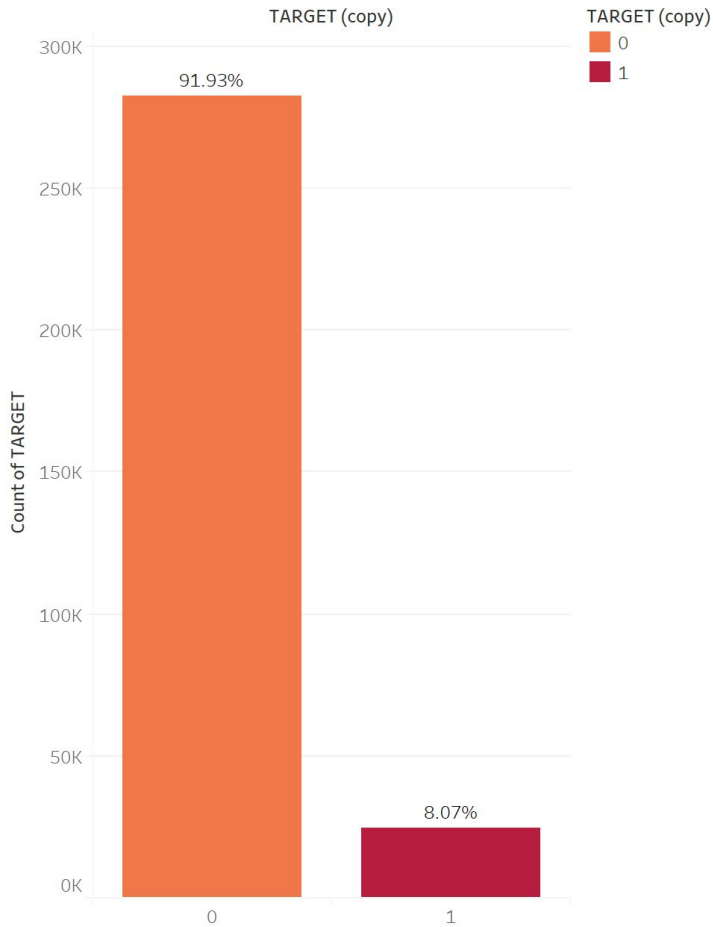


**4.**

# **Exploratory Data Analysis**

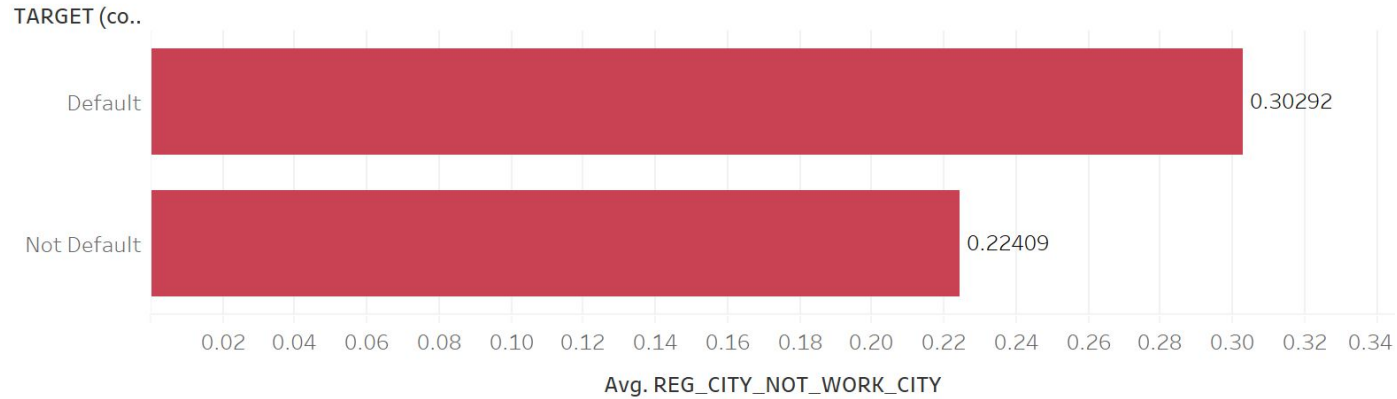


# Distribution of the Default Data



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



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.

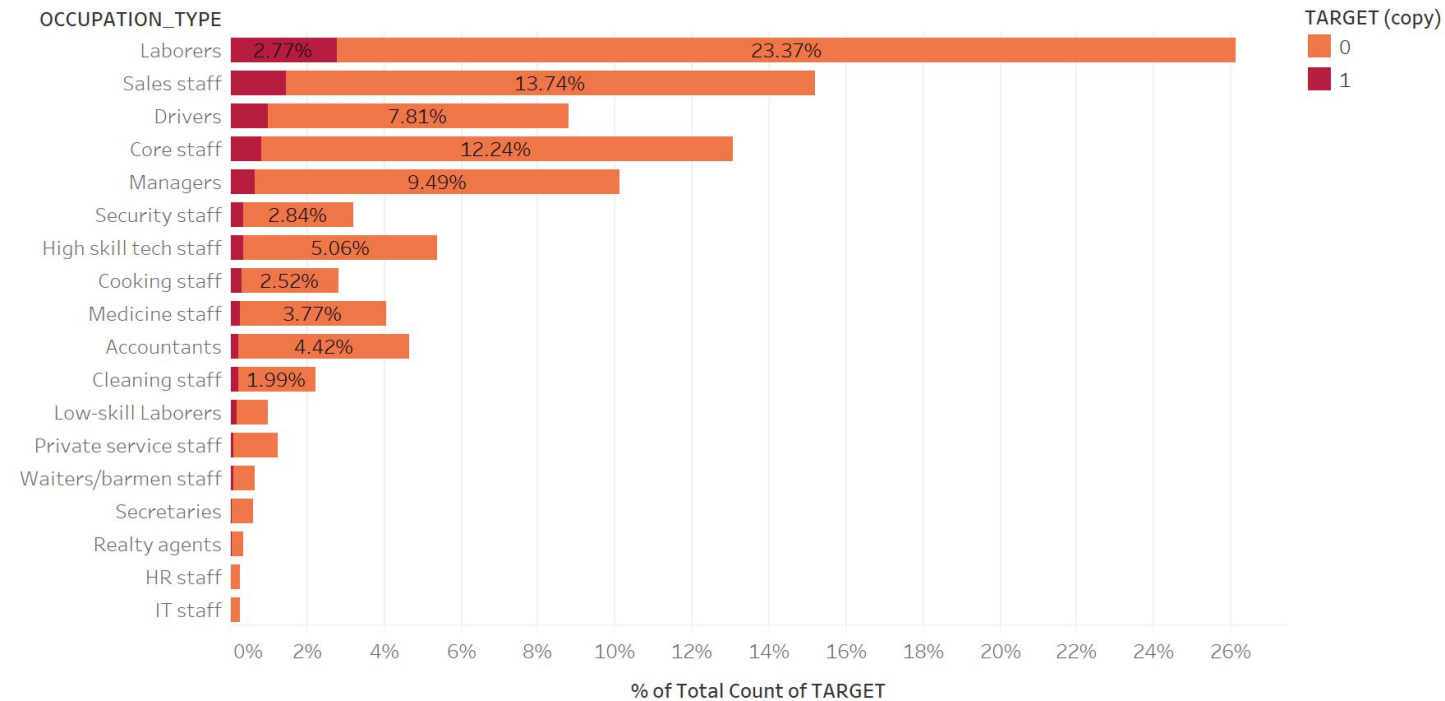
## Hypothesis 1:

Are clients with historically worse job security more likely to default on their loan payments?



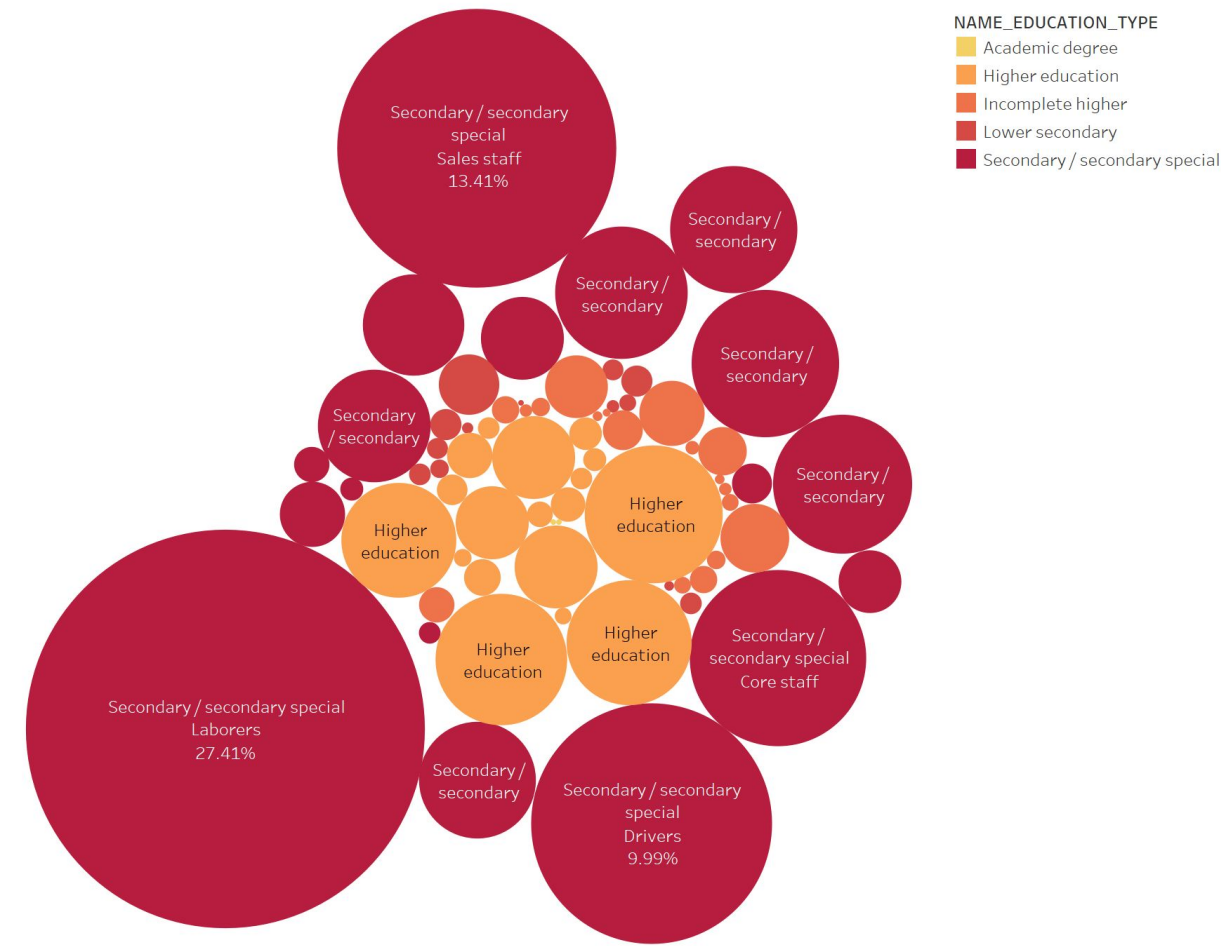


## Number of people that default on a loan based on Occupation Types



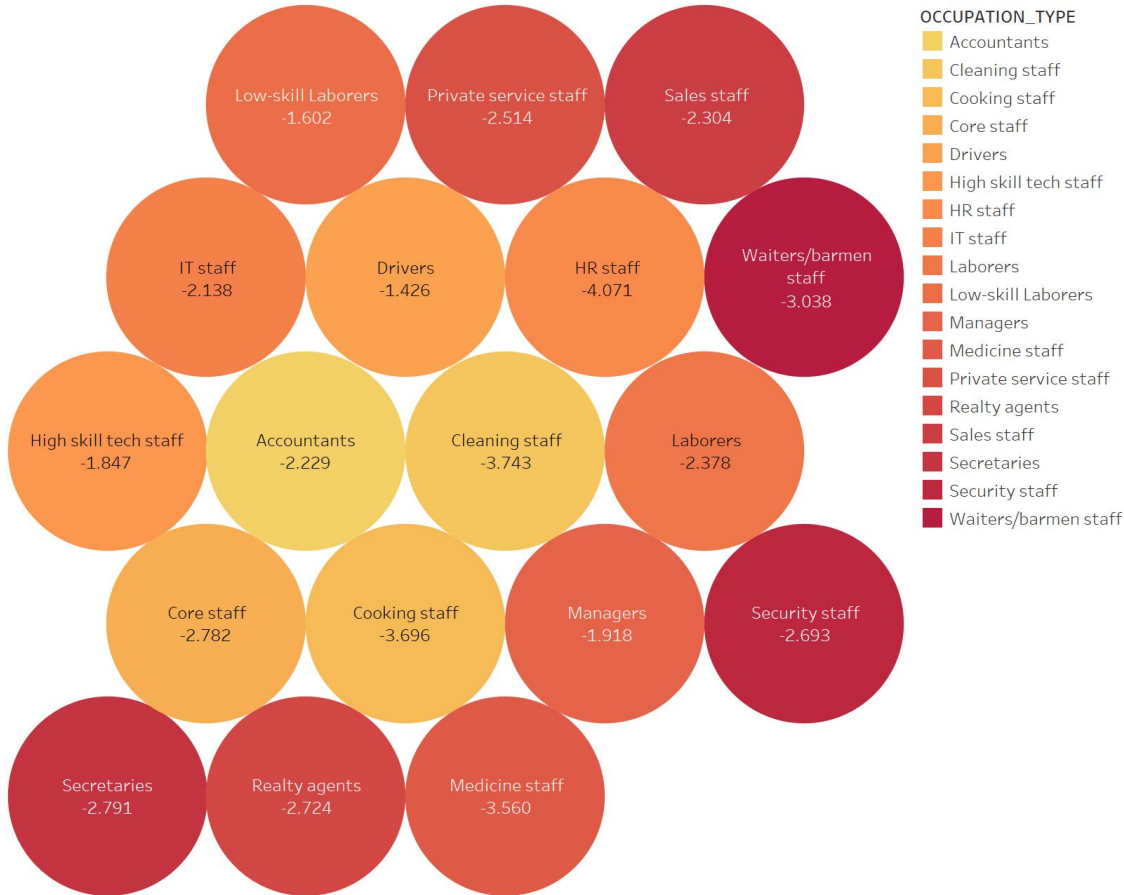
% of Total Count of TARGET for each OCCUPATION\_TYPE. Color shows details about TARGET (copy). The marks are labeled by % of Total Count of TARGET. The view is filtered on OCCUPATION\_TYPE, which excludes Null.

Number of Defaulters by Education Level and Occupation Type



NAME\_EDUCATION\_TYPE, 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 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.

NAME\_FAMILY\_STATUS

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

Married Laborers 18.91%

Married Sales staff 9.78%

Married Core staff 5.72%

Married Drivers 7.45%

Single / not married Laborers

Married High skill

Civil marriage Laborers 4.12%

Married Security

Married Managers 4.73%

Married

Civil marriage

Separated

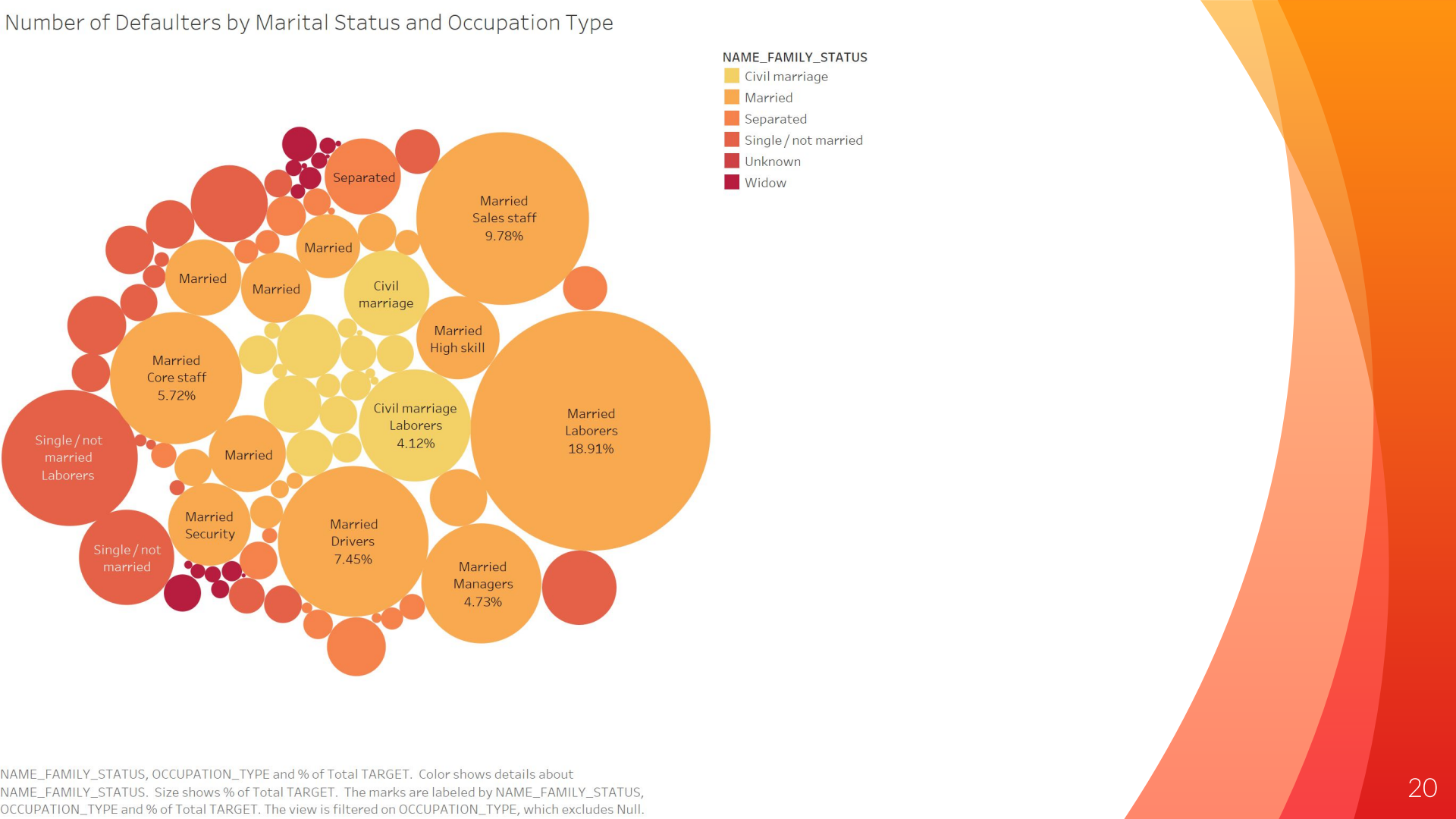
Single / not married

Unknown

Widow

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.

20



Number of Defaulters by Marital Status and Occupation Type

NAME\_FAMILY\_STATUS

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

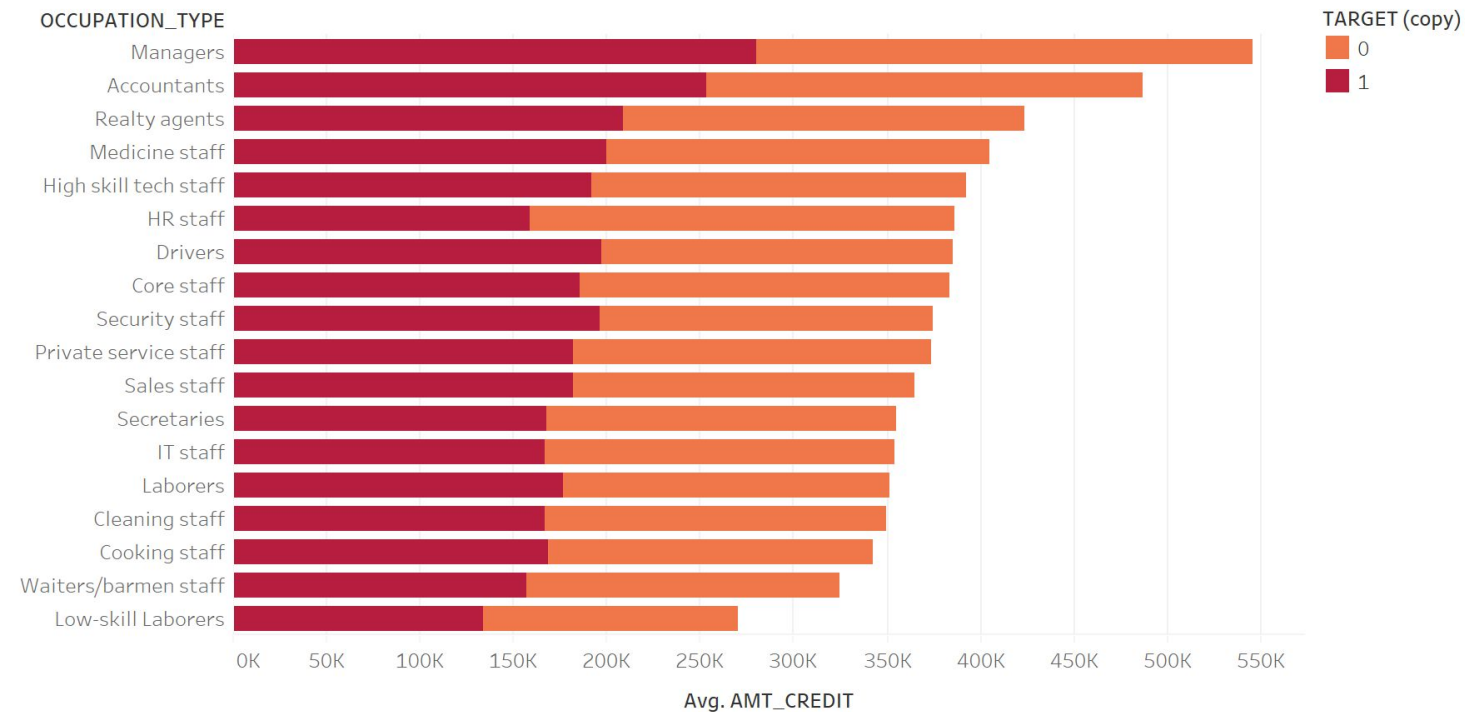
NAME_FAMILY_STATUS	OCCUPATION_TYPE	% of Total TARGET
Married	Laborers	18.91%
Married	Sales staff	9.78%
Married	Core staff	5.72%
Single / not married	Laborers	
Married	Drivers	7.45%
Civil marriage	Laborers	4.12%
Married	Managers	4.73%
Married	High skill	
Married	Security	
Single / not married		
Separated		
Widow		
Unknown		

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:

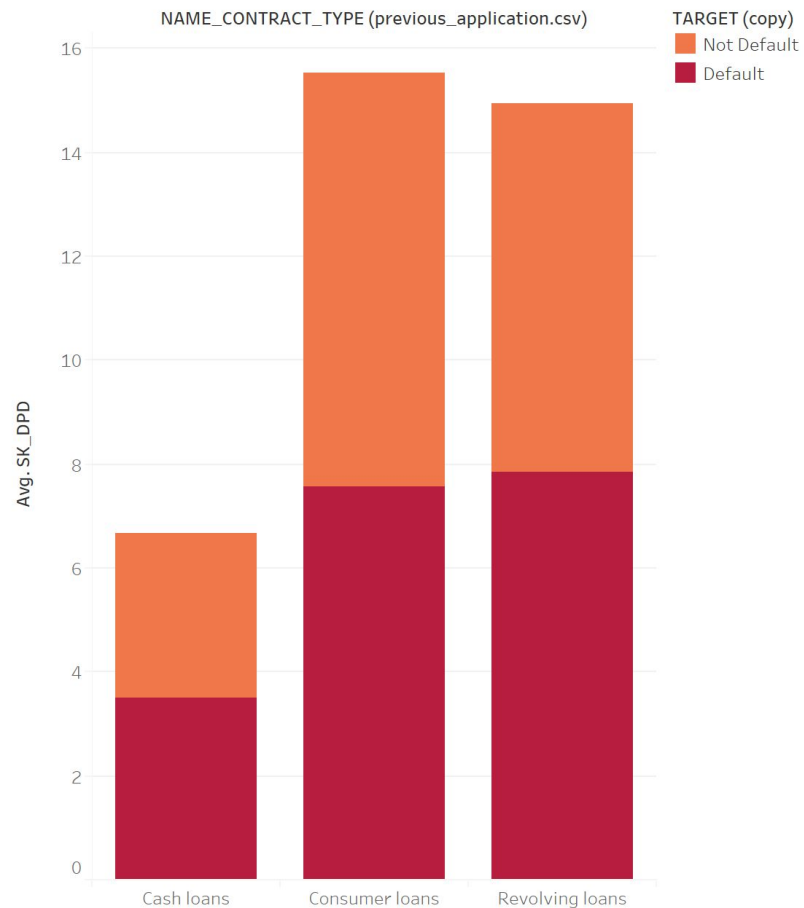
Are clients with many previous credits more likely to default on loans?

## Average Amount of Previous Credit based on Occupation Type



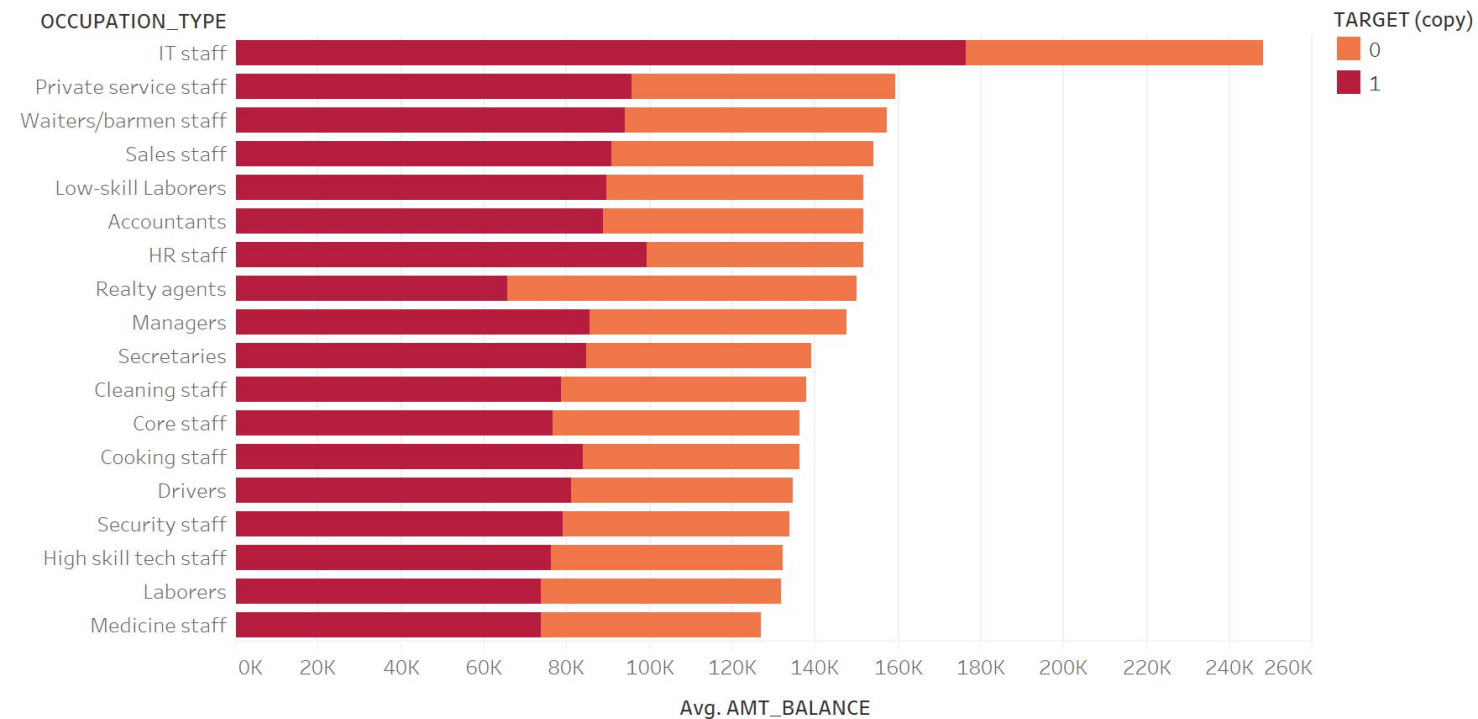
Average of AMT\_CREDIT for each OCCUPATION\_TYPE. Color shows details about TARGET (copy). The view is filtered on OCCUPATION\_TYPE and TARGET (copy). The OCCUPATION\_TYPE filter excludes Null. The TARGET (copy) filter keeps 0 and 1.

## Days past due for different loan types



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.

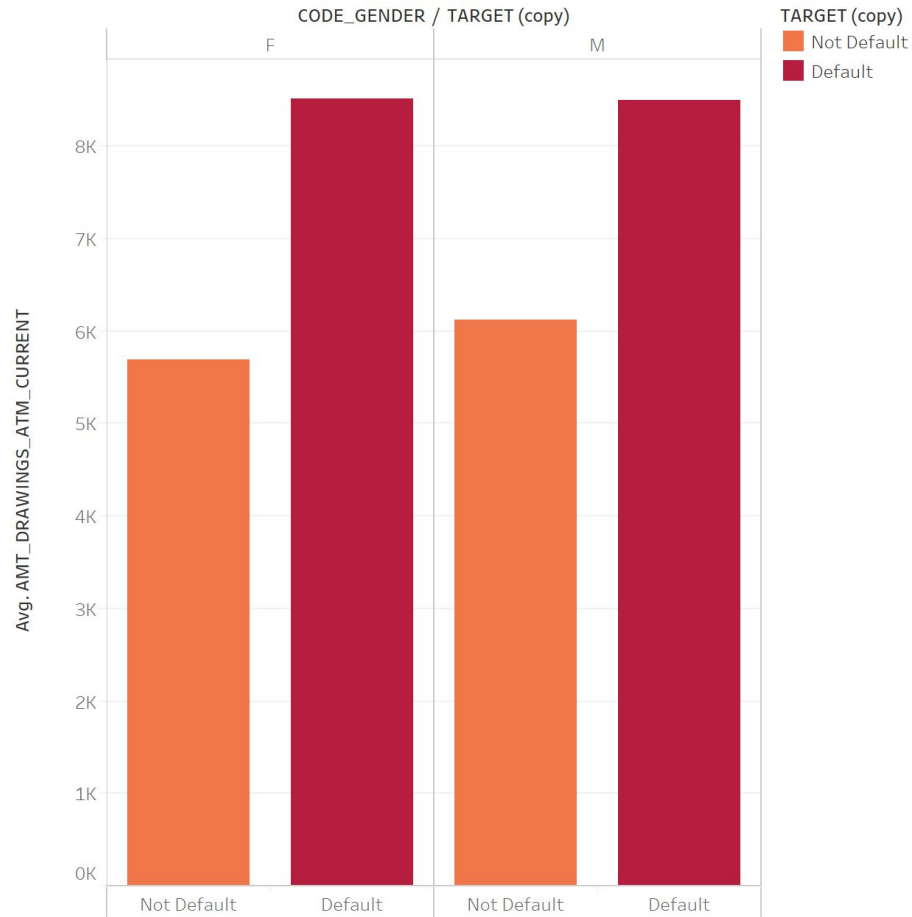
## Average Credit Card Balance by Occupation



Average of AMT\_BALANCE for each OCCUPATION\_TYPE. Color shows details about TARGET (copy). The view is filtered on OCCUPATION\_TYPE, which excludes Null.



# 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**

The background of the slide features abstract, wavy shapes in shades of orange and red. On the left side, there are overlapping orange shapes. On the right side, there are overlapping red and pink shapes that curve upwards towards the top right corner.

- '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

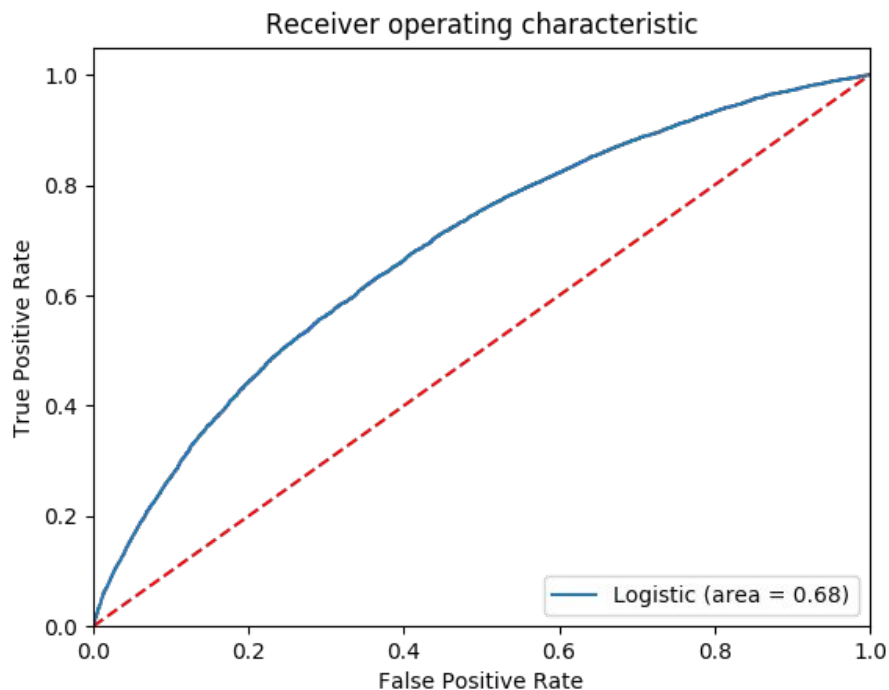
# 6. **Model Building**

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

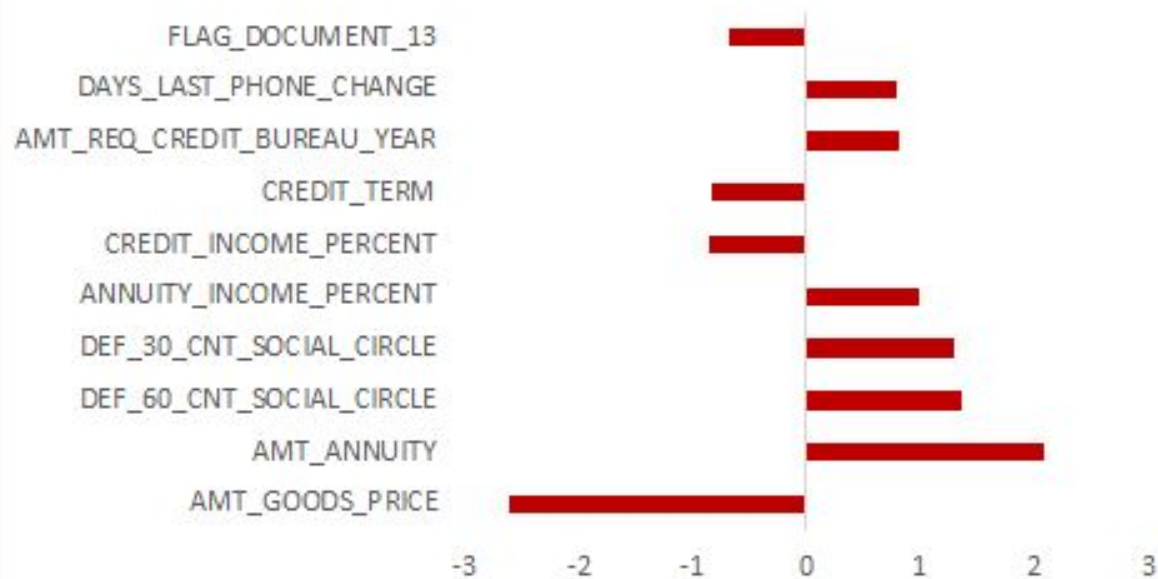


	0-Predict	1-Predict
0-Actual	55880	768
1-Actual	4574	281

## Reduce False Negatives!

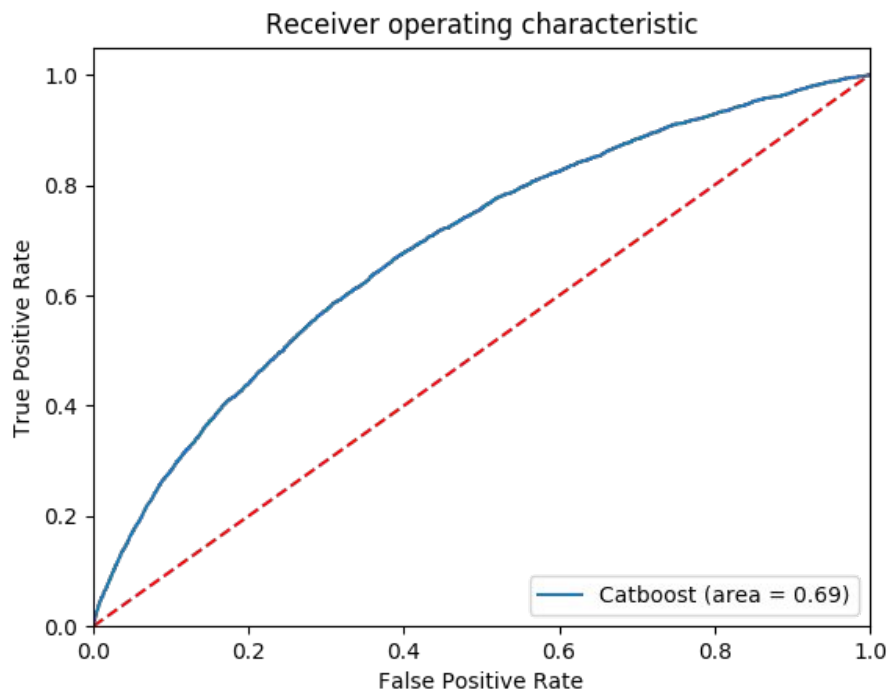
People who default but the model predicts the client won't!

## Logistic Regression: Coefficients





# Evaluation of CatBoost

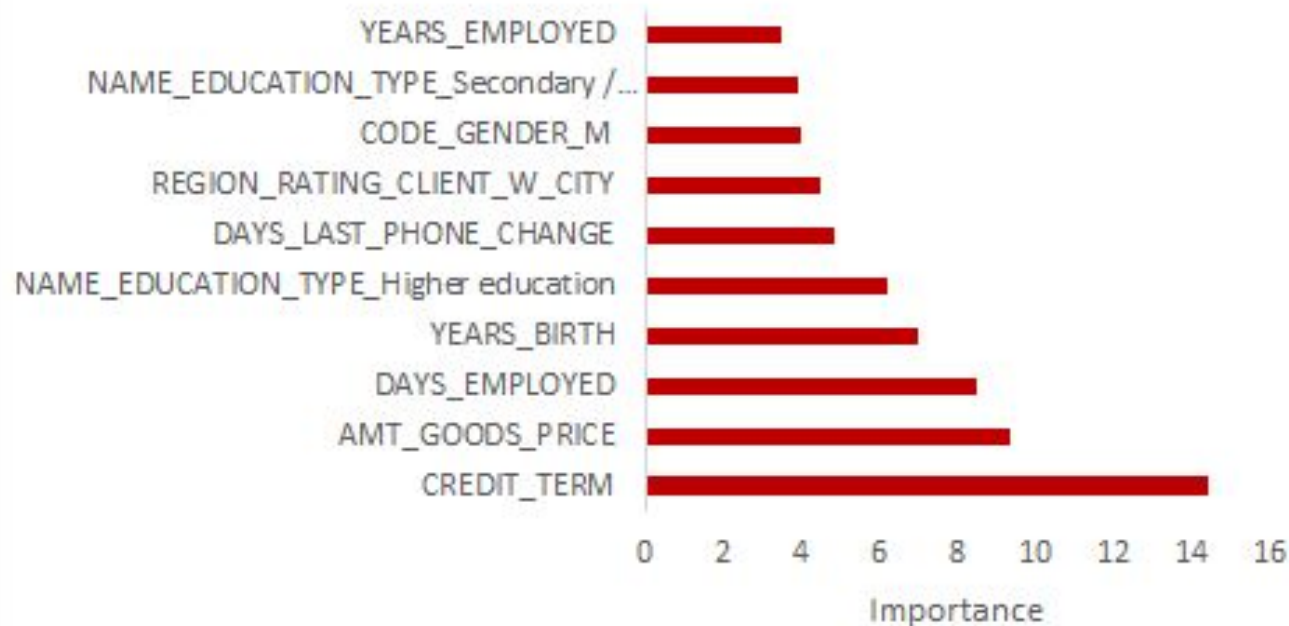


	0-Predict	1-Predict
0-Actual	54233	2415
1-Actual	4136	719

## 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?**

The background features abstract, flowing shapes in shades of orange and red. On the left, there are overlapping wavy bands of light orange and dark orange. On the right, a large, upward-sloping shape transitions from a light pink at the top to a vibrant red at the bottom, creating a sense of movement and depth.