

Loan Default Prediction

Aysan Jarvand



Problem Definition

CONTEXT

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



PROBLEM

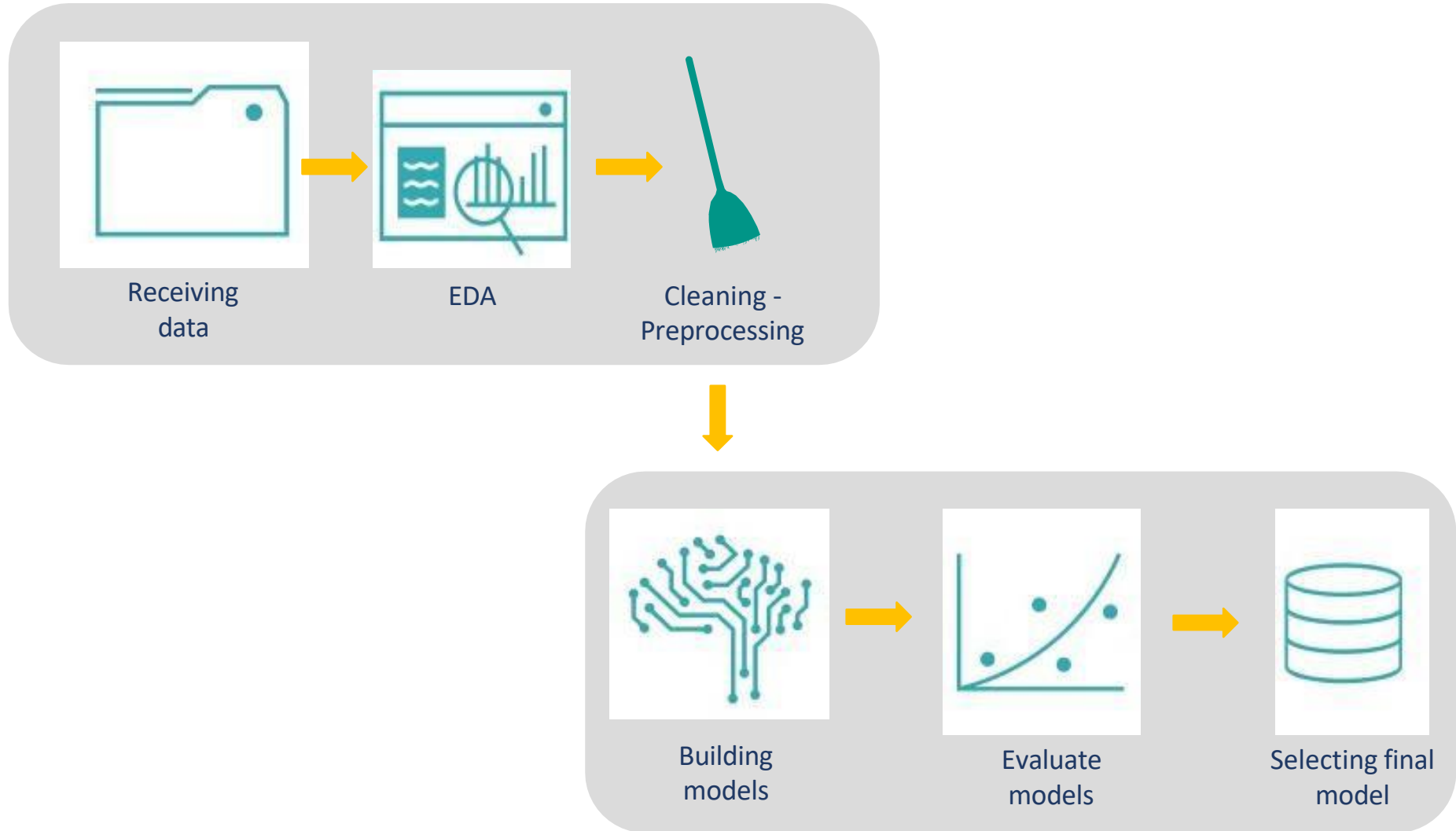
- Banks need an effective approval process.
- This process is effort-intensive and sensitive to human error and biases.



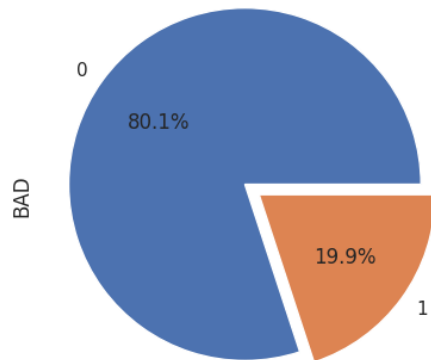
OBJECTIVES

- Predict clients who are likely to default on their loan
- Avoid the risk of misclassification of default loans predicted as non-default loans.

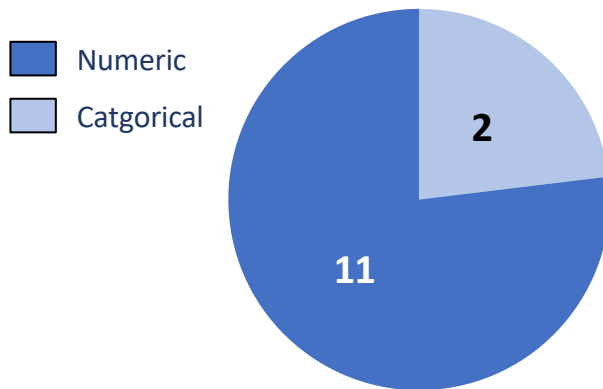
Solution approach



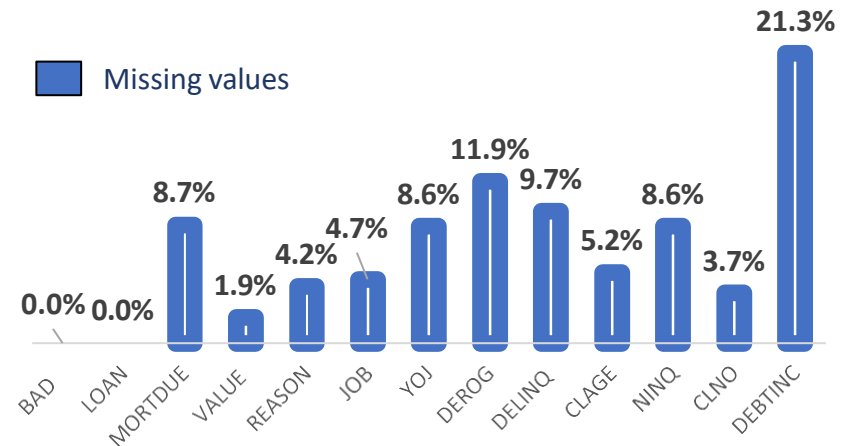
Data Insights



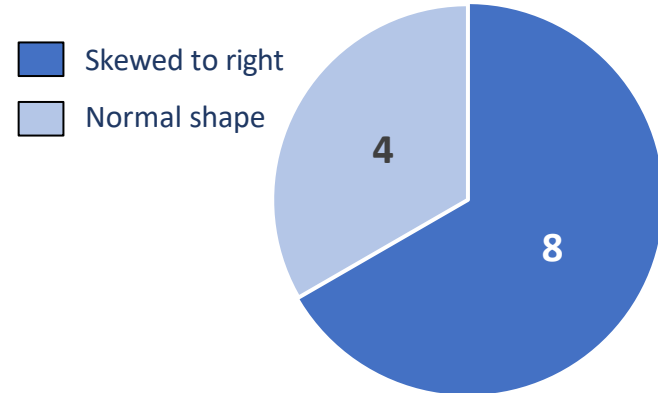
- Data contains 5,960 registers and 13 features.
- Data is unbalanced in a 80%-20% proportion.
- A balancing process of the data is needed to modeling.



- Data types are distributed in 11 numeric features and 2 categorical.
- Categorical data needs to be treated in order to use it.



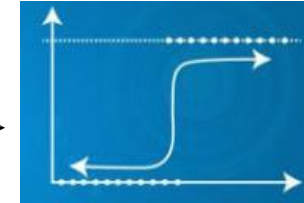
- Except Loan and BAD all features have missing values.
- A filling missing value is needed using median and mode.



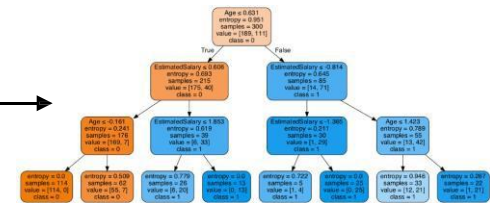
- 8 features are skewed to the right and 4 are normal shaped.
- All features have a big amount of outliers.
- Outliers need to be treated for Logistic Regression modeling

Proposed model solutions

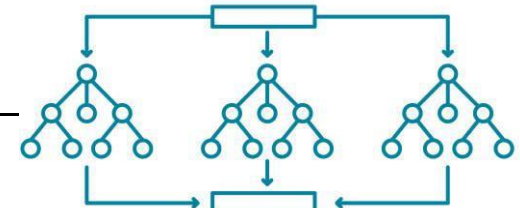
Logistic regression



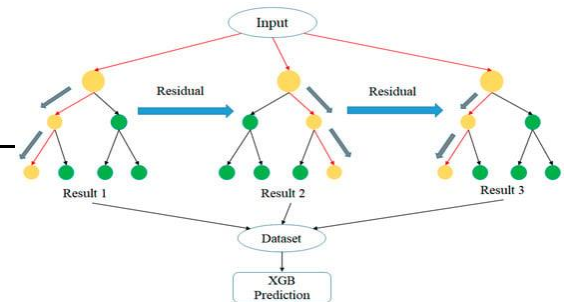
Decision Trees



Random Forest



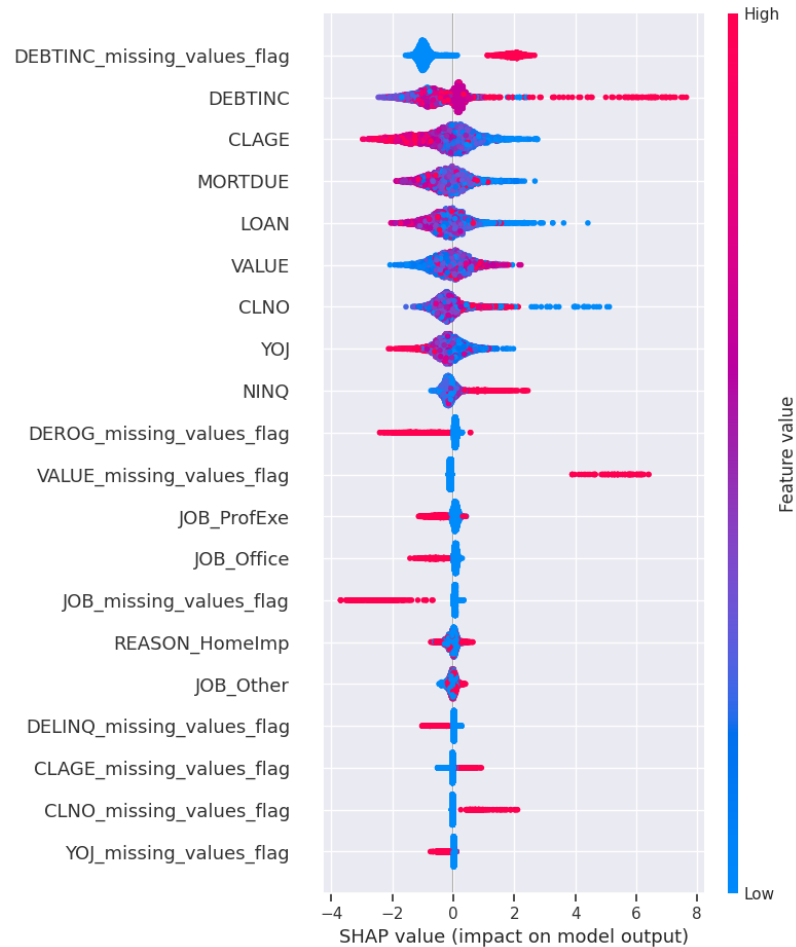
XGBoost Classifier



Comparison of techniques and performances

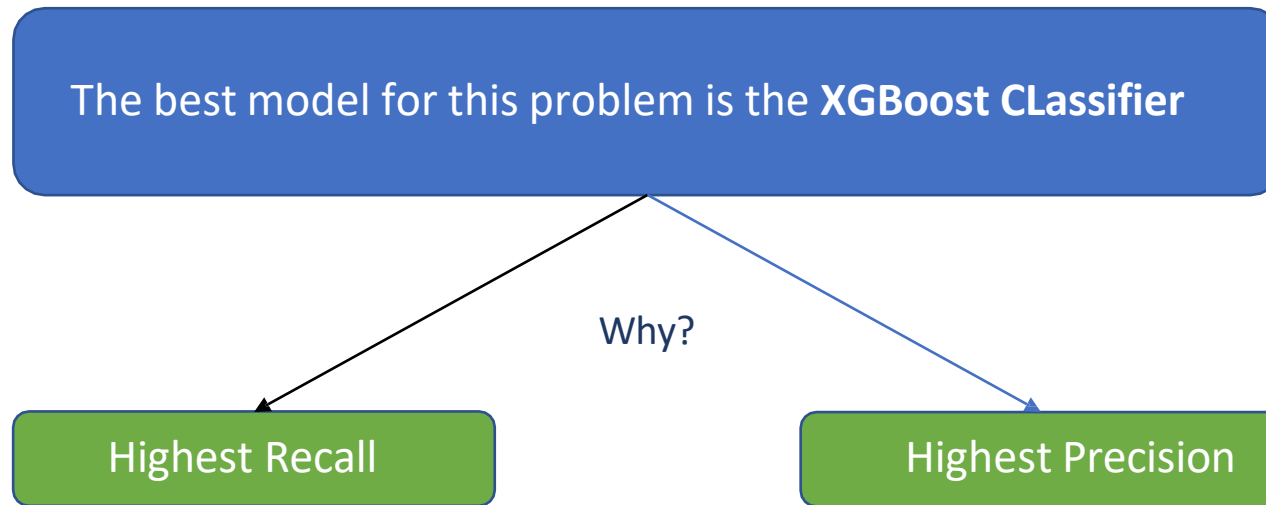
Models	Accuracy	Recall	Precision	Pros	Cons
XGBoost Classifier	.91	.82	.87	-High precision, recall & accuracy	-Non interpretable
Tuned Random Forest	0.88	0.78	0.83	-Higher precision than Decision Trees	-Non Interpretable
Random Forest	0.90	0.82	0.85	-Higher precision than Decision Trees	-Non interpretable
Tuned Decision Tree	0.86	0.76	0.78	-Highest recall -Interpretable	-Lower precision than Random Forest
Decision Tree	0.86	0.76	0.78	-Similar recall than Random Forests. -Interpretable	-Lower precision than Random Forest
Logistic Regression	.81	.54	.75	-Higher precision than Decision Trees -Interpretable	-Very low recall

Model Interpretability check with Shap values



- High DEBTINC values are on the right side, contribute positively to identify defaulter

Proposal for the Final Solution Design



- The XGBoost classifier is giving highest 0.82 recall , 0.91 accuracy and a precision of .87

Executive summary

- A XGBoost Classifier model can predict loan defaulters 82% of the time they come to ask for a home loan.
- This model has highest 82% recall for class 1 on the test data, 87% precision score, and 91% accuracy.
- The most important features to make a proper prediction are DEBINC, DELINQ and CLAGE.
- The Debt-to-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%).

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, different method for handling missing value with new algorithms like other Boosting algorithms, Support Vector Machine, KNN, neural networks.
- 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.
- Consider using balancing techniques, such as ROSE (Random Over Sampling Examples) and SMOTE (Synthetic Minority Oversampling Technique) or under sampling methods, to generate synthetic data. These methods have achieved good results on unbalanced datasets.

Risks and challenges

- The major risk associated with developing a product that meets the needs of both consumer and business.
- A big challenge is to increase the current incomes with the new model.
- There will be a emotional challenge to launch new technology while supporting old ones.
- Cost will be related to train the user with new system.

Thank You

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.

- * high recall + high precision : the class is perfectly handled by the model
- * low recall + high precision : the model can't detect the class well but is highly trustable when it does
- * high recall + low precision : the class is well detected but the model also include points of other classes in it
- * low recall + low precision : the class is poorly handled by the model

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.

