

# Machine Learning Models Based Data Quality Analysis to Detect Credit Card Frauds

<sup>1</sup>Amit Pundir, <sup>2</sup>Rajesh Pandey

<sup>1</sup>M.Tech Scholar, Shobhit Institute of Engineering & Technology, Meerut, (Uttar Pradesh), India

<sup>2</sup>Assistant Professor, Shobhit Institute of Engineering & Technology, Meerut, (Uttar Pradesh), India

<sup>1</sup>[pundirs99@gmail.com](mailto:pundirs99@gmail.com), <sup>2</sup>[rajesh@shobhituniversity.ac.in](mailto:rajesh@shobhituniversity.ac.in)

## Abstract

Considering the broader consequences and taking into account the multiple cycles, cash distortion is a serious issue in money-related business. The management of a quality data with mining of data has been successfully applied to dataset to automate the examination of vast measures of complex data. Similarly, data mining has played a major role in isolating frauds like credit card during online trades. Detection of fraud at credit card is a management of quality data issue considered under data mining, attempted for two important reasons - first, the profiles of standard and fraudulent practices change regularly, and furthermore, clarifications are needed. Second, credit card, extortion data is surprisingly sluggish. This research paper examines the performance of Decision Trees, Logistics Regression, and Random Forests that rely on deliberately heavily tilted credit cards fraud data. The dataset of credit card transaction is sourced from Kaggle (a publicly available dataset repository) with total number of transactions 284,807. These strategies are applied to values of raw data and data processing methods. Evaluation of the performance of methods depends on accuracy, sensitivity, specificity, precision, and recall. The results show ideal accuracy for decision trees, logistics regression, and random forest classifiers with 90.8%, 98.5%, and 99.1% respectively.

## Keywords

Fraud detection, Credit card fraud detection, data mining, logistics regression, decision tree, random forest, collative analysis

## 1. Introduction

Money related growth is a concern with unlimited results in public experts, corporate associations, cash business, currently high reliance on web development has given credit card based transaction in preset time, yet money extortion / digital frauds with credit cards leads the status of on/off line transactions. As money transaction over the credit card has become an inevitable installment strategy, the center has of

late computational systems to deal with this issue related to the credit card fraud. There are many distortion/misrepresentation detection tools and programming-based systems used in associations to prevent counterfeiting such as credit card, retail, online exchange, Internet business, security and enterprise. Data quality management is a famous and notable process within the data mining strategy used to address issues of fraud detection during the use of credit cards. It is difficult to be certain of the actual expectation and legality behind an application or exchange. As a general rule, the best viable option is to search for conceivable evidence of extortion from accessible information with the use of mathematical algorithms. The location of fraud in credit card is actually a way to identify exchanges that are counterfeit into the two classes titled real class and fraud class of transactions to detect the credit card frauds. Several machines learning algorithms such as genetic algorithms, artificial neural network with sequential itemset of data mining, relocation based migrating birds optimization, Logistics regression based comparative analytics, vector based support vector machine (SVM) algorithms, tree based decision tree and random forest algorithms are providing us a better version solution to resolve the problem of fraud credit card detections. Detection of frauds during transactions through credit card is a well-known problem but at the same time one of the difficult problems to deal with. First and foremost, being only a limited amount of information, a credit card tries to coordinate it with an instance for the dataset. Furthermore, many fragments in the truncated dataset of fraudsters would likewise conform to an authentic personal conduct standard. In addition, there are several imperatives in this issue. Initially, the collected dataset are not effectively open to the public in general and the results of investigations are often covered and controlled, making it difficult to access the results and, therefore, for fabricated models that trying to benchmark the existing problems. There is no reference to datasets in previous investigations with actual information in literatures. In addition, improving strategies is more troublesome because security concerns hinder business ideas and techniques in places of extortion, particularly in the detection of credit card fraud. Finally, informational collections are constantly evolving and profiling common and fraudulent practices that may be extortion in the past, actual transaction in the present or the other way around. This paper assesses three progress data mining approaches, logistics regression, decision tree, and random forest and then performs a mass test to assess to compute the comparative analysis among the modes on behalf of their performance.

A dataset that have the information about the transactions conducted by the credit cards are rarely accessible, exceptionally imbalanced and skewed. Ideal factor choice for the model, proper measurement is the single most important piece of data mining to assess the performance of strategies on the data of credit card frauds. There are some difficulties related to the discovery of credit card, in particular the fraudulent conduct profile is dynamic, i.e. false transactions will usually appear as genuine ones, detection of credit card frauds execution exceptionally used inspection approaches, determination of factors and is

affected by identity process used. At the end of this paper, decisions about the results of the classifier evaluation test are made and examined.

From the examinations, the result that is closed is that the logistic regression model performance accuracy is 99.5%, while the decision tree model shown the performance accuracy of 90.8%, yet the best results by computed by the Random forest model that show the performance accuracy that obtained with 99.1%. The results obtained along with these lines of performance accuracy with the dataset provided by ULB machine learning is the most accurate and highest accuracy in the issue of detection of frauds in credit card.

## 2. Review of Literature

Usual machine learning algorithms and strategies used into perform various tasks that are introduced in this section. A transiently illustrates of the pre-owned strategies in the research is listed in Table 1.

Research by (Akila & Reddy, 2018), present a group model sometimes called ensemble model titled the risk induced Bayesian inference bagging model abbreviated as RIBIB. This research proposes a three-enterprise approach: a dismissive design with a bounded pack construction strategy, a risk-driven Bayesian estimation technique as a basic student, and a weighted democratic combinatorial. Through the ensemble technique bagging is an interaction of consolidating multiple preparatory datasets and using them independently to produce different classifier models. They base their answer on Brazilian banking dataset and dominate the cost-limited contrast to other cutting-edge models.

Researcher (de Sá et al., 2018), proposes a converged version of Bayesian network classifier (BNC), which results in a computation called the Hyper-Heuristic Evolutionary Algorithm (HHEA). HHEA produces a custom BNC calculation, creating an overarching mix of modules that are critical to the current dataset. The dataset they have been used in their workflow is UOL PagSeguro, an online installment administration in Brazil. They assessed the results of proposed model using F1 scores and misrepresented a term they called monetary effectiveness, inferring the financial misfortune of the organization. They additionally use a method called case rechecking, while at the same time contrasting their results and other baselines. This is essentially re-examining (assigning more importance to them) false negatives (estimated as real but actually deceitful), as they are more important (in the opposite way) to installment organizations.

Researcher (Carcillo et al., 2019) consolidates guided and non-assisted procedures and presents a mixed approach to dealing with extortion. The rationale for using single learning is that deceitful conduct tends to turn up in the long run and the student needs to consider these misrepresentation designs. They indicate a way to deal with data discrepancy scores at different granularity levels. One of them is the worldwide

granularity, in which all instances of exchanges in a worldwide transport are considered. The second is neighborhood granularity, where exception scores are registered independently for each charge card. Ultimately, the group granularity lies in between the worldwide and neighborhood granularity, where client conduct, for example, is considered a measure of cash spent over the past 24 hours. To complete their classifier model, they used Balanced Random Forest (BRF) calculations. The peak precision and AUCPR (Area Under Accuracy Review Turn) measurements are used for evaluation.

<b>Researcher Publication details</b>	<b>Applied methodology algorithm/</b>	<b>Used Dataset</b>	<b>Models derived variables</b>	<b>Model Performance/ Evaluation metrics</b>
(Akila & Reddy, 2018)	Ensemble technique Risk Induced Bayesian Inference Bagging (RIBIB)	✓		Cost-Based - FPR, FNR, TNR, TPR, Recall, AUC
(de Sá et al., 2018)	Bayesian Network classifier used HHE Algorithm	✓		F1-score, and loss of economy
(Nami & Shajari, 2018)	Dynamic Random Forest (DRF), Least Threat Model with K-Nearest Neighbor (KNN)	✓	✓	Recall, Precision, F1-Score, Model Accuracy
(Deshe et al., 2018)	Logistics Regression (LR), Naïve Bayes (NB), Decision Tree (DT), Random Forest (RF), and Artificial Neural Network (ANN)	✓		Model false positive rate (FPR) and false negative rate (FNR)
(Carcillo et al., 2019)	K-means clustering, Balanced Random Forest	✓		AUC-PR, and model Precision
(Kim et al., 2019)	Logistic Regression (LR), Decision Tree (DT), Recurrent Neural Network (RNN), Convolutional Neural Network	✓		K-S Statistics, AUROC, Alert Rate, Accuracy, Recall, Cost Reduction

	(CNN)			Rate
(Dornadula & Geetha, 2019)	Least Squares Regression Model (LSM), Naïve Bayes (NB), Decision Tree (DT), Support Vector Machine (SVM)	✓	✓	Accuracy, Precision, Mathews Correlation Coefficient (MCC)
(Saia et al., 2019)	Prudential Multiple Consensus Model, Ensemble of Multi-Layer Perceptron, Random Forest (RF), Gaussian Naive Bayes (GNB), Adaptive Boosting (AB), Gradient Boosting (GB)	✓		Model sensitivity, model specificity, AUC, information miss-rate, Model Fallout
(Askari & Hussain, 2020)	Intuitionistic Fluent Logic (IFL) using Decision Tree Technique	✓	✓	Model-Sensitivity, Model-Specificity, False Negative Rate (FNR), False Positive Rate (FPR), precision value, model accuracy and F1-score
(Lucas et al., 2020)	Hidden Markov Models (HMMs) using Random Forest Classifiers technique	✓	✓	Model Precision-Recall, and AUC
(Misra et al., 2020)	Multilayer perceptron, Logistics regression (LR), k-nearest neighbor (k-NN) classification, Deep autoencoder for dimension reductions	✓		Model F1-Score

Table 1: Investigation of various actions for the discovery of fraud detection

Researcher (Kim et al., 2019) contrast two different methodologies with a place of misrepresentation: hybrid group strategies and deep learning. These exams are conducted in a system called Champion-

Challenger Investigation. The Hero model is a model that has been used for some time that uses various AI classifiers, for example, choice trees, strategic relapse, and basic neural organization. Every strategy is crafted using different examples and highlights and the best results are physically selected by experts. (Kim et al., 2019) However the Challenger system uses a late deep learning design that includes convolutional neural organization, repetitive neural organization, and their variations. This structure attempts every advanced deep learning design, determining activation capacity, dropout rate, and expenses, tracking the best hyper-parameters at that point. It uses the initial halving at that point in time as a way to skip preparation when no further corrections on the approval set have been completed. In the end, it picks the best performing models and saves the hyper-parameters and complex settings used on them, and tries to find better up-and-comers from previous tests. (Kim et al., 2019)

Various evaluation measurements are used to think of each structure: K-S Insight: Maximum value of contrast between two oscillations, AUROC: Area under receiver operating characteristics, Plot of true positive rate over false positive rate (FPR). Rate of alerts: Exchanges are alerted, generally speaking, when the client has given an alarm worth removing. Accuracy: Expected false exchange as clear (TP) in normal alarm exchange (TP + FP). Review: False exchanges are anticipated as clear (TP) on fake exchanges (TP + FN). Cost Reduction Rate: The cost of missed cheats (FN) gets added to the organization owning the exchange. It is determined by using the amount of the amount withdrawn as FN. Accordingly, the challenge system that relies on deep learning clearly outperformed the protagonist structure.

Researcher (Askari & Hussain, 2020) promoted a choice tree using intuitionistic drunken logic. He argues that it assumes the proper properties of exchanges so that real people are not caught by misrepresentation or the other way around. The motivation that narcissists use is that not as much effort is put on e-value-based fraud detection as other artificial intelligence methodologies.

The C 4.5 calculation is used in conjunction with fluffy logic and intuitionistic fluffy logic and the final computation is named IFDTC 4.5. Fluffy tests are characterized by various traits and the data acquisition ratio is determined using enrollment degree and non-participation degree. After that the available data is used to calculate an intuitionistic narcissistic argument that groups misrepresentations, exchanging common and ambiguous. To assess the final model, practically all appropriate measurements are used: sensitivity, discoverability, false negative rate, false positive rate, accuracy, and F1-score. They show that the proposed strategy outperforms existing strategies. Likewise, this calculation claims to work more productively and contrast quickly with others. (Askari & Hussain, 2020)

Researcher (Pourhabibi et al., 2020) survey diagram based anomaly space between the years 2007–18. They declare that the overall methodology is designing the right elements and setting the charts in one

element space to assemble the machine learning model. They also argue that chart-based uniqueness detection strategies have been on the rise since 2017. (Pourhabibi et al., 2020)

As it can very well be guessed, credit card exchanges are not free opportunities that have been missed; All things considered, they are a succession of exchanges. Researcher (Lucas et al., 2020) consider this asset and use the Hidden Markov Model (HMM) to plan a currency exchange for its past exchanges, to separately predict highlights, and to search for misrepresentation through those highlights. Go with a random forest classifier. Highlights made by HMM evaluate how much the succession is comparable to a single cardholder or previous group of terminals. They assessed the final model using the Precision-Recall AUC metric and showed that component designing with HMM presents a satisfactory climb in PR-AUC scores.

Researcher (Misra et al., 2020) proposes a two-stage model for the detection of credit card misrepresentation. Initially, an autoencoder is used to reduce the measurement with the goal that the exchange credits are reduced to the measurement that includes the vector. At that point, the final component vector is sent as information from a regulated classifier. Autoencoder is a type of feed-forward neural organization. It contains the same information and yield measurements as the regular one, yet a remaking stage exists in between. First, there is an encoder that changes the contribution to a smaller measure, at which time the encoder tries to grow the yield layer with the same measure as the information layer. That progress is made by the decoder. In this work, only the encoder part of the auto-encoder is used. Thus, the yield from the encoder is used as a contribution to various classifiers: multi-layer perceptron, k-nearest neighbor, and strategic relapse. The F1 score is used to assess the final classifier. It beats comparable technologies up to F1-score..(Misra et al., 2020)

Researcher (Dornadula & Geetha, 2019) check that card exchanges are not regular as in previous exchanges by a similar cardholder. In this way, the main gathering of cardholders depends on their exchange amount: the high, medium and low access portions. After a while, they use the sliding window strategy to isolate some of the additional highlights that depend on these gatherings. Then, SMOTE (Synthetic Minority Over-Sampling Technique) activity is performed on the dataset to take care of the weirdness dataset problem. The accuracy and MCC (Matthews's correlation coefficient) measures are used to assess the model. Among the different classifiers, the logistics regression, decision tree and random forest models rely heavily on the evaluation measure.(Dornadula & Geetha, 2019)

Researcher (Saia et al., 2019) Unite state-of-the-art arrangement computation with a model called Prudential Multiple Consensus. The idea is based on the fact that the results of different classifiers are not equivalent for specific exchanges. The calculation is done in two steps:

1. An exchange is seen as real if current calculations show it as real and the group's probability is higher than normal, all else being equal. Otherwise, it is seen as false.
2. After all the calculations have been done in the preliminary stage and a final choice is made, majority democracy is implemented.

Model effectiveness/sensitivity, fallout and AUC evaluation measurements are used to assess the model as a mixture of different computations, for example, multi-layer perceptron, random forest, Gaussian Naive Bayes, adaptive boosting, and gradient boosting. It performs well in terms of sensitivity and AUC.

Researcher (Nami & Shajari, 2018) offer a two-state answer to this issue. Before starting the calculation phase, they set out a few additional highlights to get an advanced understanding of the cardholders' spending conduct. At that point, at the core level, thinking that cardholders' new attitudes would be closer to their new mindset, another comparability measure dependent on exchange timing is developed. This action usually gives more weight to exchanges that are late. The subsequent step involves creating a dynamic random forest computation applying the base danger model. It is a model for determining the end result of an exchange with an expensive method. (Nami & Shajari, 2018) tried their model using a variety of measurements such as review accuracy, F-score specificity and accuracy, and showed that the base hazard approach expanded performance.

Researcher (Deshe et al., 2018) join multiple machine learning computation with client motivators. He argues that there should be alternative confirmations to obtain more accurate results. Optional checks can be applied to specific exchanges that exceed the edge value. They indicate the procedures and their conditions as to the benefits to be provided to the retailers, card backers and buyers that result in the "Mutual Benefit Win" achievement. Current methods for the most part are without anticipation (sitting idle), and are using AI processes for all exchanges. The third method is to conduct a second check with the client promoters. They explore a variety of different methods with Logistic Regression, Decision Trees, Naive Bayes, Random Forests, and Artificial Neural Network calculations.

### **3. General Challenges of data quality management**

The Credit card extortion detection problem shares some basic difficulties to consider when performing effective AI calculations: they can be aggregated as a way of defeating the imbalanced dataset problem, designing the right elements, and making the model progressively perform the conditions.

#### **3.1. Dataset imbalance**

Practically all datasets of banks or different associations contain a large number of exchanges, and each of them shares a specific problem as far as best in class machine learning computation goes: an unbalanced dataset. The issue comes out in the way that the speed of exchange of genuine misrepresentation is



evident from all exchanges. The volume of actual exchanges terminated by Tier-1 backers each day in 2017 is 5.7m, while misrepresentation exchanges in a similar classification are 1150.(Ryman-Tubb et al., 2018) This uneven information dissemination undermines the adequacy models of machine leaning.(Japkowicz & Stephen, 2002) Therefore, designing models to recognize fraudulent exchanges requires extra vigilance and thinking. The overall known methods of dealing with the imbalanced dataset problem are classified into two: sampling strategies, and cost-based techniques.(Dal Pozzolo et al., 2017) We look at how cutting-edge research handles this issue.

Research by (Fiore et al., 2019), take care of this issue by expanding the amount of "interesting yet underrepresented" opportunities in the preparation set. They accomplish this by generating valid models using Generative Adversarial Networks (GANs) that mimic the "attractive" class model, although not as much as might be expected. From the perspective of impact efficiency rate, the classifier constructed with the help of GAN gives substantial results unlike the first classifier.

Research by (Rtayli & Enneya, 2020), express that the volume of counterfeit exchanges is an exceptionally small segment of all-out exchanges, and this imbalanced dataset addresses the problem. To tackle that problem, they use random forest classifiers to select only the relevant highlights. They use this method in the area of credit risk identification and it gives accurate results based on the measurements they use; this method can also be used in the detection of credit card misrepresentation.

Research by (Zeager et al., 2017) express that common ways of dealing with the oddity of the lose class are observing the mining class (fraudulent exchanges), the lion share class (real exchanges) and the cost-contiguous spending efficiencies are reducing. They use an oversampling approach called SMOTE, which creates engineered examples of fraudulent exchanges.

Researcher (Jurgovsky et al., 2018) present an alternative way of handling and using record level handling and under sampling to conquer awkwardness. Inside and out, they label accounts that contain a fraudulent exchange, at any rate, as "trade-off", and label accounts that do not contain a fraudulent exchange as "certified". . With a probability of 0.9, they randomly choose a certified record and an undervalued account with a probability of 0.1. The interactions are repeated over and over to form the subset of dataset i.e. training set.

Researcher (Zhu et al., 2020) recommend a methodology called Weighted Extreme Learning Machine (WELM) to take care of imbalanced dataset issues. WELM is a modified version of ELM for imbalanced datasets employing different weights for different types of tests.

### **3.2. Valuable feature extraction and engineering challenge**

The adulterated exchange data separated from the association information base is very limited. Cardholder's balance, exchange time, credit limit, exchange amount are some of them. When these ready-to-use highlights are used to prepare basic AI calculations, the presentation probably won't move between them. To make a distinction, precise element designing turns into an indisputable necessity. We'll take a look at some of the investigations involved in designing the detection of fraud through credit card transaction.

Researcher (Zhang et al., 2019) form an element designing strategy that is subject to uniformity-based conduct examination, expressing that conduct checks should be conducted at particular gatherings of exchanges with similar exchange properties. These properties can be extracted from the data for time, geographic location, exchange amount, and exchange repetition. For each trademark discovered, two techniques are devised to incorporate designing: transaction aggregate and rules-based systems.

Researcher (Roy et al., 2018) broaden benchmark highlights and add new highlights to their model: frequency of exchanges per month, filling in missing information spurious factors, most extreme, average approval amount over an 8-month time frame, New Factors to Show When in a predefined area such as eateries, corner stores, and so forth, another variable determines whether the exchange amount at a retailer is signed by more than 10% of the standard deviation of the mean of the actual exchanges that retailer.

Researcher (Chouiekh & Haj, 2018) use convolutional neural networks in misrepresentation recognition tests and argue that since deep learning computation uses deep engineering inside, it can be performed with search layers from base to top. In different level paths naturally separate your highlights. Around there, one element is designing interactions that avoid the burning of time and assets.(Chouiekh & Haj, 2018)

Researcher (Wu et al., 2019) are eager for visas as opposed to exchanges to include designing. They are essentially centered on the Visa cash-out issue. This is an extortion method that burns through all the cutoff points on Credit card. The exam incorporates additional highlights into the model by acknowledging data from industry experts, tips, reports, and news shared by fraudsters on the web. The number of completed points that appear in the exam is 521, which forms a pool to contain the selection check. The classifier model built using these capabilities increases the accuracy performance by 4.6%-8.1%.

### **3.3. Issues during working on a real time environment**

Since exchanges designed to be consistent frameworks are unnecessary, and the practices of cardholders and fraudsters can change quickly, the order model must be rehabilitated habitually. This highlighted the

theme of how effective the models created are. We find that some investigations were done that attempted to drive a productive framework to work sustainably.

Researcher (Carcillo et al., 2018) use open source giant information tools, for example, Apache Spark, Kafka and Cassandra, to create persistent misrepresentation locators called Scalable Real-Time Fraud Finders (SCARFF). They underline that the framework has been widely tested in terms of adaptability, effectiveness and accuracy; And it measures 200 exchanges per second, which he argues is significantly more than his ally, with a speed of 2.4 exchanges per second.(Carcillo et al., 2018)

Researcher (Patil et al., 2018) recommend a misrepresentation recognition framework on Credit card that is constantly scrutinizing incoming exchanges. It uses Hadoop technology that used to encode information in HDFS design and the SAS framework transforms records into raw information. Crude information is taken into practical models to fabricate the information model. That cycle assists the framework with learning the model without the help of anyone else in an adaptable and ongoing manner.

#### 4. Proposed Methodology

The methods proposed in this paper have been used to identify extortion in the credit card framework. Tested for various machine learning algorithm like Logistic Regression, Decision Tree, Random Forest to find which calculation is most appropriate and can be adjusted with Credit card merchants to recognize extortion. Figure 1 shows the engineering chart for addressing the structure of the common structure.

The preparation steps are examined in Table 1 to identify the best count for the given dataset.

---

**Algorithm:**

---

Step 1	Read the dataset.
Step 2	To accommodate this random sampling is performed on the informative index named balancing the dataset values.
Step 3	Split the dataset into two sections, i.e. Train dataset and the Test dataset.
Step 4	Apply the feature options to the proposed models.
Step 5	Accuracy and performance measurements have been determined to know the effectiveness of various calculations.
Step 6	Then retrieve the best count dependent on the effectiveness for the given dataset.

---

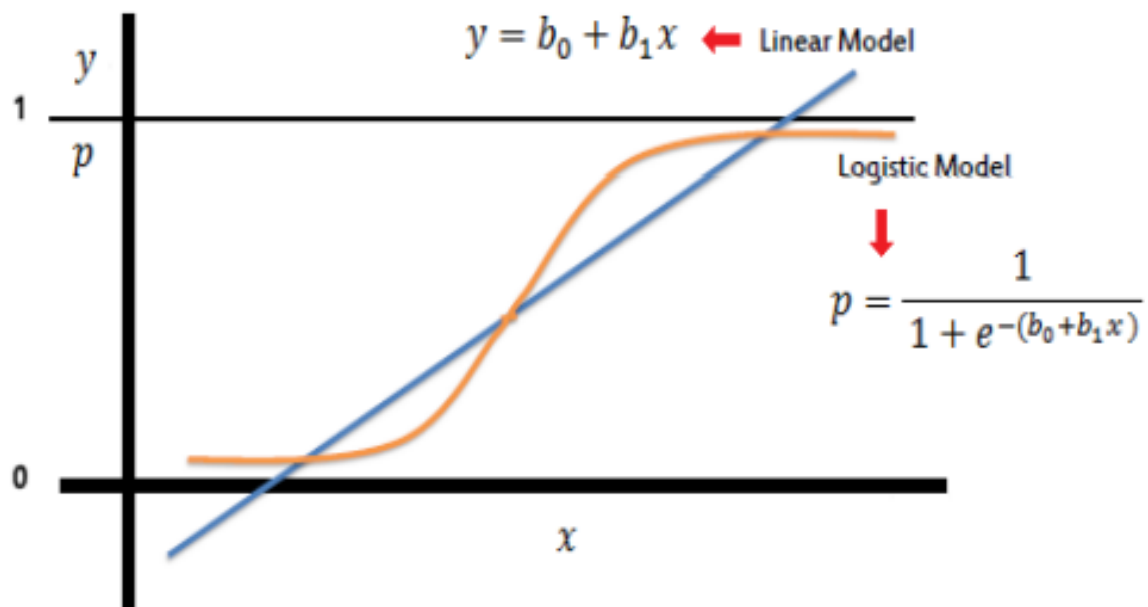
#### 4.1. Logistics Regression Model

A famous supervised learning classification technique named Logistic regression (LR) that takes advantage of the probability of parallel ward variables anticipated from the free factorization of the dataset, strategic recursion predicts the probability of an outcome that has two attributes either zero or one, yes or no, and bogus or valid. Calculated recursion has a similarity to direct recursion, although as indirect recursion acquires a straight line, strategic recursion shows a turning point. The use of one or a few indicators or the free factor depending on what the forecast is based on; Strategic relapse produces calculated turns that plot property somewhere in the range of one and the other.

LR is a regression model where the dependent variable is eliminated and the relationship between the various autonomic factors is broken. There are several types of computed relapse models such as the dual strategic model, the various computed models, and the binomial calculated model. Binary logistic regression models are used to assess the probability of a double response dependent on at least one indicator.

$$p = \frac{1}{1 + e^{-(b_0 + b_1 x)}}$$

The above situation addresses the calculated mathematical formula of logistics regression.



**Figure 1: Logistic curve**

Graph in figure 1 shows the difference between a linear regression and a logistics regression, where the logistics regression shows a turning point, although the linear regression refers to a straight line.

## 4.2. Decision Tree Model

A decision tree is an algorithm that uses a tree-like chart or model of options and their possible outcomes to estimate the final option, this calculation uses a restrictive control explanation. A decision tree is a computation for moving towards discrete-honed objective abilities, in which the choice tree is represented by a learned ability. For inductive learning, these types of calculations are highly acclaimed and have been effectively applied to a wide range of tasks. We give another exchange a mark as to whether it is genuine or misrepresentation for which the class name is ambiguous and a subsequent attempt is made to honor the exchange against the tree of choice, and then yield from the root hub / that the class name is done away with for exchange.

The principles of choice decide the outcome of the substance of the leaf hub. Overall the guidelines have the type of 'if condition 1 and condition 2 but not condition 3'. A choice tree helps improve understanding and understanding of what is most feared, best, and expected for different situations, and allows for the expansion of new possible situations.

$$Entropy(S) = \sum_{i=1}^c -p_i \log_2 p_i$$

$$Gain(S, A) = Entropy(S) - \sum_{v \in Value(A)} \frac{|S_v|}{|S|} Entropy(S_v)$$

The step to settle on a decision tree is to calculate the entropy of each feature using the dataset in issue right off the bat, splitting the dataset into subsets using the feature for which the gain is most extreme OR entropy is then at least an decision tree containing that property at the hub and in conclusion, recursion is performed on the subset using the remaining properties to form an decision tree.

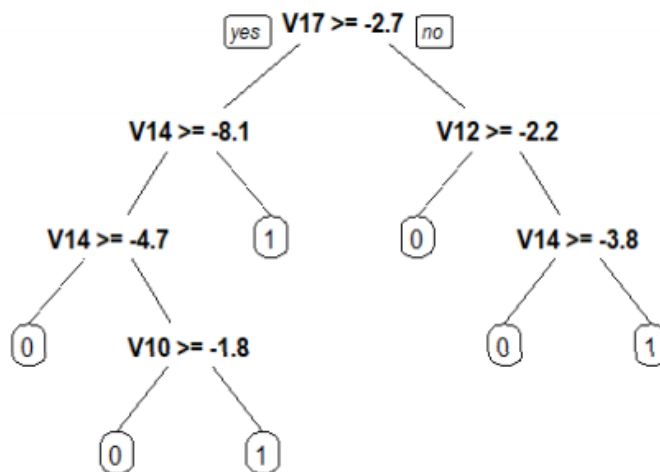


Figure 2: Decision Tree

### 4.3. Random Forest

There is a algorithm for random forest arrangement and relaxation. Immediately, this choice is the classification of the tree classifier. Arbitrary woodland enjoys an upper hand over the tree of choice because it accommodates the tendency for over-fitting for their preparation set. A subset of the preparation set is arbitrarily inspected to construct each tree and subsequently construct a choice tree; each hub at that point on a component is chosen from an arbitrary subset of the full list of capabilities. In any event, preparation is incredibly fast in unregulated woodland and because each tree is prepared independently of the others, for vast informational indexes with many highlights and information opportunities. Random forest calculations have been found to give a good gauge of speculative error and to be impervious to over-fitting.

$$P(c|x) = \frac{P(x|c)P(c)}{P(x)}$$

Where,  $P(c|x)$ : represents the posterior probability,  $P(x|c)$ : represents the likelihood,  $P(c)$ : represents the class of prior probability,  $P(x)$ : represents the predictor prior probability.

$$P(c|X) = P(x_1|c) \times P(x_2|c) \times P(x_3|c) \times \dots \times P(x_n|c) \times P(c)$$

Random forest has the importance of factors in the issue of relapse or arrangement should be made possible by the random forest in a specific way.

## 5. Experiments & Results

To begin with, the credit card dataset is taken from the source, and cleanup and approval are performed on the dataset which includes additional extractions, capturing the blanks in segments, transforming critical variables into components or sections, then splitting the information into two segments. is split, one is preparing the dataset and the other is the test informational index. Currently, K overlap cross-approval is performed which is the first instance that is randomly divided into k equivalent estimated sub-samples. Out of k sub-samples, a single sub-sample is kept as approval information for testing the model, and the remaining k-1 sub-samples are used as information preparation, model logistic For relapse, decision tree, SVM, random forest and subsequent accuracy, impact potential, explicator, accuracy are determined and a correlation is made.

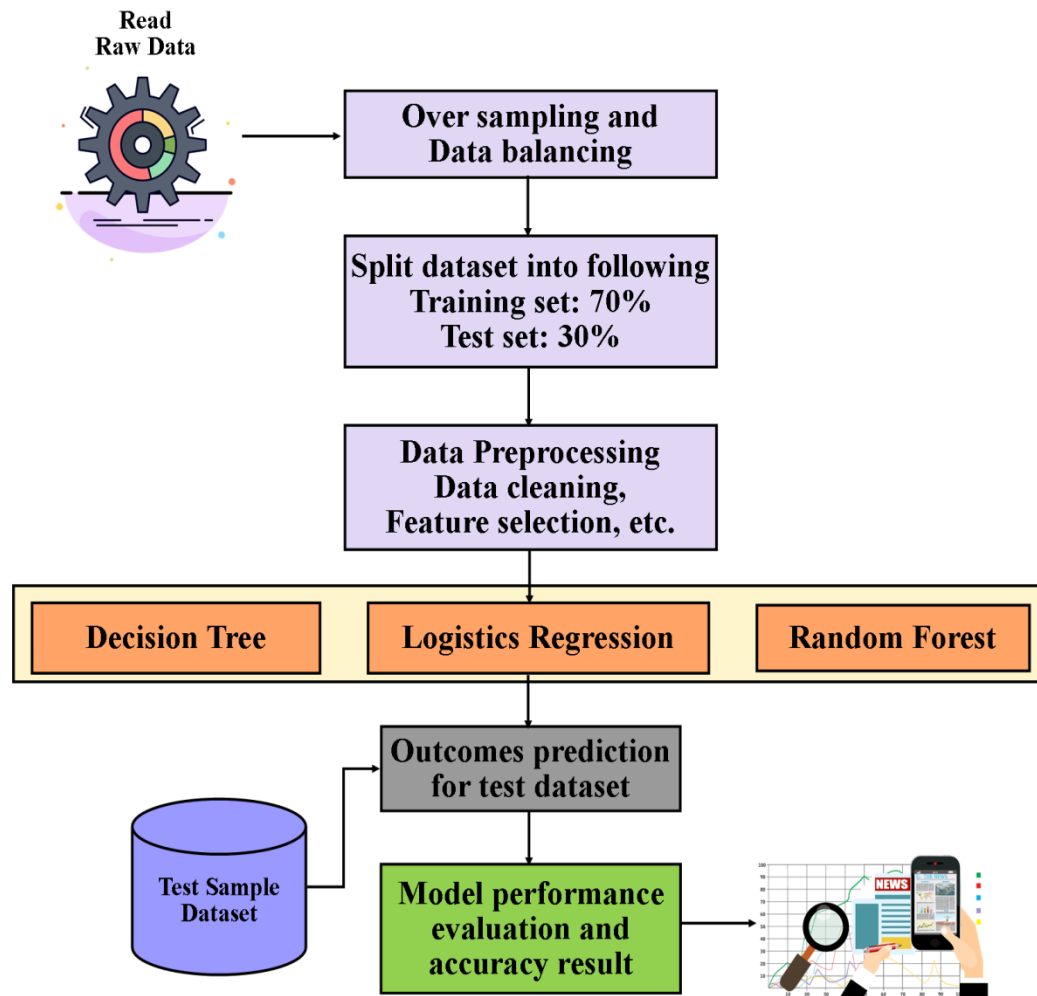


Figure 3: System Architecture

### 5.1. Performance metrics

There is an assortment of measures for different calculations and these actions are designed to measure completely different things. So it should be the standard for evaluating other proposed techniques. Bogus Positive (FP), False Negative (FN), True Positive (TP), True Negative (TN), and the relationship between them are commonly accepted with credit card misrepresentation location experts to see the accuracy of the various methods. The meanings of the limits referred to are presented below:

- **True Positive (TP):** The true positive rate addresses the bits of fraudulent exchanges effectively known as false exchanges.

$$\text{True positive} = \frac{TP}{(TP + FN)}$$

- **True Negative (TN):** The true negative rate addresses the segment of general exchanges accurately assigned to specific exchanges.

$$\text{True negative} = \frac{TN}{(TN + FP)}$$

- **Fake Positive (FP):** The false positive rate reflects the portion of non-fake exchanges falsely named fake exchanges.

$$\text{False positive} = \frac{FP}{(FP + TN)}$$

- **False Negative (FN):** The false negative rate refers to the portion of non-fraudulent exchanges incorrectly named specific exchanges.

$$\text{False negative} = \frac{FN}{(FN + TP)}$$

- **Confusion matrix:** A confusion matrix gives more understanding in the exposition of a retrospective model, yet also which classes are being predicted exactly, which is wrong, and what kind of mistakes are being made. With negative and positive squares, the least complicated dislocation grid accounts for the two-class ordering issue. In this type of disorder lattice, each cell in the table has a special and precisely known name.

Predicted	Positive	Negative
Positive	TP	FN
Negative	FP	TN

- **Accuracy:** Accuracy is the level of effectively ordered opportunities. This is probably the most commonly used characterization performance measure.

$$\text{Accuracy} = \frac{\text{Number of correct predictions}}{\text{Total Number of predictions}}$$

Or for the two fold order model. Accuracy can be formulated as:

$$\text{Accuracy} = \frac{(TP + TN)}{(TP + TN + FP + FN)}$$

- **Precision:** Precision is the amount of grouped positive or fake opportunities that are actually positive examples.

$$\text{Precision} = \frac{TP}{(TP + FP)}$$

- **Recall:** Recall is a metric that measures the amount of correct fixed expectations from all of the forecasts that may have been made. Unlike precision that a solitary commentary on some of the true expectations out of every certain forecast, reviews lack positive expectations. Reviews are



determined as the number of true positives by separating them from the total number of true positives and false negatives.

$$Recall = \frac{TP}{(TP + FN)}$$

- **F1-score:** F1-score is the weighted average of Precision and Recall. Therefore, this score takes both false positives and false negatives into account.

$$F1\ Score = \frac{2 \times (Recall \times Precision)}{(Recall + Precision)}$$

- **Support:** This performance metrics of model represents the number of tests of actual feedback that are out there. The backing is the number of actual occurrences of the class in the predefined dataset. Unbalanced help in preparation information may show primary deficiencies in the classifier's manifest scores and may require separate oversight or rebalancing. Support does not change between models but rather analyzes the evaluation cycle.

## 5.2. Results and Discussion

### 5.2.1. Importing Dataset

The dataset we will use is publically available on Kaggle titled Credit Card Fraud Detection dataset<sup>1</sup>. It contains highlights from V1 to V28 which are the main sections obtained by PCA. We will disregard the time that there is no use for assembling models. The remaining highlights include 'Sum' which contains the total of cash to be executed and 'Class' highlights in which the exchange is an extortion case.

In figure 4, we do import information using the 'read\_csv' method of panda's library of python and print the information to check it in python.

```

      V1      V2      V3      V4      V5      V6      V7  \
0 -1.359807 -0.072781  2.536347  1.378155 -0.338321  0.462388  0.239599
1  1.191857  0.266151  0.166480  0.448154  0.060018 -0.082361 -0.078803
2 -1.358354 -1.340163  1.773209  0.379780 -0.503198  1.800499  0.791461
3 -0.966272 -0.185226  1.792993 -0.863291 -0.010309  1.247203  0.237609
4 -1.158233  0.877737  1.548718  0.403034 -0.407193  0.095921  0.592941

      V8      V9      V10  ...      V21      V22      V23      V24  \
0  0.098698  0.363787  0.090794  ... -0.018307  0.277838 -0.110474  0.066928
1  0.085102 -0.255425 -0.166974  ... -0.225775 -0.638672  0.101288 -0.339846
2  0.247676 -1.514654  0.207643  ...  0.247998  0.771679  0.909412 -0.689281
3  0.377436 -1.387024 -0.054952  ... -0.108300  0.005274 -0.190321 -1.175575
4 -0.270533  0.817739  0.753074  ... -0.009431  0.798278 -0.137458  0.141267

      V25      V26      V27      V28  Amount  Class
0  0.128539 -0.189115  0.133558 -0.021053   149.62      0
1  0.167170  0.125895 -0.008983  0.014724     2.69      0
2 -0.327642 -0.139097 -0.055353 -0.059752   378.66      0
3  0.647376 -0.221929  0.062723  0.061458   123.50      0
4 -0.206010  0.502292  0.219422  0.215153    69.99      0

```

Figure 4: Dataset importing and variable exploration

### 5.2.2. Data exploration and visualization process

This part will examine the information and decide the best way to deal with the preprocessing, which will be done in the accompanying segment.

**Features/variables:** Our information in the dataset contains 28 anonymous features/variables, each of float64 information types. The dataset additionally contains information on the exchange amount, time, and class.

```
>>>
===== RESTART: C:\AMIT_PUNDIR\M.TECH_THESIS\FINAL_CODE.py =====
The dataset contains 560 rows and 31 columns.
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 560 entries, 0 to 559
Data columns (total 31 columns):
Time                560 non-null int64
V1                  560 non-null float64
V2                  560 non-null float64
V3                  560 non-null float64
V4                  560 non-null float64
V5                  560 non-null float64
V6                  560 non-null float64
V7                  560 non-null float64
V8                  560 non-null float64
V9                  560 non-null float64
V10                 560 non-null float64
V11                 560 non-null float64
V12                 560 non-null float64
V13                 560 non-null float64
V14                 560 non-null float64
V15                 560 non-null float64
V16                 560 non-null float64
V17                 560 non-null float64
V18                 560 non-null float64
V19                 560 non-null float64
V20                 560 non-null float64
V21                 560 non-null float64
V22                 560 non-null float64
V23                 560 non-null float64
V24                 560 non-null float64
V25                 560 non-null float64
V26                 560 non-null float64
V27                 560 non-null float64
V28                 560 non-null float64
Amount              560 non-null float64
Class               560 non-null int64
dtypes: float64(29), int64(2)
memory usage: 135.8 KB
Normal transactions count: 332
Fraudulent transactions count: 228
Original dataset shape Counter({1: 224, 0: 151})
Resampled dataset shape Counter({0: 227, 1: 224})
```

Figure 5: Dataset variables summary exploration

**Selection of anonymous features:** The first dataset used in this venture has also undergone dimensionality reduction with PCE in unnamed highlights. This was probably done both to anonymize the information and to make the information easier to work with.

**Time:** Time is addressed as the number of seconds since the hour of the primary exchange in the dataset.

**Class Variable:** Marks the class variable exchange as false (1) or non-cheat (0). This will fill in as our objective. We will use the highlights in the dataset to foresee the class.

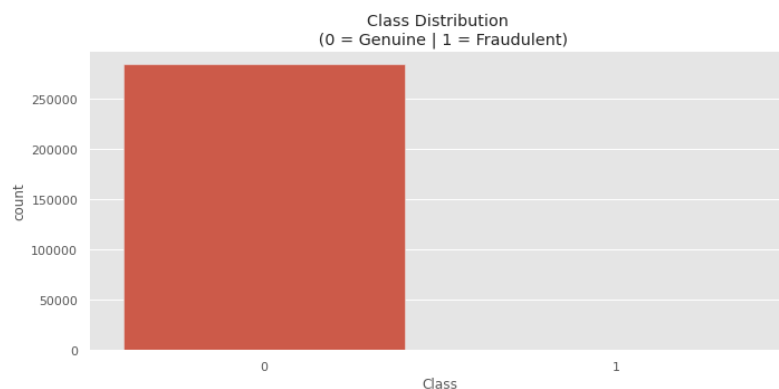
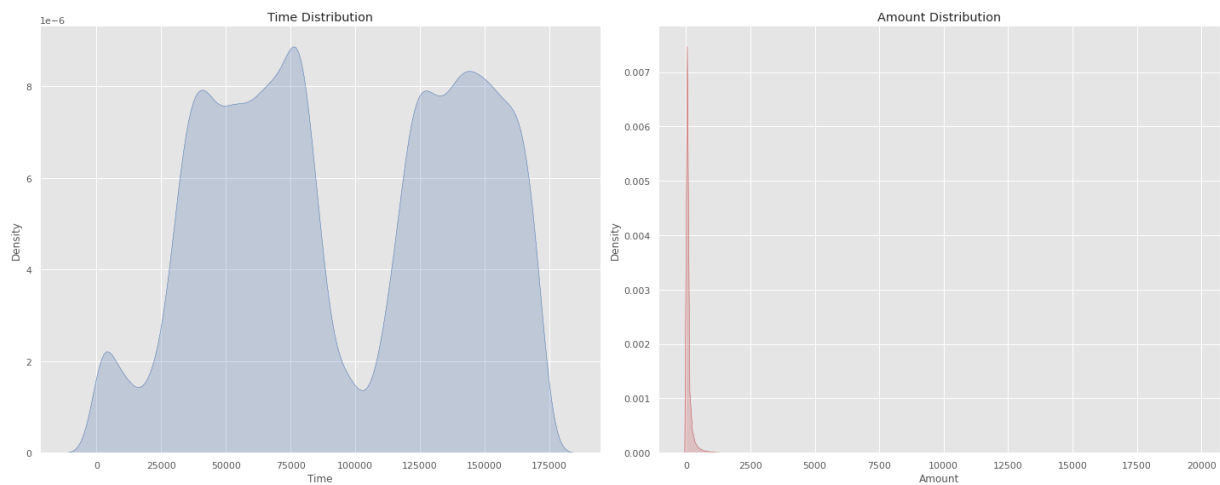


Figure 6: Dataset variable classification

### Time and amount value of transaction

The sum and the time factor are the only two highlights that have not been changed or scaled in the first dataset. Allow us to examine these highlights in the present.



**Figure 7: Distribution of dataset quality among time and amount**

Both of these variables/features have attractive appropriations.

Amount section is tilted very correctly, and we can see that although the typical exchange amount is \$88, the maximum value is over \$25,000! There are many exceptions to this component section, and we will look at expanding this element in subsequent advances.

We also see that there is an intriguing appropriation of time; however, there is no exception to this. Once again, we will zero in on scaling the variables/features after managing the skewed objective variables by inspecting/testing our dataset.

### Exploration of anonymous features

Finally, we should take a look at our 28 unnamed variables/features. Even though we can't know exactly what these variables/features mean because of the dimensionality reduction before, we can see how they are expressed.

As should be clear, each highlight has an alternate appropriation, and many have discrepancies. They all appear to be centered at 0, yet the scope of the properties differs in the variables/features.

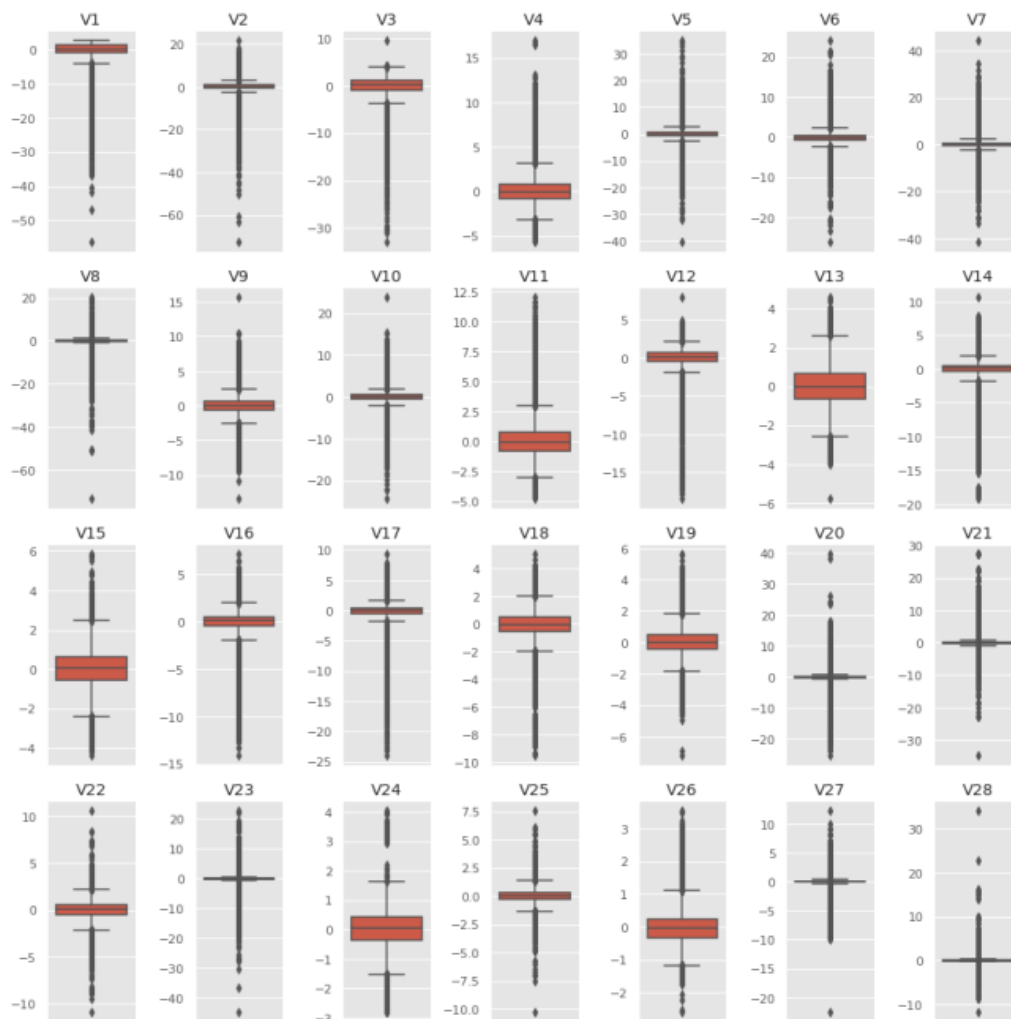


Figure 8: Dataset individual variables visualization for value exploration

### 5.2.3. Preprocessing of data

Since we have visualized our information, we will preprocess the information before boosting our algorithms of machine learning.

#### Missing values in the dataset

##### Missing Values

```
>>>
===== RESTART: C:\AMIT_PUNDIR\M.TECH_THESIS\FINAL_CODE.py =====
# Check for null values
train.isnull().values.any()
> FALSE
```

**Figure 9: Python environment dataset exploration to computing any missing value in occur in dataset**

There are no missing value properties in our dataset, so there is no need for additional activity on this

### Scale out the features

As we've seen, a significant number of our variables/features are tilted, and they have vastly different areas of properties. To address this, we will measure our information.

While the anonymous variables/features have been effectively reduced via PCA and accordingly scaling has started before PCA, we will scale them in this progression. The rationale for doing this is to move all of our variables/features to comparable scales, even if they have been scaled before. So our main concern is to guarantee that our model predictions can rely on the data conveyed by our variables/features rather than on their reach.

*shape of independent training dataset ( $X_{train}$ ): (227845, 30)*

*shape of dependent training dataset ( $y_{train}$ ): (227845,)*

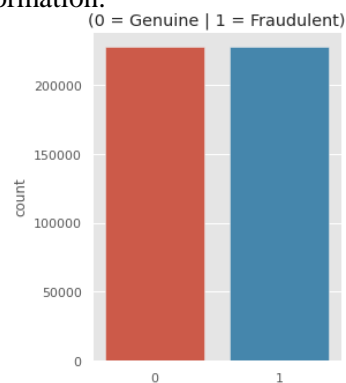
*shape of independent test training dataset ( $X_{test}$ ): (56962, 30)*

*shape of dependent test dataset ( $y_{test}$ ): (56962,)*

### Controlling the imbalance of dataset

As the objective variable is sharply tilted, presenter models may be inclined to expect that experiments are real. In any case, we need our models the option to decide whether the new exchanges are fraudulent, looking at the information variables/features. If we don't fix the weirdness of the information one way or another, our models will over-fit the actually exchanged information and are probably going to be ruthless to fraudulent exchanges.

In this investigative work, we are looking at ways to deal with amendments to unilateralism. Undersampling involves testing an unregulated subset of information with the end goal being addressed to proven and false exchanges. The fundamental issue with this method is that for our situation, we will only have ~800 information completely concentrated as there are about 400 false exchanges. This will waste a lot of our knowledge and result in models that are completely less unreliable than those that have been given a lot of preparation information.



**Figure 10: Distribution of target variable after under-sampling**

#### 5.2.4. Model evaluation results:

```

===== Decision Tree =====
Cross Validation Mean Score: 90.5%

Model Accuracy: 90.5%

Confusion Matrix:
[[227  0]
 [ 43 181]]

Classification Report:
      precision    recall  f1-score   support

     0       0.84      1.00      0.91       227
     1       1.00      0.81      0.89       224

 accuracy          0.90          451
  macro avg          0.92      0.90      0.90          451
 weighted avg          0.92      0.90      0.90          451

```

**Figure 11: Decision tree model evaluation report**

```

===== LogisticRegression =====
Cross Validation Mean Score: 97.7%

Model Accuracy: 98.8%

Confusion Matrix:
[[225  2]
 [ 3 220]]

Classification Report:
      precision    recall  f1-score   support

     0       1.00      1.00      1.00       227
     1       1.00      1.00      1.00       224

 accuracy          1.00          451
  macro avg          1.00      1.00      1.00          451
 weighted avg          1.00      1.00      1.00          451

```

**Figure 12: Logistics Regression model evaluation report**

```

===== RandomForest Classifier =====
Cross Validation Mean Score: 98.5%

Model Accuracy: 99.8%

Confusion Matrix:
[[227  2]
 [ 2 222]]

Classification Report:
      precision    recall  f1-score   support

     0       1.00      1.00      1.00       227
     1       1.00      1.00      1.00       224

 accuracy          1.00          451
  macro avg          1.00      1.00      1.00          451
 weighted avg          1.00      1.00      1.00          451

```

**Figure 11: Random Forest model evaluation report**

### 5.2.5. Model Test Results and discussion

#### Decision Tree Model Test Results

==== Decision Tree =====

Model Accuracy: 90.8%

Confusion Matrix:

```
[[77  0]
 [17 91]]
```

Classification Report:

	precision	recall	f1-score	support
0	0.82	1.00	0.90	77
1	1.00	0.84	0.91	108
accuracy			0.91	185
macro avg	0.91	0.92	0.91	185
weighted avg	0.92	0.91	0.91	185

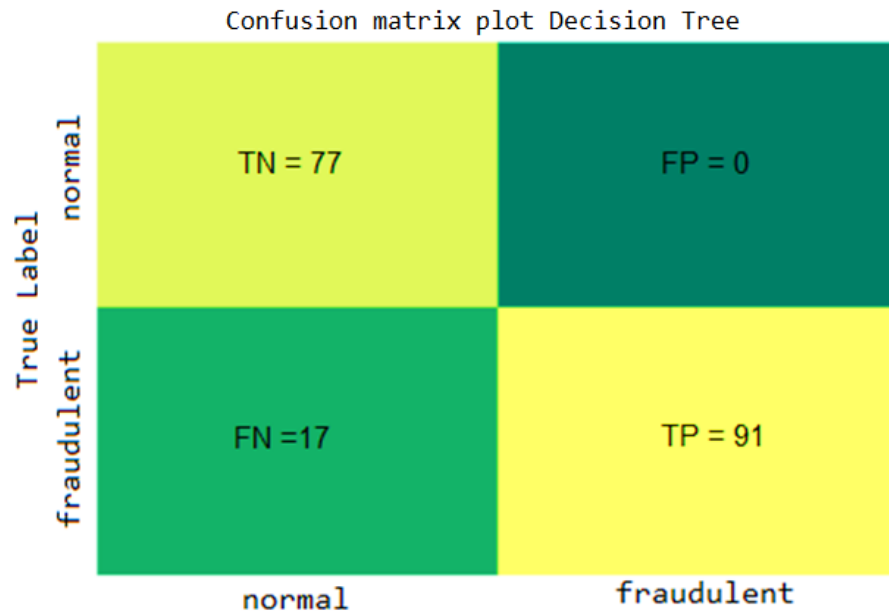


Figure 12: Confusion matrix plot of decision tree model test results

## Logistics Regression Model Test Results

```
=== LogisticRegression ===
```

```
Model Accuracy: 98.5%
```

```
Confusion Matrix:
```

```
[[ 75  2 ]
 [  2 106]]
```

```
Classification Report:
```

	precision	recall	f1-score	support
0	1.00	0.99	0.99	77
1	0.99	1.00	1.00	108
accuracy			0.99	185
macro avg	1.00	0.99	0.99	185
weighted avg	0.99	0.99	0.99	185

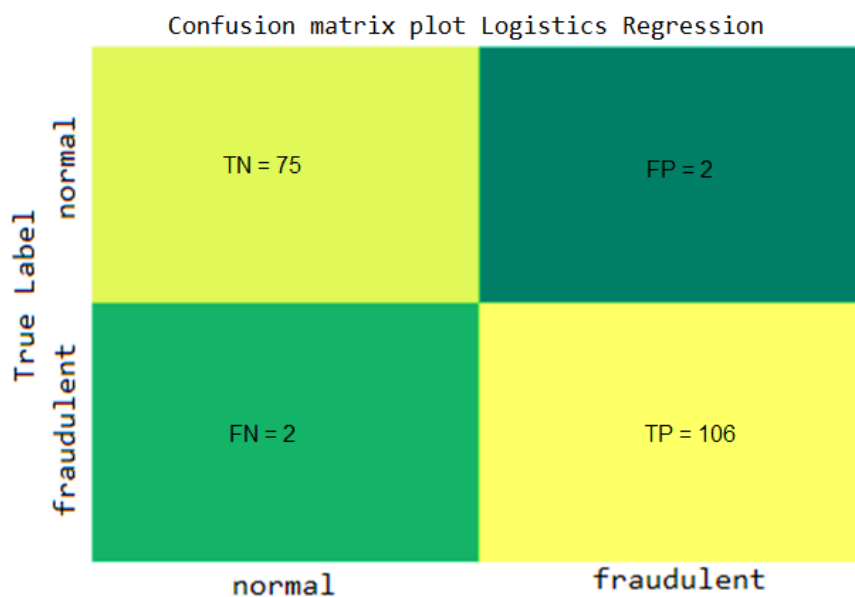


Figure 13: Confusion matrix plot of Logistics regression tree model test results



## Random Forest Model Test Results

```
===== Model Test Results =====
```

```
=== RandomForest Classifier ===
```

```
Model Accuracy: 99.1%
```

```
Confusion Matrix:
```

```
[[ 75  2 ]
 [  1 107]]
```

```
Classification Report:
```

	precision	recall	f1-score	support
0	1.00	1.00	1.00	77
1	1.00	1.00	1.00	108
accuracy			1.00	185
macro avg	1.00	1.00	1.00	185
weighted avg	1.00	1.00	1.00	185

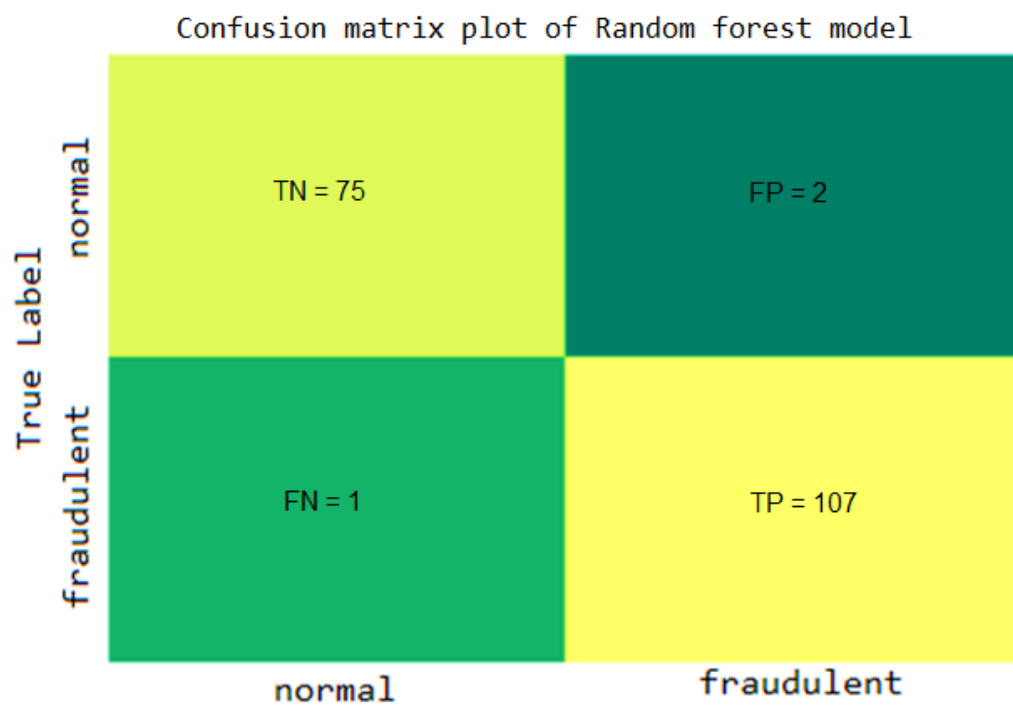


Figure 14: Confusion matrix plot of Random Forest model test results

## ROC curve of test models

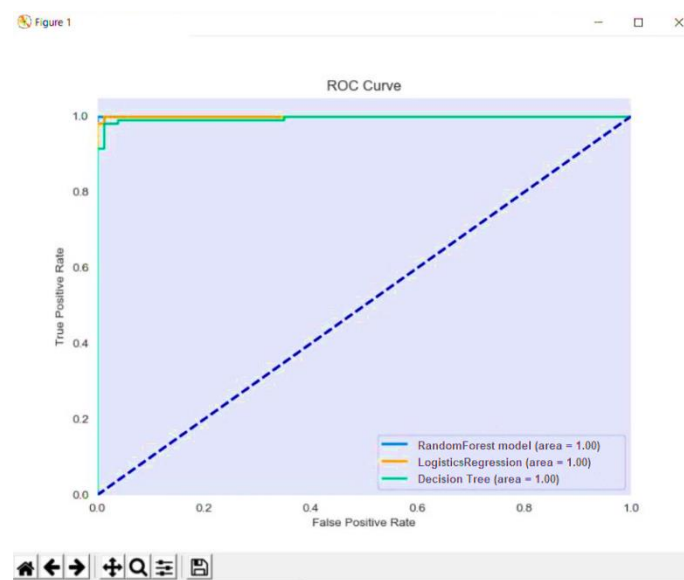


Figure 15: ROC curve of all three models

## 6. Conclusion

In this research paper, we considered the use of machine learning algorithms such as Logistics regression, Decision Trees, Random forest to enhance data quality by implementing information pre-processing steps and showed that it is effective in cutting out spurious exchange and fake alarms. Supervised learning algorithm are one of the novel and significant way to as pre the literature review in terms of their arability and covering applications areas. If these algorithms are applied to the Bank credit card fraud detection and prediction for security framework, the possibility of fraud during money transmission after credit card exchanges can be predicted earlier to secure the customer accounts and services. In addition, advances in anti-fraud techniques can be adopted in a hostile manner to protect banks from unrelenting misfortunes and reduce misclassification cost as well. The objective of this research aimed specifically at contrasting general order issues in which we had a variable misclassification cost. Model performance metrics including accuracy, precision, recall, F1-score, support, and accuracy are used to assess the presence of the proposed framework. In contrast to each of the three strategies, we found that the random forest classifier with the boosting strategy is superior to the logistic regression and decision tree model.

## References

- [1] Akila, S., & Reddy, U. S. (2018). Cost-sensitive Risk Induced Bayesian Inference Bagging (RIBIB) for credit card fraud detection. *Journal of Computational Science*, 27, 247–254.
- [2] Askari, S. M. S., & Hussain, M. A. (2020). IFDTC4. 5: Intuitionistic fuzzy logic based decision tree for E-transactional fraud detection. *Journal of Information Security and Applications*, 52, 102469.

- [3] Carcillo, F., Dal Pozzolo, A., Le Borgne, Y.-A., Caelen, O., Mazzer, Y., & Bontempi, G. (2018). Scarff: a scalable framework for streaming credit card fraud detection with spark. *Information Fusion*, 41, 182–194.
- [4] Carcillo, F., Le Borgne, Y.-A., Caelen, O., Kessaci, Y., Oblé, F., & Bontempi, G. (2019). Combining unsupervised and supervised learning in credit card fraud detection. *Information Sciences*.
- [5] Chouiekh, A., & Haj, E. L. H. I. E. L. (2018). Convnets for fraud detection analysis. *Procedia Computer Science*, 127, 133–138.
- [6] Dal Pozzolo, A., Boracchi, G., Caelen, O., Alippi, C., & Bontempi, G. (2017). Credit card fraud detection: a realistic modeling and a novel learning strategy. *IEEE Transactions on Neural Networks and Learning Systems*, 29(8), 3784–3797.
- [7] de Sá, A. G. C., Pereira, A. C. M., & Pappa, G. L. (2018). A customized classification algorithm for credit card fraud detection. *Engineering Applications of Artificial Intelligence*, 72, 21–29.
- [8] Deshe, W., Chen, B., & Chen, J. (2018). *Credit Card Fraud Detection Strategies with Consumer Incentives*. ELSEVIER.
- [9] Dornadula, V. N., & Geetha, S. (2019). Credit card fraud detection using machine learning algorithms. *Procedia Computer Science*, 165, 631–641.
- [10] Fiore, U., De Santis, A., Perla, F., Zanetti, P., & Palmieri, F. (2019). Using generative adversarial networks for improving classification effectiveness in credit card fraud detection. *Information Sciences*, 479, 448–455.
- [11] Japkowicz, N., & Stephen, S. (2002). The class imbalance problem: A systematic study. *Intelligent Data Analysis*, 6(5), 429–449.
- [12] Jurgovsky, J., Granitzer, M., Ziegler, K., Calabretto, S., Portier, P.-E., He-Guelton, L., & Caelen, O. (2018). Sequence classification for credit-card fraud detection. *Expert Systems with Applications*, 100, 234–245.
- [13] Kim, E., Lee, J., Shin, H., Yang, H., Cho, S., Nam, S., Song, Y., Yoon, J., & Kim, J. (2019). Champion-challenger analysis for credit card fraud detection: Hybrid ensemble and deep learning. *Expert Systems with Applications*, 128, 214–224.
- [14] Lucas, Y., Portier, P.-E., Laporte, L., He-Guelton, L., Caelen, O., Granitzer, M., & Calabretto, S. (2020). Towards automated feature engineering for credit card fraud detection using multi-perspective HMMs. *Future Generation Computer Systems*, 102, 393–402.
- [15] Misra, S., Thakur, S., Ghosh, M., & Saha, S. K. (2020). An autoencoder based model for detecting fraudulent credit card transaction. *Procedia Computer Science*, 167, 254–262.

- [16] Nami, S., & Shajari, M. (2018). Cost-sensitive payment card fraud detection based on dynamic random forest and k-nearest neighbors. *Expert Systems with Applications*, 110, 381–392.
- [17] Patil, S., Nemade, V., & Soni, P. K. (2018). Predictive modelling for credit card fraud detection using data analytics. *Procedia Computer Science*, 132, 385–395.
- [18] Pourhabibi, T., Ong, K.-L., Kam, B. H., & Boo, Y. L. (2020). Fraud detection: A systematic literature review of graph-based anomaly detection approaches. *Decision Support Systems*, 133, 113303.
- [19] Roy, A., Sun, J., Mahoney, R., Alonzi, L., Adams, S., & Beling, P. (2018). Deep learning detecting fraud in credit card transactions. *2018 Systems and Information Engineering Design Symposium (SIEDS)*, 129–134.
- [20] Rtayli, N., & Enneya, N. (2020). Selection features and support vector machine for credit card risk identification. *Procedia Manufacturing*, 46, 941–948.
- [21] Ryman-Tubb, N. F., Krause, P., & Garn, W. (2018). How Artificial Intelligence and machine learning research impacts payment card fraud detection: A survey and industry benchmark. *Engineering Applications of Artificial Intelligence*, 76, 130–157.
- [22] Saia, R., Carta, S., Reforgiato Recupero, D., & Fenu, G. (2019). Fraud Detection for E-commerce Transactions by Employing a Prudential Multiple Consensus Model. *Journal of Information Security and Applications*, 46. <https://doi.org/10.1016/j.jisa.2019.02.007>
- [23] Wu, Y., Xu, Y., & Li, J. (2019). Feature construction for fraudulent credit card cash-out detection. *Decision Support Systems*, 127, 113155.
- [24] Zeager, M. F., Sridhar, A., Fogal, N., Adams, S., Brown, D. E., & Beling, P. A. (2017). Adversarial learning in credit card fraud detection. *2017 Systems and Information Engineering Design Symposium (SIEDS)*, 112–116.
- [25] Zhang, X., Han, Y., Xu, W., & Wang, Q. (2019). HOBA: A novel feature engineering methodology for credit card fraud detection with a deep learning architecture. *Information Sciences*.
- [26] Zhu, H., Liu, G., Zhou, M., Xie, Y., Abusorrah, A., & Kang, Q. (2020). Optimizing Weighted Extreme Learning Machines for imbalanced classification and application to credit card fraud detection. *Neurocomputing*, 407, 50–62.