

# HOME CREDIT DEFAULT RISK

# INTRODUCTION

The focus of this ADS (Automated Decision System) is to examine the Home Credit default risk machine learning model, which utilizes historical loan application information to determine whether a loan applicant will be capable of repaying the loan

This is a standard supervised classification task:

- Supervised: The labels are included in the training data and the goal is to train a model to learn to predict the labels from the features
- Classification: The label is a binary variable, 0 (will repay loan on time), 1 (will have difficulty repaying loan)

Additionally, we assessed the model's potential bias by examining the protected characteristic of gender, with females assigned the value 1 and males assigned the value 0.

# OBJECTIVE

This ADS model's objective is to forecast a borrower's ability to repay a loan using past information from loan applications. Accurate predictions will help Home Credit expand financial inclusion for the unbanked population. However, if a loan application is wrongfully denied, it can have a negative impact on the individual borrower. To analyze the algorithm for biases and fairness, a multitude of evaluation metrics will be used.

# DATA

- 7 different data sets available.
- But the ADS utilized only Application Train to evaluate the model.

## Application Train Head

	SK_ID_CURR	TARGET	NAME_CONTRACT_TYPE	CODE_GENDER	FLAG_OWN_CAR	FLAG_OWN_REALTY	CNT_CHILDREN	AMT_INCOME_TOTAL	AMT_CREDIT	AMT_ANNUITY
0	100002	1	Cash loans	M	N	Y	0	202500.0	406597.5	24700.5
1	100003	0	Cash loans	F	N	N	0	270000.0	1293502.5	35698.5
2	100004	0	Revolving loans	M	Y	Y	0	67500.0	135000.0	6750.0
3	100006	0	Cash loans	F	N	Y	0	135000.0	312682.5	29686.5
4	100007	0	Cash loans	M	N	Y	0	121500.0	513000.0	21865.5

## Key Categories

- **SK\_ID\_CURR:** Customer ID
- **TARGET:** 0-will repay loan on time, 1-will have difficulty repaying loan
- **Code Gender:** M- Male, F- Female

# DATA PROCESS



These steps have been taken for the data cleaning and any other pre-processing.

- Data Cleaning: Invalid characters removed to ensure compatibility with model interpretation methods
- Data Shuffling: Shuffled training data to ensure that the order of the data does not affect the learning of the model.

The original code only did a few other data cleaning and preprocessing steps. These include handling missing values and normalizing numerical features.

# SYSTEM IMPLEMENTATION

## These steps have been taken for Model

1. Dataset Loading: The application\_train dataset is a set of loan applications with features indicating difficulty in repaying the loan.
2. Data Balancing: The dataset has a significant imbalance between classes, so a random subset of the majority class is selected to match the number of samples in the minority class. (This is a form of undersampling to balance the dataset and reduce the potential for the model to be biased toward the majority class.)
3. Handling Categorical Features: Factorize() from pandas converts categorical features into numerical features, which is necessary for many machine learning models as they can only handle numerical input.
4. Data Splitting: The application\_train was split into a training and testing set, with 80% of the data used for training and 20% for testing.
5. Model Training: The LightGBM model is trained on a training dataset with parameters defined.
6. Feature Importance Assessment: Feature importance is calculated and visualized to identify which features are most influential in predicting the target variable.
7. Model Prediction: The trained LightGBM model predicts the target variable on the test set, resulting in a data frame.

# RESULT



	SK_ID_CURR	TARGET
285628	430803	0.348510
156868	281818	0.188688
83433	196772	0.590860
2145	102520	0.237021
248579	387603	0.562758
...	...	...
48099	155703	0.778788
181437	310290	0.267582
176837	304920	0.823233
191368	321888	0.708272
204044	336556	0.371301

[9930 rows x 2 columns]

Output of 9,930 customers' probability of difficulty of repaying the loan (Target).

# ADS VALIDATION

Train-Test Split: The code divides the dataset into training and testing. Holdout validation is a type of validation in which the model is trained on one set of data and then validated on another set of data that it did not see during training.

Besides Train-Test Split, the original ADS has no other validation methodology.



# CONFUSION MATRIX & ACCURACY METRICS

Confusion Matrix for Female Subpopulation (1):

```
[[2485  841]
 [1039 1798]]
```

Confusion Matrix for Male Subpopulation (0):

```
[[ 979  660]
 [ 494 1634]]
```

Metrics for Female Subpopulation (1):

Accuracy: 0.6950

Precision: 0.7471

Recall (Sensitivity): 0.7052

F1-score: 0.7255

Metrics for Male Subpopulation (0):

Accuracy: 0.6937

Precision: 0.5973

Recall (Sensitivity): 0.6646

F1-score: 0.6292

From these matrices, we can see that the model seems to have a slightly higher False Negative rate for the female subpopulation and a slightly higher False Positive rate for the male subpopulation.

Overall, the model seems to perform slightly better on the female subpopulation in terms of precision, recall, and F1-score. However, the accuracy is marginally better for the male subpopulation.

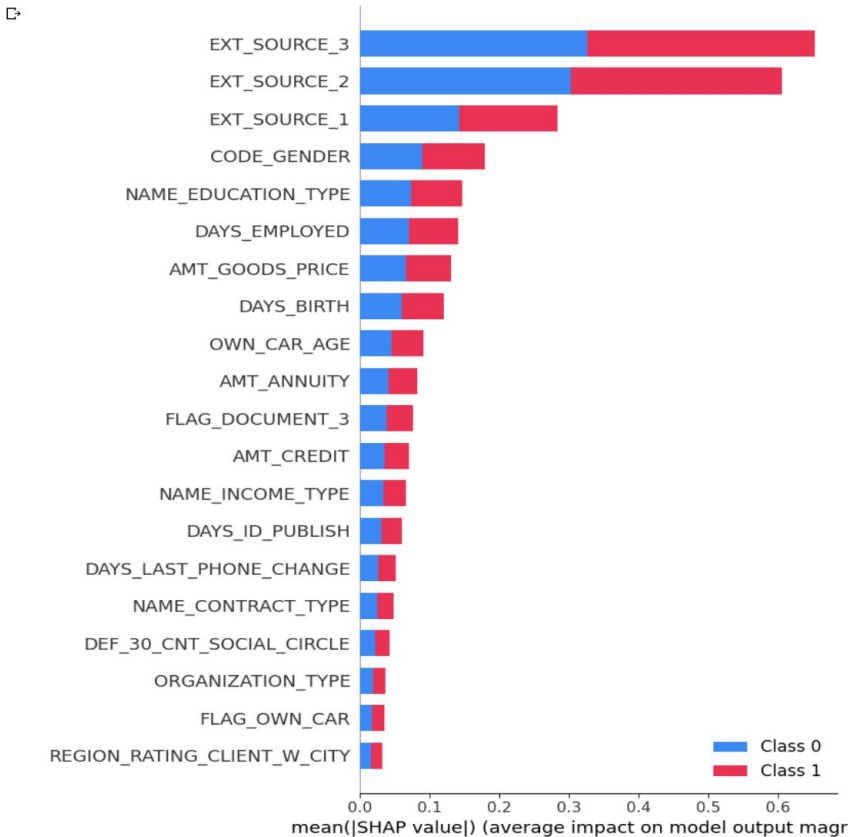
# FAIRNESS METRICS

Demographic Disparity Difference (Selection Rate): 0.1808  
Demographic Disparity Ratio (Selection Rate): 0.6839  
Equalized Odds (TPR): 0.1498  
Equalized Odds (FPR): 0.1341  
Equal Opportunity: 0.1498  
Type I Parity (Equalized False Discovery Rate): 0.0405  
Type I Parity (Equalized False Positive Rate): 0.1341  
Type II Parity (Equalized False Omission Rate): 0.0310  
Type II Parity (Equalized False Negative Rate): 0.1498

These fairness metrics suggest there are disparities in how the model performs for different subpopulations (male vs. female).

# SHAP

SHAP VALUE



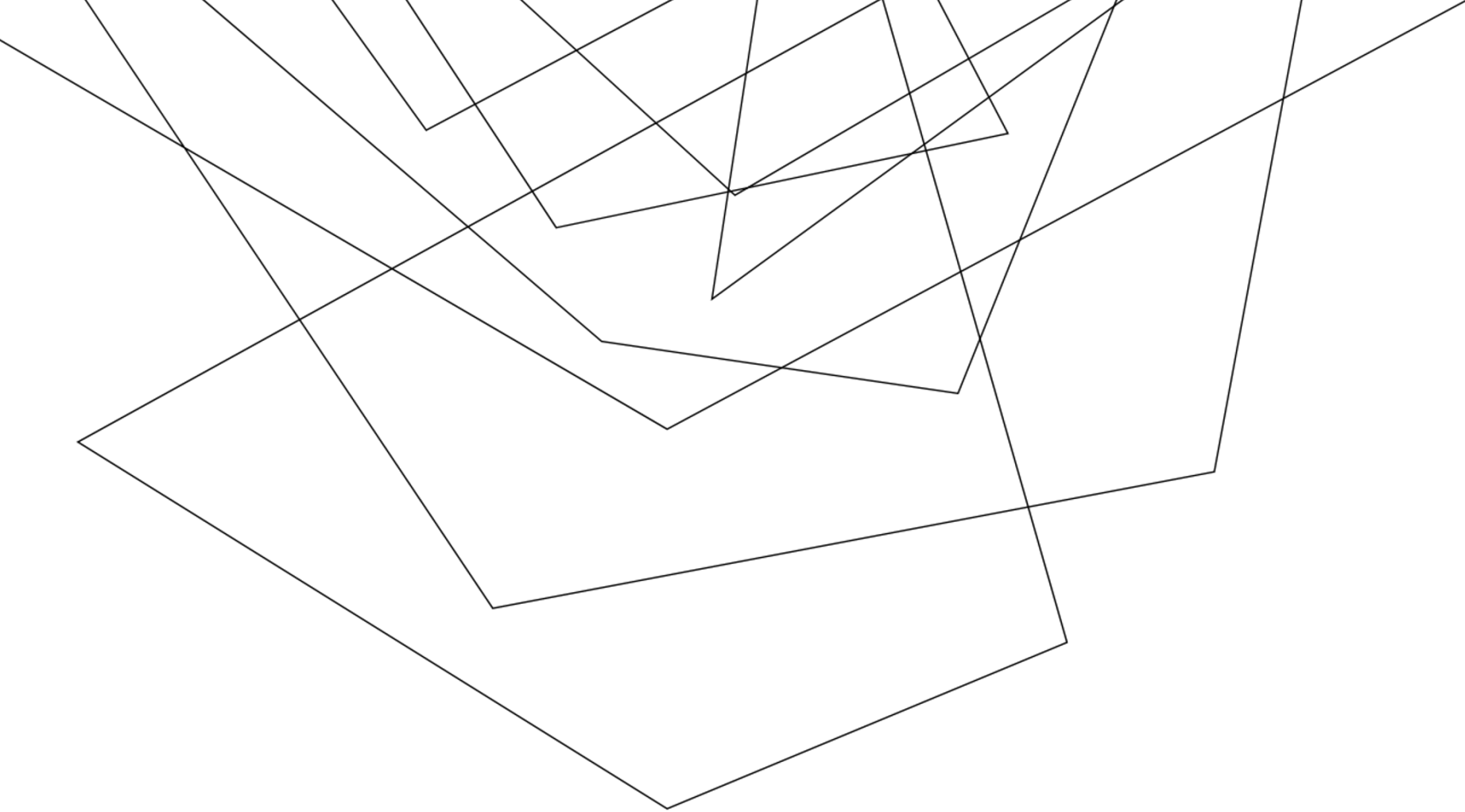
These are SHAP values plot and summary plot.



# SUMMARY



- The ADS provides a good starting point
- Further validation, testing, and potential adjustment before deployment in the public sector or industry.
- A few improvements could be made to improve the data collection, processing, and analysis methodology.



# HOME CREDIT DEFAULT RISK