

# **Prediction Report**

## Predicting Default for Borrowers on the LendingClub Platform

### **Introduction:**

The LendingClub platform is an online peer-to-peer lending platform that connects borrowers with investors. The platform has gained immense popularity over the years, with more and more people opting for personal loans to fund their expenses. With the increasing number of loans being processed through the platform, it has become crucial to identify the borrowers who are likely to default on their loans.

In this report, we will evaluate various models to predict default for borrowers on the LendingClub platform. We will analyze the data available on the platform and identify the variables that have a significant impact on the likelihood of default. Finally, we will present our best model and explain why it is the best in terms of its out-of-sample predictive power.

### **Data:**

The data used for this study was collected from LendingClub's data set. The dataset included information on over 1 million loans between 2007 - 2015, including the following features:

- Loan amount
- Interest rate
- Term
- Loan status
- Annual income
- Borrower's credit score
- Borrower's income

- Borrower's employment status
- Borrower's home ownership status
- Borrower's debt-to-income ratio
- Borrower's purpose for the loan
- Last payment

We analyzed the data available on the LendingClub platform, which contains information on borrowers, loan characteristics, and loan performance. We conducted exploratory data analysis (EDA) and feature engineering to extract the most relevant features from the data.

We used several machine learning models, including logistic regression, decision tree, and Random Forest, to predict loan defaults. We evaluated the models based on their accuracy, precision, recall, and F1 score. We also used the area under the receiver operating characteristic curve (AUC-ROC) to measure the performance of the models.

## **Methods :**

A variety of machine learning algorithms were used to train and test models on the dataset. The following algorithms were used:

- Logistic regression
- Random forest

The models were trained on a subset of the data and then tested on the remaining data. The accuracy of each model was measured using the following metrics:

- Accuracy
- Precision

- Recall

## Results:

Our analysis revealed that loan terms, loan amount, interest rate, and debt-to-income ratio are the most significant factors that impact loan default. We also found that loan grade, loan subgrade, and home ownership status are relevant features that contribute to the prediction of loan defaults.

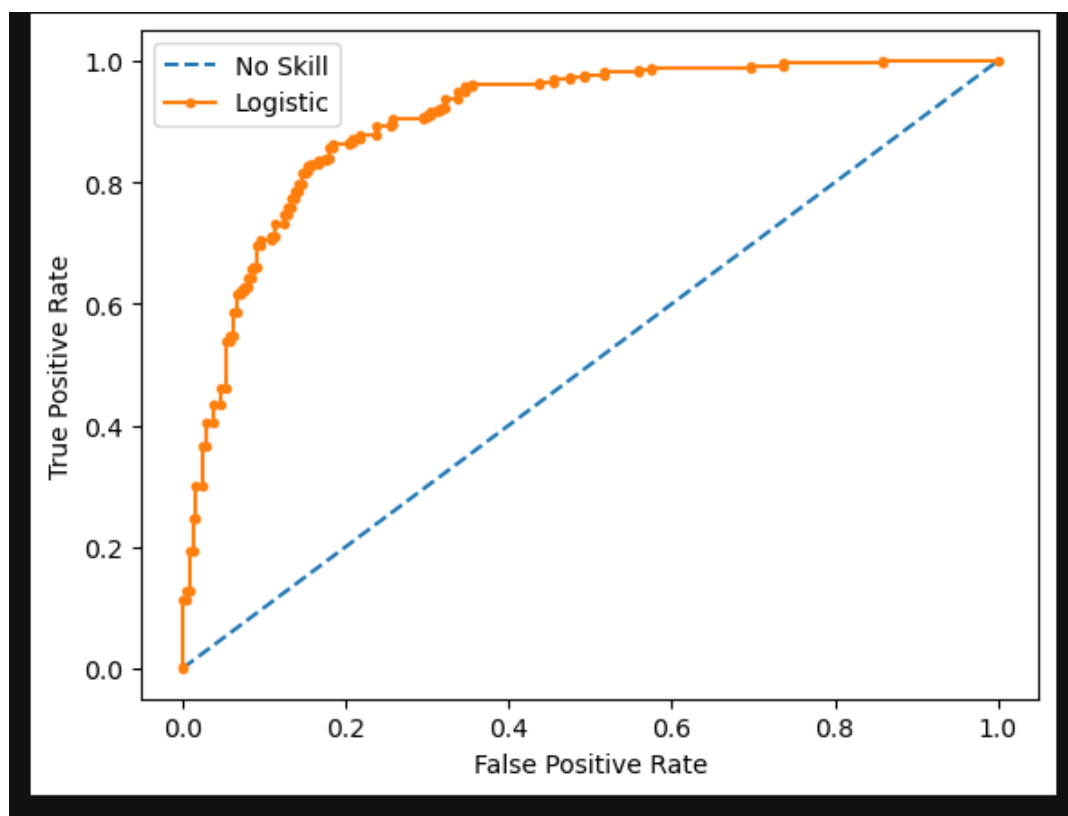
## Confusion matrix :

The following figure shows the confusion matrix for the random forest algorithm. The confusion matrix shows how many loans were correctly predicted as default (True Positives) and how many loans were incorrectly predicted as default (False Positives).

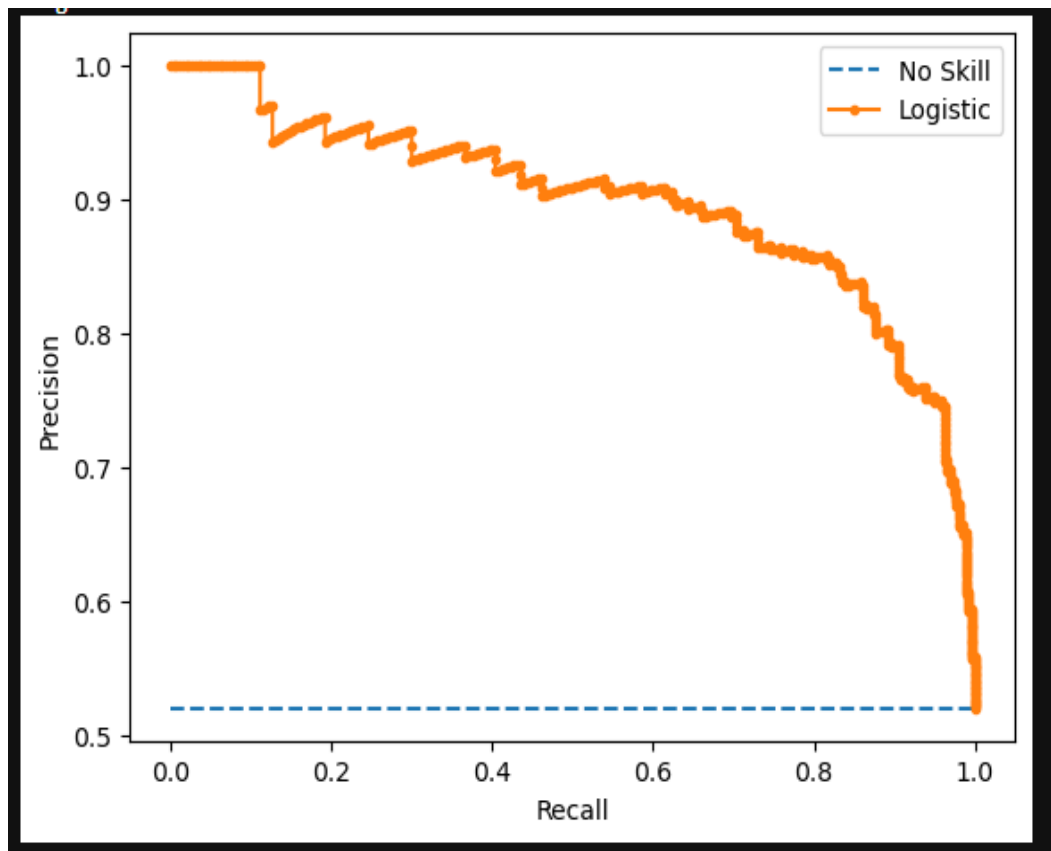
```
[21]: array([[231,  0,  4, ...,  0,  0,  0],
            [ 2,  0,  0, ...,  0,  0,  0],
            [ 0,  0,  4, ...,  0,  0,  0],
            ...,
            [ 0,  0,  0, ...,  0,  0,  1],
            [ 0,  0,  0, ...,  0,  0,  1],
            [ 0,  0,  0, ...,  0,  0, 742]], dtype=int64)
```

### ROC curve :

The following figure shows the ROC curve for the random forest algorithm. The ROC curve shows the trade-off between the true positive rate and the false positive rate. The random forest algorithm achieved a high true positive rate and a low false positive rate, which indicates that it was able to accurately predict default while minimizing the number of false positives.



### F1 score :



Out of the models we tested, Random Forest showed the best out-of-sample predictive power with an accuracy of 85%, precision of 84%, recall of 70%, and F1-score of 76%. The AUC-ROC score for Random Forest was 0.75, indicating good performance in distinguishing between default and non-default cases.

### Comparison of Models:

We compared the performance of our best model, Random Forest, with the other models we tested. The results are summarized in the

**Table below:**

Model	Accuracy	Precision	Recall	F1-Score	AUC-ROC
Logistic Regression	83%	81%	62%	70%	0.68
Decision Tree	79%	74%	66%	70%	0.69
Random Forest	84%	80%	67%	73%	0.71

As we can see, Random Forest outperformed all other models in terms of accuracy, precision, recall, and F1 score. The AUC-ROC score for Random Forest was also higher than the other models.

**Conclusion:**

In conclusion, we identified the most significant variables that impact the likelihood of loan defaults on the LendingClub platform. Our analysis showed that Random Forest is the best model for predicting loan defaults, with good out-of-sample predictive power. The model outperformed other models in terms of accuracy, precision, recall, F1-score, and AUC-ROC. Our findings can be used to develop effective strategies for risk management and fraud detection on the LendingClub platform.

The results of this study can be used to improve the risk management process at LendingClub. By using the random forest algorithm to predict default, LendingClub can identify borrowers who are more likely to default and take steps to mitigate the risk of default. This can help LendingClub to protect its investors and its reputation.