



Luddy School of Informatics, Computing and Engineering

Home Credit Default Risk

Team Members:

1. Akash Gangadharan (akganga@iu.edu)
2. Laya Harwin (lharwin@iu.edu)
3. Dhairya Shah (shahds@iu.edu)
4. Shobhit Sinha(shosinha@iu.edu)



INDIANA UNIVERSITY BLOOMINGTON

Agenda

01

Project Objective

Understanding the business objective and data

02

EDA and Data Preparation

Analyzing the historical data and generating insights out of it

03

Modelling Pipelines

Exploring and preprocessing the data

04

Comparing results and Deriving Insights

Comparing results from different model and evaluating

05

Future Scope

Future work to be done



SECTION 1

Project Objective

What are we aiming to solve?

Who?

Home Credit Group



What?

The goal of this Project is to predict whether a loan applicant is likely to default on a loan, given information about the applicant and their previous loan applications.

Why?

To make informed decision that can lead to less default on the loan and thereby reducing the loss for the organization

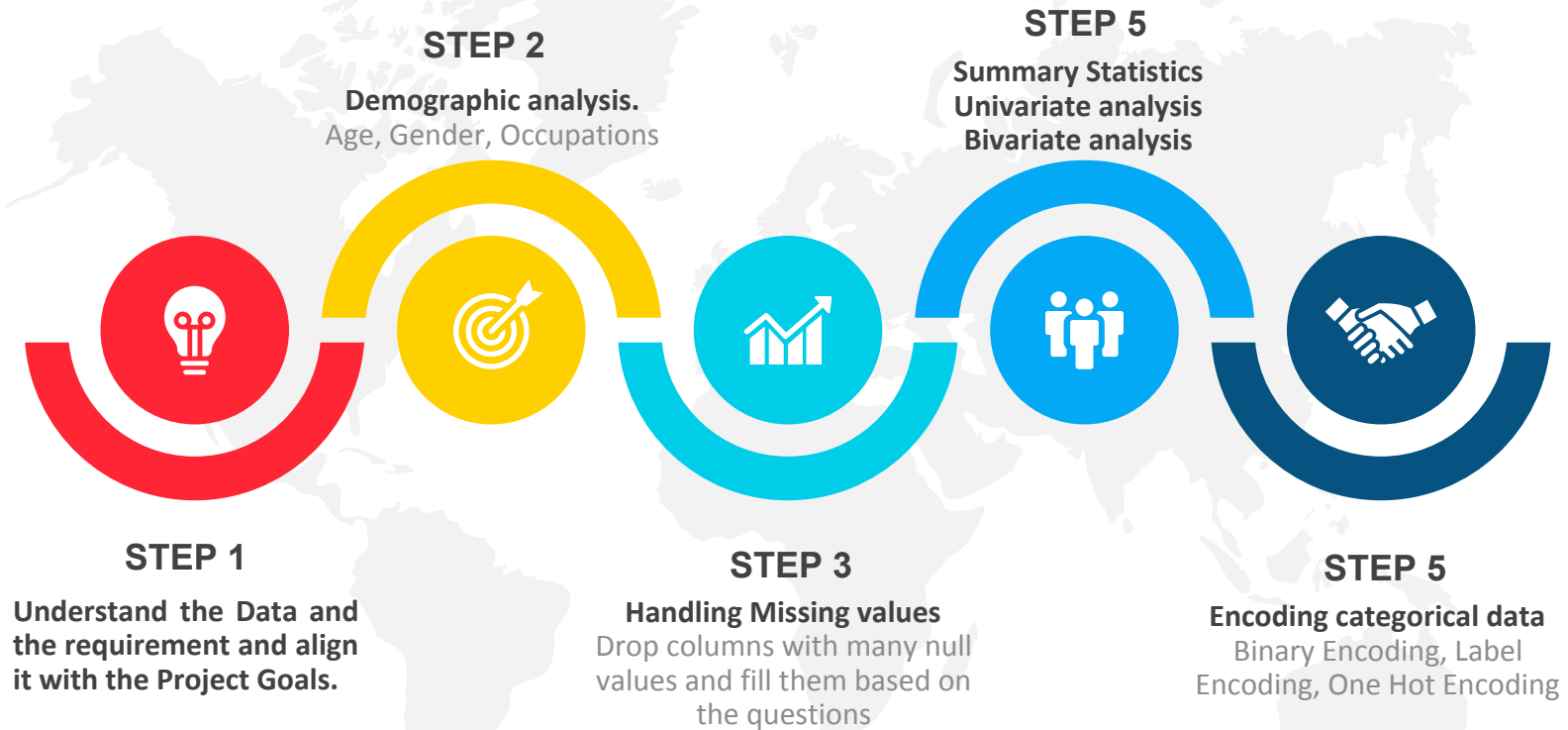
How?

Creating machine learning models using different demographic features and other usage data of the customer to predict whether the customer will default or not.

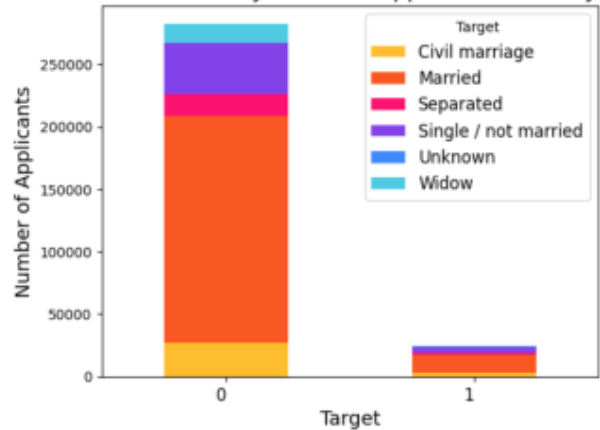
SECTION 2

EDA and Data Preparation

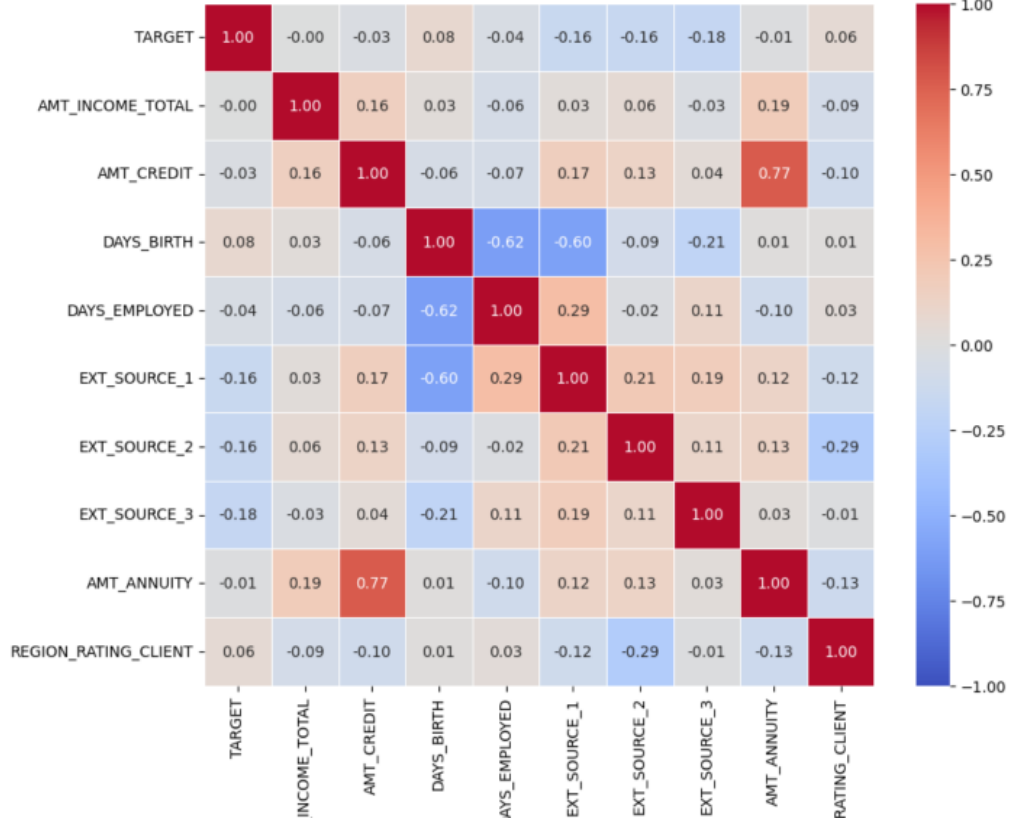
EDA and Data Preparation



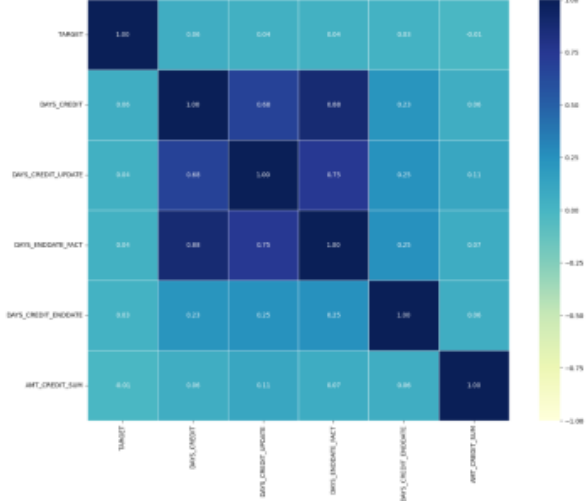
Distribution of Family Status in Application Train by Target



Correlation Matrix for Selected Variables

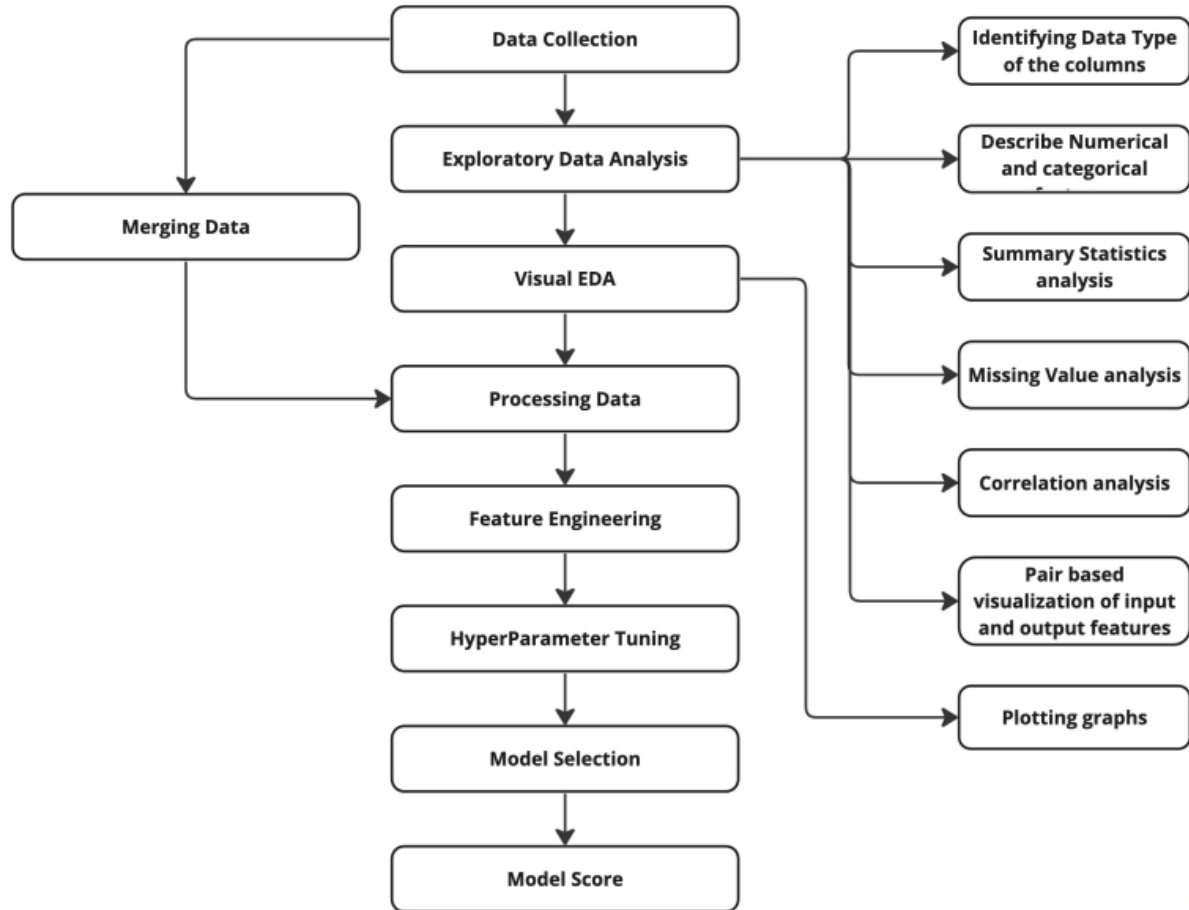


Correlation Matrix of top 5 features with Target



SECTION 3

Modelling Pipeline



SECTION 4

Comparing results

	Pipeline	Parameters	TrainAcc	ValidAcc	TestAcc	Train Time(s)	Test Time(s)
1	Baseline Pipeline(steps=[('dt', DecisionTreeClassifier())]) with 320 inputs	{'dt__max_depth': [5, 10]}	93.83%	93.75%	93.46%	0.196276	0.934601
2	Baseline Pipeline(steps=[('lr', LogisticRegression())]) with 320 inputs	{'lr__C': [0.01], 'lr__penalty': ['l1', 'l2']}	91.90%	92.25%	91.93%	0.078855	0.919290
3	Baseline Pipeline(steps=[('ada', AdaBoostClassifier(base_estimator=DecisionTreeClassifier()))]) with 320 inputs	{'ada__n_estimators': [50], 'ada__learning_rate': [0.01, 0.1], 'ada__base_estimator__max_depth': [1, 5]}	92.60%	92.69%	92.46%	10.076050	0.924596

ROC AUC score

Logistic Regression : 0.7470648772667208

Decision Tree : 0.8074479725505205

Random Forest : 0.9999999934281516

ADA Boost : 0.8129134572402118



Home Credit Default Risk

Can you predict how capable each applicant is of repaying a loan?

\$70,000

Prize Money



Home Credit Group · 7,176 teams · 5 years ago

[Overview](#) [Data](#) [Code](#) [Discussion](#) [Leaderboard](#) [Rules](#) [Team](#)

[Submissions](#)

[Late Submission](#)

...

Submissions

You selected 0 of 2 submissions to be evaluated for your final leaderboard score. Since you selected less than 2 submission, Kaggle auto-selected up to 2 submissions from among your public best-scoring unselected submissions for evaluation. The evaluated submission with the best Private Score is used for your final score.

0/2

■ Submissions evaluated for final score

All

Successful

Selected

Errors

Recent ▼

Submission and Description

Private Score ⓘ

Public Score ⓘ

Selected



submission.csv

Complete (after deadline) · now · Group15 AML Logistic Regression

0.73315

0.73762



submission.csv

Complete (after deadline) · 7h ago · group 15 AML

0.66457

0.65961



submission.csv

Complete (after deadline) · 7h ago · Submission from Group 15 AML

0.65139

0.66622



SECTION 5

Future Scope

To conclude, in this phase we cleaned the data and found the most relevant features to the Target variable and prediction.

Further we are planning to improvise the feature engineering, perform hyperparameter tuning for our models alongside using K-Fold cross validation and GridSearchCV, we might also use some advanced gradient boosting models so that we could get as close to the best accuracy as we can.





INDIANA UNIVERSITY BLOOMINGTON