

Home Credit Default Risk

Capstone Project

Omar Alqasem
November 3rd, 2018

I. Definition

Project Overview

Credit risk is the risk of not meeting the legal obligations from a borrower not making the payments required. Banks use strict credit risk management to minimize lending borrowers who fail to repay the obligated payments upon agreed terms based on applicant credit history. Therefore, many applicants with insufficient or non-existing credit history struggle to get loans. As a result, these applicants are taken advantage of by untrustworthy lenders.

Home Credit Group (HCG) goal is to help applicants with non-existing or insufficient credit history with cash or revolving loans. The project solution shall help HCG to take the right decision whether to give a loan to an applicant based on their application information. The input data of applicants' information has been given to the public in a competition hosted by Kaggle to select the best solution with the highest evaluation score. This project will use the same dataset given for the competition.

Problem Statement

Home Credit Group's (HCG) goal is to open their doors for applicants with insufficient or non-existing credit history to apply for loans and receive a positive loan experience. At the same time, HCG wants to eliminate applicants who fail to repay their obligated payments. Therefore, in order to help eliminating rejections for capable borrowers of repaying their loans, Machine Learning will be applied to this problem to predict new applicants' ability to repay their loans. This can be achieved by training a classification model on historical data of past applicants. The dataset of past applicants is taken from Kaggle website that's part of a previous competitions. The end trained model shall classify new applicants' ability to repay for the request loan.

Metrics

The dataset given by Home Credit Group (HCG) is an imbalanced dataset. This means that the dataset contains much more records applicants that were capable to repay the loan than applicants whom defaulted. Therefore, accuracy measure shall not be reliable to evaluate the classification models since the majority of the dataset belongs to only one class. In this project, we will mainly use Area Under the Receiver Operating Characteristic Curve (ROC AUC) score for evaluating the model performance since it calculates the sensitivity and specificity across a continuum of cutoffs.

II. Analysis

Data Exploration

The dataset used for this project has been taken from one of Kaggle's past competitions in CSV format. It constructs of 122 features and over 300K records. These records represent historical loan applications and their target labels of whether the applicant has repaid the full loan installments. Like any other real-world scenario, the number of positive class (default applicants) weighs only 8% of the dataset and the rest is marked as negative class (paid successfully). To get a sense of how the dataset looks like, Table II.I below shows the first 5 records of the dataset along.

	SK_ID_CURR	TARGET	NAME_CONTRACT_TYPE	CODE_GENDER	FLAG_OWN_CAR	FLAG_OWN_REALTY	CNT_CHILDREN	AMT_INCOME_TOTAL	AMT_CRED
0	100002	1	Cash loans	M	N	Y	0	202500.0	406597
1	100003	0	Cash loans	F	N	N	0	270000.0	1293502
2	100004	0	Revolving loans	M	Y	Y	0	67500.0	135000
3	100006	0	Cash loans	F	N	Y	0	135000.0	312682
4	100007	0	Cash loans	M	N	Y	0	121500.0	513000

Table II.I: First 5 records of Home Credit Default Risk taken from Kaggle competitions

The dataset contains many features related to the loan application and applicant's personal information. Features related to the loan application include loan type whether cash or revolving loan, credit total requested by client, final credit total given for client, down payment ...etc. Likewise, features related to applicant's personal information include gender, number of children, education level, family status, duration of employment ...etc.

During the data exploration, it has been encountered some abnormalities within the dataset. As an example, one of the applicant's employment duration exceeded 1000 years. Likewise, another applicant's gender is marked as XNA. These abnormalities have been mutated accordingly. Furthermore, the dataset suffers from lots of missing values for some features where they exceeded 60% of the dataset missing. Those missing features have been addressed to mutate with the average values. Finally, categorical variables do exist in the dataset and they have been one-hot encoded for the classification models to understand.

Exploratory Visualization

It's a common practice to visualize some of the characteristics in the dataset used for any machine learning problem in order to make sense out of it. As for this project, we begin with the question, how many applications have paid their loan successfully. Figure 2.1 shows that more than 250K loan applications have been successfully paid while less than 30K applications have defaulted. This figure indicates that HCG data is an imbalanced dataset.

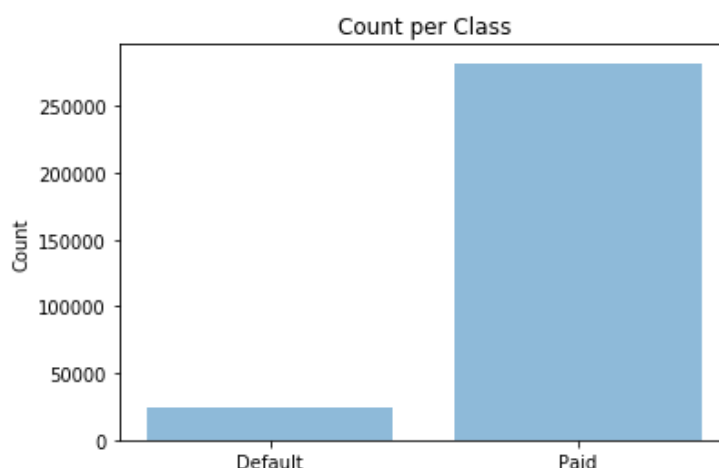


Figure 2.1: Number of loan applications that have been paid versus loan applications that have been defaulted

Next, categorical features are good starting point to find any relevance with the target variable. We begin with features related to applicants' personal information such as gender, family status, children count etc. This will provide good insights whether these features have any relationship with the ability of repaying the loan installments. As shown in Figure 2.2, we can identify that cash loans applications happen to have more default than revolving loans.

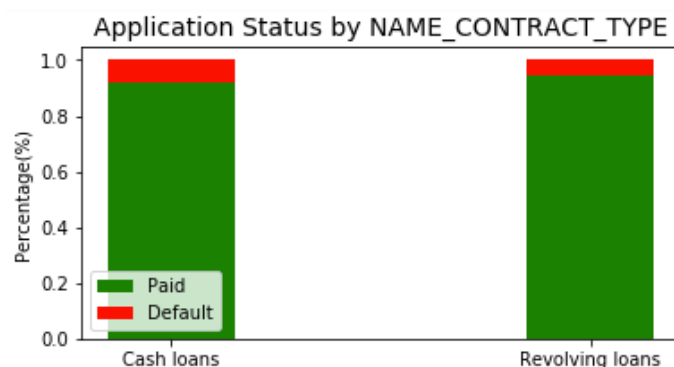


Figure 2.2: % of paid application versus % of default application by contract type.

Furthermore, Figure 2.3 shows that 85% of male applicants have successfully paid the full loan installments and 90% of female applicant have successfully paid the full

loan installments. This shows that male applicants are more likely to default than female applicants.

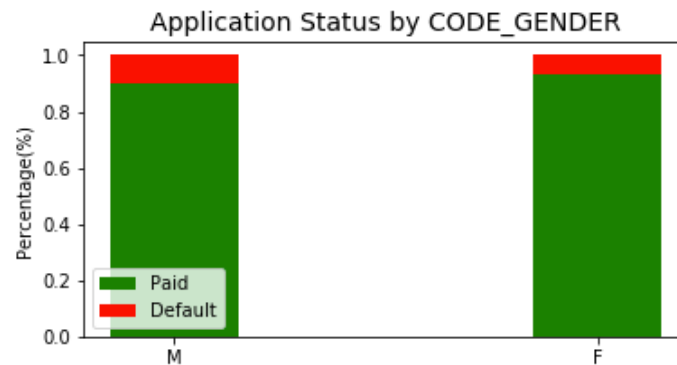


Figure 2.3: % of paid application versus % of default application by applicant gender.

In addition, income type is an important factor in identifying whether an applicant is willing to repay the loan. Figure 2.4 shows that students and businessmen are more likely to successfully pay the loan installments. On the other hand, unemployed and maternity leave applicants have a higher risk to default than any other applicants.

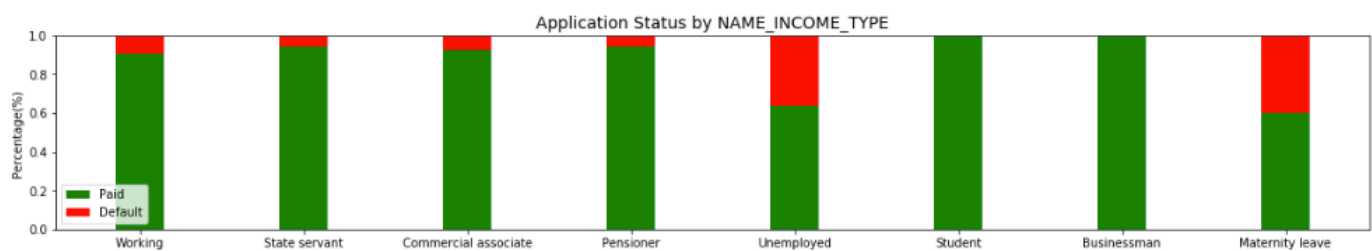


Figure 2.4: % of paid application versus % of default application by applicant income type.

Besides categorical features, numerical features provide meaningful insights as well. One of the important numerical feature in identifying applicant's ability to repay the full loan amount is the applicant age. Figure 2.5 shows that younger applicants are more likely to default than older applicants.

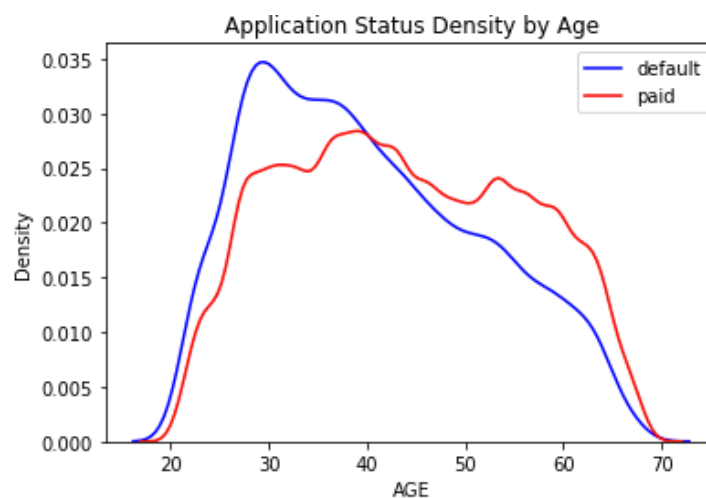


Figure 2.5: Density of paid applications versus Density of default application by applicant age.

Algorithms and Techniques

As a starting point, it is unclear which classification model will best solve the problem of predicting home credit default risk. Therefore, we will use different classification models for training and testing then a final classification model with the best evaluation metric will be selected. In this project, we will use classification algorithms such as Logistic Regression, Adaptive Boosting, Naïve Bayes, XGBoost and Light GBM classifiers for predicting home credit default risk.

Logistic Regression is a type of classification algorithm involving a linear discriminant. logistic regression outputs a probability that the given input point belongs to a certain class.

Naive Bayes Classifiers are a family of simple "probabilistic classifiers" based on applying Bayes' theorem with strong (naive) independence assumptions between the features. [1]

Boosting algorithms such as Adaboost and XGBoost fit a sequence of weak learners on different weighted training data. It starts by predicting original data set and gives equal weight to each observation. If prediction is incorrect using the first learner, then it gives higher weight to observation which have been predicted incorrectly. Being an iterative process, it continues to add learner(s) until a limit is reached in the number of models or accuracy. [2]

LightGBM grows tree vertically while other algorithm grows trees horizontally meaning that Light GBM grows tree leaf-wise while other algorithm grows level-wise. It will choose the leaf with max delta loss to grow. When growing the same leaf, Leaf-wise algorithm can reduce more loss than a level-wise algorithm [3].

Recall that the dataset used for this project is imbalanced where the number of applications that have been successfully paid weighs more than 90% of the dataset. Therefore, we will scale the positive class to control the balance of positive class and negative class in the initial run for applicable classification models such as Logistic Regression, XGBoost and LightGBM models.

[1] https://en.wikipedia.org/wiki/Naive_Bayes_classifier

[2] <https://www.analyticsvidhya.com/blog/2015/11/quick-introduction-boosting-algorithms-machine-learning/>

[3] <https://medium.com/@pushkarmandot/https-medium-com-pushkarmandot-what-is-lightgbm-how-to-implement-it-how-to-fine-tune-the-parameters-60347819b7fc>

Benchmark

There are many existing research studies that predict whether a loan will default using different datasets. In this project, we will use one of these studies' final AUC score as a benchmark for our classification model. One of these studies is The LendingClub's Loan [1] and their best AUC score was 0.711. The aim for this project is to produce a model that results with a better AUC score than 0.711.

III. Methodology

Data Preprocessing

Prior to training a classification model, it is important to do some preprocessing required to avoid problems during model training. The preprocessing steps for this project include handling missing values in specific features, dealing with abnormalities in the dataset, converting categorical features to binary format, and normalizing the numerical features.

Missing values in the dataset's features cause runtime errors during model training. The dataset contains 128 features and 67 out of these features have missing values. The number of missing values ranges between less than 1% of the dataset and up to 70% for some of the features. For categorical features, missing values have been mutated with the most repetitive feature. For example, NAME_TYPE_SUITE is a categorical feature that's missing 42% of the dataset. Figure 3.1 shows that "Unaccompanied" is the most repetitive feature value. Missing values were mutated with Unaccompanied value.

```
Unaccompanied    248526
Family           40149
Spouse, partner   11370
Children         3267
Other_B          1770
Other_A          866
Group of people   271
Name: NAME_TYPE_SUITE, dtype: int64
```

Figure 3.1: NAME_TYPE_SUITE feature values count.

For numerical features, a different approach has been taken to handle missing values. Feature like OWN_CAR_AGE has more than 65% missing values. This feature is interrelated with FLAG_OWN_CAR in a way such that if the applicant does not own a car, OWN_CAR_AGE will not show a value. Therefore, we replaced missing data with a zero if FLAG_OWN_CAR is true. For other numerical missing values, we replaced them with -999 to represent a missing value instead of NaN.

Abnormalities do exist in the dataset as discussed in previous sections. CODE_GENDER that represents the applicant's gender contain values either man or female. The dataset happens to have a third value of XNA which has been replaced with Female as it's the most common value for that feature. Another abnormality occurs in the dataset is DAYS_EMPLOYED which represents the duration of employment in days. One record has 365243 as a value for DAYS_EMPLOYED which is about 1000 years. This is not a realistic value and therefore it has been replaced with a zero.

Many machine learning algorithms cannot understand categorical directly. They require input data to be in numerical format. To handle this matter, one-hot encoding

is performed on all 16 categorical features to convert to binary format. Furthermore, it is a common practice to perform logarithmic transformation on highly skewed numerical features so that large and small data do not negatively affect the performance of the learning algorithm. We applied the logarithmic transformation to numerical features in the dataset with care to logarithm of 0 to avoid undefined values. Finally, in addition to performing logarithmic transformation, we normalized the numerical features in the dataset to limit the values between 0 and 1. This is to ensure that all features are treated equally when applying a supervised learner.

Implementation

AdaBoostClassifier, GaussianNB, and LogisticRegression classification models are implemented using scikit-learn Python module, while XGBoost and LightGBM are implemented using XGBoost and LightGBM module respectively. The dataset has been split into training and testing dataset with a test size of 30% of the dataset. The classification algorithms are initialized with the default parameters with a random state of '48'. Each model is trained on the training dataset by using .fit function. All classification models accept a Pandas object for the dataset except LightGBM which requires an LGBM dataset by using the below line of code.

```
lgb.Dataset(data=X_train,label=y_train)
```

Since the home credit default risk dataset is an imbalanced dataset, all classification algorithms predict a negative class for each input data. Therefore, we addressed this issue by scaling the positive class to equal the total number of the training dataset divided by the total number of positive class of the training dataset as illustrated below.

```
pos_class = target.count() // np.sum(target) # positive class scaler
```

After scaling the positive class, the classification models started to automatically adjust weights inversely proportional to class frequencies in the input data. After performing training and testing, the resulted confusion matrix and evaluation metric is captured below for each classification model. Note that all classification models return a binary prediction for each test data input except LightGBM which returns a probability. LightGBM probability will be converted to 1 if the probability is greater than or equal to 0.5, and 0 if the probability is less than 0.5. This is required since confusion matrix, accuracy, recall, and F1 evaluation metric functions expect binary inputs.

Model	TP	FP	TN	FN
AdaBoostClassifier	4	9	84596	7557
GaussianNB	7266	78774	5831	295
LogisticRegression	4680	30370	54235	2881
XGBClassifier	5261	28051	56554	2300
LightGBMClassifier	2361	7588	77017	5200

Figure 3.2: Confusion matrix for each classification model

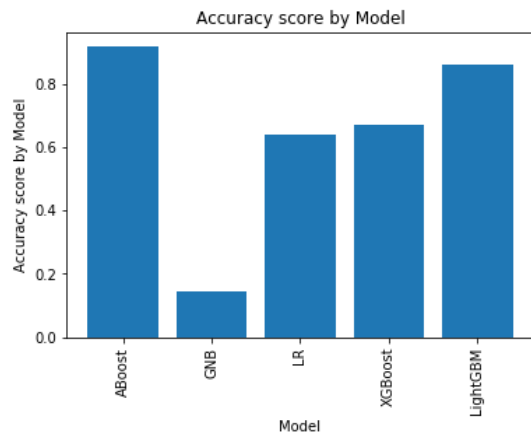


Figure 3.3: Accuracy score by classification model

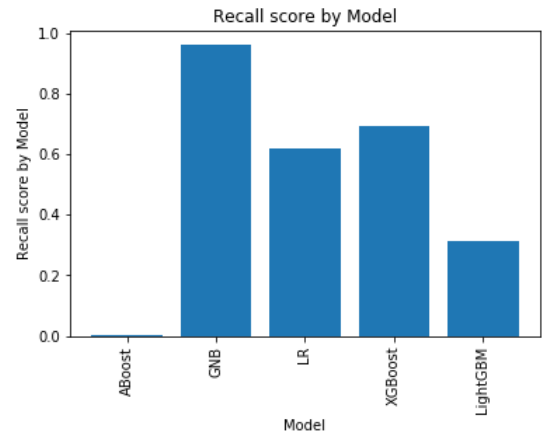


Figure 3.4 Recall score by classification model

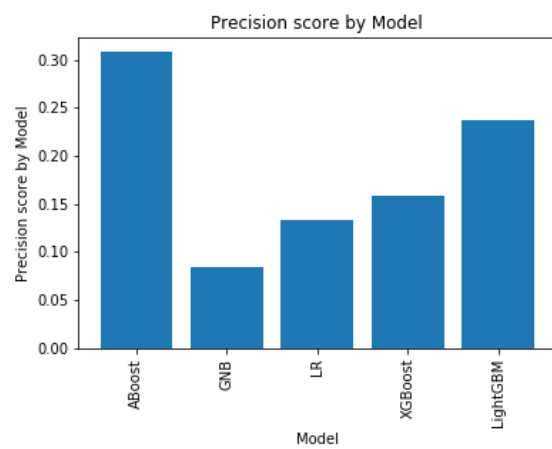


Figure 3.5: Precision score by classification model

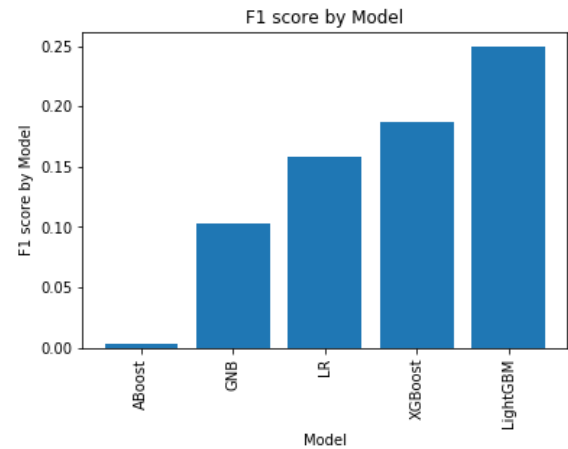


Figure 3.6: F1 score by classification model

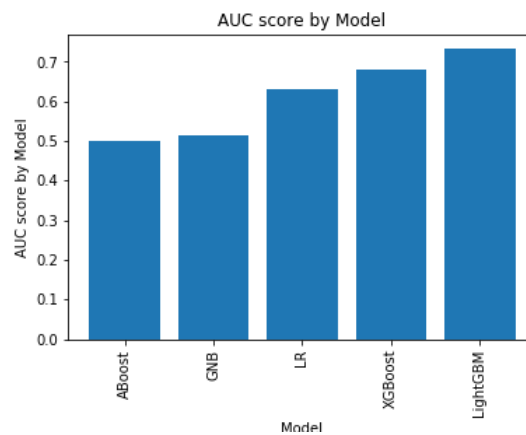


Figure 3.7: AUC score by classification model

After visualizing the evaluation metrics, we can clearly see that LightGBM has the highest F1 and AUC evaluation metrics. Therefore, we choose LightGBM for further improvements which will be detailed in the next section.

Refinement

In order to maximize optimization of classification models, hyper-parameter tuning is performed to achieve this goal. Since LightGBM performed the best compared to the other selected models, we will try to find its best hyper-parameters that can optimally predict home credit default risk on the given dataset. LightGBM has many fine parameters to be tuned which are shown below.

```
params = {
    'boosting_type': 'gbdt', # np.random.choice(['dart', 'gbdt']),
    'objective': 'binary',
    'metric': ['auc', 'f1'],

    'learning_rate': 0.1,

    'num_leaves': np.random.randint(64, 128),
    'max_depth': np.random.randint(6, 12),
    'min_data_in_leaf': int(2 ** (np.random.rand()*3.5 + 9)),

    'feature_fraction': np.random.rand()*0.35+0.65,
    'bagging_fraction': np.random.rand()*0.35+0.65,
    'bagging_freq': 1,

    'lambda_l1': 10 ** (np.random.rand() * 4),
    'lambda_l2': 10 ** (np.random.rand() * 3 + 2),
    'min_gain_to_split': 0.0,
    'min_sum_hessian_in_leaf': 0.1,

    'num_threads': 16,
    'verbose': 0,
    'is_training_metric': 'True',
    'scale_pos_weight': pos_class
}
```

We have set the number of boosting iterations to 5000 and an early stopping rounds 250 so that the model will train until the validation score stops improving.

```
evals_result = {} # dict to store evaluation results of all the items in `valid_sets`
dev_X, val_X, dev_y, val_y = train_test_split(X_train, y_train, test_size = 0.2, random_state = 42)
pos_class = dev_y.count()/np.sum(dev_y) # recalculate positive class weight after CV split

# para searching
```

```
lgb_train = lgb.Dataset(data=dev_X, label=dev_y)
lgb_val = lgb.Dataset(data=val_X, label=val_y)
lgb_test = lgb.Dataset(data=X_test, label=y_test)

gbm = lgb.train(params, lgb_train, 5000, valid_sets=[lgb_train, lgb_val], \
                valid_names = ['train', 'test'], evals_result=evals_result, \
                early_stopping_rounds=250, verbose_eval=100, \
                categorical_feature='auto')
predictions_test_tuned = gbm.predict(X_test)
```

During training, the model has improved AUC score to 0.76 on the cross-validation dataset as shown below.

```

Training until validation scores don't improve for 250 rounds.
[100] cv's auc: 0.753214
[200] cv's auc: 0.758614
[300] cv's auc: 0.760037
[400] cv's auc: 0.760226
[500] cv's auc: 0.760255
[600] cv's auc: 0.760183
[700] cv's auc: 0.76004
Early stopping, best iteration is:
[510] cv's auc: 0.760331

```

Model	TP	FP	TN	FN
LightGBMClassifier	2394	7364	77074	5062

Table 3.8: Confusion Matrix before hyper-parameter tuning

Model	TP	FP	TN	FN
LightGBMClassifier	5177	25976	58734	2279

Table 3.9: Confusion Matrix after hyper-parameter tuning

	Before Hyper-Parameters Tuning		After Hyper-Parameters Tuning	
	Train	Test	Train	Test
Accuracy	0.863	0.862	0.863	0.693
Precision	0.244	0.238	0.244	0.166
Recall	0.330	0.321	0.330	0.694
F1	0.257	0.251	0.257	0.195
AUC	0.748	0.736	0.748	0.759

We can clearly see in table 3.10 a comparison of the evaluation metrics on LightGBM classifier before and after hyper-parameter tuning. We can clearly see an excellent improvement on AUC score from 0.736 to 0.759 for the test dataset.

IV. Results

(approx. 2-3 pages)

Model Evaluation and Validation

Initially, we selected many classification models to predict home credit default risk on a given dataset which are AdaBoostClassifier, Naïve Bayes Classifier, LogisticRegression, XGBoostClassifier, and LightGBM. Since the dataset is an imbalanced dataset, accuracy evaluation metric does not really tell us how the model is performing. Instead, the primary evaluation metric used for this problem is AUC score. After training and testing the dataset, we have collected each model's evaluation metrics and plotted them in below figures.

Model	TP	FP	TN	FN
-------	----	----	----	----

AdaBoostClassifier	4	9	84596	7557
GaussianNB	7266	78774	5831	295
LogisticRegression	4680	30370	54235	2881
XGBClassifier	5261	28051	56554	2300
LightGBMClassifier	2394	7364	77074	5062

Figure 4.1: Confusion matrix for each classification model

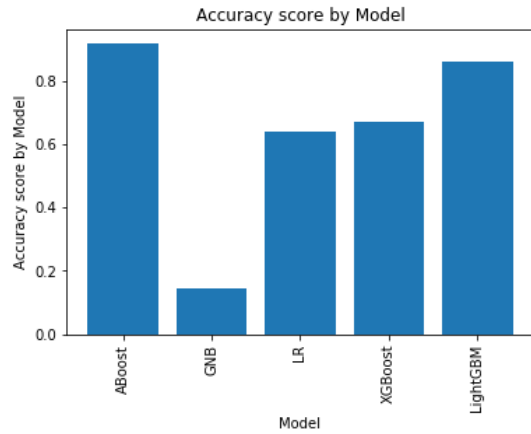


Figure 4.2: Accuracy score by classification model

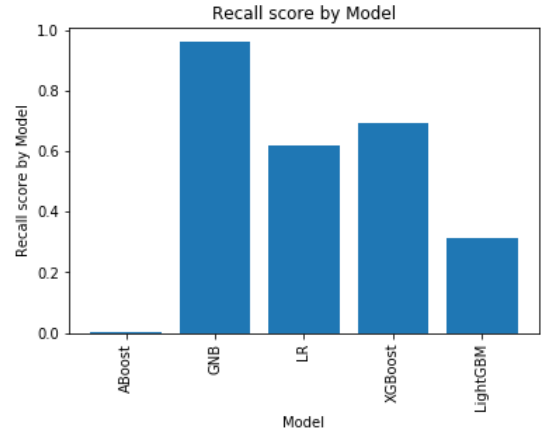


Figure 4.3 Recall score by classification model

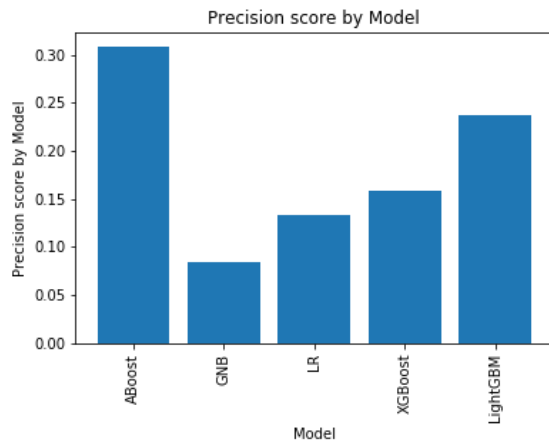


Figure 4.4: Precision score by classification model

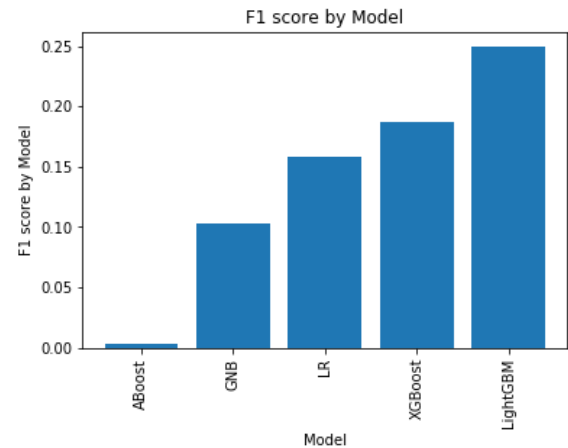


Figure 4.5: F1 score by classification model

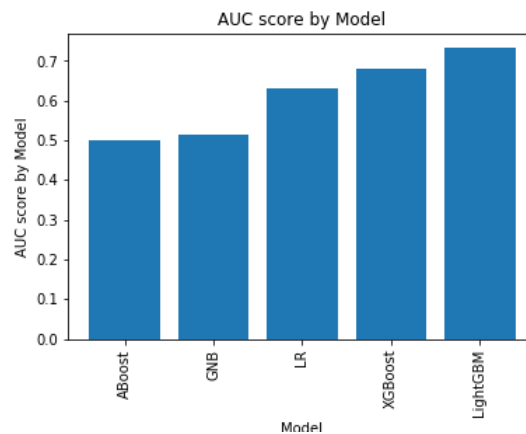


Figure 3.7: AUC score by classification model

It turned out that LightGBM has the highest AUC score of 0.731 and therefore been selected for further improvements. After tuning LightGBM classifier hyper-parameters, we gained an improvement of AUC score to 0.759. The parameters that increased model evaluation is shown below.

```
Best Parameters that improved our final model LightGBM:
{'boosting_type': 'gbdt', 'objective': 'binary', 'metric': ['auc', 'f1'], 'learning_rate': 0.1, 'num_leaves': 116, 'max_depth': 7, 'min_data_in_leaf': 2394, 'feature_fraction': 0.8480511816917773, 'bagging_fraction': 0.7682161532143479, 'bagging_freq': 1, 'lambda_l1': 134.29783711684772, 'lambda_l2': 3105.4418845715327, 'min_gain_to_split': 0.0, 'min_sum_hessian_in_leaf': 0.1, 'num_threads': 16, 'verbose': 0, 'is_training_metric': 'True', 'scale_pos_weight': 12}
```

Model	TP	FP	TN	FN
LightGBMClassifier	5177	25976	58734	2279

Table 3.9: Confusion Matrix after hyper-parameter tuning

	Before Hyper-Parameters Tuning		After Hyper-Parameters Tuning	
	Train	Test	Train	Test
Accuracy	0.863	0.862	0.863	0.693
Precision	0.244	0.238	0.244	0.166
Recall	0.330	0.321	0.330	0.694
F1	0.257	0.251	0.257	0.195
AUC	0.748	0.736	0.748	0.759

Finally, after random parameter search process, the model has gained an improvement of AUC score to 0.76 on the cross-validation test set and 0.759 on the final test dataset. This suggests that the model is sufficiently robust for its intended application and does not show overfitting.

Justification

Recall that the aim is to benchmark our model with existing studies on similar application of predicting loan default on LendingClub dataset. Comparing our model evaluation metric with the benchmark, we can see that our model received an AUC score of 0.759 while the benchmark AUC score is 0.713.

While this performance comparison highlights the advantage of using a complex model over a simple mode in this particular application, it is good to assess the final model on its standalone metrics and consider its limitation of a low precision score of 0.166 is low for a precision-optimized model which highlights a difficulty of minimizing false positive cases.

Despite the limitation of this model, it still can serve its purpose of identifying home credit default risks. This model provides a higher level of confidence in identifying applicants whom will result in failing to repay their loan installments than a naïve prediction. The model results in minimizing home credit default risk.

In this section, your model's final solution and its results should be compared to the benchmark you established earlier in the project using some type of statistical analysis. You should also justify whether these results and the solution are

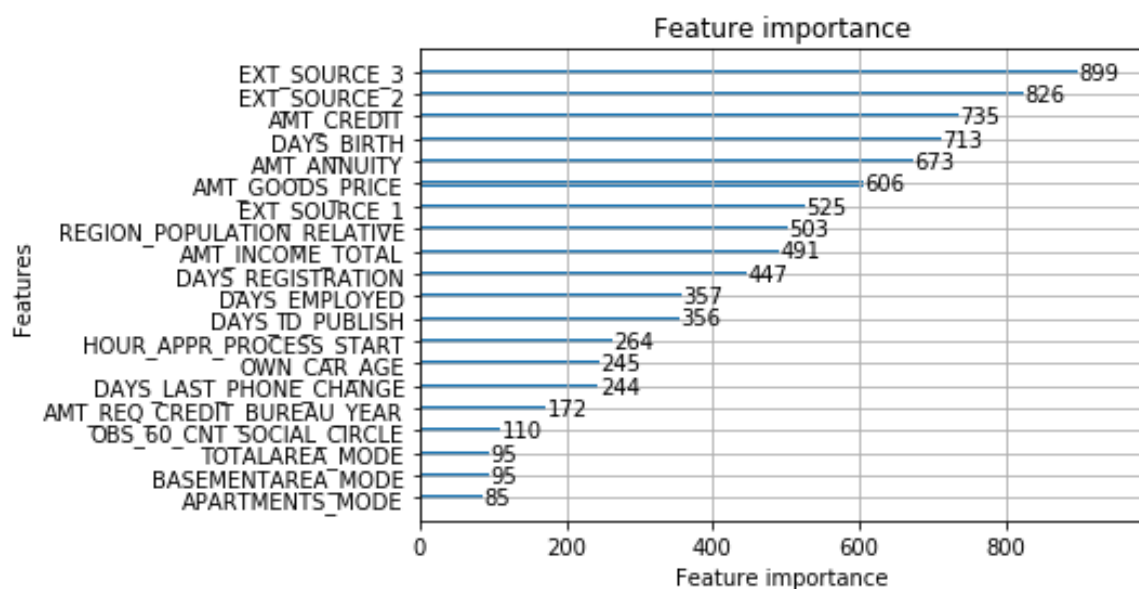
significant enough to have solved the problem posed in the project. Questions to ask yourself when writing this section:

V. Conclusion

(approx. 1-2 pages)

Free-Form Visualization

After model training, it is good to know the top contributing features that help to predict whether an applicant is likely to default on home credit repayment. These features importance have the highest weight to the final prediction on each data input. The feature importance figure below details the top 20 important features for the given dataset of predicting home credit default risk



Reflection

The most difficult task faced in implementing this project was data cleansing and preprocessing. The dataset has many missing values which I was not sure whether to remove or mutate record with missing feature values. I started off with removing as this is the most convenient approach however, I ended up having a poor model performance. Another challenge was dealing with the imbalanced dataset. I noticed that all models were predicting applicants to repay the loan since over 90% of the dataset is negative class. I tried to balance the dataset myself which turned out to be a greater challenge which led me to scale the positive class instead that worked out successfully.

Improvement

Although the final model's evaluation metrics are acceptable, further improvements can be done to gain higher metrics score. There is a high potential of improvement in performing feature engineering of the home credit default risk dataset such as credit to income ratio, credit to annuity ratio, income per child ratio ...etc. Feature engineering also include reduction of unnecessary features that do not add a value in constructing a final classification model should be removed to minimize model complexity.

Another potential improvement can be made by removing outliers from the dataset. Although the aim is to encounter any applicants whom at default risk, the model can greatly benefit from removing these outliers as they may affect the prediction of such risks. Finally, although scaling positive class in an imbalanced dataset helps to avoid naïve prediction, adding more positive class data in the dataset can greatly help the model to predict default risks.