



Predicting Credit Card Default

BT2101 – Decision making tools and method - Group Project

Contents

Chapter 1:	Project Introduction	3
Chapter 2:	Literature review about Machine Learning models group by Tribe	3
Part 1:	Tribes machine learning:	3
Part 2:	Tribe #1: Symbolists	4
Part 3:	Tribe #2: Bayesians	5
Part 4:	Tribe #3: Analogizers.....	6
Part 5:	Tribe #4: Connectionists	7
Chapter 3:	Exploratory data analysis.....	8
Part 1:	Describing data	8
Part 2:	Distribution plot for each column according to “default” values.....	8
Part 3:	Correlation analysis.....	10
Chapter 4:	Data pre-processing.....	11
Chapter 5:	Feature selection.....	12
Chapter 6:	Splitting data and Imbalanced Class treatment.....	13
Part 1:	Data snooping prevention:.....	13
Part 2:	Imbalanced class treatment:	13
Chapter 7:	Model Selection and evaluation	14
Part 1:	Models and datasets	14
Part 2:	Model evaluation results	15
Part 3:	Experiential lessons in this project:.....	19
Chapter 8:	Future plan for Model Improvement.....	20
Part 1:	Overall picture	20
Part 2:	Improvement on existing models	20
Part 3:	Potential models have not been implemented:.....	20
Chapter 9:	References.....	20

Chapter 1: Project Introduction

BACKGROUND:

Credit card is considered as a flexible tool by which you can use bank's money for a short period of time. If you accept a credit card, you agree to pay your bills by the due date listed on your credit card statement. Otherwise, the credit card will be defaulted. When a customer is not able to pay back the loan by the due date and the bank is totally certain that they are not able to collect the payment, it will usually try to sell the loan. After that, if the bank recognizes that they are not able to collect the payment, they will write it off. This is called a charge-off. This results in significant financial losses to the bank on top of the damaged credit rating of the customer and thus it is an important problem to be tackled. [1]

Predicting accurately which customers are most probable to default represents significant business opportunity for all banks. Bank cards are widely used over the world, which emphasizes the impact of risk prediction to both consumers and banks. Hence, risk prediction is essential for predicting business performance or individual customers' credit risk and to reduce the damage and uncertainty. [1]

PROJECT DESIGN:

Motivation: The fundamental objective is to implement a model to predict default to help bank identify and act on customers with high probability of defaulting to improve their bottom line.

Even though plenty of solutions to the default prediction have been previously done, the scope of this project extends beyond that. Not only we aim to create an effective model to predict default, discussion between different models are emphasized and become one of the top goals in our project.

Dataset: The dataset is obtained from UCI repository.

It contains payment information of 30,000 credit card holders obtained from a bank in Taiwan. Each data sample is described by 23 feature attributes.

In this dataset, we aim to perform binary classification for target feature to be predicted is binary valued 0 (=not default) or 1 (=default).

Chapter 2: Literature review about Machine Learning models group by Tribe

Part 1: Tribes machine learning:

Models in Machine learning can be categorized into different tribes

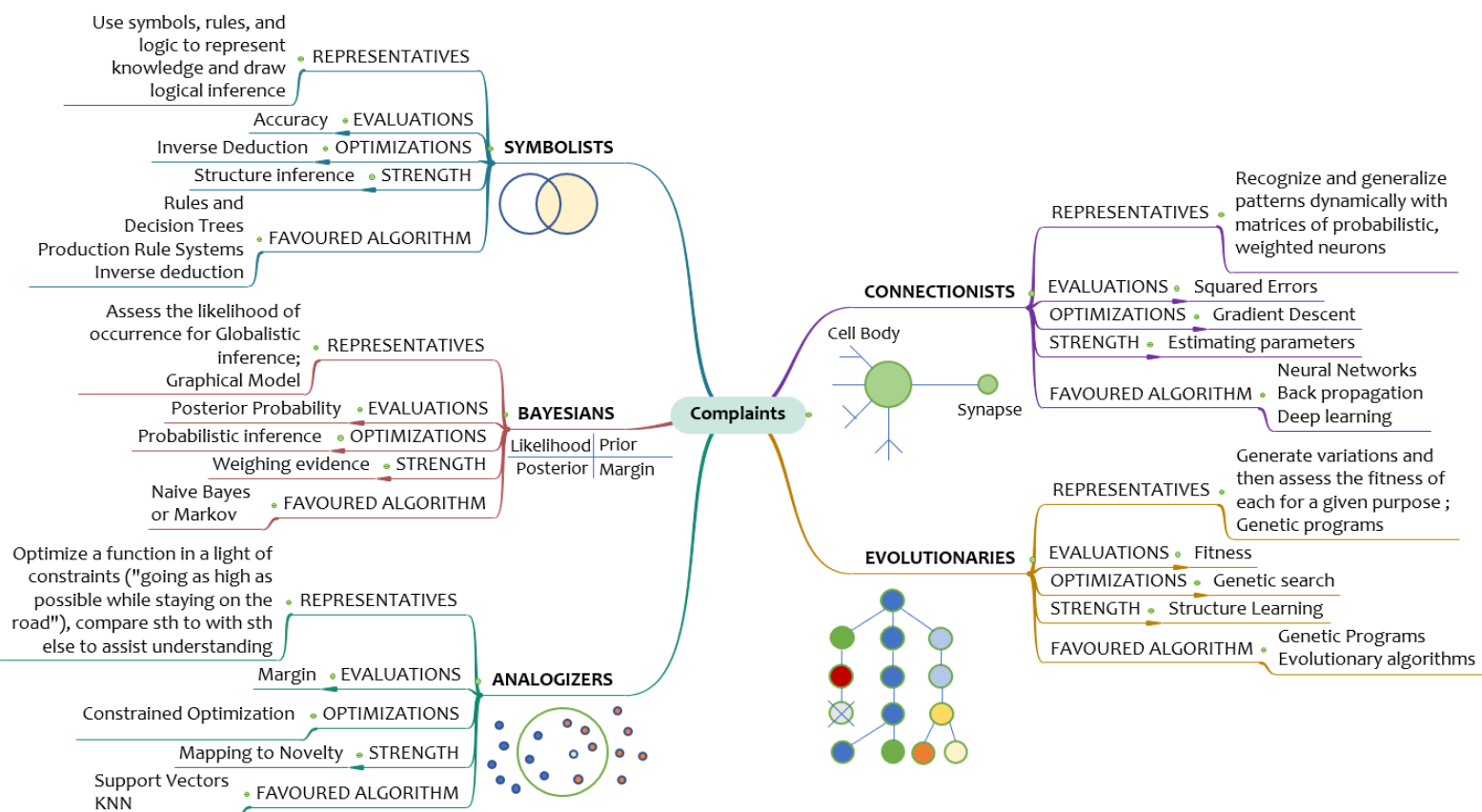


Figure 1: Five different tribes in Machine Learning

Part 2: Tribe #1: Symbolists

DECISION TREE (DT):

A tree can be “learned” by splitting training dataset into subsets based on an features value test.

—Each internal node represents a “test” on a feature resulting on split of current sample.

—At each step algorithm selects feature and a cut-off value that maximises a given metric.

Different metrics exist for regression tree (target is continuous) or classification tree (the target is qualitative).

—This process is repeated on each derived subset in a recursive manner called recursive partitioning. Recursion is completed when subset at a node has all same value of target variable, or when splitting no longer adds value to predictions. This general principle is implemented by many recursive partitioning tree algorithms.

—Decision trees are simple to understand and interpret however they tend to overfit data.

—However, decision trees tend to overfit training set. Leo Breiman propose random forest to deal with this issue.

Decision tree learning algorithm:

Aim: find a small tree minimize J_{train}

Idea: (recursively) choose most significant attribute as root of tree

function *DTL*(*examples*, *attribute*, *default*) return decision tree

if *examples* is **empty**, **return** *default*

else if *attributes* is **empty**, **return** *MODE*(*examples*) **else**

best ← *CHOOSE.ATTRIBUTE*(*attributes*, *examples*)

tree ← a **new decision tree with root test** *best*

for each value v_i of *best* **do**:

examples_i ← {element of *examples* with *best* = v_i }

subtree ← *DTL*(*examples_i*, *attributes* – *best*, *MODE*(*example*))

add a branch to *tree* with label v_i **and** *subtree*

subtree

return *tree*

Picking an attribute via entropy

(*CHOOSE.ATTRIBUTE* function):

Idea: a good attribute split training set into subsets that are (ideally) “all positive” or “all negative”, hence, purity of subproblems.

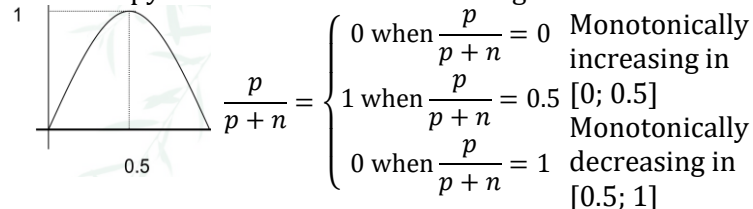
To implement *CHOOSE.ATTRIBUTE* in decision tree learning enumerated feature sets with C possible values, we first need a concept of purity. We will use entropy:

$$H(X) = - \sum_{i=0}^C p \log_2(p)$$

For training set having p positive examples and n negative examples:

$$H\left(\frac{p}{p+n}, \frac{n}{p+n}\right) = -\frac{p}{p+n} \log_2\left(\frac{p}{p+n}\right) - \frac{n}{p+n} \log_2\left(\frac{n}{p+n}\right)$$

For entropy curve: When there are 2 target classes



RANDOM FOREST (RF):

It is a meta estimator that fits several decision tree learners on subsamples of dataset and use averaging to improve predictive accuracy and control over-fitting. A tree-based technique that uses a high number of decision trees built out of randomly selected set of features. Contrary to simple decision tree, it is highly uninterpretable but it's generally good performance. Strong learners composed of multiple trees can be called “forests”. Trees that compose a forest can be chosen to be either shallow (few depths) or deep (lot of depths, if not fully grown). [2]

Shallow trees have less variance but higher bias, hence, better choice for sequential methods. Deep trees, on other side, have low bias but high variance, hence, bagging method that is mainly focused at reducing variance. [2]

The random forest approach is a bagging method where deep trees, fitted on bootstrap samples, are combined to produce an output with lower variance. Random forests also use another trick to make multiple fitted trees a bit less correlated with each other's: when growing each tree, instead of only sampling over observations in dataset to generate a bootstrap sample, we also sample over features and keep only a random subset of them to build tree. [2]

Algorithm: The steps are as follows:

Randomly select a subset of data from training set with replacement.

Generate a decision tree from selected data samples. At each node of tree:

- Randomly select K features without replacement.
- Split node by searching among K features best cut to maximize information gain.

Repeat steps 1 and 2 N times to obtain N trees.

Final prediction is achieved by majority voting.

GRADIENT BOOSTING:

In gradient boosting, ensemble model we try to build is a weighted sum of weak learners

$$s_L(.) = \sum_{l=1}^L c_l \times w_l(.); \text{ } c_l \text{ are coefficient; } w_l \text{ weak learners}$$

Gradient boosting casts problem into a gradient descent one: at each iteration we fit a weak learner to opposite of gradient of current fitting error with respect to current ensemble model. [2]

First, theoretical gradient descent process over ensemble model can be written $s_l(.) = s_{l-1}(.) + c_l \times \nabla_{s_{l-1}} E(s_{l-1})(.)$

where $E(.)$ is fitting error of given model, c_l is a coefficient corresponding to step size and $-\nabla_{s_{l-1}} E(s_{l-1})(.)$ is opposite of gradient of fitting error with respect to ensemble model at step $l - 1$. This (abstract) opposite of gradient is a function that can only be evaluated for observations in training dataset (for which we know inputs and outputs): these evaluations are called **pseudo-residuals** attached to each observation. [2]

Moreover, even if we know for observations values of these pseudo-residuals, we don't want to add to our ensemble model any kind of function: we only want to add a new instance of weak model. So, natural thing to do is to fit a weak learner to pseudo-residuals computed for each observation. [2]

Finally, coefficient c_l is computed following a one-dimensional optimisation process (line-search to obtain best step size c_l). [2]

Illustration of boosting process:

At very beginning of algorithm (first model of sequence), pseudo-residuals are set equal to observation values. Then, we repeat L times (for L models of sequence):

- Fit best possible weak learner to pseudo-residuals (approximate opposite of gradient with respect to current strong learner)

- Compute value of optimal step size that defines by how much we update ensemble model in direction of new weak learner

- Update ensemble model by adding new weak learner multiplied by step size (make a step of gradient descent)

- Compute new pseudo-residuals that indicate, for each observation, in which direction we would like to update next ensemble model predictions

Repeating these steps, we have then build sequentially our L models and aggregate them following a gradient descent approach.

STACKING:

Stacking mainly differ from boosting on two points.

- Stacking uses heterogeneous weak learners (different learning algorithms are combined) whereas bagging and boosting consider mainly homogeneous weak learners.

- Stacking learns to combine base models using a meta-model whereas bagging and boosting combine weak learners following deterministic algorithms. [2]

Idea of stacking is to learn several different weak learners and combine them by training a meta-model to output predictions based on multiple predictions returned by these weak models. So, we need to define two things in order to build our stacking model: L learners we want to fit and meta-model that combines them. [2]

Classification problem example: We can choose as weak learners a KNN classifier, a logistic regression and a SVM, and decide to learn a neural network as meta-model. Then, neural network will take as inputs outputs of our three weak learners and will learn to return final predictions based on it.

Illustration of stacking process: We want to fit a stacking ensemble composed of L weak learners. Then we must follow steps thereafter:

- split training data in two folds

- choose L weak learners and fit them to data of first fold

- for each of L weak learners, make predictions for observations in second fold fit meta-model on second fold, using predictions made by weak learners as inputs

In previous steps, we split dataset in two folds because predictions on data that have been used for training of weak learners are not relevant for training of meta-model.

Thus, an obvious drawback of this split of our dataset in two parts is that we only have half of data to train base models and half of data to train meta-model.

In order to overcome this limitation, we can however follow some kind of “k-fold cross-training” approach (similar to what is done in k-fold cross-validation) such that all observations can be used to train meta-model: for any observation, prediction of weak learners are done with instances of these weak learners trained on $k - 1$ folds that do not contain considered observation.

Part 3: Tribe #2: Bayesians

NAÏVE BAYES CLASSIFICATION:

Naive Bayes models are a group of extremely fast and simple classification algorithms that are often suitable for very high dimensional datasets. Because they are so fast and have so few tuneable parameters, they end up being very useful as a quick and dirty baseline for a classification problem. This section will focus on an intuitive explanation of how naive Bayes classifiers work, followed by a couple examples of them in action on some datasets. [3]

Naive Bayes classifiers are built on Bayesian classification methods. These rely on Bayes’ theorem, which is an equation describing the relationship of conditional probabilities of statistical quantities. In Bayesian classification, we’re interested in finding the probability of a label given some observed features, which we can write as $P(L|features)$.

Bayes’ theorem tells us how to express this in terms of quantities we can compute more directly:

$$P(L|features) = \frac{P(features|L)P(L)}{P(features)}$$

If we are trying to decide between two labels—let’s call them L_1 and L_2 then one way to make this decision is to compute the ratio of the posterior probabilities for each label:

$$\frac{P(L_1|features)}{P(L_1|features)} = \frac{P(features|L_1)P(L_1)}{P(features|L_2)P(L_2)}$$

Gaussian Naïve Bayes: in this classifier, the assumption is data from each label is drawn from a simple Gaussian distribution.

Multinomial Naïve Bayes: classifier where the features are assumed to be generated from a simple multinomial distribution. The multinomial distribution describes the probability of observing counts among a number of categories, and thus multinomial naive Bayes is most appropriate for features that represent counts or count rates.

The idea is precisely the same as before, except that instead of modeling the data distribution with the best-fit Gaussian, we model the data distribution with a best-fit multinomial distribution.

Special cases for Naïve Bayes well performance:

Naïve Bayes tends to perform especially well in one of the following situations:

- When the Naïve assumption actually match the data (very rare in practice)

- For very well-separated categories, when model complexity is less important

- For very high-dimensional data, when model complexity is less important

The last two points seem distinct, but they actually are related: as the dimension of a dataset grows, it is much less likely for any two points to be found close together (after all, they must be close in every single dimension to be close overall). This means that clusters in high dimensions tend to be more separated, on average, than clusters in low dimensions, assuming the new dimensions actually add information. For this reason, simplistic classifiers like naïve Bayes tend to work as well or better than more complicated classifiers as the dimensionality grows: once you have enough data, even a simple model can be very powerful. [3]

Part 4: Tribe #3: Analogizers

SUPPORT VECTOR MACHINE (LINEAR SUPPORT VECTOR CLASSIFIER)

Support vector machines (SVMs) are a particularly powerful and flexible class of supervised algorithms for both classification and regression

In classification problem, we will try to find a hyperplane to separate data samples.

In order to achieve that, we will find a hyperplane whose distances from points nearest to it have maximum values.

On the overall picture, we will find hyperplane that maximize the geometric margin to itself.

We can pose following optimization problem:

$$\max_{\gamma, w, b} \gamma \text{ such that}$$

$$y_i(w^T x_i + b) \geq \gamma; i = 1 \dots m; \|w\| = 1$$

*We want to maximise γ subject to each training example having functional margin at least γ . The $\|w\| = 1$ constraint moreover ensures that functional margin equals to geometric margin, so we guaranteed that all geometric margins are at least γ , hence, solving problem will result in (w, b) with largest possible geometric margin with respect to training set.

*If we could solve optimization problem above, we'd be done. But " $\|w\| = 1$ " constraint is a nasty (nonconvex) 1 and this problem certainly isn't in any format that we can plug into standard optimization software to solve, hence, try transforming problem into a nicer one.:

$$\max_{\hat{\gamma}, w, b} \frac{\hat{\gamma}}{\|w\|} \text{ such that } y_i(w^T x_i + b)$$

*We can add an arbitrary scaling constraint on w and b without changing anything. We will introduce scaling constraint that functional margin of w, b with respect to training set must be 1: $\hat{\gamma} = 1$.

*Since multiplying w and b by some constant results in functional margin being multiplied by that same constant, this is indeed a **scaling constraint** and can be satisfied by rescaling w, b, hence, $\max_{\hat{\gamma}, w, b} \hat{\gamma} / \|w\| = \max 1 / \|w\|$

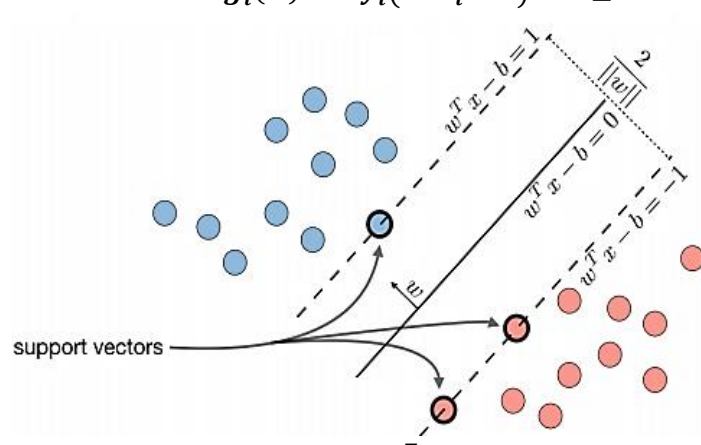
*The **optimal margin classifier h** is such that:

$$h(x) = \text{sign}(w^T x - b)$$

where (w, b) is solution of above optimization problem:

$$\min_{\gamma, w, b} \frac{1}{2} \|w\|^2 \text{ such that } y_i(w^T x_i - b) \geq 1, i = 1 \dots m$$

We can write as: $g_i(w) = -y_i(w^T x_i + b) + 1 \leq 0$



Remarks: line is defined as $(w^T x - b)$

To make algorithm work for **non-linearly separable** datasets as well as be less sensitive to outliers, we reformulate our optimization (ℓ_1 regularization) as follows:

$$\min_{\gamma, w, b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i \text{ such that}$$

$$y_i(w^T x_i + b) \geq 1 - \xi_i \text{ and } \xi_i \geq 0; i = 1 \dots m;$$

Parameter C controls relative weighting between twin goals of making $\|w\|^2$ small and of ensuring that most examples have functional margin ≥ 1

As before we can form Lagrangian:

$$\mathcal{L}(w, b, \xi, \alpha, r) = \frac{1}{2} w^T w + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m \alpha_i [y_i(w^T x_i + b) - 1 + \xi_i] - \sum_{i=1}^m r_i \xi_i$$

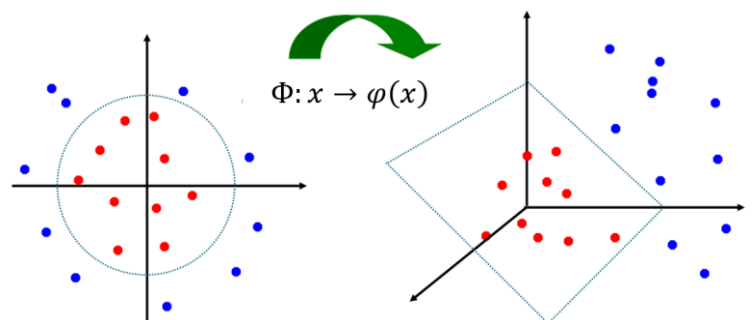
where α_i 's and r_i 's are our Lagrange multipliers (constrained to be ≥ 0). After setting derivatives of dual with respect to w and b to zero as before, substituting them back in and simplifying we obtain

DUAL OPTIMIZATION PROBLEM:

$$\max_{\alpha} W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y_i y_j \alpha_i \alpha_j \langle x_i, x_j \rangle$$

$$\text{such that } 0 \leq \alpha_i \leq C, i = 1 \dots m \text{ and } \sum_{i=1}^m \alpha_i y_i = 0$$

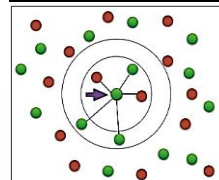
KERNEL TRICKS: General idea: original input space can be mapped to some higher dimensional feature space where training set is separable



Common types of inner product kernel for SVM $K(x^T x_i)$

Type	Kernel	Comment
Linear	$x^T x_i$	None
Polynomial learning machine	$(x^T x_i + 1)^p$	p is specified a priori by user.
Radial basis function network (Gaussian kernel)	$\exp\left(-\frac{1}{\sigma^2} \ x - x_i\ ^2\right)$	the width σ^2 is specified a priori by user
Two-layer perceptron (sigmoid)	$\tanh(\beta_0 x^T x_i + \beta_1)$	only for some values of β_0 and β_1

K - NEAREST NEIGHBOURS (K-NN):



This model will look at 'k' closest labelled data points and then we apply majority vote
Finds records in database have similar numerical values of set of predictor variables

*Measure Euclidean distance between records in training data set. Nearest neighbour to record in training data set is 1 that that has smallest distance from it.

*If $k = 1$, hence, $1 - NN$ rule classifies record in same category as its nearest neighbour.

* $k - NN$ rule finds k-Nearest Neighbours in training data set to each record we want to classify; hence, assigns classification as classification of majority of k nearest neighbours.

*Typically, various values of k are used; then results inspected to determine which is best.

MODEL COMPLEXITY: Small value of K means that noise will have higher influence on result → more complex model → can lead to *overfitting*

Large value of K makes it computationally expensive; defeats basic idea behind KNN (that points that are near might have similar classes) → Smoother decision boundary, hence, less complex model

HOW TO CHOOSE VALUE K?

Higher k, hence, higher bias, hence, Complex surface

Lower k, hence, lower variance, hence, Smoother surface

Selecting value of K in K-nearest neighbour is most critical prob.

⇒ Simple approach to select k is $k = n^{1/2}$.

Part 5: Tribe #4: Connectionists

LOGISTIC REGRESSION

Logistic regression is a statistical model widely used in binary classification problem.

The output of the problem is a probability.

Assumptions made in Logistic regression: [4]

To make **inference** on model parameters. In other words, infer properties about unknown population coefficients β_0 and β_1 from sample $(X_1, Y_1), \dots (X_n, Y_n)$

—Linearity (of relationship between Y and X):

Residual vs. fitted – *Find straight horizontal line*

$$\mathbb{E}[Y|X = x] = \beta_0 + \beta_1 x$$

—Normality of errors = Errors (e; residuals) are normally distributed: Normal Q-Q plot – Look for linear relationship

$$\varepsilon_i \sim \mathcal{N}(0, \sigma^2) \text{ for } i = 1 \dots N$$

—Homoscedasticity = Constant / Equal variance of errors (e) for all values of X = Impact of X on Y is same for all X values: Residual vs. fitted; Scale-location – Look for straight horizontal line

$$\text{Var}[\varepsilon_i] = \sigma^2 \text{ with } \sigma^2 \text{ constant for } i = 1 \dots n$$

—Independence of errors = There is no correlation between errors (e) calculated from regression model – Need additional plot/test

- Panel/time-series data need to check
- Cross-sectional data – data is collected only once, from different individuals/entities
- Panel/time-series data – data is collected multiple times from each individual/entity

$\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$ are independent (uncorrelated)

$$\mathbb{E}[\varepsilon_i \varepsilon_j] = 0 ; i \neq j \rightarrow \text{assumed to be normal}$$

Good-to-1-line summary of model (independence assume) $Y|X = x \sim \mathcal{N}(\beta_0 + \beta_1 x ; \sigma^2)$

Algorithm:

1. Initialise weights at $t = 0$ to $\theta(0)$

2. Do:

Compute gradient for cost function: $\nabla(t) = \nabla J(\theta(t))$

$$= -\frac{1}{m} X^T (h_{\theta}(X) - y) = -\frac{1}{m} \sum_{j=1}^m \frac{y^{(j)} x^{(j)}}{1 + \exp(y^{(j)} \theta(t)^T x^{(j)})}$$

$\nabla J(\theta(t))$ is gradient and should be n – dimensional vector, where i – th element of $\nabla J(\theta(t))$ is partial derivative of loss function $J(\theta)$ with respect to θ_i

// Move in direction $v(t) = -\nabla(t)$

Update weights $\theta(t+1) = \theta(t) - \alpha \nabla_{\theta(t)} J(\theta(t))$ where $\theta(t)$ represents weights in t – th iteration and α denotes learning rate

Continue next iteration, until it is time to stop

3. Return final weights θ^*

STOCHASTIC GRADIENT DESCENT (SGD)

Gradient Descent models are used to reduce the loss, or residuals. However, it takes very slow when dataset is in huge size, the amount of computations is impressively high.

Say we have 10,000 data points and 10 features. The sum of squared residuals consists of as many terms as there are data points, so 10000 terms in our case. We need to compute the derivative of this function with respect to each of the features, so in effect we will be doing $10000 * 10 = 100,000$ computations per iteration. It is common to take 1000 iterations; in effect we have $100,000 * 1000 = 100,000,000$ computations to complete the algorithm. That is pretty much an overhead and hence gradient descent is slow on huge data.

“Stochastic” gradient descent model comes to rescue.

It will randomly picks one data point from the whole data set at each iteration to reduce the computations enormously. [5]

Chapter 3: Exploratory data analysis

Part 1: Describing data

The data consists of 30,000 customers and 24 columns of variables. Each sample corresponds to a single customer. The columns consist of the following variables:

- **DEFAULT** (Yes or no) as a binary response variable – Categorical variable
- **BALANCE LIMIT** – Continuous variable
- **SEX** encoded as 1 or 2 – Categorical variable
- **EDUCATION** encoded as 0 → 6
- **MARRIAGE** shows Marital status encoded as 0 → 3 – Categorical variable
- **AGE** (Years) – Discrete variable
- **PAY(PAY_0, PAY_1, ...)** shows Payment status – Categorical variable
- **BILL_AMT(BILL_AMT1, BILL_AMT2, ...)** shows Bill amount – Continuous variable
- **PAY_AMT(PAY_AMT1, PAY_AMT2, ...)** shows Payment amount

We renamed last column to “default” in short

The variables Balance limit, Age, Sex, Education, Marital status, Employer, and Location are defined as demographic variables, since they describe a demography of customers and are available for new customers.

	count	mean	std	min	25%	50%	75%	max
LIMIT_BAL	30000.0	167484.322667	129747.661567	10000.0	50000.00	140000.0	240000.00	1000000.0
SEX	30000.0	1.603733	0.489129	1.0	1.00	2.0	2.00	2.0
EDUCATION	30000.0	1.853133	0.790349	0.0	1.00	2.0	2.00	6.0
MARRIAGE	30000.0	1.551867	0.521970	0.0	1.00	2.0	2.00	3.0
AGE	30000.0	35.485500	9.217904	21.0	28.00	34.0	41.00	79.0
PAY_0	30000.0	-0.016700	1.123802	-2.0	-1.00	0.0	0.00	8.0
PAY_2	30000.0	-0.133767	1.197186	-2.0	-1.00	0.0	0.00	8.0
PAY_3	30000.0	-0.166200	1.196868	-2.0	-1.00	0.0	0.00	8.0
PAY_4	30000.0	-0.220667	1.169139	-2.0	-1.00	0.0	0.00	8.0
PAY_5	30000.0	-0.266200	1.133187	-2.0	-1.00	0.0	0.00	8.0
PAY_6	30000.0	-0.291100	1.149988	-2.0	-1.00	0.0	0.00	8.0
BILL_AMT1	30000.0	51223.330900	73635.860576	-165580.0	3558.75	22381.5	67091.00	964511.0
BILL_AMT2	30000.0	49179.075167	71173.768783	-69777.0	2984.75	21200.0	64006.25	983931.0
BILL_AMT3	30000.0	47013.154800	69349.387427	-157264.0	2666.25	20088.5	60164.75	1664089.0
BILL_AMT4	30000.0	43262.948967	64332.856134	-170000.0	2326.75	19052.0	54506.00	891586.0
BILL_AMT5	30000.0	40311.400967	60797.155770	-81334.0	1763.00	18104.5	50190.50	927171.0
BILL_AMT6	30000.0	38871.760400	59554.107537	-339603.0	1256.00	17071.0	49198.25	961664.0
PAY_AMT1	30000.0	5663.580500	16563.280354	0.0	1000.00	2100.0	5006.00	873552.0
PAY_AMT2	30000.0	5921.163500	23040.870402	0.0	833.00	2009.0	5000.00	1684259.0
PAY_AMT3	30000.0	5225.681500	17606.961470	0.0	390.00	1800.0	4505.00	896040.0
PAY_AMT4	30000.0	4826.076867	15666.159744	0.0	296.00	1500.0	4013.25	621000.0
PAY_AMT5	30000.0	4799.387633	15278.305679	0.0	252.50	1500.0	4031.50	426529.0
PAY_AMT6	30000.0	5215.502567	17777.465775	0.0	117.75	1500.0	4000.00	528666.0
default	30000.0	0.221200	0.415062	0.0	0.00	0.0	0.00	1.0

Figure 2: Description of Dataset

Part 2: Distribution plot for each column according to “default” values

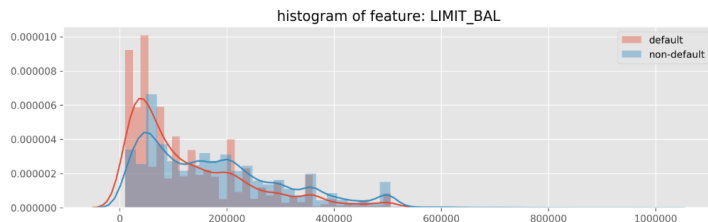


Figure 3: Distribution plot of LIMIT_BAL

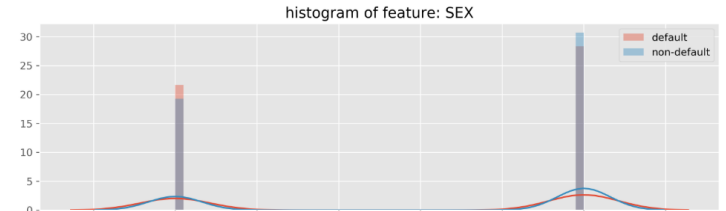


Figure 4: Distribution plot of SEX

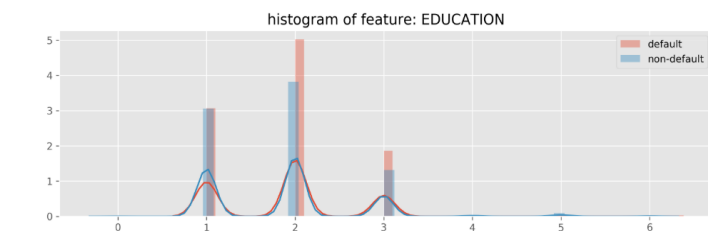


Figure 5: Distribution plot of EDUCATION

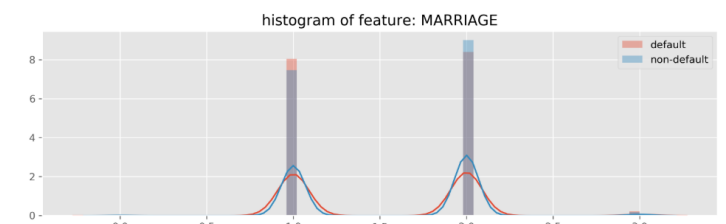


Figure 6: Distribution plot of MARRIAGE

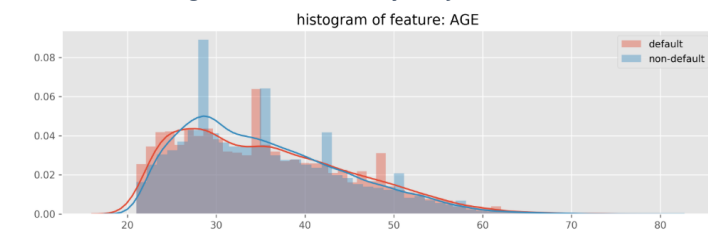


Figure 7: Distribution plot of AGE

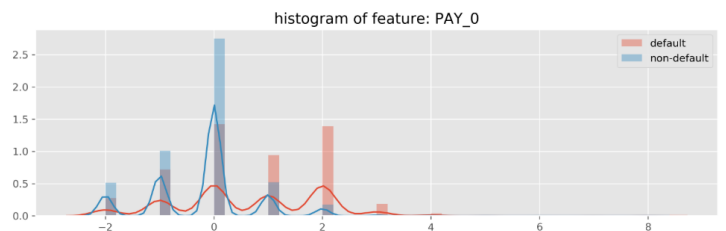


Figure 8: Distribution plot of PAY_0

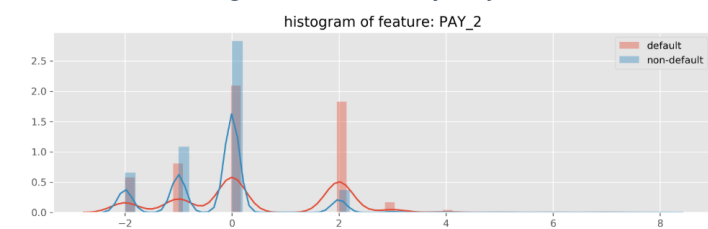


Figure 9: Distribution plot of PAY_2

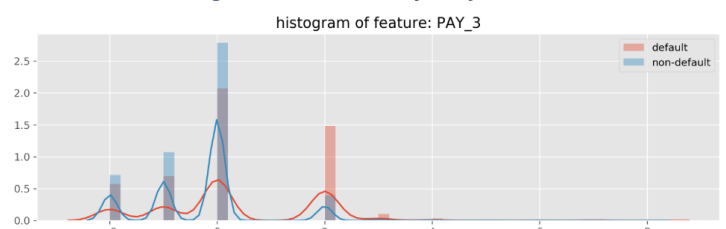


Figure 10: Distribution plot of PAY_5

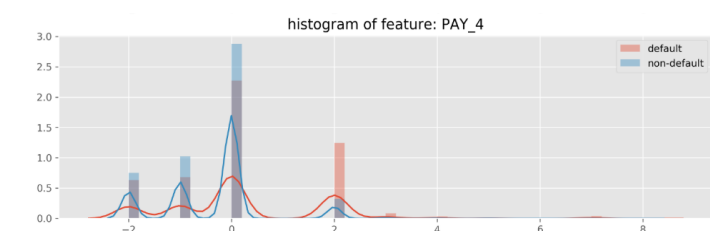


Figure 11: Distribution plot of PAY_4

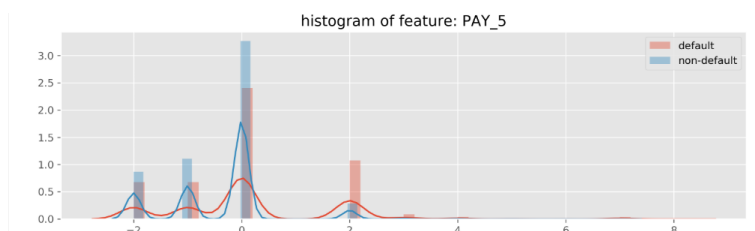


Figure 12: Distribution plot of PAY_5

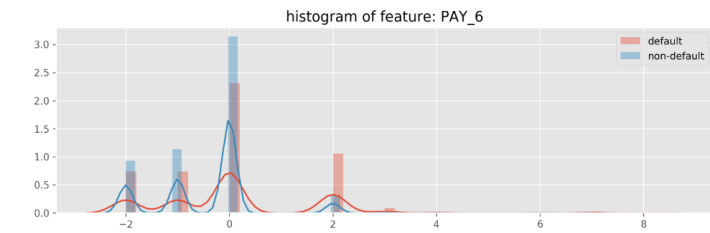


Figure 13: Distribution plot of PAY_6

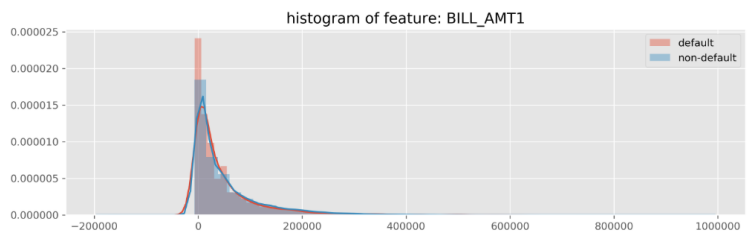


Figure 14: Distribution plot of BILL_AMT1

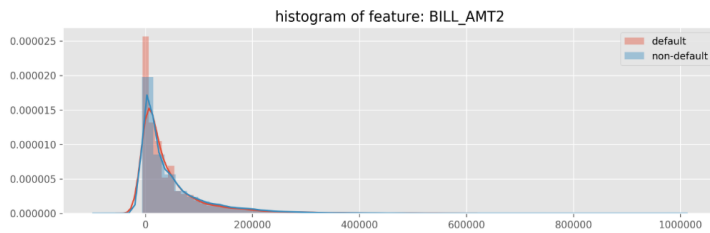


Figure 15: Distribution plot of BILL_AMT2

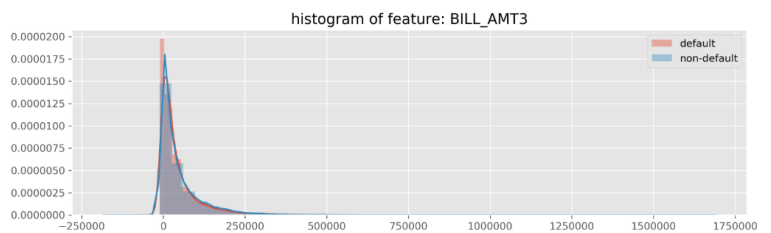


Figure 16: Distribution plot of BILL_AMT3

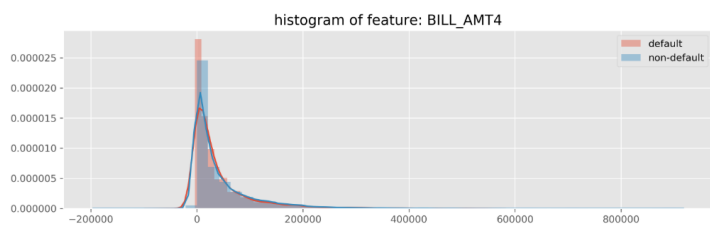


Figure 17: Distribution plot of BILL_AMT4

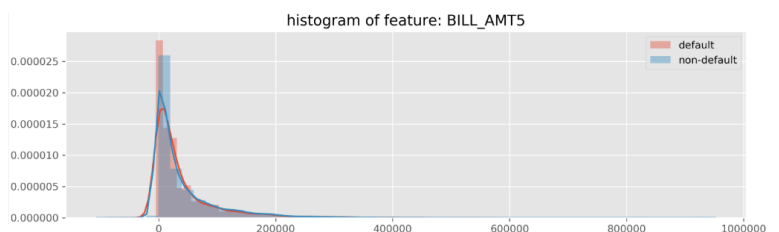


Figure 18: Distribution plot of PAY_AMT1

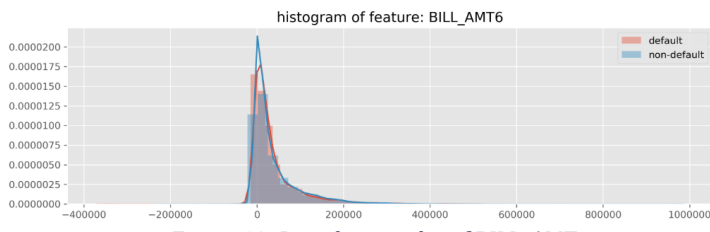


Figure 19: Distribution plot of BILL_AMT6

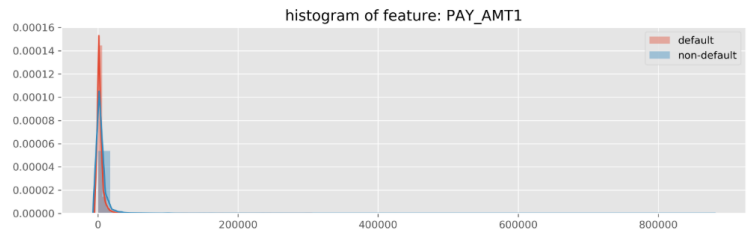


Figure 20: Distribution plot of PAY_AMT1

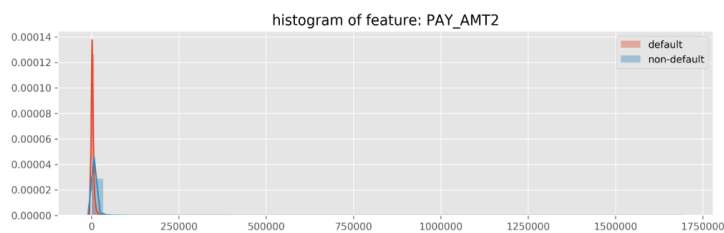


Figure 21: Distribution plot of PAY_AMT2

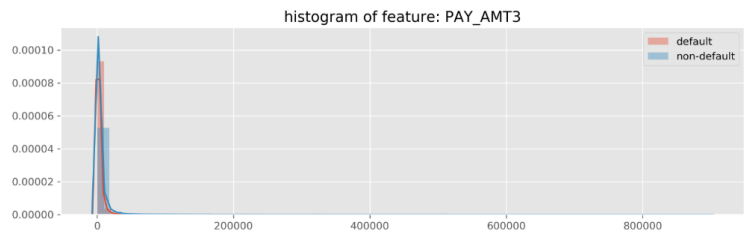


Figure 22: Distribution plot of PAY_AMT3

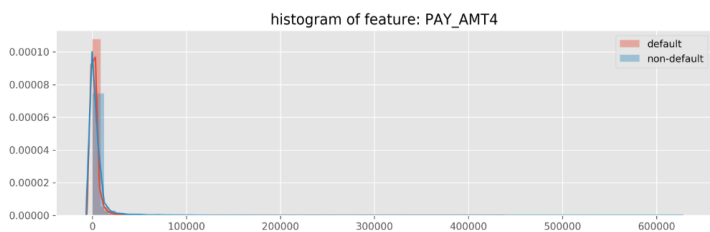


Figure 23: Distribution plot of PAY_AMT4

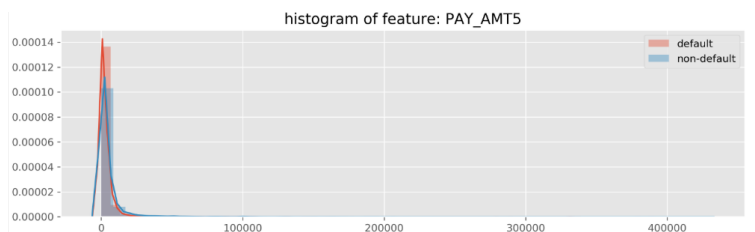


Figure 24: Distribution plot of PAY_AMT5

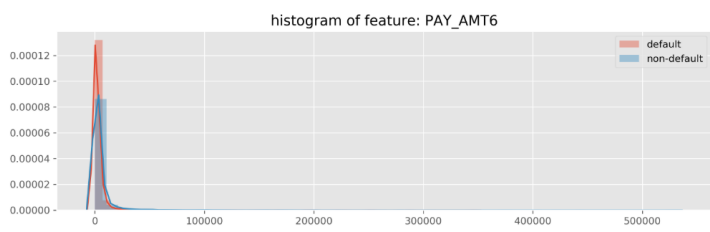


Figure 25: Distribution plot of PAY_AMT6

Part 3: Correlation analysis

F1	AGE	Bill Amt1	Bill Amt2	Bill Amt3	Bill Amt4	Bill Amt5	Bill Amt6	Education	Limit Bal	Marriage	Pay 0	Pay 2	Pay 3	Pay 4	Pay 5	Pay 6	Pay Amt1	Pay Amt2	Pay Amt3	Pay Amt4	Pay Amt5	Pay Amt6	SEX	Default
AGE	1.000	0.056	0.054	0.054	0.051	0.049	0.048	0.175	0.145	-0.414	-0.039	-0.050	-0.053	-0.050	-0.054	-0.049	0.026	0.022	0.029	0.021	0.023	0.019	-0.091	0.014
BILL_AMT1	0.056	1.000	0.951	0.892	0.860	0.830	0.803	0.024	0.285	-0.023	0.187	0.235	0.208	0.203	0.207	0.207	0.140	0.099	0.157	0.158	0.167	0.179	-0.034	-0.020
BILL_AMT2	0.054	0.951	1.000	0.928	0.892	0.860	0.832	0.019	0.278	-0.022	0.190	0.235	0.237	0.226	0.227	0.227	0.280	0.101	0.151	0.147	0.158	0.174	-0.031	-0.014
BILL_AMT3	0.054	0.892	0.928	1.000	0.924	0.884	0.853	0.013	0.283	-0.025	0.180	0.224	0.227	0.245	0.243	0.241	0.244	0.317	0.130	0.143	0.180	0.182	-0.025	-0.014
BILL_AMT4	0.051	0.860	0.892	0.924	1.000	0.940	0.901	0.000	0.294	-0.023	0.179	0.222	0.227	0.246	0.272	0.266	0.233	0.208	0.300	0.130	0.160	0.178	-0.022	-0.010
BILL_AMT5	0.049	0.830	0.860	0.884	0.940	1.000	0.946	-0.008	0.296	-0.025	0.181	0.221	0.225	0.243	0.270	0.291	0.217	0.181	0.252	0.293	0.142	0.164	-0.017	-0.007
BILL_AMT6	0.048	0.803	0.832	0.853	0.901	0.946	1.000	-0.009	0.290	-0.021	0.177	0.219	0.222	0.239	0.263	0.285	0.200	0.173	0.234	0.250	0.308	0.115	-0.017	-0.005
EDUCATION	0.175	0.024	0.019	0.013	0.000	-0.008	-0.009	1.000	-0.219	-0.143	0.105	0.122	0.114	0.109	0.098	0.082	-0.037	-0.030	-0.040	-0.038	-0.040	-0.037	0.014	0.028
LIMIT_BAL	0.145	0.285	0.278	0.283	0.294	0.296	0.290	-0.219	1.000	-0.108	-0.271	-0.296	-0.286	-0.267	-0.249	-0.235	0.195	0.178	0.210	0.203	0.217	0.220	0.025	-0.154
MARRIAGE	-0.414	-0.023	-0.022	-0.025	-0.023	-0.025	-0.021	-0.143	-0.108	1.000	0.020	0.024	0.033	0.033	0.036	0.034	-0.006	-0.008	-0.004	-0.013	-0.001	-0.007	-0.031	-0.024
PAY_0	-0.039	0.187	0.190	0.180	0.179	0.181	0.177	0.105	-0.271	0.020	1.000	0.672	0.574	0.539	0.509	0.475	-0.079	-0.070	-0.071	-0.064	-0.058	-0.059	-0.058	0.325
PAY_2	-0.050	0.235	0.235	0.224	0.222	0.221	0.219	0.122	-0.296	0.024	0.672	1.000	0.767	0.662	0.623	0.576	-0.081	-0.059	-0.056	-0.047	-0.037	-0.037	-0.071	0.264
PAY_3	-0.053	0.208	0.237	0.227	0.227	0.225	0.222	0.114	-0.286	0.033	0.574	0.767	1.000	0.777	0.687	0.633	0.001	-0.067	-0.053	-0.046	-0.036	-0.036	-0.066	0.235
PAY_4	-0.050	0.203	0.226	0.245	0.246	0.243	0.239	0.109	-0.267	0.033	0.539	0.662	0.777	1.000	0.820	0.716	-0.009	-0.002	-0.069	-0.043	-0.034	-0.027	-0.060	0.217
PAY_5	-0.054	0.207	0.227	0.243	0.272	0.270	0.263	0.098	-0.249	0.036	0.509	0.623	0.687	0.820	1.000	0.817	-0.006	-0.003	0.009	-0.058	-0.033	-0.023	-0.055	0.204
PAY_6	-0.049	0.207	0.227	0.241	0.266	0.291	0.285	0.082	-0.235	0.034	0.475	0.576	0.633	0.716	0.817	1.000	-0.001	-0.005	0.006	0.019	-0.046	-0.025	-0.044	0.187
PAY_AMT1	0.026	0.140	0.280	0.244	0.233	0.217	0.200	-0.037	0.195	-0.006	-0.079	-0.081	0.001	-0.009	-0.006	-0.001	1.000	0.286	0.252	0.200	0.148	0.186	0.000	-0.073
PAY_AMT2	0.022	0.099	0.101	0.317	0.208	0.181	0.173	-0.030	0.178	-0.008	-0.070	-0.059	-0.067	-0.002	-0.003	-0.005	0.286	1.000	0.245	0.180	0.181	0.158	-0.001	-0.059
PAY_AMT3	0.029	0.157	0.151	0.130	0.300	0.252	0.234	-0.040	0.210	-0.004	-0.071	-0.056	-0.053	-0.069	0.009	0.006	0.252	0.245	1.000	0.216	0.159	0.163	-0.009	-0.056
PAY_AMT4	0.021	0.158	0.147	0.143	0.130	0.293	0.250	-0.038	0.203	-0.013	-0.064	-0.047	-0.046	-0.043	-0.058	0.019	0.200	0.180	0.216	1.000	0.152	0.158	-0.002	-0.057
PAY_AMT5	0.023	0.167	0.158	0.180	0.160	0.142	0.308	-0.040	0.217	-0.001	-0.058	-0.037	-0.036	-0.034	-0.033	-0.046	0.148	0.181	0.159	0.152	1.000	0.155	-0.002	-0.055
PAY_AMT6	0.019	0.179	0.174	0.182	0.178	0.164	0.115	-0.037	0.220	-0.007	-0.059	-0.037	-0.036	-0.027	-0.023	-0.025	0.186	0.158	0.163	0.158	0.155	1.000	-0.003	-0.053
SEX	-0.091	-0.034	-0.031	-0.025	-0.022	-0.017	-0.017	0.014	0.025	-0.031	-0.058	-0.071	-0.066	-0.060	-0.055	-0.044	0.000	-0.001	-0.009	-0.002	-0.002	-0.003	1.000	-0.040
default	0.014	-0.020	-0.014	-0.014	-0.010	-0.007	-0.005	0.028	-0.154	-0.024	0.325	0.264	0.235	0.217	0.204	0.187	-0.073	-0.059	-0.056	-0.057	-0.055	-0.053	-0.040	1.000

Figure 26: Correlation matrix of the entire dataset featuring all variables.

Second one shows the correlation between columns by showing the correlation between 2 columns is the size of square.

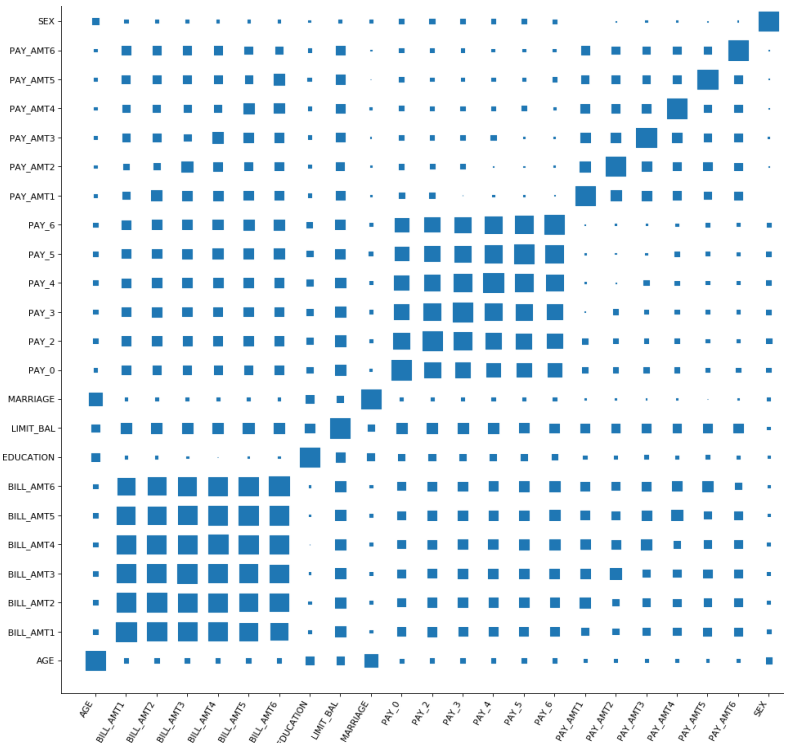


Figure 27: Correlation matrix shown by size of square cells

Most importantly, the only variable with a notable correlation with default is payment status (PAY columns)

The next highest correlation with default occurs with balance limit at -0.15, which itself is somewhat negatively correlated with payment statuses. These measures indicate that customers with lower balance limits have more delays in payments and are more likely to default. This argument is also supported by Figure 3, which you can see defaulting customers have low balance limits, and payment status (PAY columns) are correlated with default.

There is a positive correlation coefficient of 0.14 between limit balance and age. And there is a negative -0.23 correlation between limit balance and education. Age seems do not have much correlation with default, but based on figure 7, defaulting customer groups tends to grow with age.

Education is negatively correlated with Balance limit. Bills amount and payment amounts are moderately positively correlation with balance limit, as expected, since customer have higher credit are more likely to spend more.

Chapter 4: Data pre-processing

1. SPECIAL TREATMENT FOR AGE COLUMN:

For the age columns, we will replace the Age column with $\log_{10}(\text{age} + 1)$

2. DATA PROCESSING USING MIN-MAX SCALER WITH RANGE (0, 1):

Before we apply any method, such as PCA, or feature selection, data should be scaled and non-negative:

	LIMIT_BAL	SEX	EDUCATION	MARRIAGE	AGE	PAY_0	PAY_2	PAY_3	PAY_4	PAY_5	PAY_6	BILL_AMT1	BILL_AMT2	BILL_AMT3	BILL_AMT4
0	0.010101	1.0	0.333333	0.333333	0.051724	0.4	0.4	0.1	0.1	0.0	0.0	0.149982	0.069164	0.086723	0.160138
1	0.111111	1.0	0.333333	0.666667	0.086207	0.1	0.4	0.2	0.2	0.2	0.4	0.148892	0.067858	0.087817	0.163220
2	0.080808	1.0	0.333333	0.666667	0.224138	0.2	0.2	0.2	0.2	0.2	0.2	0.172392	0.079532	0.093789	0.173637
3	0.040404	1.0	0.333333	0.333333	0.275862	0.2	0.2	0.2	0.2	0.2	0.2	0.188100	0.111995	0.113407	0.186809
4	0.040404	0.0	0.333333	0.333333	0.620690	0.1	0.2	0.1	0.2	0.2	0.2	0.154144	0.071601	0.106020	0.179863

BILL_AMT5	BILL_AMT6	PAY_AMT1	PAY_AMT2	PAY_AMT3	PAY_AMT4	PAY_AMT5	PAY_AMT6
0.080648	0.260979	0.000000	0.000409	0.000000	0.000000	0.000000	0.000000
0.084074	0.263485	0.000000	0.000594	0.001116	0.001610	0.000000	0.003783
0.095470	0.272928	0.001738	0.000891	0.001116	0.001610	0.002345	0.009458
0.109363	0.283685	0.002290	0.001199	0.001339	0.001771	0.002506	0.001892
0.099633	0.275681	0.002290	0.021779	0.011160	0.014493	0.001615	0.001284

3. DATA PROCESSING WITH PRINCIPLE COMPONENT ANALYSIS: PCA (Principal Component Analysis) mainly using to reduce the size of the feature space while retaining as much of the information as possible. In order to decide how many components, we should reduce to, we will plot a graph for Cumulative Summation of the Explained Variance

Choosing n components when performing PCA (model did not exclude any columns):

```
n_component = 1 → 0.5528528428073621
n_component = 2 → 0.6964800622705274
n_component = 3 → 0.7904022739803037
n_component = 4 → 0.8455618207700792
n_component = 5 → 0.8839247959378522
n_component = 6 → 0.9166721899058589
n_component = 7 → 0.9380700545764724
n_component = 8 → 0.9562071457419538
n_component = 9 → 0.9687710243399839
n_component = 10 → 0.9771386310509677
n_component = 11 → 0.9833384310033771
n_component = 12 → 0.9875115316576303
n_component = 13 → 0.9906108232655144
n_component = 14 → 0.9931429460461323
n_component = 15 → 0.9954486759841078
n_component = 16 → 0.9968243940331725
n_component = 17 → 0.9978175525224455
n_component = 18 → 0.9986091155659746
n_component = 19 → 0.9991613442997643
n_component = 20 → 0.9994913120632744
n_component = 21 → 0.9997417197558053
n_component = 22 → 0.999912543541501
n_component = 23 → 1.0
```

Using this graph, we can decide to set n components = 13 in PCA

After applying PCA with $n_{\text{components}} = 13$, below are dimensional reduction for first few samples:

	#1	#2	#3	#4	#5	#6	#7	#8	#9	#10	#11	#12	#13	default
0	-0.404232	-0.008005	-0.002597	-0.212566	-0.211858	0.142288	-0.353981	-0.055854	-0.037792	0.047794	0.098949	0.061357	0.008969	NaN
1	-0.391889	0.209224	-0.196988	-0.063890	-0.009020	0.000373	0.104670	-0.040404	0.036161	0.249256	0.028769	0.093556	-0.011375	1.0
2	-0.392642	0.084095	-0.126567	-0.077123	0.080733	-0.050136	0.021148	0.006008	0.000249	0.006144	0.001706	0.000031	-0.002720	1.0
3	-0.394038	0.068885	0.154030	-0.108658	-0.085688	0.014440	0.002314	0.077498	0.005925	0.001265	0.001589	-0.000635	-0.005089	0.0
4	0.610813	-0.089679	0.334445	-0.141113	0.074264	-0.165212	0.070620	0.116729	0.010839	0.046725	0.077771	0.028073	0.007249	0.0

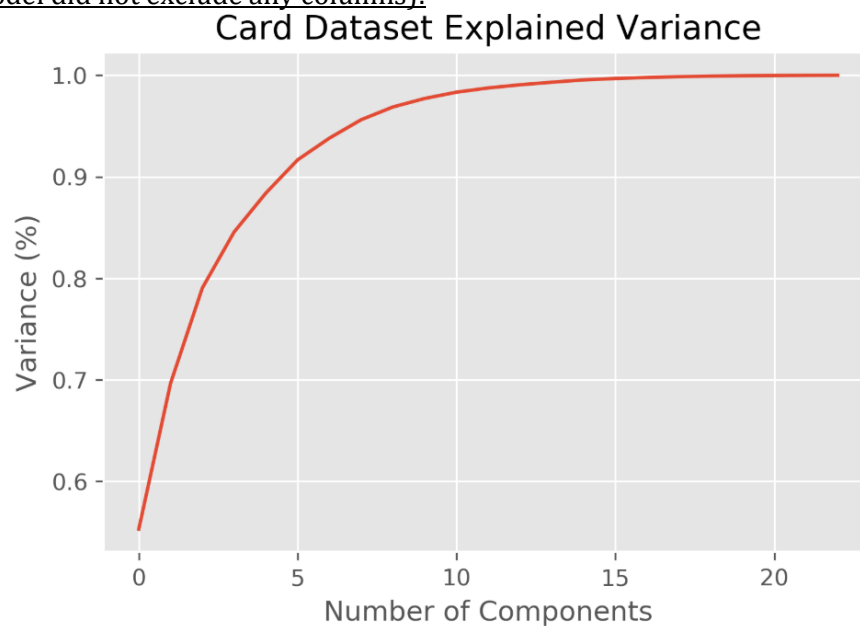


Figure 28: Number of components in PCA vs. Explained Variance of the dataset

Chapter 5: Feature selection

Some popular feature selection techniques such as univariate selection using Chi-squared test of independence, feature importance with tree classifiers and variance thresholds were ran to understand whether a feature is important or noise. These techniques were chosen to represent different aspects of the dataset.

Chi-squared test of independence is a statistical tool used to understand the level of dependence between a feature and the target variable. If the target variable is highly dependent on a feature, it means that changing the feature is likely to induce a change in the target variable thus making it desirable as a predictor.

$$\text{Entropy}, H(X) = - \sum_{i=0}^c p_i \log_2(p_i) ; \text{Information gain}, IG(S, x_i) = H(S) - \sum_{j=1}^{c_i} \frac{|S_j|}{|S|} H(S_j)$$

Decision tree classifier was also used to compute feature importance as a cross reference to support our decision. Decision tree classifiers are good feature selectors as at every node in the tree, a feature is selected for that node and the information gain for that node is calculated based on entropy reduction. Entropy is a metric for purity of a sample. The feature with the most information gain is then used for the node. Above formula shows the mathematical formula of calculating entropy and the information gain from entropy. As shown in the formula, maximising information gain also maximises the loss of entropy which makes it a good metric for feature importance.

METHOD 1: UNIVARIATE SELECTION:

Univariate selection using Chi2 to compute Score

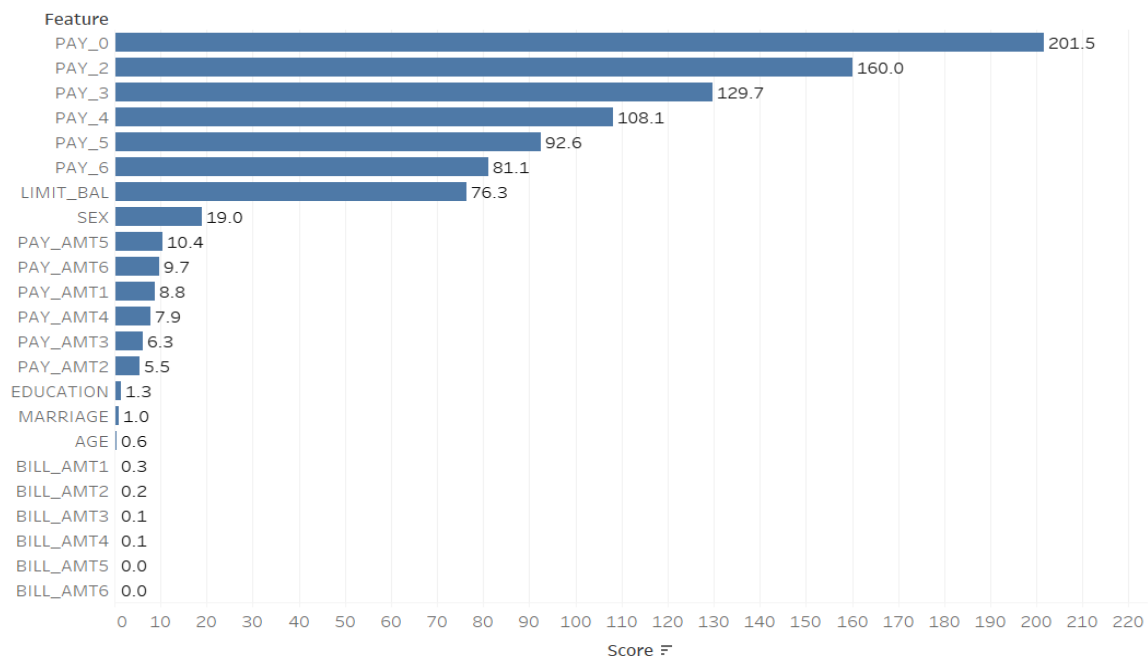


Figure 29: Univariate Selection using Chi2 to compute score

METHOD 2: FEATURE IMPORTANCE BY ExtraTreesClassifier:

Feature Importance using ExtraTreesClassifier

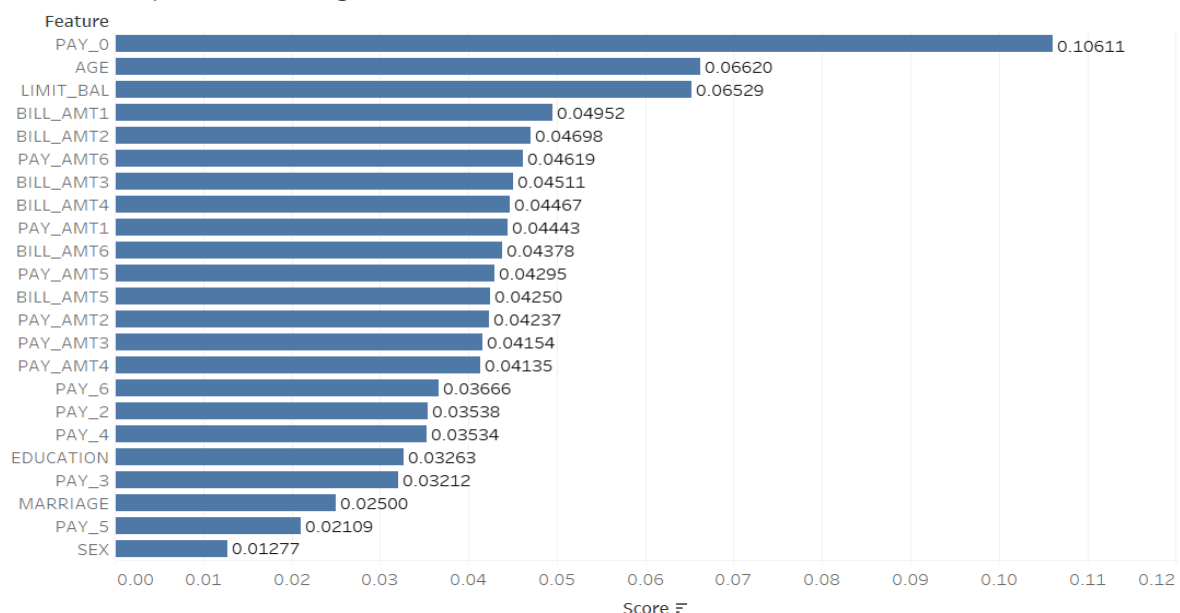


Figure 30: Feature importance using Extra Tree Classifier

Chapter 6: Splitting data and Imbalanced Class treatment

Part 1: Data snooping prevention:

Before doing any further research or training models, to avoid data snooping, we must split the dataset into 2 sub parts, naming training and testing dataset with ratio: 7: 3

Number transactions training dataset: 21000

Number transactions testing dataset: 9000

Total number of transactions: 30000

Part 2: Imbalanced class treatment:

As you can see, there are many records have “default payment next month” = 0 while less in 1 (there are 22.12 % of no-default and 77.88 % default next month), hence, this phenomenon is called “Imbalanced class” The reason behind why we should find a treatment for Imbalanced class is because conventional methods such as Decision Tree and Logistic Regression have a bias towards classes which have number of instances. They tend to only predict the majority class data. The features of the minority class are treated as noise and are often ignored. Thus, there is a high probability of misclassification of the minority class as compared to the majority class.

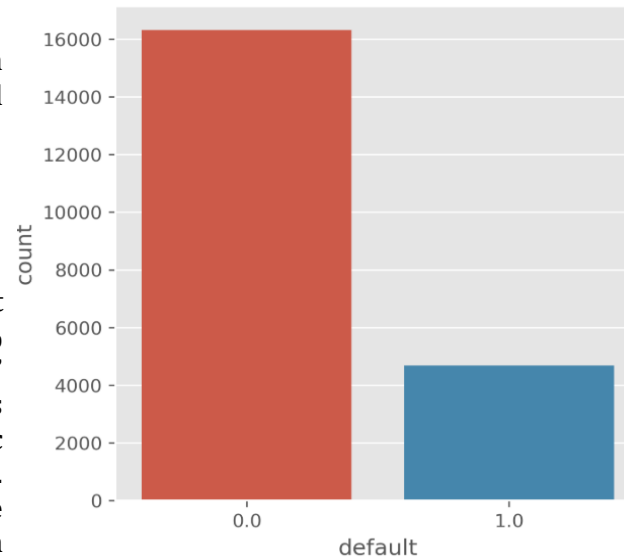


Figure 31: Count-plot for target column: "default payment next month"

RESAMPLING TECHNIQUES:

Under sampling	Over sampling
When you randomly eliminate instances from the majority class of a dataset and assign it to the minority class (without filling out the void created in majority class). The void that gets created in the majority dataset for this makes the process <i>random</i> .	Taking any help from the majority class, you increase the instances corresponding to the minority class by replicating them up to a constant degree. In this case, you do not decrease the number of instances assigned to the majority class.
It can help improve the runtime of the model and solve the memory problems by reducing the number of training data samples when the training data set is enormous.	Unlike undersampling, this method leads to no information loss.
It can discard useful information about the data itself which could be necessary for building rule-based classifiers such as Random Forests. The sample chosen by random undersampling may be a biased sample. And it will not be an accurate representation of the population in that case. Therefore, it can cause the classifier to perform poorly on real unseen data.	It increases the likelihood of overfitting since it replicates the minority class events.

INBUILT OVERSAMPLING AND UNDERSAMPLING METHODS:

Synthetic Minority Over-sampling Technique (SMOTE)

SMOTE is a technique based on nearest neighbor judged by Euclidean Distance between data points in feature space. There is a percentage of Over-Sampling which indicates the number of synthetic samples to be created and this percentage parameter of Over-sampling is always a multiple of 100. If the percentage of Over-sampling is 100, then for each instance, a new sample will be created. Hence, the number of minority class instances will get doubled. Similarly, if the percentage of Over-sampling is 200, then the total number of minority class samples will get tripled.

In *SMOTE*, for each minority instance, k number of nearest neighbours are found such that they also belong to the same class where, $k = (SMOTE\%)/100$

The difference between the feature vector of the considered instance and the feature vectors of the k nearest neighbor are found. So, k number of difference vectors are obtained.

The k difference vectors are each multiplied by a random number between 0 and 1 (excluding 0 and 1).

Now, the difference vectors, after being multiplied by random numbers, are added to the feature vector of the considered instance (original minority instance) at each iteration. [6]

Adaptive synthetic sampling method (ADASYN) is an improved version of *SMOTE*. There is a small improvement. After creating some, it adds a small value to the points thus making it more realistic. In order words, instead of all sampling being linearly correlated to the parent, they have a little more variance in them, more scattered. [7]

Random under sample: perform random under sampling by undersampling majority class by randomly picking samples with or without replacement [8]

Combination of undersampling and oversampling into a hybrid strategy SMOTE + ENN: Combination we will focus on this project is SMOTE and Edited Nearest Neighbours (ENN). [9]

Chapter 7: Model Selection and evaluation

Part 1: Models and datasets

MODELS USING ORIGIN DATASET WITH OPTIONS "CLASS WEIGHT" IN SKLEARN MODELS

In order to overcome Imbalanced class, there is an option we can use while avoiding spending time to resample dataset. In `sklearn` models, the option `class_weight` is available for several models such as:

- Random Forest classifier
- Logistic regression
- Linear Support Vector Classifier
- Stochastic gradient descent classifier
- Decision tree classifier,...

The default value for `class_weight` is `None`. When we want to use it, define it as "balanced" value:

- Weights associated with classes in the form `{class_label: weight}`. If not given, all classes are supposed to have weight one.
- The "balanced" mode uses the values of `y` to automatically adjust weights inversely proportional to class frequencies in the input data as $n_samples / (n_classes * np.bincount(y))$.
- Note that these weights will be multiplied with `sample_weight` (passed through the fit method) if `sample_weight` is specified.

For binary classification problem, if class 0 samples are much less than class 1, more weight will be given to class 0 samples. Following this approach, we will achieve a balanced weight model eventually. Therefore, this is also an effective method to overcome Imbalanced class.

MODEL SYSTEM STRUCTURE:

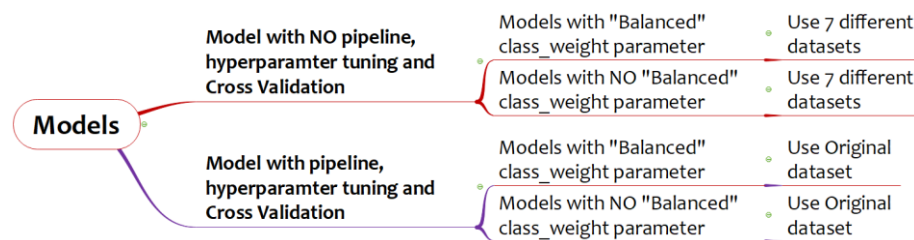


Figure 32: Model system structure

We have tried several models and tested them both with their default setting or with cross validated versions.

Also, for each model created, we will test them with different datasets, these datasets are created with resampling treatment methods shown above.

There are **7 total datasets**, namely: Original dataset (no resampling treatment), SMOTE dataset, ADASYN dataset, EMOTE + ENN dataset, Random Under sample dataset, Oversampling and Undersampling dataset (we did design our own oversample, under sample methods)

There are different models, such as: Naïve Bayes, Support Vector Machine, Logistic regression, Multilayer perceptron, Random Forest, Xtreme Gradient Boosting, Stochastic gradient descent, K-nearest neighbors, ...

ENSEMBLE MODEL: STACKING

We did try to improve model performance by using different pipeline or scaler and hyperparameter tuning. Every model has its own strength and weaknesses. Eventually, we came up with an idea of combining different models together to cover one model's weakness by using another model's strength. There is a term to define this process, it is called Stacking. The explanation of Stacking is shown in Chapter 2, part 2. Below is our stacking model:

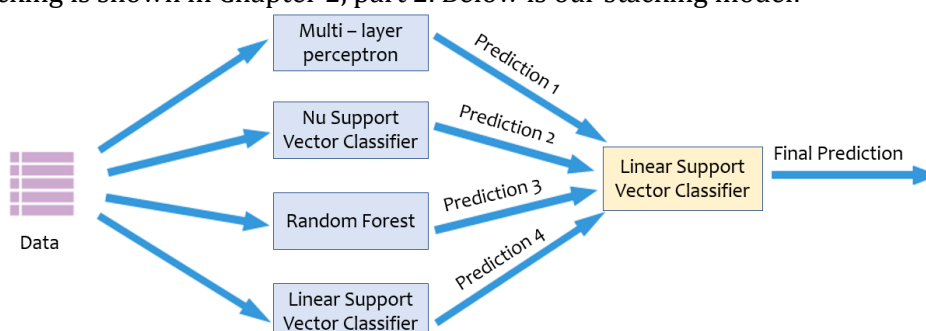


Figure 33: Our project Stacking model structure

MODELS EVALUATION BENCHMARK

Various models along with different pre-processing methods were tested and the performance of each classifier based on their F1 score, balanced accuracy, precision, recall, sensitivity, specificity. Balanced accuracy or F1 score can be a benchmark for the evaluation as they are more robust to imbalanced dataset. The strategy used was to start from simple models and gradually increase their complexity. Thus, Logistic Regression was chosen as a baseline since it is a relatively simple linear model and thus easier to understand compared to other learning models. Also, training time was factored in when making this decision. After tuning hyperparameters of Logistic Regression models, a F1 score of **0.4733** was achieved, which became our baseline of comparison.

Part 2: Model evaluation results

Details of all models with all available training datasets:

- Explanation for model type:

 - *Original* means model with NO pipeline or Cross validation.
 - *Original Stack* means original model used as a classifier in Stacking, but now is it evaluated solely.
 - *GSCV* means model only use Grid Search Cross Validation, no pipeline
 - *Pipeline_GSCV* means model with pipeline use Grid Search Cross Validation
 - *Pipeline_RSCV* means model with pipeline use Randomized Search Cross Validation
 - *Multiple SVM* means a Support Vector Machine model with different kernels as Hyperparamter.
- Explanation of Model Name:

 - *NB* stands for Naïve Bayes
 - *SVM* stands for Support Vector Machine
 - *SVC* stands for Support Vector classifier
 - *RBF* stands for Radial Basis Function (a kernel of Support Vector Machine)
 - *SVC* stands for Support Vector Classifier (a linear kernel Support Vector Machine)
 - *SGD* stands for Stochastic Gradient Descent
 - *XGB* stands for Xtreme Gradient Boosting
 - *RF* stands for Random Forest
 - *MLP* stands for Multi-layer perceptron

Name	Data set	Type	Class weight	Balanced Accuracy	F1 score	Accuracy	Sensitivity	Specificity	Precision	Recall
Bernoulli NB	ADASYN	Original	None	<div></div> 0.6936	<div></div> 0.5007	<div></div> 0.7649	<div></div> 0.5413	<div></div> 0.8459	<div></div> 0.4658	<div></div> 0.5413
	Origin	Original	None	<div></div> 0.6780	<div></div> 0.4827	<div></div> 0.7680	<div></div> 0.4969	<div></div> 0.8591	<div></div> 0.4692	<div></div> 0.4969
	Over sample	Original	None	<div></div> 0.7018	<div></div> 0.4868	<div></div> 0.7369	<div></div> 0.5730	<div></div> 0.8307	<div></div> 0.4231	<div></div> 0.5730
	Random US	Original	None	<div></div> 0.7022	<div></div> 0.4849	<div></div> 0.7340	<div></div> 0.5750	<div></div> 0.8294	<div></div> 0.4193	<div></div> 0.5750
	SMOTE	Original	None	<div></div> 0.6686	<div></div> 0.4787	<div></div> 0.7779	<div></div> 0.4684	<div></div> 0.8689	<div></div> 0.4896	<div></div> 0.4684
	SMOTE + ENN	Original	None	<div></div> 0.7217	<div></div> 0.5076	<div></div> 0.7336	<div></div> 0.6306	<div></div> 0.8128	<div></div> 0.4247	<div></div> 0.6306
	Under sample	Original	None	<div></div> 0.6998	<div></div> 0.4909	<div></div> 0.7449	<div></div> 0.5648	<div></div> 0.8349	<div></div> 0.4341	<div></div> 0.5648
Decision Tree	ADASYN	Original	None	<div></div> 0.6635	<div></div> 0.3947	<div></div> 0.6783	<div></div> 0.4816	<div></div> 0.8454	<div></div> 0.3344	<div></div> 0.4816
	Origin	Original	Balanced	<div></div> 0.6361	<div></div> 0.3867	<div></div> 0.7324	<div></div> 0.3872	<div></div> 0.8849	<div></div> 0.3861	<div></div> 0.3872
			None	<div></div> 0.6424	<div></div> 0.3932	<div></div> 0.7267	<div></div> 0.4066	<div></div> 0.8781	<div></div> 0.3806	<div></div> 0.4066
	Over sample	Original	None	<div></div> 0.5788	<div></div> 0.2318	<div></div> 0.6604	<div></div> 0.2352	<div></div> 0.9224	<div></div> 0.2284	<div></div> 0.2352
	Random US	Original	None	<div></div> 0.7023	<div></div> 0.4171	<div></div> 0.6204	<div></div> 0.6235	<div></div> 0.7812	<div></div> 0.3133	<div></div> 0.6235
	SMOTE	Original	None	<div></div> 0.6610	<div></div> 0.3990	<div></div> 0.6918	<div></div> 0.4699	<div></div> 0.8521	<div></div> 0.3468	<div></div> 0.4699
	SMOTE + ENN	Original	None	<div></div> 0.7064	<div></div> 0.4560	<div></div> 0.6856	<div></div> 0.6051	<div></div> 0.8078	<div></div> 0.3658	<div></div> 0.6051
Gaussian NB	Under sample	Original	None	<div></div> 0.6400	<div></div> 0.3022	<div></div> 0.5038	<div></div> 0.4934	<div></div> 0.7867	<div></div> 0.2178	<div></div> 0.4934
	ADASYN	Original	None	<div></div> 0.6782	<div></div> 0.3821	<div></div> 0.3572	<div></div> 0.9128	<div></div> 0.4435	<div></div> 0.2417	<div></div> 0.9128
	Origin	Original	None	<div></div> 0.7251	<div></div> 0.4815	<div></div> 0.6918	<div></div> 0.6571	<div></div> 0.7931	<div></div> 0.3799	<div></div> 0.6571
	Over sample	Original	None	<div></div> 0.6143	<div></div> 0.3648	<div></div> 0.2879	<div></div> 0.9388	<div></div> 0.2898	<div></div> 0.2264	<div></div> 0.9388
	Random US	Original	None	<div></div> 0.7368	<div></div> 0.4339	<div></div> 0.5592	<div></div> 0.7755	<div></div> 0.6980	<div></div> 0.3012	<div></div> 0.7755
	SMOTE	Original	None	<div></div> 0.6992	<div></div> 0.3908	<div></div> 0.3927	<div></div> 0.8944	<div></div> 0.5040	<div></div> 0.2500	<div></div> 0.8944
	SMOTE + ENN	Original	None	<div></div> 0.7022	<div></div> 0.3920	<div></div> 0.3996	<div></div> 0.8888	<div></div> 0.5156	<div></div> 0.2514	<div></div> 0.8888
K-Nearest Neighbors	Under sample	Original	None	<div></div> 0.6874	<div></div> 0.3723	<div></div> 0.4057	<div></div> 0.8092	<div></div> 0.5656	<div></div> 0.2417	<div></div> 0.8092
	ADASYN	Original	None	<div></div> 0.7074	<div></div> 0.4271	<div></div> 0.6296	<div></div> 0.6342	<div></div> 0.7806	<div></div> 0.3220	<div></div> 0.6342
	Origin	Original	None	<div></div> 0.6255	<div></div> 0.4211	<div></div> 0.7929	<div></div> 0.3459	<div></div> 0.9050	<div></div> 0.5381	<div></div> 0.3459
	Over sample	Original	None	<div></div> 0.6290	<div></div> 0.2924	<div></div> 0.5493	<div></div> 0.4276	<div></div> 0.8305	<div></div> 0.2222	<div></div> 0.4276
	Random US	Original	None	<div></div> 0.7192	<div></div> 0.4585	<div></div> 0.6646	<div></div> 0.6520	<div></div> 0.7863	<div></div> 0.3535	<div></div> 0.6520
	SMOTE	Original	None	<div></div> 0.7018	<div></div> 0.4377	<div></div> 0.6640	<div></div> 0.6005	<div></div> 0.8030	<div></div> 0.3444	<div></div> 0.6005
	SMOTE + ENN	Original	None	<div></div> 0.7224	<div></div> 0.4377	<div></div> 0.6147	<div></div> 0.6888	<div></div> 0.7560	<div></div> 0.3208	<div></div> 0.6888
Multi Layer Perceptron	Under sample	Original	None	<div></div> 0.6426	<div></div> 0.3056	<div></div> 0.4996	<div></div> 0.5056	<div></div> 0.7796	<div></div> 0.2190	<div></div> 0.5056
	ADASYN	Original Stack	None	<div></div> 0.6865	<div></div> 0.4968	<div></div> 0.7707	<div></div> 0.5199	<div></div> 0.8531	<div></div> 0.4757	<div></div> 0.5199
	Origin	Original Stack	None	<div></div> 0.6187	<div></div> 0.4381	<div></div> 0.8193	<div></div> 0.3235	<div></div> 0.9140	<div></div> 0.6788	<div></div> 0.3235
	Over sample	Original Stack	None	<div></div> 0.6229	<div></div> 0.3739	<div></div> 0.7470	<div></div> 0.3469	<div></div> 0.8989	<div></div> 0.4055	<div></div> 0.3469
	Random US	Original Stack	None	<div></div> 0.6842	<div></div> 0.5061	<div></div> 0.7831	<div></div> 0.5102	<div></div> 0.8581	<div></div> 0.5020	<div></div> 0.5102
	SMOTE	Original Stack	None	<div></div> 0.6811	<div></div> 0.5026	<div></div> 0.7838	<div></div> 0.5015	<div></div> 0.8606	<div></div> 0.5036	<div></div> 0.5015
	SMOTE + ENN	Original Stack	None	<div></div> 0.7115	<div></div> 0.4782	<div></div> 0.7100	<div></div> 0.6102	<div></div> 0.8128	<div></div> 0.3932	<div></div> 0.6102
	Under sample	Original Stack	None	<div></div> 0.6761	<div></div> 0.3777	<div></div> 0.6053	<div></div> 0.5500	<div></div> 0.8021	<div></div> 0.2876	<div></div> 0.5500

Logistic Regression	ADASYN	Original	None	<div><div></div></div> 0.7261	<div><div></div></div> 0.4518	<div><div></div></div> 0.6368	<div><div></div></div> 0.6872	<div><div></div></div> 0.7650	<div><div></div></div> 0.3365	<div><div></div></div> 0.6872
		Pipeline GSCV	Balanced	<div><div></div></div> 0.7297	<div><div></div></div> 0.4531	<div><div></div></div> 0.6312	<div><div></div></div> 0.7015	<div><div></div></div> 0.7580	<div><div></div></div> 0.3346	<div><div></div></div> 0.7015
			None	<div><div></div></div> 0.7278	<div><div></div></div> 0.4550	<div><div></div></div> 0.6396	<div><div></div></div> 0.6908	<div><div></div></div> 0.7648	<div><div></div></div> 0.3392	<div><div></div></div> 0.6908
	Origin	Original	Balanced	<div><div></div></div> 0.7159	<div><div></div></div> 0.4696	<div><div></div></div> 0.6896	<div><div></div></div> 0.6311	<div><div></div></div> 0.8007	<div><div></div></div> 0.3739	<div><div></div></div> 0.6311
			None	<div><div></div></div> 0.5867	<div><div></div></div> 0.3525	<div><div></div></div> 0.8102	<div><div></div></div> 0.2372	<div><div></div></div> 0.9362	<div><div></div></div> 0.6858	<div><div></div></div> 0.2372
		Pipeline GSCV	Balanced	<div><div></div></div> 0.7173	<div><div></div></div> 0.4733	<div><div></div></div> 0.6929	<div><div></div></div> 0.6337	<div><div></div></div> 0.8008	<div><div></div></div> 0.3777	<div><div></div></div> 0.6337
			None	<div><div></div></div> 0.5850	<div><div></div></div> 0.3470	<div><div></div></div> 0.8093	<div><div></div></div> 0.2327	<div><div></div></div> 0.9374	<div><div></div></div> 0.6826	<div><div></div></div> 0.2327
	Over sample	Original	None	<div><div></div></div> 0.6885	<div><div></div></div> 0.4194	<div><div></div></div> 0.6610	<div><div></div></div> 0.5622	<div><div></div></div> 0.8148	<div><div></div></div> 0.3344	<div><div></div></div> 0.5622
		Pipeline GSCV	None	<div><div></div></div> 0.5000	<div><div></div></div> 0.3577	<div><div></div></div> 0.2178	<div><div></div></div> 1.0000	<div><div></div></div> 0.0000	<div><div></div></div> 0.2178	<div><div></div></div> 1.0000
	Random US	Original	None	<div><div></div></div> 0.7132	<div><div></div></div> 0.4695	<div><div></div></div> 0.6941	<div><div></div></div> 0.6214	<div><div></div></div> 0.8050	<div><div></div></div> 0.3772	<div><div></div></div> 0.6214
		Pipeline GSCV	Balanced	<div><div></div></div> 0.7189	<div><div></div></div> 0.4687	<div><div></div></div> 0.6826	<div><div></div></div> 0.6429	<div><div></div></div> 0.7949	<div><div></div></div> 0.3687	<div><div></div></div> 0.6429
			None	<div><div></div></div> 0.7176	<div><div></div></div> 0.4728	<div><div></div></div> 0.6916	<div><div></div></div> 0.6352	<div><div></div></div> 0.8000	<div><div></div></div> 0.3766	<div><div></div></div> 0.6352
	SMOTE	Original	None	<div><div></div></div> 0.7157	<div><div></div></div> 0.4685	<div><div></div></div> 0.6881	<div><div></div></div> 0.6311	<div><div></div></div> 0.8003	<div><div></div></div> 0.3725	<div><div></div></div> 0.6311
		Pipeline GSCV	Balanced	<div><div></div></div> 0.7199	<div><div></div></div> 0.4677	<div><div></div></div> 0.6791	<div><div></div></div> 0.6474	<div><div></div></div> 0.7924	<div><div></div></div> 0.3661	<div><div></div></div> 0.6474
			None	<div><div></div></div> 0.7186	<div><div></div></div> 0.4679	<div><div></div></div> 0.6818	<div><div></div></div> 0.6423	<div><div></div></div> 0.7948	<div><div></div></div> 0.3679	<div><div></div></div> 0.6423
	SMOTE + ENN	Original	None	<div><div></div></div> 0.7286	<div><div></div></div> 0.4297	<div><div></div></div> 0.5741	<div><div></div></div> 0.7367	<div><div></div></div> 0.7205	<div><div></div></div> 0.3033	<div><div></div></div> 0.7367
		Pipeline GSCV	Balanced	<div><div></div></div> 0.7148	<div><div></div></div> 0.4706	<div><div></div></div> 0.6930	<div><div></div></div> 0.6265	<div><div></div></div> 0.8031	<div><div></div></div> 0.3768	<div><div></div></div> 0.6265
			None	<div><div></div></div> 0.7313	<div><div></div></div> 0.4399	<div><div></div></div> 0.5957	<div><div></div></div> 0.7291	<div><div></div></div> 0.7334	<div><div></div></div> 0.3150	<div><div></div></div> 0.7291
	Under sample	Original	None	<div><div></div></div> 0.6192	<div><div></div></div> 0.3156	<div><div></div></div> 0.6650	<div><div></div></div> 0.3546	<div><div></div></div> 0.8839	<div><div></div></div> 0.2843	<div><div></div></div> 0.3546
Multinomial NB	ADASYN	Original	None	<div><div></div></div> 0.7018	<div><div></div></div> 0.4469	<div><div></div></div> 0.6799	<div><div></div></div> 0.5939	<div><div></div></div> 0.8098	<div><div></div></div> 0.3583	<div><div></div></div> 0.5939
	Over sample	Original	None	<div><div></div></div> 0.5018	<div><div></div></div> 0.0101	<div><div></div></div> 0.7820	<div><div></div></div> 0.0051	<div><div></div></div> 0.9986	<div><div></div></div> 0.4545	<div><div></div></div> 0.0051
	Random US	Original	None	<div><div></div></div> 0.7092	<div><div></div></div> 0.4391	<div><div></div></div> 0.6500	<div><div></div></div> 0.6291	<div><div></div></div> 0.7892	<div><div></div></div> 0.3373	<div><div></div></div> 0.6291
	SMOTE	Original	None	<div><div></div></div> 0.7082	<div><div></div></div> 0.4386	<div><div></div></div> 0.6513	<div><div></div></div> 0.6255	<div><div></div></div> 0.7909	<div><div></div></div> 0.3377	<div><div></div></div> 0.6255
	SMOTE + ENN	Original	None	<div><div></div></div> 0.6758	<div><div></div></div> 0.3767	<div><div></div></div> 0.3607	<div><div></div></div> 0.8872	<div><div></div></div> 0.4643	<div><div></div></div> 0.2391	<div><div></div></div> 0.8872
	Under sample	Original	None	<div><div></div></div> 0.6236	<div><div></div></div> 0.4174	<div><div></div></div> 0.7928	<div><div></div></div> 0.3408	<div><div></div></div> 0.9064	<div><div></div></div> 0.5383	<div><div></div></div> 0.3408
Multiple SVM	ADASYN	Pipeline GSCV	None	<div><div></div></div> 0.6398	<div><div></div></div> 0.3643	<div><div></div></div> 0.6894	<div><div></div></div> 0.4087	<div><div></div></div> 0.8709	<div><div></div></div> 0.3287	<div><div></div></div> 0.4087
	Origin	Pipeline GSCV	None	<div><div></div></div> 0.5960	<div><div></div></div> 0.3698	<div><div></div></div> 0.8046	<div><div></div></div> 0.2633	<div><div></div></div> 0.9287	<div><div></div></div> 0.6209	<div><div></div></div> 0.2633
	Over sample	Pipeline GSCV	None	<div><div></div></div> 0.5913	<div><div></div></div> 0.2415	<div><div></div></div> 0.6127	<div><div></div></div> 0.2832	<div><div></div></div> 0.8993	<div><div></div></div> 0.2105	<div><div></div></div> 0.2832
	Random US	Pipeline GSCV	None	<div><div></div></div> 0.6909	<div><div></div></div> 0.5034	<div><div></div></div> 0.7713	<div><div></div></div> 0.5321	<div><div></div></div> 0.8498	<div><div></div></div> 0.4776	<div><div></div></div> 0.5321
	SMOTE	Pipeline GSCV	None	<div><div></div></div> 0.6429	<div><div></div></div> 0.3753	<div><div></div></div> 0.6989	<div><div></div></div> 0.4153	<div><div></div></div> 0.8706	<div><div></div></div> 0.3423	<div><div></div></div> 0.4153
	SMOTE + ENN	Pipeline GSCV	None	<div><div></div></div> 0.6972	<div><div></div></div> 0.4336	<div><div></div></div> 0.6668	<div><div></div></div> 0.5857	<div><div></div></div> 0.8087	<div><div></div></div> 0.3442	<div><div></div></div> 0.5857
	Under sample	Pipeline GSCV	None	<div><div></div></div> 0.6390	<div><div></div></div> 0.3043	<div><div></div></div> 0.5306	<div><div></div></div> 0.4714	<div><div></div></div> 0.8065	<div><div></div></div> 0.2247	<div><div></div></div> 0.4714
Nu SVC	ADASYN	Original Stack	None	<div><div></div></div> 0.6745	<div><div></div></div> 0.3612	<div><div></div></div> 0.5602	<div><div></div></div> 0.5709	<div><div></div></div> 0.7781	<div><div></div></div> 0.2642	<div><div></div></div> 0.5709
	Origin	Original Stack	None	<div><div></div></div> 0.6040	<div><div></div></div> 0.2491	<div><div></div></div> 0.5452	<div><div></div></div> 0.3464	<div><div></div></div> 0.8616	<div><div></div></div> 0.1945	<div><div></div></div> 0.3464
	Over sample	Original Stack	None	<div><div></div></div> 0.6020	<div><div></div></div> 0.2477	<div><div></div></div> 0.5567	<div><div></div></div> 0.3352	<div><div></div></div> 0.8689	<div><div></div></div> 0.1965	<div><div></div></div> 0.3352
	Random US	Original Stack	None	<div><div></div></div> 0.6342	<div><div></div></div> 0.2943	<div><div></div></div> 0.5129	<div><div></div></div> 0.4663	<div><div></div></div> 0.8020	<div><div></div></div> 0.2150	<div><div></div></div> 0.4663
	SMOTE	Original Stack	None	<div><div></div></div> 0.6436	<div><div></div></div> 0.3116	<div><div></div></div> 0.5321	<div><div></div></div> 0.4862	<div><div></div></div> 0.8010	<div><div></div></div> 0.2293	<div><div></div></div> 0.4862
	SMOTE + ENN	Original Stack	None	<div><div></div></div> 0.6996	<div><div></div></div> 0.3871	<div><div></div></div> 0.5367	<div><div></div></div> 0.6719	<div><div></div></div> 0.7273	<div><div></div></div> 0.2719	<div><div></div></div> 0.6719
	Under sample	Original Stack	None	<div><div></div></div> 0.6014	<div><div></div></div> 0.2442	<div><div></div></div> 0.5439	<div><div></div></div> 0.3383	<div><div></div></div> 0.8646	<div><div></div></div> 0.1910	<div><div></div></div> 0.3383
SVM (RBF)	ADASYN	Pipeline GSCV	None	<div><div></div></div> 0.7132	<div><div></div></div> 0.4463	<div><div></div></div> 0.6548	<div><div></div></div> 0.6388	<div><div></div></div> 0.7875	<div><div></div></div> 0.3429	<div><div></div></div> 0.6388
	Origin	Pipeline GSCV	Balanced	<div><div></div></div> 0.6597	<div><div></div></div> 0.4781	<div><div></div></div> 0.7904	<div><div></div></div> 0.4408	<div><div></div></div> 0.8785	<div><div></div></div> 0.5224	<div><div></div></div> 0.4408
	Random US	Pipeline GSCV	None	<div><div></div></div> 0.6907	<div><div></div></div> 0.4990	<div><div></div></div> 0.7671	<div><div></div></div> 0.5327	<div><div></div></div> 0.8488	<div><div></div></div> 0.4694	<div><div></div></div> 0.5327
	SMOTE	Pipeline GSCV	None	<div><div></div></div> 0.6965	<div><div></div></div> 0.4609	<div><div></div></div> 0.7114	<div><div></div></div> 0.5663	<div><div></div></div> 0.8266	<div><div></div></div> 0.3885	<div><div></div></div> 0.5663
	SMOTE + ENN	Pipeline GSCV	None	<div><div></div></div> 0.7186	<div><div></div></div> 0.4437	<div><div></div></div> 0.6372	<div><div></div></div> 0.6643	<div><div></div></div> 0.7730	<div><div></div></div> 0.3331	<div><div></div></div> 0.6643
	Under sample	Pipeline GSCV	None	<div><div></div></div> 0.6402	<div><div></div></div> 0.3320	<div><div></div></div> 0.6217	<div><div></div></div> 0.4316	<div><div></div></div> 0.8488	<div><div></div></div> 0.2697	<div><div></div></div> 0.4316
Stacking	Origin	Pipeline GSCV	None	<div><div></div></div> 0.6165	<div><div></div></div> 0.4341	<div><div></div></div> 0.8198	<div><div></div></div> 0.3173	<div><div></div></div> 0.9157	<div><div></div></div> 0.6865	<div><div></div></div> 0.3173
	Random US	Pipeline GSCV	None	<div><div></div></div> 0.6866	<div><div></div></div> 0.5125	<div><div></div></div> 0.7861	<div><div></div></div> 0.5163	<div><div></div></div> 0.8570	<div><div></div></div> 0.5088	<div><div></div></div> 0.5163
	SMOTE + ENN	Pipeline GSCV	None	<div><div></div></div> 0.7291	<div><div></div></div> 0.4936	<div><div></div></div> 0.7033	<div><div></div></div> 0.6638	<div><div></div></div> 0.7945	<div><div></div></div> 0.3928	<div><div></div></div> 0.6638
XGB	ADASYN	Pipeline GSCV	None	<div><div></div></div> 0.6417	<div><div></div></div> 0.4686	<div><div></div></div> 0.8084	<div><div></div></div> 0.3878	<div><div></div></div> 0.8955	<div><div></div></div> 0.5919	<div><div></div></div> 0.3878
	Origin	Original	None	<div><div></div></div> 0.6237	<div><div></div></div> 0.4499	<div><div></div></div> 0.8207	<div><div></div></div> 0.3367	<div><div></div></div> 0.9106	<div><div></div></div> 0.6776	<div><div></div></div> 0.3367
		Pipeline GSCV	None	<div><div></div></div> 0.6320	<div><div></div></div> 0.4664	<div><div></div></div> 0.8210	<div><div></div></div> 0.3592	<div><div></div></div> 0.9047	<div><div></div></div> 0.6648	<div><div></div></div> 0.3592
	Over sample	Original	None	<div><div></div></div> 0.6386	<div><div></div></div> 0.3192	<div><div></div></div> 0.5929	<div><div></div></div> 0.4383	<div><div></div></div> 0.8390	<div><div></div></div> 0.2510	<div><div></div></div> 0.4383
	Random US	Pipeline GSCV	None	<div><div></div></div> 0.7241	<div><div></div></div> 0.5190	<div><div></div></div> 0.7442	<div><div></div></div> 0.6337	<div><div></div></div> 0.8146	<div><div></div></div> 0.4395	<div><div></div></div> 0.6337
	SMOTE	Pipeline GSCV	None	<div><div></div></div> 0.6420	<div><div></div></div> 0.4685	<div><div></div></div> 0.8079	<div><div></div></div> 0.3888	<div><div></div></div> 0.8952	<div><div></div></div> 0.5893	<div><div></div></div> 0.3888
	SMOTE + ENN	Pipeline GSCV	None	<div><div></div></div> 0.6954	<div><div></div></div> 0.5148	<div><div></div></div> 0.7771	<div><div></div></div> 0.5429	<div><div></div></div> 0.8479	<div><div></div></div> 0.4894	<div><div></div></div> 0.5429
	Under sample	Original	None	<div><div></div></div> 0.6403	<div><div></div></div> 0.3055	<div><div></div></div> 0.5251	<div><div></div></div> 0.4796	<div><div></div></div> 0.8011	<div><div></div></div> 0.2241	<div><div></div></div> 0.4796
		Pipeline GSCV	None	<div><div></div></div> 0.6387	<div><div></div></div> 0.2982	<div><div></div></div> 0.4618	<div><div></div></div> 0.5250	<div><div></div></div> 0.7524	<div><div></div></div> 0.2082	<div><div></div></div> 0.5250

Random Forest	ADASYN	Original	None	<div><div></div></div> 0.6645	<div><div></div></div> 0.4615	<div><div></div></div> 0.7663	<div><div></div></div> 0.4597	<div><div></div></div> 0.8694	<div><div></div></div> 0.4632	<div><div></div></div> 0.4597	
		Original Stack	None	<div><div></div></div> 0.7137	<div><div></div></div> 0.5021	<div><div></div></div> 0.7382	<div><div></div></div> 0.6061	<div><div></div></div> 0.8212	<div><div></div></div> 0.4286	<div><div></div></div> 0.6061	
		Pipeline GSCV	Balanced	<div><div></div></div> 0.7219	<div><div></div></div> 0.5017	<div><div></div></div> 0.7257	<div><div></div></div> 0.6342	<div><div></div></div> 0.8097	<div><div></div></div> 0.4150	<div><div></div></div> 0.6342	
			None	<div><div></div></div> 0.7183	<div><div></div></div> 0.4997	<div><div></div></div> 0.7283	<div><div></div></div> 0.6230	<div><div></div></div> 0.8137	<div><div></div></div> 0.4172	<div><div></div></div> 0.6230	
		Pipeline RSCV	Balanced	<div><div></div></div> 0.7201	<div><div></div></div> 0.4906	<div><div></div></div> 0.7134	<div><div></div></div> 0.6337	<div><div></div></div> 0.8066	<div><div></div></div> 0.4003	<div><div></div></div> 0.6337	
			None	<div><div></div></div> 0.7320	<div><div></div></div> 0.5037	<div><div></div></div> 0.7133	<div><div></div></div> 0.6679	<div><div></div></div> 0.7961	<div><div></div></div> 0.4043	<div><div></div></div> 0.6679	
	Origin	Original	Balanced	<div><div></div></div> 0.6089	<div><div></div></div> 0.3977	<div><div></div></div> 0.8028	<div><div></div></div> 0.2990	<div><div></div></div> 0.9189	<div><div></div></div> 0.5937	<div><div></div></div> 0.2990	
			None	<div><div></div></div> 0.6120	<div><div></div></div> 0.4228	<div><div></div></div> 0.8186	<div><div></div></div> 0.3051	<div><div></div></div> 0.9188	<div><div></div></div> 0.6881	<div><div></div></div> 0.3051	
		Pipeline GSCV	Balanced	<div><div></div></div> 0.6973	<div><div></div></div> 0.5125	<div><div></div></div> 0.7723	<div><div></div></div> 0.5495	<div><div></div></div> 0.8451	<div><div></div></div> 0.4802	<div><div></div></div> 0.5495	
			None	<div><div></div></div> 0.6171	<div><div></div></div> 0.4355	<div><div></div></div> 0.8200	<div><div></div></div> 0.3189	<div><div></div></div> 0.9153	<div><div></div></div> 0.6868	<div><div></div></div> 0.3189	
		Pipeline RSCV	Balanced	<div><div></div></div> 0.6856	<div><div></div></div> 0.5039	<div><div></div></div> 0.7790	<div><div></div></div> 0.5153	<div><div></div></div> 0.8559	<div><div></div></div> 0.4929	<div><div></div></div> 0.5153	
	Over sample	Original	None	<div><div></div></div> 0.5248	<div><div></div></div> 0.1074	<div><div></div></div> 0.7469	<div><div></div></div> 0.0699	<div><div></div></div> 0.9796	<div><div></div></div> 0.2314	<div><div></div></div> 0.0699	
		Original Stack	None	<div><div></div></div> 0.6119	<div><div></div></div> 0.2865	<div><div></div></div> 0.6287	<div><div></div></div> 0.3423	<div><div></div></div> 0.8814	<div><div></div></div> 0.2463	<div><div></div></div> 0.3423	
	Random US	Original	None	<div><div></div></div> 0.7175	<div><div></div></div> 0.5157	<div><div></div></div> 0.7491	<div><div></div></div> 0.6133	<div><div></div></div> 0.8217	<div><div></div></div> 0.4449	<div><div></div></div> 0.6133	
		Original Stack	None	<div><div></div></div> 0.7203	<div><div></div></div> 0.5160	<div><div></div></div> 0.7457	<div><div></div></div> 0.6224	<div><div></div></div> 0.8182	<div><div></div></div> 0.4406	<div><div></div></div> 0.6224	
		Pipeline GSCV	Balanced	<div><div></div></div> 0.7088	<div><div></div></div> 0.5220	<div><div></div></div> 0.7674	<div><div></div></div> 0.5832	<div><div></div></div> 0.8345	<div><div></div></div> 0.4725	<div><div></div></div> 0.5832	
			None	<div><div></div></div> 0.7171	<div><div></div></div> 0.5147	<div><div></div></div> 0.7486	<div><div></div></div> 0.6122	<div><div></div></div> 0.8219	<div><div></div></div> 0.4440	<div><div></div></div> 0.6122	
		Pipeline RSCV	Balanced	<div><div></div></div> 0.7151	<div><div></div></div> 0.5237	<div><div></div></div> 0.7613	<div><div></div></div> 0.6026	<div><div></div></div> 0.8276	<div><div></div></div> 0.4631	<div><div></div></div> 0.6026	
			None	<div><div></div></div> 0.7252	<div><div></div></div> 0.5211	<div><div></div></div> 0.7453	<div><div></div></div> 0.6362	<div><div></div></div> 0.8141	<div><div></div></div> 0.4413	<div><div></div></div> 0.6362	
	SMOTE	Original	None	<div><div></div></div> 0.6579	<div><div></div></div> 0.4615	<div><div></div></div> 0.7770	<div><div></div></div> 0.4388	<div><div></div></div> 0.8770	<div><div></div></div> 0.4867	<div><div></div></div> 0.4388	
		Original Stack	None	<div><div></div></div> 0.7030	<div><div></div></div> 0.5044	<div><div></div></div> 0.7559	<div><div></div></div> 0.5704	<div><div></div></div> 0.8357	<div><div></div></div> 0.4521	<div><div></div></div> 0.5704	
		Pipeline GSCV	Balanced	<div><div></div></div> 0.7184	<div><div></div></div> 0.5158	<div><div></div></div> 0.7480	<div><div></div></div> 0.6163	<div><div></div></div> 0.8206	<div><div></div></div> 0.4435	<div><div></div></div> 0.6163	
			None	<div><div></div></div> 0.7152	<div><div></div></div> 0.5171	<div><div></div></div> 0.7539	<div><div></div></div> 0.6051	<div><div></div></div> 0.8252	<div><div></div></div> 0.4515	<div><div></div></div> 0.6051	
		Pipeline RSCV	Balanced	<div><div></div></div> 0.7141	<div><div></div></div> 0.5217	<div><div></div></div> 0.7604	<div><div></div></div> 0.6000	<div><div></div></div> 0.8282	<div><div></div></div> 0.4615	<div><div></div></div> 0.6000	
			None	<div><div></div></div> 0.7177	<div><div></div></div> 0.5198	<div><div></div></div> 0.7537	<div><div></div></div> 0.6122	<div><div></div></div> 0.8231	<div><div></div></div> 0.4516	<div><div></div></div> 0.6122	
	SMOTE + ENN	Original	None	<div><div></div></div> 0.7090	<div><div></div></div> 0.5055	<div><div></div></div> 0.7489	<div><div></div></div> 0.5893	<div><div></div></div> 0.8286	<div><div></div></div> 0.4425	<div><div></div></div> 0.5893	
		Original Stack	None	<div><div></div></div> 0.7286	<div><div></div></div> 0.4885	<div><div></div></div> 0.6966	<div><div></div></div> 0.6653	<div><div></div></div> 0.7920	<div><div></div></div> 0.3859	<div><div></div></div> 0.6653	
		Pipeline GSCV	Balanced	<div><div></div></div> 0.7121	<div><div></div></div> 0.5134	<div><div></div></div> 0.7538	<div><div></div></div> 0.5964	<div><div></div></div> 0.8277	<div><div></div></div> 0.4507	<div><div></div></div> 0.5964	
			None	<div><div></div></div> 0.7178	<div><div></div></div> 0.5147	<div><div></div></div> 0.7476	<div><div></div></div> 0.6148	<div><div></div></div> 0.8209	<div><div></div></div> 0.4427	<div><div></div></div> 0.6148	
		Pipeline RSCV	Balanced	<div><div></div></div> 0.7119	<div><div></div></div> 0.5137	<div><div></div></div> 0.7543	<div><div></div></div> 0.5959	<div><div></div></div> 0.8280	<div><div></div></div> 0.4515	<div><div></div></div> 0.5959	
			None	<div><div></div></div> 0.7418	<div><div></div></div> 0.4980	<div><div></div></div> 0.6898	<div><div></div></div> 0.7066	<div><div></div></div> 0.7769	<div><div></div></div> 0.3845	<div><div></div></div> 0.7066	
	Under sample	Original	None	<div><div></div></div> 0.6134	<div><div></div></div> 0.2679	<div><div></div></div> 0.5568	<div><div></div></div> 0.3724	<div><div></div></div> 0.8543	<div><div></div></div> 0.2092	<div><div></div></div> 0.3724	
		Original Stack	None	<div><div></div></div> 0.6418	<div><div></div></div> 0.3057	<div><div></div></div> 0.5110	<div><div></div></div> 0.4944	<div><div></div></div> 0.7893	<div><div></div></div> 0.2213	<div><div></div></div> 0.4944	
SGD	ADASYN	Original	None	<div><div></div></div> 0.6952	<div><div></div></div> 0.4803	<div><div></div></div> 0.7391	<div><div></div></div> 0.5536	<div><div></div></div> 0.8369	<div><div></div></div> 0.4242	<div><div></div></div> 0.5536	
		Pipeline GSCV	Balanced	<div><div></div></div> 0.7299	<div><div></div></div> 0.4324	<div><div></div></div> 0.5781	<div><div></div></div> 0.7378	<div><div></div></div> 0.7221	<div><div></div></div> 0.3058	<div><div></div></div> 0.7378	
			None	<div><div></div></div> 0.7152	<div><div></div></div> 0.4694	<div><div></div></div> 0.6906	<div><div></div></div> 0.6286	<div><div></div></div> 0.8018	<div><div></div></div> 0.3746	<div><div></div></div> 0.6286	
	Origin	Original	Balanced	<div><div></div></div> 0.6885	<div><div></div></div> 0.4925	<div><div></div></div> 0.7632	<div><div></div></div> 0.5276	<div><div></div></div> 0.8495	<div><div></div></div> 0.4618	<div><div></div></div> 0.5276	
			None	<div><div></div></div> 0.5011	<div><div></div></div> 0.0061	<div><div></div></div> 0.7824	<div><div></div></div> 0.0031	<div><div></div></div> 0.9991	<div><div></div></div> 0.6000	<div><div></div></div> 0.0031	
		Pipeline GSCV	Balanced	<div><div></div></div> 0.6045	<div><div></div></div> 0.3919	<div><div></div></div> 0.8066	<div><div></div></div> 0.2862	<div><div></div></div> 0.9227	<div><div></div></div> 0.6213	<div><div></div></div> 0.2862	
			None	<div><div></div></div> 0.5686	<div><div></div></div> 0.2941	<div><div></div></div> 0.8032	<div><div></div></div> 0.1883	<div><div></div></div> 0.9490	<div><div></div></div> 0.6721	<div><div></div></div> 0.1883	
	Over sample	Original	None	<div><div></div></div> 0.5287	<div><div></div></div> 0.1385	<div><div></div></div> 0.7844	<div><div></div></div> 0.0796	<div><div></div></div> 0.9779	<div><div></div></div> 0.5342	<div><div></div></div> 0.0796	
		Pipeline GSCV	None	<div><div></div></div> 0.5838	<div><div></div></div> 0.3642	<div><div></div></div> 0.2633	<div><div></div></div> 0.9689	<div><div></div></div> 0.1987	<div><div></div></div> 0.2243	<div><div></div></div> 0.9689	
	Oversample	Pipeline GSCV	Balanced	<div><div></div></div> 0.6792	<div><div></div></div> 0.3763	<div><div></div></div> 0.3693	<div><div></div></div> 0.8735	<div><div></div></div> 0.4850	<div><div></div></div> 0.2398	<div><div></div></div> 0.8735	
	Random US	Original	None	<div><div></div></div> 0.6899	<div><div></div></div> 0.4900	<div><div></div></div> 0.7586	<div><div></div></div> 0.5327	<div><div></div></div> 0.8471	<div><div></div></div> 0.4537	<div><div></div></div> 0.5327	
		Pipeline GSCV	Balanced	<div><div></div></div> 0.6902	<div><div></div></div> 0.4971	<div><div></div></div> 0.7658	<div><div></div></div> 0.5316	<div><div></div></div> 0.8488	<div><div></div></div> 0.4668	<div><div></div></div> 0.5316	
			None	<div><div></div></div> 0.6885	<div><div></div></div> 0.4973	<div><div></div></div> 0.7684	<div><div></div></div> 0.5260	<div><div></div></div> 0.8509	<div><div></div></div> 0.4716	<div><div></div></div> 0.5260	
	SMOTE	Original	None	<div><div></div></div> 0.6763	<div><div></div></div> 0.4959	<div><div></div></div> 0.7838	<div><div></div></div> 0.4883	<div><div></div></div> 0.8643	<div><div></div></div> 0.5037	<div><div></div></div> 0.4883	
		Pipeline GSCV	Balanced	<div><div></div></div> 0.6922	<div><div></div></div> 0.4979	<div><div></div></div> 0.7638	<div><div></div></div> 0.5378	<div><div></div></div> 0.8467	<div><div></div></div> 0.4635	<div><div></div></div> 0.5378	
			None	<div><div></div></div> 0.6904	<div><div></div></div> 0.4974	<div><div></div></div> 0.7658	<div><div></div></div> 0.5321	<div><div></div></div> 0.8487	<div><div></div></div> 0.4669	<div><div></div></div> 0.5321	
	SMOTE + ENN	Original	None	<div><div></div></div> 0.7284	<div><div></div></div> 0.4504	<div><div></div></div> 0.6286	<div><div></div></div> 0.6990	<div><div></div></div> 0.7578	<div><div></div></div> 0.3323	<div><div></div></div> 0.6990	
		Pipeline GSCV	Balanced	<div><div></div></div> 0.7117	<div><div></div></div> 0.4702	<div><div></div></div> 0.6980	<div><div></div></div> 0.6153	<div><div></div></div> 0.8080	<div><div></div></div> 0.3804	<div><div></div></div> 0.6153	
			None	<div><div></div></div> 0.7218	<div><div></div></div> 0.4549	<div><div></div></div> 0.6526	<div><div></div></div> 0.6658	<div><div></div></div> 0.7778	<div><div></div></div> 0.3455	<div><div></div></div> 0.6658	
	Under sample	Pipeline GSCV	None	<div><div></div></div> 0.5918	<div><div></div></div> 0.3644	<div><div></div></div> 0.2691	<div><div></div></div> 0.9622	<div><div></div></div> 0.2213	<div><div></div></div> 0.2248	<div><div></div></div> 0.9622	
	Undersample	Pipeline GSCV	Balanced	<div><div></div></div> 0.5906	<div><div></div></div> 0.3647	<div><div></div></div> 0.2681	<div><div></div></div> 0.9648	<div><div></div></div> 0.2163	<div><div></div></div> 0.2249	<div><div></div></div> 0.9648	

SVM (Linear)	ADASYN	Original	None	<div></div> 0.7267	<div></div> 0.4511	<div></div> 0.6339	<div></div> 0.6908	<div></div> 0.7627	<div></div> 0.3349	<div></div> 0.6908	
		Original Stack	None	<div></div> 0.6855	<div></div> 0.4931	<div></div> 0.7681	<div></div> 0.5179	<div></div> 0.8532	<div></div> 0.4706	<div></div> 0.5179	
		Pipeline GSCV	Balanced	<div></div> 0.7299	<div></div> 0.4536	<div></div> 0.6320	<div></div> 0.7015	<div></div> 0.7583	<div></div> 0.3352	<div></div> 0.7015	
			None	<div></div> 0.7276	<div></div> 0.4553	<div></div> 0.6406	<div></div> 0.6898	<div></div> 0.7655	<div></div> 0.3398	<div></div> 0.6898	
	Origin	Original	Balanced	<div></div> 0.7149	<div></div> 0.4724	<div></div> 0.6958	<div></div> 0.6255	<div></div> 0.8042	<div></div> 0.3796	<div></div> 0.6255	
			None	<div></div> 0.5417	<div></div> 0.1951	<div></div> 0.7938	<div></div> 0.1148	<div></div> 0.9685	<div></div> 0.6503	<div></div> 0.1148	
		Original Stack	None	<div></div> 0.6065	<div></div> 0.4086	<div></div> 0.8170	<div></div> 0.2903	<div></div> 0.9226	<div></div> 0.6897	<div></div> 0.2903	
		Pipeline GSCV	Balanced	<div></div> 0.7156	<div></div> 0.4762	<div></div> 0.7001	<div></div> 0.6260	<div></div> 0.8053	<div></div> 0.3843	<div></div> 0.6260	
			None	<div></div> 0.5602	<div></div> 0.2660	<div></div> 0.8013	<div></div> 0.1653	<div></div> 0.9551	<div></div> 0.6807	<div></div> 0.1653	
	Over sample	Original	None	<div></div> 0.6925	<div></div> 0.4141	<div></div> 0.6407	<div></div> 0.5832	<div></div> 0.8018	<div></div> 0.3211	<div></div> 0.5832	
		Original Stack	None	<div></div> 0.6979	<div></div> 0.3782	<div></div> 0.4450	<div></div> 0.7750	<div></div> 0.6207	<div></div> 0.2501	<div></div> 0.7750	
	Random US	Original	None	<div></div> 0.7106	<div></div> 0.4724	<div></div> 0.7032	<div></div> 0.6102	<div></div> 0.8110	<div></div> 0.3854	<div></div> 0.6102	
			Original Stack	<div></div> 0.6804	<div></div> 0.5000	<div></div> 0.7822	<div></div> 0.5000	<div></div> 0.8608	<div></div> 0.5000	<div></div> 0.5000	
		Pipeline GSCV	Balanced	<div></div> 0.7168	<div></div> 0.4706	<div></div> 0.6896	<div></div> 0.6337	<div></div> 0.7999	<div></div> 0.3743	<div></div> 0.6337	
			None	<div></div> 0.7149	<div></div> 0.4734	<div></div> 0.6972	<div></div> 0.6250	<div></div> 0.8048	<div></div> 0.3810	<div></div> 0.6250	
	SMOTE	Original	None	<div></div> 0.7150	<div></div> 0.4700	<div></div> 0.6918	<div></div> 0.6276	<div></div> 0.8024	<div></div> 0.3757	<div></div> 0.6276	
			Original Stack	<div></div> 0.6693	<div></div> 0.5001	<div></div> 0.7972	<div></div> 0.4658	<div></div> 0.8728	<div></div> 0.5399	<div></div> 0.4658	
		Pipeline GSCV	Balanced	<div></div> 0.7182	<div></div> 0.4701	<div></div> 0.6861	<div></div> 0.6393	<div></div> 0.7971	<div></div> 0.3717	<div></div> 0.6393	
			None	<div></div> 0.7181	<div></div> 0.4714	<div></div> 0.6882	<div></div> 0.6383	<div></div> 0.7980	<div></div> 0.3737	<div></div> 0.6383	
	SMOTE + ENN	Original	None	<div></div> 0.7278	<div></div> 0.4303	<div></div> 0.5787	<div></div> 0.7306	<div></div> 0.7250	<div></div> 0.3049	<div></div> 0.7306	
			Original Stack	<div></div> 0.7272	<div></div> 0.4697	<div></div> 0.6686	<div></div> 0.6740	<div></div> 0.7805	<div></div> 0.3604	<div></div> 0.6740	
		Pipeline GSCV	Balanced	<div></div> 0.7119	<div></div> 0.4751	<div></div> 0.7049	<div></div> 0.6133	<div></div> 0.8105	<div></div> 0.3877	<div></div> 0.6133	
			None	<div></div> 0.7306	<div></div> 0.4410	<div></div> 0.6002	<div></div> 0.7240	<div></div> 0.7373	<div></div> 0.3170	<div></div> 0.7240	
	Under sample	Original	None	<div></div> 0.6125	<div></div> 0.2897	<div></div> 0.6338	<div></div> 0.3429	<div></div> 0.8822	<div></div> 0.2507	<div></div> 0.3429	
		Original Stack	None	<div></div> 0.6978	<div></div> 0.3905	<div></div> 0.5598	<div></div> 0.6474	<div></div> 0.7481	<div></div> 0.2795	<div></div> 0.6474	

Figure 34: Evaluation metrics of all models for all dataset.

It is observed that most of the models using Pipeline have F1 score improved (F1 Score). Hence, we do have better model.

Best performance in each model (sorted by F1 Score):

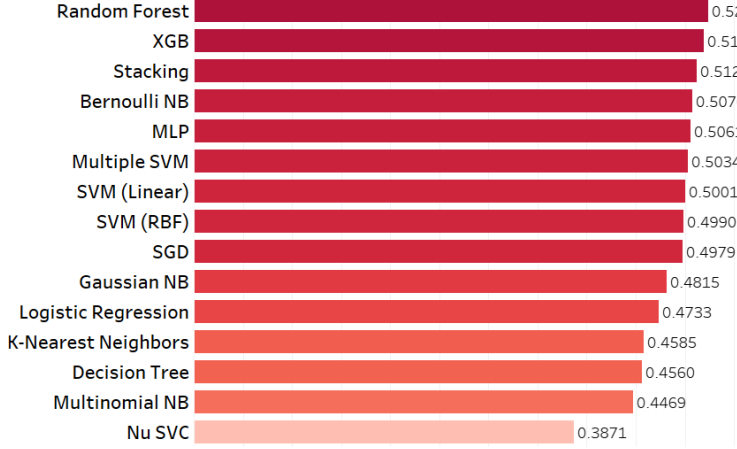


Figure 35: Maximum F1 Score for each model

The difference between F1 score of top models are insignificant.

The most 3 impressive ones are Random Forest, Xtreme Gradient Boosting and Stacking.

Although there is not much different in performance of these 3 models. Stacking seems to have vey big potential that can explored in future.

Worth note taking potential of Stacking model:



Figure 36: Maximum F1 Score in each classifier component of Stacking model

Stacking is the combination of four different models, namely Random Forest, Multi-layer perceptron, Linear Support Vector Machine and Nu SVC (shown in figure 33). Although the performance of Nu SVC is not high. And three others are high, the performance of Stacking model is not affected much. This proves how powerful of this model.

Part 3: Experiential lessons in this project:

For **linear models**, the performances of the models converge around the F1 score of around 0.47 with no sign of improvement no matter how much time is spent to tune the hyper-parameters of these models. This is likely caused by the lack of robustness of the models as they are not designed to accurately this non-linear dataset.

For **Naive Bayes**, the training time and the testing time of this model is the fastest out of all the models implemented in this project. And the F1 score of Bernoulli Naïve Bayes turns out to be impressive. This is a rare case as the dataset matches to Bernoulli distribution. However, Naïve Bayes are not on our preferable model lists as this performance are not consistent if we have a new set of data that does not follow Bernoulli Distribution, then our performance is very bad.

For **KNN**, the model is trained very quickly, but it takes a very long time for it to make its prediction. The more data used for training this model, the longer it takes for it to find the nearest neighbours of the data that it would like to classify. This model learn by memorising the data, which caused it to be easily overfitted. The F1 score converges to 0.46

For **SVM**, training of the data was not able converge easily as a high number of iterations is observed even though features have been normalized and the weights of the classes have been balanced. When number of iterations is limited, the F1 score observed during cross validation was not ideal. This shows that data was not able to fit well as it was hard to find a suitable separating hyperplane and computationally intensive to tune the hyperparameters. [10]

Extreme Gradient Boosting trees are an efficient implementation of traditional Gradient Boosting trees, having clever penalisation of trees, proportional shrinking of leaf nodes, extra randomisation parameters which helped increase both the performance and the speed of the classifiers [11]. The intuition behind gradient boosting trees and the other ensemble methods used in this paper is the idea that crowds generally perform better than individuals.

In **Random Forests classifiers**, multiple weak classifiers are trained in parallel using techniques such as bootstrapping on both data points and features, before aggregating them to reduce the variance. On the other hand, gradient boosting trees allow the results of each weak classifier to affect the next weak classifier. Since each individual classifier is made deliberately weak, it is unable to explore the full parameter space, thus making it unlikely to overfit. Also, with the aggregation of multiple weak learners, the resulting model can also prevent overfitting, which explains why Extreme Gradient Boosting trees, an advanced implementation with heuristics performed much better than the other models.

Similar to the idea of deep networks, **stacking** classifiers follow the idea that hierarchical approach to learning might be better in terms of variance and bias as it is hard for models to learn multiple things at once without being too large and resulting in overfitting from the curse of dimensionality phenomenon [12]. In stacking, by splitting models into 2 levels in this paper (the lower level classifiers with an additional layer as meta-classifiers), lower level models can explore different parts of the problem space and higher-level meta-classifiers can combine their predictions in meaningful ways. This means that when these 4 classifiers are combined together, with their predictions as input to the higher-level meta-classifier, the meta-classifier can prioritise certain classifiers when presented with certain types of data points, thus achieving a slightly better overall result. To achieve a better result, the weights of the meta-classifier must also be adjusted which presents a new set of hyperparameters to tune. As such, the complexity of the model has increased not only in the number of models but also the number of hyperparameters, increasing the difficulty of tuning such models. In this paper, only 4 models were used for stacking, but the number of models could be increased significantly, by not only varying the model type, but also the pre-processing steps or hyperparameters for model. This presents a tradeoff between model performance and understandability of the model, which must be duly considered when deciding to use similar techniques.

Part 1: Overall picture

Predicting credit card default is a complex problem with many features and the exploratory description analysis shows this complexity in heatmap and distribution plots.

Although there are multiple reports showing the solution for this problem. In this project, we aim for a higher goal than that. We have explored a stacking model, which has not been implemented in articles we have read. The performance turns out to be optimistic and we would like to improve this model further.

Regretfully, some of the features in this dataset were obfuscated, possibly for privacy reasons and the contextual meaning was lost. In the future, new datasets could be created to include more features as well as replace some of the obfuscated features in this dataset with more relevant ones. This could help improve the exploratory description analysis part, so that enable us to improve model for higher performance.

Also, the focus in this paper was only on predicting whether a customer will default next month or not, but the goal can be extended into predict a customer will default within a range of period, not only one month.

Part 2: Improvement on existing models

Feature selection: In models, we have chosen some columns that can be removed in data-preprocessing step based on Feature Importance using Extra Tree Classifier, such as:

`"SEX", "PAY_5", "MARRIAGE", "PAY_3", "EDUCATION", "PAY_2", "PAY_4", "PAY_6"`

It is suggesting that we should try other feature selection method as well, namely Chi2 Univariate selection.

Potential models: In our opinion, one of the top potential models are **Ensemble cross validated Stacking**. First, it can combine different strong learners models. If we do in the right way, the Stacking model is expected to outperform its sub-models. However, the choice for meta-learner models has to be as simple as possible. Otherwise, overfitting will likely occur, resulting in worse performance. Another criterion to take note when we do Stacking is timing. In this project, we did put Support Vector Machine model as a meta-classifier in the end, and hyperparameter tuning for the desire of achieving higher score. However, efficiency is not achieved as it takes very long to train a dataset (more than 5 hours). In the future, we planned to consider Logistic regression as our meta-classifier instead as it is faster to train, and it will prevent overfitting by combining with regularizers.

Part 3: Potential models have not been implemented

Neural network model is highly anticipated to have a potential outperform our existing models. However, if we want to tune Neural Network from scratch to build an efficient model will consume a lot of effort and energy. Therefore, we will recommend utilizing pre-built model.

Chapter 9: References

- [1] M. Merikoski, "Predicting and Preventing Credit Card Default," 2018.
- [2] J. Rocca, "Ensemble methods: bagging, boosting and stacking," 23 04 2019. [Online]. Available: <https://towardsdatascience.com/ensemble-methods-bagging-boosting-and-stacking-c9214a10a205>.
- [3] J. VanderPlas, Python Data Science Handbook.
- [4] Book Down, "Lab notes for Statistics for Social Sciences II: Multivariate Techniques," [Online]. Available: <https://bookdown.org/egarpor/SSS2-UC3M/>.
- [5] A. V. Srinivasan, "Stochastic Gradient Descent — Clearly Explained," 07 09 2019. [Online]. Available: <https://towardsdatascience.com/stochastic-gradient-descent-clearly-explained-53d239905d31>.
- [6] S. A. Rahim, "SMOTE AND NEAR MISS IN PYTHON: MACHINE LEARNING IN IMBALANCED DATASETS," 4 06 2018. [Online]. Available: <https://medium.com/@saeedAR/smote-and-near-miss-in-python-machine-learning-in-imbalanced-datasets-b7976d9a7a79>.
- [7] K. Vala, "ADASYN: Adaptive Synthetic Sampling Method for Imbalanced Data," 10 09 2019. [Online]. Available: <https://towardsdatascience.com/adasyn-adaptive-synthetic-sampling-method-for-imbalanced-data-602a3673ba16>.
- [8] imbalanced-learn.readthedocs, "Documentation for imblearn.under_sampling," [Online]. Available: https://imbalanced-learn.readthedocs.io/en/stable/generated/imblearn.under_sampling.RandomUnderSampler.html.
- [9] imbalanced-learn.readthedocs, "Documentation for imbalanced-learn SMOTE ENN," [Online]. Available: <https://imbalanced-learn.readthedocs.io/en/stable/generated/imblearn.combine.SMOTEENN.html>.
- [10] V. M. Y. Cherkassky, "Practical selection of SVM parameters and noise estimation for SVM regression," 2004.
- [11] R. Gandhi, "Gradient Boosting and XGBoost," 2018. [Online]. Available: <https://hackernoon.com/gradient-boosting-and-xgboost-90862daa6c77>.
- [12] D. H. Wolpert, "Stacked generalization," 1992. [Online]. Available: <https://www.sciencedirect.com/science/article/abs/pii/S0893608005800231>.