

Loan default prediction using Machine Learning techniques

Daniel Jordan

04/08/2022



Problem definition

CONTEXT

- Retail banks offer home loans to obtain profits.
- Loans are borrowed by bank customers.
- Banks are rigorous while approving loans.

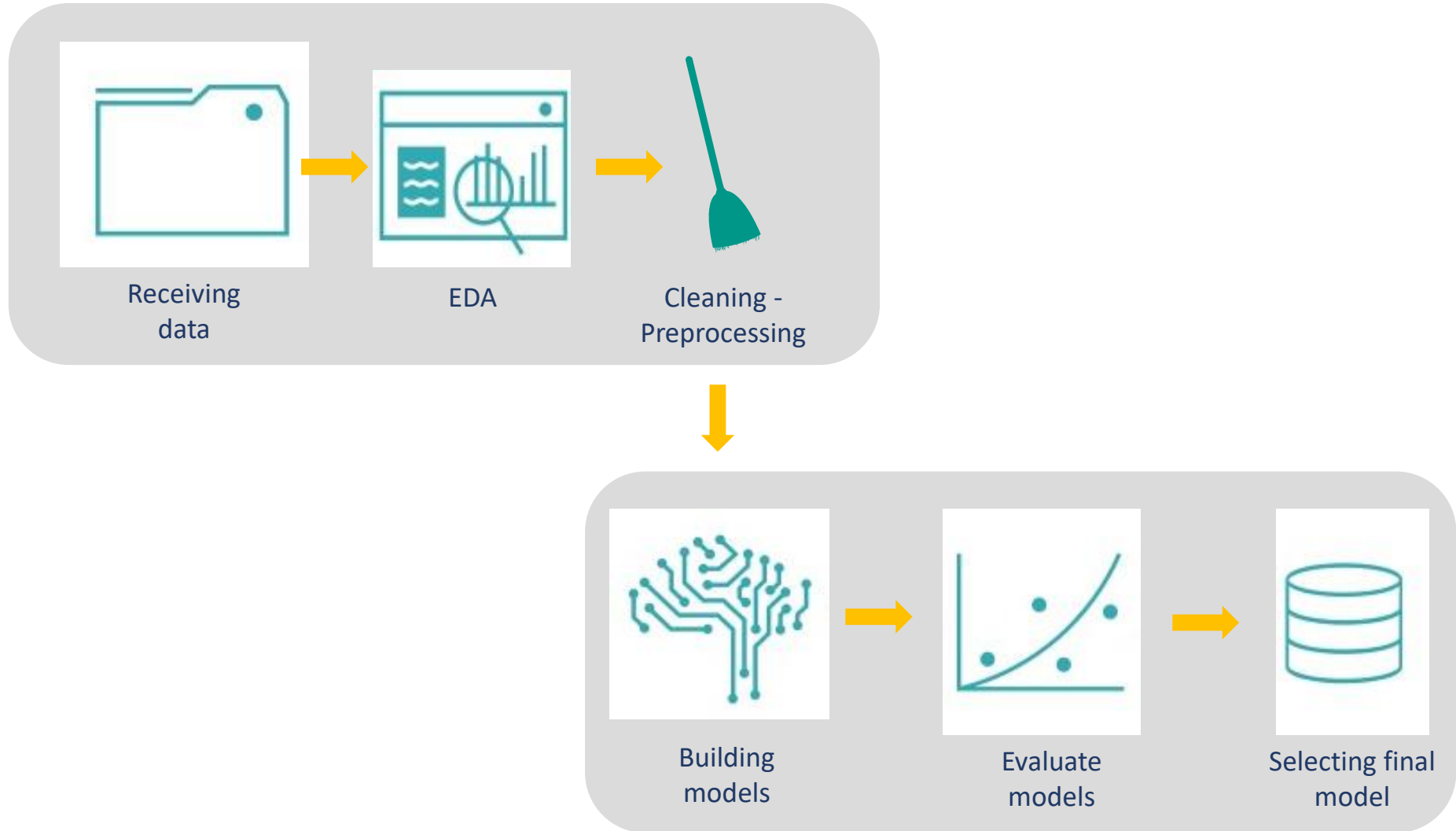
PROBLEM

- Banks need an effective approval process.
- This process is effort-intensive and sensitive to human error and biases.

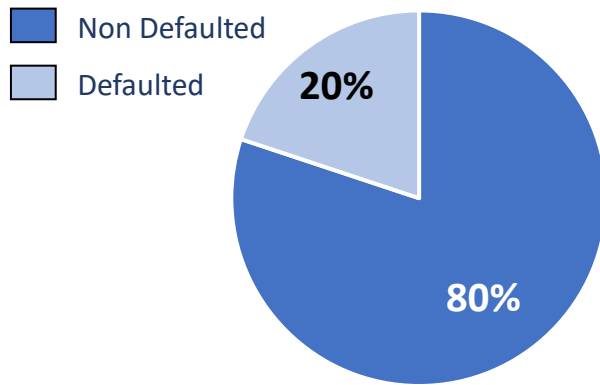
OBJECTIVE

Build a Machine Learning model free of biases and more efficient.

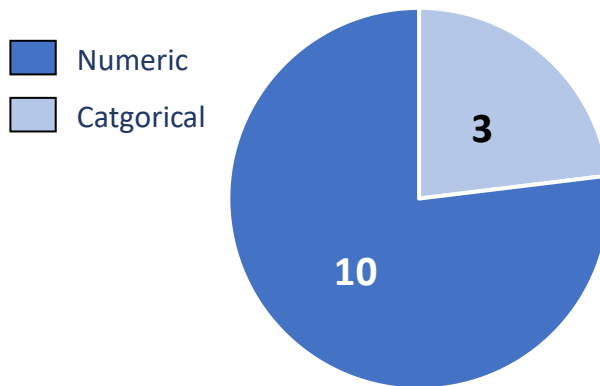
Solution approach



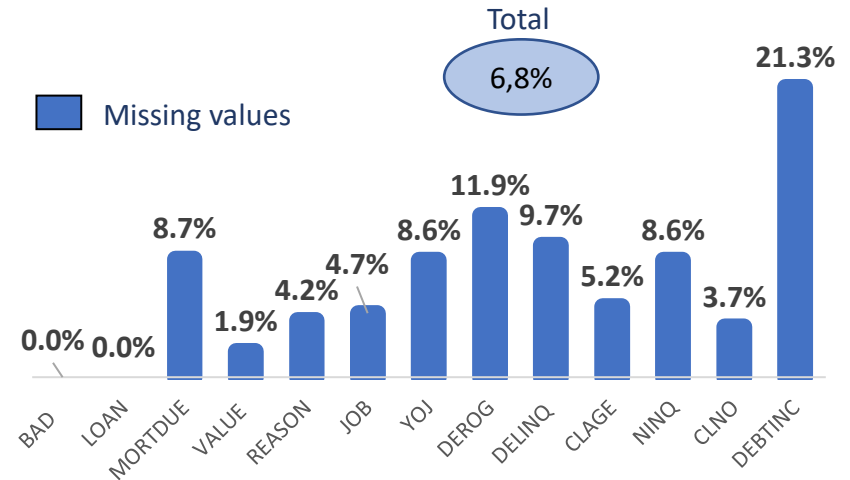
Data Insights



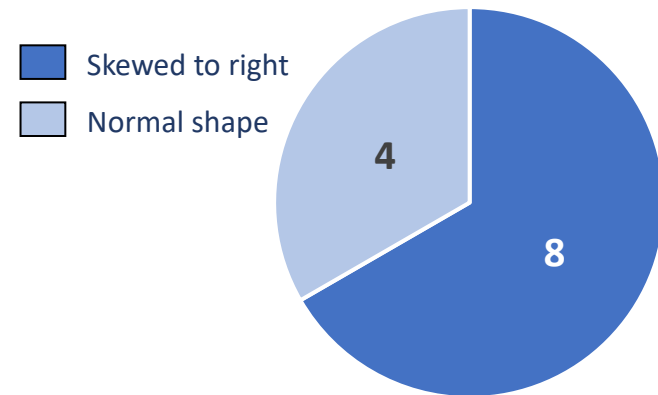
- Data contains 5,960 registers and 13 features.
- Data is unbalanced in a 80%-20% proportion.
- A balancing process of the data is needed to modeling.



- Data types are distributed in 10 numeric features and 3 categorical.
- Categorical data needs to be treated in order to use it.



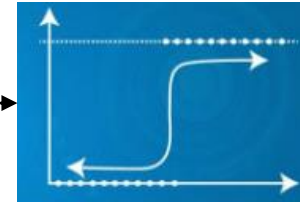
- 6,8% of the data has missing values.
- All features have missing values.
- A filling data process is needed using median and mode.



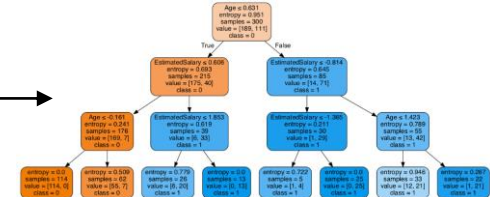
- 8 features are skewed to the right and 4 are normal shaped.
- All features have a big amount of outliers.
- Outliers need to be treated for Logistic Regression modeling

Proposed model solutions

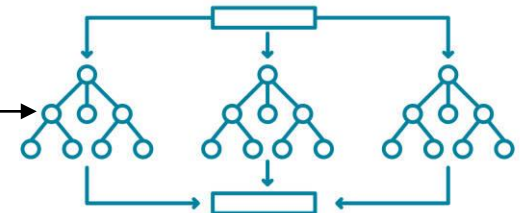
Logistic regression



Decision Trees



Random Forest



Model performances

Baseline model

Tuned model

Logistic
Regression

AC	RC	PR	F1
.81	.08	.68	.15

AC

PR

1417	14
327	30

AC	RC	PR	F1
.76	.41	.39	.40

AC

PR

1208	223
212	145

Decision
tree

AC	RC	PR	F1
.88	.62	.73	.67

(Overfitting)

AC

PR

1349	82
134	223

AC	RC	PR	F1
.86	.74	.62	.68

AC

PR

1270	161
93	264

Random
Forest

AC	RC	PR	F1
.91	.68	.84	.75

(Overfitting)

AC

PR

1345	46
115	242

AC	RC	PR	F1
.91	.68	.82	.75

AC

PR

1379	52
113	244

Comparison of techniques and performances

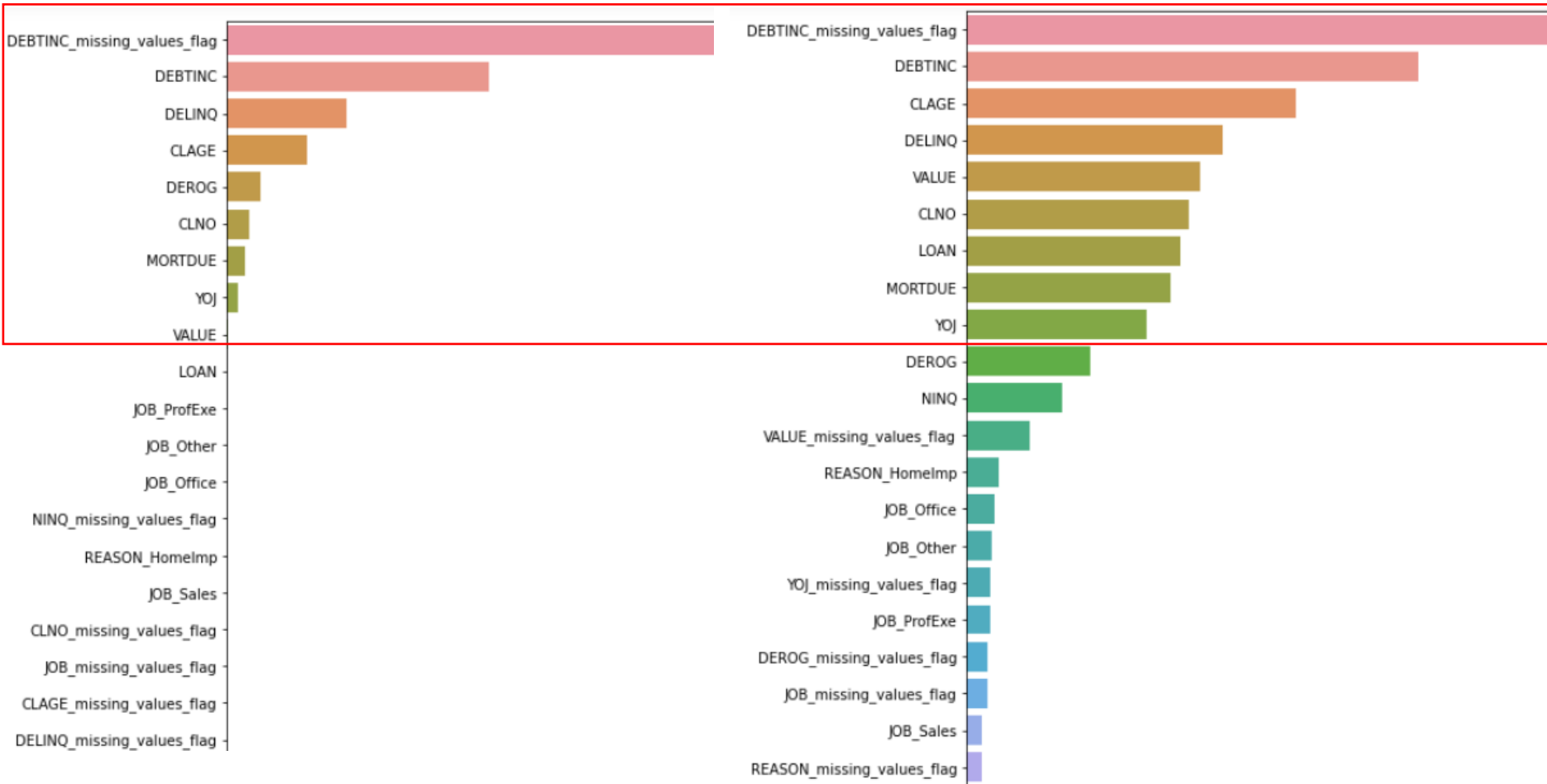
Models	Accuracy	Recall	Precision	Pros	Cons
Tuned Random Forest	0.91	0.68	0.82	-Higher precision than Decision Trees	-Non interpretable
Random Forest	0.90	0.63	0.85	-Higher precision than Decision Trees	-Non interpretable
Tuned Decision Tree	0.86	0.74	0.62	-Highest recall -Interpretable	-Lower precision than Random Forest
Decision Tree	0.88	0.62	0.73	-Similar recall than Random Forests. -Interpretable	-Lower precision than Random Forest

- * Data observed from test sets
- ** Logistic Regression wasn't included due to its very low recall

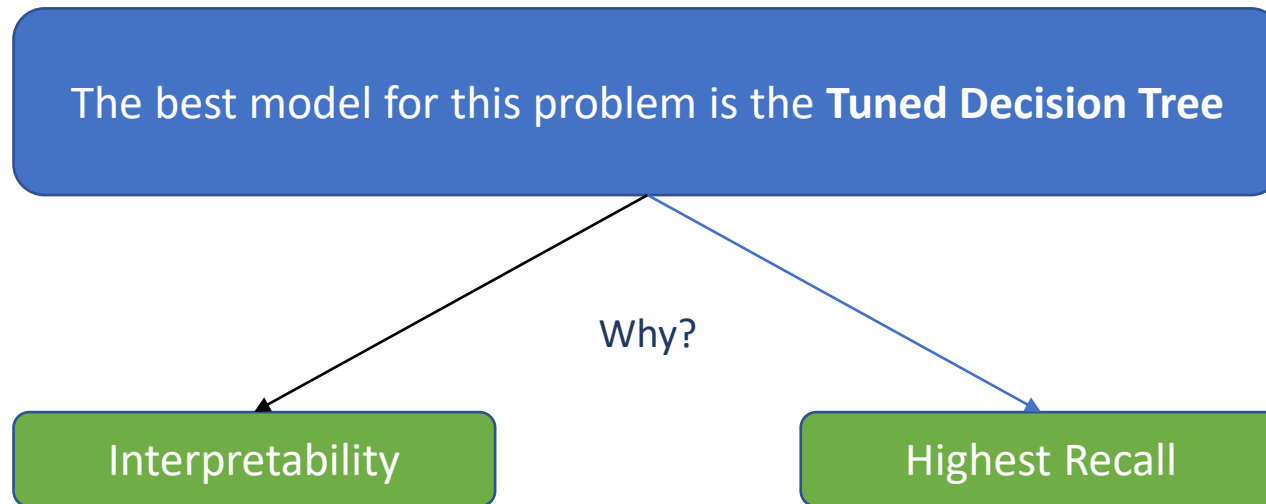
Importance features

Tuned decisión tree

Tuned random forest



Proposal for the Final Solution Design



- The Tuned Decision Tree is giving 0.74 recall (highest), 0.86 accuracy and a precision of .62
- A decision tree is by nature very interpretable
- The most important features used by it are the same than the tuned random forest.

Executive summary

- A decision tree model can predict loan defaulters 74% of the time they come to ask for a home loan.
- This model is highly interpretable.
- The most important features to make a proper prediction are DEBINC, DELINQ and CLAGE.
- The Debt/Income ratio is the most important feature but also the one with the most missing data (21.3%) which is similar to the proportion of defaulted customers (20%).
- It is recommendable to explore the possibility to create an alternative business process to manage and take decisions on those clients with no Debt/Income ratio available.

Recommendations and next steps

- Check the possibility to create an alternative business process to manage and take decisions on those clients with no Debt/Income ratio available.
- Explore other machine learning techniques such as engineering features, dropping columns, support vector machine, neural networks, among others.
- Create a pilot test with the new model and compare the results with the current manual process before completing the transition to the new model.
- Check if there is a way to complete the missing values in the data.

Risks and challenges

- The major risk associated with this project is to underperform versus the current and manual process.
- A big challenge will be changing the internal culture of the bank to adapt the new model to it.
- Another challenge is to exceed the current incomes with the new model.

Appendix

Why is recall important to this project?

The model can make two types of wrong predictions:

1. Predicting a client will pay his loan when the client can't pay.
2. Predicting a client won't pay his loan when the client can pay.

Which case is more important?

• **Predicting that a client will pay but the client can't**, i.e., losing money immediately. This would be considered the most important case of wrong predictions because bad loans (NPA) usually eat up a major chunk of the bank profits.

How to reduce this loss i.e the need to reduce False Negatives?

• **The bank would want the RECALL to be maximized**, the greater the Recall, the higher the chances of minimizing false negatives. Hence, the focus should be on increasing the Recall (minimizing the false negatives) or, in other words, identifying the true positives (i.e. Class 1) This would help in increasing the bank profit that comes from interests in the form of home loans.

