# Credit Risk Predictive Modelling

## ID/X Partners – Rakamin Academy
## Virtual Internship Experience
## Final Task

**Supported by:**
**Rakamin Academy**
Career Acceleration School
www.rakamin.com

**Created by:**
**Kenneth Wahyudi, S.Si.**

kennethwahyudi48@gmail.com
https://www.linkedin.com/in/kenneth-wahyudi-80b886209/

I am an aspiring data scientist with interest in machine learning and paddling through the data lake. I have experiences as a research intern in BATAN, and several data science and machine learning projects under my belt. Specifically intermediate to advanced knowledge in Exploratory Data Analysis, Data Pre-processing, Supervised & Unsupervised Learning, as well as Data Visualization and Storytelling. And I am always looking forward to exploring new and uncharted territories that might allow myself to grow and improve.
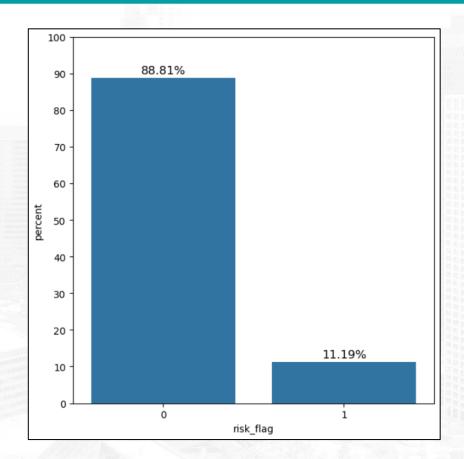
Rakamin Academy

Welcome to a captivating journey into the world of modern finance. In the next few moments, we will explore the art and science behind developing a credit risk predictive model, offering insights to enhance your financial acumen. Fasten your seatbelts as we delve into the complexities of credit risk modeling, where numbers meet narratives, and the future of finance comes into sharper focus. The dataset contains 466285 rows, and customer's loan information from 2007-2014. The objective is to create a predictive model that would solve the problem of defaulting customers.

All results are from outputs of codes written in Python and its packages. JupyterLab was the GUI used in this project.

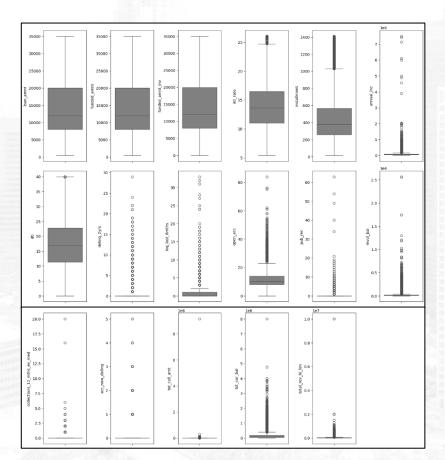The dataset used in this project can be found in the following link.

The full code in .ipynb file format can be found in the following link. And in .py file format can be found in the following link.
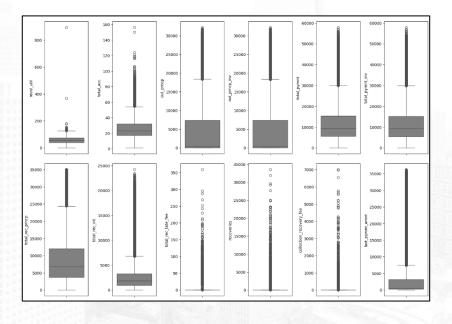
As we can see from the countplot to the left, the target classes are imbalanced. With at risk customers being only about 11.2% of the total dataset, and non-risk customers being an overwhelming 88.8% of the dataset. This means that we need to handle these imbalances later on before modelling.
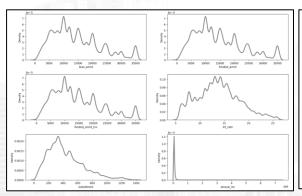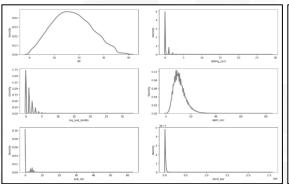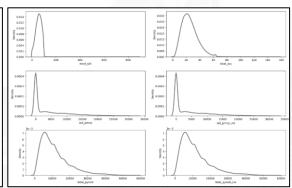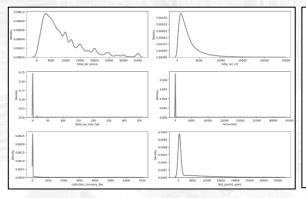
As we can see from the boxplots above, there are a lot of columns that have a lot of outliers within them. We might need to use tree-based algorithms later on since they are robust towards outliers, unlike linear models.
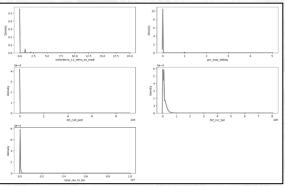
# Exploratory Data Analysis



From the kdeplots above, we can see that the distributions of most of the columns are very skewed due to outliers. Most are skewed normal distributions, although a few are multi-modal.

# Exploratory Data Analysis



From the countplots above, it is clear that there is no column that has balanced classes within them.

That being said, the imbalances are still somewhat manageable and not too severe. There a few columns that has way too many unique values in them, and this needs to be handled in the preprocessing later on.

As we can see from the correlation heatmap on the left, there are a few features that have quite strong correlation with the target. And there are quite a lot of features that have little to no correlation to the target. These features must be considered to be dropped. There are also a couple of features that correlate to each other very well. One of those features also must be dropped later on.

From the correlation heatmap on the left, we can see that there are few features that correlate to the target. But quite a lot of features that have little to no correlation to the target. Again, these features must be considered to be dropped. There are also a couple of features that correlate to each other very well. One of those features also must be dropped later on.
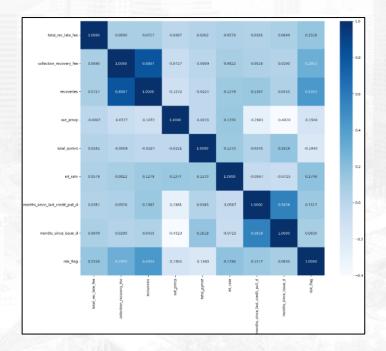
Since there are so many columns, we risk of hitting the curse of dimensionality and expenses of computing. Therefore, we are going to do a feature selection where features that have less than 0.1 correlation to the target will be dropped. And we're also going to drop one of the features that have a correlation of 0.9 or higher with another feature, since one of those features is redundant.



Since we're not going to use linear/distance based models, we don't have to scale the features, and can keep them as is. Now, the preprocessing is done. We have now narrowed the dataset to only 8 features and 1 target variable.

It is time to split the dataset. We are going to split the dataset into 3 parts. One training dataset, one test/validation dataset, and one validation set that will not be involved in the modelling (let's call it the back test set). The back test set will only be used in the final confusion matrix, so therefore the model will not be influenced by it whatsoever. Out of 466285 datapoints, we're going to use 66285 as the back test dataset. Leaving 400000 as the dataset that we will use for the modelling. We will also use random under sampler as a means of handling the dataset imbalance.

# Modelling

Since our target is quite heavily imbalanced, we cannot solely use the accuracy metric alone. We need to account for other metrics such as ROC-AUC. And since it is imperative that we suppress the false negatives (since we don't want to give loans to defaulting customers more than we don't want to lose potential customers), we therefore should also look at the recall score. So as a summary, the metrics that we will look to maximize will be the recall and ROC-AUC metric. The rest are there as complimentary metrics.

```python
from sklearn.tree import DecisionTreeClassifier
dt = DecisionTreeClassifier(random_state=42)
dt.fit(X_train_over, y_train_over)
eval_classification(dt)
```
```
Accuracy (Test Set): 0.87
Accuracy (Train Set): 0.90

Recall (Test Set): 0.86
Recall (Train Set): 1.00

Precision (Test Set): 0.46
Precision (Train Set): 0.53

F1-Score (Test Set): 0.60
F1-Score (Train Set): 0.69

ROC-AUC (Test-proba Set): 0.87
ROC-AUC (Train-proba Set): 0.94
ROC-AUC (Cross-Val Test):  0.869236890659138
ROC-AUC (Cross-Val Train): 0.9447592624626346
```

```python
from sklearn.ensemble import RandomForestClassifier
rf = RandomForestClassifier(n_jobs=-1, random_state=42)
rf.fit(X_train_over, y_train_over)
eval_classification(rf)
```
```
Accuracy (Test Set): 0.92
Accuracy (Train Set): 0.94

Recall (Test Set): 0.86
Recall (Train Set): 1.00

Precision (Test Set): 0.60
Precision (Train Set): 0.66

F1-Score (Test Set): 0.70
F1-Score (Train Set): 0.80

ROC-AUC (Test-proba Set): 0.96
ROC-AUC (Train-proba Set): 0.99
ROC-AUC (Cross-Val Test):  0.9597836868150725
ROC-AUC (Cross-Val Train): 0.9927604510009539
```

```python
from xgboost import XGBClassifier
xgb = XGBClassifier(n_jobs=-1)
xgb.fit(X_train_over, y_train_over)
eval_classification(xgb)
```
```
Accuracy (Test Set): 0.92
Accuracy (Train Set): 0.93

Recall (Test Set): 0.85
Recall (Train Set): 0.87

Precision (Test Set): 0.62
Precision (Train Set): 0.62

F1-Score (Test Set): 0.72
F1-Score (Train Set): 0.73

ROC-AUC (Test-proba Set): 0.96
ROC-AUC (Train-proba Set): 0.97
ROC-AUC (Cross-Val Test):  0.9590515596344485
ROC-AUC (Cross-Val Train): 0.9683867908419227
```

From the results above, we can see that the best model is XGBoost, since it is the only model that does not overfit. We can now move on to the hyperparameter tuning.
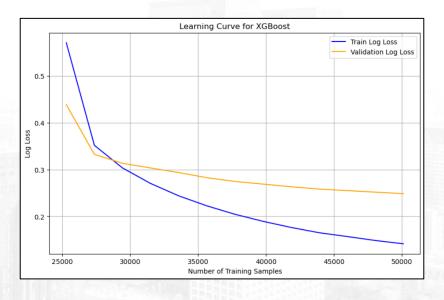
```
from xgboost import XGBClassifier
xgb = XGBClassifier(subsample=0.6666666666666666, n_estimators=357, min_child_weight=13, max_depth=20, eta=0.08102040816326529, scale_pos_weight=2, n_jobs=-1)
xgb.fit(X_train_over, y_train_over)
eval_classification(xgb)

Accuracy (Test Set): 0.90
Accuracy (Train Set): 0.91

Recall (Test Set): 0.91
Recall (Train Set): 0.98

Precision (Test Set): 0.52
Precision (Train Set): 0.55

F1-Score (Test Set): 0.66
F1-Score (Train Set): 0.70

ROC-AUC (Test-proba Set): 0.97
ROC-AUC (Train-proba Set): 0.99
ROC-AUC (Cross-Val Test):  0.9675382745740625
ROC-AUC (Cross-Val Train): 0.9857633600798033
```

As we can see above, the recall has improved significantly from 0.85 to 0.91 in the test set, whilst also maintaining a good fit to the model, and not sacrificing too much on accuracy.
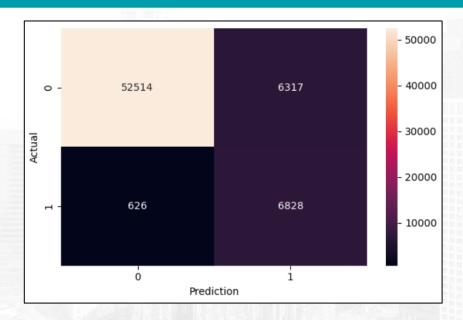
We can also check the learning curve to be sure that the model does not overfit.

From the learning curve above, we can see that the train log loss does not differ by a lot compared to the validation log loss. This means that the model is a good fit. We can now proceed to the confusion matrix.
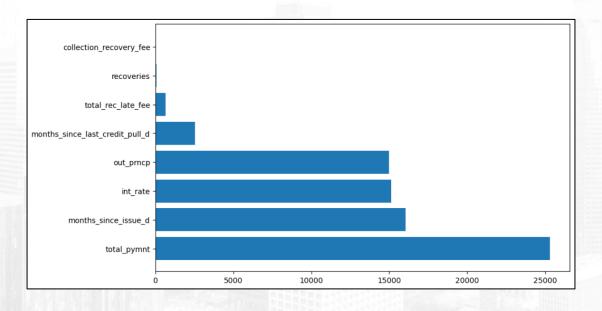
For the confusion matrix, we are going to use the back test dataset, just like we mentioned earlier. The back test dataset is a dataset that has not been touched in the modelling process, crucial for a fair and transparant review of the model's performance. From the confusion matrix above, we can calculate the accuracy, recall, and precision. And they are as follows :

1. **Accuracy** : (True Positive + True Negative) / Total. (6828 + 52514) / 66285 = **0.895**
2. **Recall** : True Positive / (True Positive + False Negative). 6828 / (6828 + 626) = **0.916**
3. **Precision** : True Positive / (True Positive + False Positive). 6828 / (6828 + 6317) = **0.519**

Let us now look at the feature importances of the model.

From the feature importances above, we can see that the 'total_pymnt' feature is the most impactful, with the recovery related features not bringing as much impact compared to the others.

Before we did the modelling, the default rate in the dataset was 11.19%. That means that out of all of the customers whose loan application were accepted, 11.19% of them eventually defaulted on their loan. However, with the predictive model, less than 1% (0.94%) defaulted on their loans. For further simulation, we need to make assumptions. Let's say, that on average, every loan is 1,000 dollars in amount, and the interest that we are going to get from that is 100 dollars. And let's say, that on average, when a customer defaults, we lose 50% or half of the amount that we loaned to them & its interest as well. Now that we get that out of the way, it's time to simulate.

Before the modelling, there are 466285 customers, of which 52186 defaulted. That means that out of all the loan that we gave (466285 x 1000) = **466,285,000 dollars**, we receive back (414099 x 1100) + (52186 x 550) = **484,211,200 dollars** which is a **3.8% profit**.

Now after the modelling, looking at the confusion matrix in the previous section, out of 66285 applications, we only gave out loans to 53140 customers, of which 626 defaulted. That means that out of all the loan that we gave (53140 x 1000) = **53,140,000 dollars**, we receive back (52514 x 1100) + (626 x 550) = **58,109,700 dollars** which is a **9.35% profit**.

We also need to address the fact that there are opportunity losses, since we are maximizing the recall score. If we do the first scenario (giving loans to all customers that applied) on the second back test dataset, than out of all the loan that we gave (66285 x 1000) = **66,285,000 dollars**, we receive back ((52514 + 6317) x 1100) + ((626 + 6828) x 550) = **68,813,800 dollars**. That's a profit ratio of **3.8%**, and a total profit of **2,528,800 dollars**, whilst the second scenario previously had a total profit of **4,969,700 dollars**.

Therefore in conclusion, there is a **146.05%** increase in profit margins after implementing the predictive model. And despite the opportunity loss, we still have a total profit that is higher, **96.52%** higher to be exact.