Project Title: **Credit Card Fraud Detection System**

**Information Security Management**
**CSE3502**
**J Component**
**Review - 3**
**Submitted by:**
**Hrithik Ahuja   (18BCE2154)**
**Ayush Rana (18BCE2305)**
**Rakshith Sachdev (18BCI0109)**
**Rohan Allen (18BCI0247)**

Under the Guidance of:
**Prof. AMUTHA PRABHAKAR M.**

**TITLE: CREDIT CARD FRAUD DETECTION SYSTEM**

**DEMO VIDEO LINK:**

https://drive.google.com/file/d/1cI3vY6sR9snj_7c0qqLlhyjkivy0Pm3P/view

**TABLE OF CONTENTS:**

| 3. | **BRIEF ON ISOLATION FOREST ALGORITHM** |
|---|---|

| | |
|---|---|
| 4. | **BRIEF ON LOF ALGORITHM** |
| 5. | **LITERATURE SURVEY** |
| 6. | **PROPOSED METHODOLOGY** |
| 7. | **IMPLEMENTATION** |
| 8. | **PROJECT OUTCOME** |

## Abstract:

The main purpose of this project is understand and implement the distinct approach of Isolation forest algorithm and LOF algorithm to identify fraudulent transactions in a database instead of using generic Random Forest approach. The model will be able to identify the transactions with greater accuracy and by comparing both the approaches we will be working towards a more optimal solution. The Credit Card Fraud Detection Problem includes modelling past credit card transactions with the knowledge of the ones that turned out to be fraud. This model is then used to identify whether a new transaction is fraudulent or not. The aim of the project here is to detect 100% of the fraudulent transactions while minimizing the incorrect fraud classifications.

## INTRODUCTION:

The following are reasons why we need to develop a robust system to detect fraudulent transactions. The following challenges have to be overcome to make sure that the final product is resilient.

The challenge is to recognize fraudulent credit card transactions so that the customers of credit card companies are not charged for items that they did not purchase. Main challenges involved in credit card fraud detection are:

- Enormous Data is processed every day and the model build must be fast enough to respond to the scam in time.
- Imbalanced Data i.e. most of the transactions (99.8%) are not fraudulent which makes it really hard for detecting the fraudulent ones Data availability as the data is mostly private.
- Misclassified Data can be another major issue, as not every fraudulent transaction is caught and reported.
- Adaptive techniques used against the model by the scammers.

## BRIEF ON ISOLATION FOREST ALGORITHM:

- Isolation forest is an unsupervised learning algorithm for anomaly detection that works on the principle of isolating anomalies, instead of the most common techniques of profiling normal points.
- In statistics, an anomaly (a.k.a. outlier) is an observation or event that deviates so much from other events to arouse suspicion it was generated by a different mean.
- Anomalies in a big dataset may follow very complicated patterns, which are difficult to detect "by eye" in the great majority of cases. This is the reason why the field of anomaly detection is well suited for the application of Machine Learning techniques.
- The most common techniques employed for anomaly detection are based on the construction of a profile of what is "normal": anomalies are reported as those instances in the dataset that do not conform to the normal profile.
- Isolation Forest uses a different approach: instead of trying to build a model of normal instances, it explicitly isolates anomalous points in the dataset. The main advantage of this approach is the possibility of exploiting sampling techniques to an extent that is not allowed to the profile-based methods, creating a very fast algorithm with a low memory demand

## BRIEF ON LOF (Local Outlier Factor Algorithm) ALGORITHM:

The Local Outlier Factor (LOF) algorithm is an unsupervised anomaly detection method which computes the local density deviation of a given data point with respect to its neighbours. It considers as outliers the samples that have a substantially lower density than their neighbours.

The number of neighbors considered (parameter n_neighbors) is typically greater than the minimum number of samples a cluster has to contain, so that other samples can be local outliers relative to this cluster, and 2) smaller than the maximum number of close by samples that can potentially be local outliers. In practice, such informations are generally not available, and taking n_neighbors=20 appears to work well in general.

**ALGORITHMS USED TO SOLVE THE PROBLEM OF ANAMOLY DETECTION UNTIL NOW:**

1. Distributed Data Mining
2. Artificial Neural Networks
3. Neural Data mining
4. Hidden Markov Models
5. Bayesian neural networks
6. Computational Intelligence
7. Dempster–Shafer theory and Bayesian learning
8. Parallel granular neural networks
9. Phone usage patterns
10. Transaction aggregation methodology
11. Feature engineering

**ALGORITHMS USED TO SOLVE THIS PROBLEM IN THIS PROJECT:**

1. Isolation forest algorithm
2. Local Outlier Factor

**LITERATURE SURVEY:**

1. **Chan, P. K., Fan, W., Prodromidis, A. L., & Stolfo, S. J. (1999). Distributed data mining in credit card fraud detection.** *IEEE Intelligent Systems and Their Applications*, *14*(6)

   Large scale data-mining techniques can improve the state of the art in commercial practice. Scalable techniques to analyze massive amounts of transaction data that efficiently compute fraud detectors in a timely manner is an important problem, especially for e-commerce. Besides scalability and efficiency, the fraud-detection task exhibits technical problems that include skewed distributions of training data and non- uniform cost per error, both of which have not been widely studied in the knowledge- discovery and data mining community. The proposed methods of combining multiple

learned fraud detectors under a "cost model" are general and demonstrably useful; the empirical results demonstrate that we can significantly reduce loss due to fraud through distributed data mining of fraud models.

2. **Ghosh, S., & Reilly, D. L. (1994, January). Credit card fraud detection with a neural- network. In *System Sciences, 1994. Proceedings of the Twenty-Seventh Hawaii International Conference on* (Vol. 3, pp. 621-630). IEEE**

   Using data from a credit card issuer, a neural network based fraud detection system was trained on a large sample of labelled credit card account transactions and tested on a holdout data set that consisted of all account activity over a subsequent two-month period of time. The neural network was trained on examples of fraud due to lost cards, stolen cards, application fraud, counterfeit fraud, mail-order fraud and NRI (non- received issue) fraud. The network detected significantly more fraud accounts (an order of magnitude more) with significantly fewer false positives (reduced by a factor of 20) over rule-based fraud detection procedures. We discuss the performance of the network on this data set in terms of detection accuracy and earliness of fraud detection.

3. **Raj, S. B. E., & Portia, A. A. (2011, March). Analysis on credit card fraud detection methods. In *2011 International Conference on Computer, Communication and Electrical Technology (ICCCET)* (pp. 152-156). IEEE**

   The most commonly used fraud detection methods are Neural Network (NN), rule-induction techniques, fuzzy system, decision trees, Support Vector Machines (SVM), Artificial Immune System (AIS), genetic algorithms, K-Nearest Neighbour algorithms. These techniques can be used alone or in collaboration using ensemble or meta-learning techniques to build classifiers. This paper presents a survey of various techniques used in credit card fraud detection and evaluates each methodology based on certain design criteria.

4. **Brause, R., Langsdorf, T., & Hepp, M. (1999, November). Neural data mining for credit card fraud detection. In *Proceedings 11th International Conference on Tools with Artificial Intelligence* (pp. 103-106)**

   The prevention of credit card fraud is an important application for prediction techniques. One major obstacle for using neural network training techniques is the high necessary diagnostic quality: Since only one financial transaction of a thousand is invalid no prediction success less than 99.9% is acceptable. Due to these credit card transaction

proportions complete new concepts had to be developed and tested on real credit card data. This paper shows how advanced data mining techniques and neural network algorithm can be combined successfully to obtain a high fraud coverage combined with a low false alarm rate.

5. **Srivastava, A., Kundu, A., Sural, S., & Majumdar, A. (2008). Credit card fraud detection using hidden Markov model.** *IEEE Transactions on dependable and secure computing*, *5*(1)

   Due to a rapid advancement in the electronic commerce technology, the use of credit cards has dramatically increased. As credit card becomes the most popular mode of payment for both online as well as regular purchase, cases of fraud associated with it are also rising. In this paper, the sequence of operations in credit card transaction processing are modelled using a hidden Markov model (HMM) and show how it can be used for the detection of frauds. An HMM is initially trained with the normal behaviour of a cardholder. If an incoming credit card transaction is not accepted by the trained HMM with sufficiently high probability, it is considered to be fraudulent.

6. *Maes, S., Tuyls, K., Vanschoenwinkel, B., & Manderick, B. (2002, January).* **Credit card fraud detection using Bayesian and neural networks. In** *Proceedings of the 1st international naiso congress on neuro fuzzy technologies*

   This paper discusses automated credit card fraud detection by means of machine learning. In an era of digitalization, credit card fraud detection is of great importance to financial institutions. Two machine learning techniques suited for reasoning under uncertainty: artificial neural networks and Bayesian belief networks to the problem and show their significant results on real world financial data. Finally, future directions are indicated to improve both techniques and results.

7. **Aleskerov, E., Freisleben, B., & Rao, B. (1997, March). Cardwatch: A neural network based database mining system for credit card fraud detection. In** *Proceedings of the IEEE/IAFE 1997 computational intelligence for financial engineering (CIFEr)* **(pp. 220- 226). IEEE**
   CARDWATCH, a database mining system used for credit card fraud detection, is presented. The system is based on a neural network learning module, provides an interface to a variety of commercial databases and has a comfortable graphical user interface. Test results obtained for synthetically generated credit card data and an auto associative neural network model show very successful fraud detection rates.

8.  **Quah, J. T., & Sriganesh, M. (2008). Real-time credit card fraud detection using computational intelligence.** *Expert systems with applications*, *35***(4)**

    Online banking and e-commerce have been experiencing rapid growth over the past few years and show tremendous promise of growth even in the future. This has made it easier for fraudsters to indulge in new and abstruse ways of committing credit card fraud over the Internet. This paper focuses on real-time fraud detection and presents a new and innovative approach in understanding spending patterns to decipher potential fraud cases. It makes use of self-organization map to decipher, filter and analyze customer behaviour for detection of fraud.

9.  **Ogwueleka, F. N. (2011). Data mining application in credit card fraud detection system.** *Journal of Engineering Science and Technology*, *6***(3)**

    Data mining is popularly used to combat frauds because of its effectiveness. It is a well- defined procedure that takes data as input and produces models or patterns as output. Neural network, a data mining technique was used in this study. The design of the neural network (NN) architecture for the credit card detection system was based on unsupervised method, which was applied to the transactions data to generate four clusters of low, high, risky and high-risk clusters. The self-organizing map neural network (SOMNN) technique was used for solving the problem of carrying out optimal classification of each transaction into its associated group, since a prior output is unknown. The receiver-operating curve (ROC) for credit card fraud (CCF) detection watch detected over 95% of fraud cases without causing false alarms unlike other statistical models and the two-stage clusters. This shows that the performance of CCF detection watch is in agreement with other detection software, but performs better.

10. *Ogwueleka, F. N. (2011).* **Data mining application in credit card fraud detection system.** *Journal of Engineering Science and Technology*

    The fraud detection system (FDS) consists of four components, namely, rule-based filter, Dempster–Shafer adder, transaction history database and Bayesian learner. In the rule- based component, we determine the suspicion level of each incoming transaction based on the extent of its deviation from good pattern. Dempster–Shafer's theory is used to combine multiple such evidences and an initial belief is computed. The transaction is classified as normal, abnormal or suspicious depending on this initial belief. Once a transaction is found to be suspicious, belief is further strengthened or weakened

according to its similarity with fraudulent or genuine transaction history using Bayesian learning. Extensive simulation with stochastic models shows that fusion of different evidences has a very high positive impact on the performance of a credit card fraud detection system as compared to other methods.

11. **Syeda, M., Zhang, Y. Q., & Pan, Y. (2002, May). Parallel granular neural networks for fast credit card fraud detection. In** *2002 IEEE World Congress on Computational Intelligence*

This paper proposes an intelligent credit card fraud detection model for detecting fraud from highly imbalanced and anonymous credit card transaction datasets. The class imbalance problem is handled by finding legal as well as fraud transaction patterns for each customer by using frequent item set mining. A matching algorithm is also proposed to find to which pattern (legal or fraud) the incoming transaction of a particular customer is closer and a decision is made accordingly. In order to handle the anonymous nature of the data, no preference is given to any of the attributes and each attribute is considered equally for finding the patterns.

12. **Dal Pozzolo, A., Caelen, O., Le Borgne, Y. A., Waterschoot, S., & Bontempi, G. (2014). Learned lessons in credit card fraud detection from a practitioner perspective.** *Expert systems with applications*

Billions of dollars of loss are caused every year due to fraudulent credit card transactions. The design of efficient fraud detection algorithms is key for reducing these losses, and more and more algorithms rely on advanced machine learning techniques to assist fraud investigators. The design of fraud detection algorithms is however particularly challenging due to non stationary distribution of the data, highly imbalanced classes distributions and continuous streams of transactions. At the same time public data are scarcely available for confidentiality issues, leaving unanswered many questions about which is the best strategy to deal with them. In this paper they provided some answers from the practitioner's perspective by focusing on three crucial issues: non-balancedness, non-stationary and assessment. The analysis is made possible by a real credit card dataset provided by our industrial partner.

13. **Stolfo, S., Fan, D. W., Lee, W., Prodromidis, A., & Chan, P. (1997, July). Credit card fraud detection using meta-learning: Issues and initial results.**

The paper describes initial experiments using meta-learning techniques to learn models of fraudulent credit card transactions. The experiments reported here are the first step towards a better understanding of the advantages and limitations of current meta- learning strategies on real-world data. They argue that, for the fraud detection domain, fraud catching rate (True Positive rate) and false alarm rate (False Positive rate) are better metrics than the overall accuracy when evaluating the learned fraud classifiers. They showed that given a skewed distribution in the original data, artificially more balanced training data leads to better classifiers. They demonstrated how meta-learning can be used to combine different classifiers and maintain, and in some cases, improve the performance of the best classifier.

14. **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*

Due to the growing volume of electronic payments, the monetary strain of credit-card fraud is turning into a substantial challenge for financial institutions and service providers, thus forcing them to continuously improve their fraud detection systems. However, modern data-driven and learning-based methods, despite their popularity in other domains, only slowly find their way into business applications. In this paper, we phrase the fraud detection problem as a sequence classification task and employ Long Short-Term Memory (LSTM) networks to incorporate transaction sequences. We also integrate state-of-the-art feature aggregation strategies and report our results by means of traditional retrieval metrics. A comparison to a baseline random forest (RF) classifier showed that the LSTM improves detection accuracy on offline transactions where the card-holder is physically present at a merchant. Both the sequential and non-sequential learning approaches benefit strongly from manual feature aggregation strategies. A subsequent analysis of true positives revealed that both approaches tend to detect different frauds, which suggests a combination of the two. We conclude our study with a discussion on both practical and scientific challenges that remain unsolved.

15. **Carneiro, N., Figueira, G., & Costa, M. (2017). A data mining based system for credit- card fraud detection in e-tail.** *Decision Support Systems*

Credit-card fraud leads to billions of dollars in losses for online merchants. With the development of machine learning algorithms, researchers have been finding increasingly sophisticated ways to detect fraud, but practical implementations are rarely reported. We describe the development and deployment of a fraud detection system in

a large e-tail merchant. The paper explores the combination of manual and automatic classification, gives insights into the complete development process and compares different machine learning methods. The paper can thus help researchers and practitioners to design and implement data mining based systems for fraud detection or similar problems. This project has contributed not only with an automatic system, but also with insights to the fraud analysts for improving their manual revision process, which resulted in an overall superior performance.

16. **Zhu, Kevin. "Phone usage pattern as credit card fraud detection trigger.**

A credit card fraud reduction system is disclosed. The system comprises a computer system and an application. The application, when executed on the computer system, applies increased credit card fraud prevention procedures to use of an electronic credit card application in a mobile electronic device, based on a changed communication usage pattern of the mobile electronic device.

17. *Shen, A., Tong, R., & Deng, Y. (2007, June). Application of classification models on credit card fraud detection. In *2007 International conference on service systems and service management*

Along with the great increase in credit card transactions, credit card fraud has become increasingly rampant in recent years. This study investigates the efficacy of applying classification models to credit card fraud detection problems. Three different classification methods, i.e. decision tree, neural networks and logistic regression are tested for their applicability in fraud detections. This paper provides a useful framework to choose the best model to recognize the credit card fraud risk.

18. *Whitrow, C., Hand, D. J., Juszczak, P., Weston, D., & Adams, N. M. (2009). Transaction aggregation as a strategy for credit card fraud detection. *Data mining and knowledge discovery*

The problem of pre-processing transaction data for supervised fraud classification is considered. It is impractical to present an entire series of transactions to a fraud detection system, partly because of the very high dimensionality of such data but also because of the heterogeneity of the transactions. Hence, a framework for transaction aggregation is considered and its effectiveness is evaluated against transaction-level detection, using a variety of classification methods and a realistic cost-based

performance measure. These methods are applied in two case studies using real data. Transaction aggregation is found to be advantageous in many but not all circumstances. Also, the length of the aggregation period has a large impact upon performance. Aggregation seems particularly effective when a random forest is used for classification. Moreover, random forests were found to perform better than other classification methods, including SVMs, logistic regression and KNN. Aggregation also has the advantage of not requiring precisely labelled data and may be more robust to the effects of population drift.

19. **Bahnsen, A. C., Aouada, D., Stojanovic, A., & Ottersten, B. (2016). Feature engineering strategies for credit card fraud detection.** *Expert Systems with Applications*

In this paper they expanded the transaction aggregation strategy, and proposed to create a new set of features based on analyzing the periodic behaviour of the time of a transaction using the von Mises distribution. Then, using a real credit card fraud dataset provided by a large European card processing company, we compare state-of-the-art credit card fraud detection models, and evaluate how the different sets of features have an impact on the results. By including the proposed periodic features into the methods, the results show an average increase in savings of 13%.

20. **Patidar, R., & Sharma, L. (2011). Credit card fraud detection using neural network.**

In this paper they tried to detect fraudulent transaction through the neural network along with the genetic algorithm. As we will see that artificial neural network when trained properly can work as a human brain, though it is impossible for the artificial neural network to imitate the human brain to the extent at which brain work, yet neural network and brain, depend for there working on the neurons, which is the small functional unit in brain as well as ANN. Genetic algorithm are used for making the decision about the network topology, number of hidden layers, number of nodes that will be used in the design of neural network for our problem of credit card fraud detection. For the learning purpose of artificial neural network we will use supervised learning feed forward back propagation algorithm.

### EXISITING METHODOLODY

1. Heavy emphasis on neural networks
2. Tiny emphasis on pre-processing phase
3. Direct use of dataset without any cleansing
4. No use of feature scaling to more existing features more computational.

### PROPOSED METHODOLOGY

**NOVEL APPROACH TO BE FOLLOWED**

1. Problem solving without neural networks. Though neural networks can prove effective in this domain, yet they can never solve the problem of time constraints. By using a simplified approach we will be able to produce better throughput within limited time frame.
2. Using algorithms which have not been used to solve problems in this domain until 2020, like(LOF and Isolation Forest)
3. Feature scaling to adjust features in accordance with other features.
4. Dimensionality reduction using PCA. PCA technique will be used to reduce the dimensions of the existing data which will help yielding better results.

**ISOLATION FOREST METHODOLOGY**

- Isolation forest is an unsupervised learning algorithm for anomaly detection that works on the principle of isolating anomalies, instead of the most common techniques of profiling normal points.
- In statistics, an anomaly (a.k.a. outlier) is an observation or event that deviates so much from other events to arouse suspicion it was generated by a different mean.
- Anomalies in a big dataset may follow very complicated patterns, which are difficult to detect "by eye" in the great majority of cases. This is the reason why the field of anomaly detection is well suited for the application of Machine Learning techniques.
- The most common techniques employed for anomaly detection are based on the construction of a profile of what is "normal": anomalies are reported as those instances in the dataset that do not conform to the normal profile.
- Isolation Forest uses a different approach: instead of trying to build a model of normal instances, it explicitly isolates anomalous points in the dataset. The main advantage of this approach is the possibility of exploiting sampling techniques to an extent that is not allowed to the profile-based methods, creating a very fast algorithm with a low memory demand

**LOCAL OUTLIER FACTOR METHODOLOGY**

- The local outlier factor is based on a concept of a local density, where locality is given by k nearest neighbors, whose distance is used to estimate the density.
- By comparing the local density of an object to the local densities of its neighbors, one can identify regions of similar density, and points that have a substantially lower density than their neighbors. These are considered to be outliers.
- The local density is estimated by the typical distance at which a point can be "reached" from its neighbors.

- The definition of "reachability distance" used in LOF is an additional measure to produce more stable results within clusters.
- The "reachability distance" used by LOF has some subtle details that are often found incorrect in secondary sources.

## IMPLEMENTATION:

## Dataset after pre-processing, feature scaling and PCA:



## CODE :

```python
In [2]: import numpy as np
        import pandas as pd
        import matplotlib.pyplot as plt
        import seaborn as sns
```

```python
In [3]: dataset = pd.read_csv('creditcard.csv')
```

```python
In [4]: print(dataset.columns)
```

```
Index(['Time', 'V1', 'V2', 'V3', 'V4', 'V5', 'V6', 'V7', 'V8', 'V9', 'V10',
       'V11', 'V12', 'V13', 'V14', 'V15', 'V16', 'V17', 'V18', 'V19', 'V20',
       'V21', 'V22', 'V23', 'V24', 'V25', 'V26', 'V27', 'V28', 'Amount',
       'Class'],
      dtype='object')
```

```python
In [5]: print(dataset.shape)
```

```
(284807, 31)
```

```python
In [6]: dataset = dataset.sample(frac = 0.1,random_state=1)
```

```python
In [7]: print(dataset.shape)
```

```
(28481, 31)
```

```python
In [8]: dataset.hist(figsize=(20,20))
        plt.show()
```



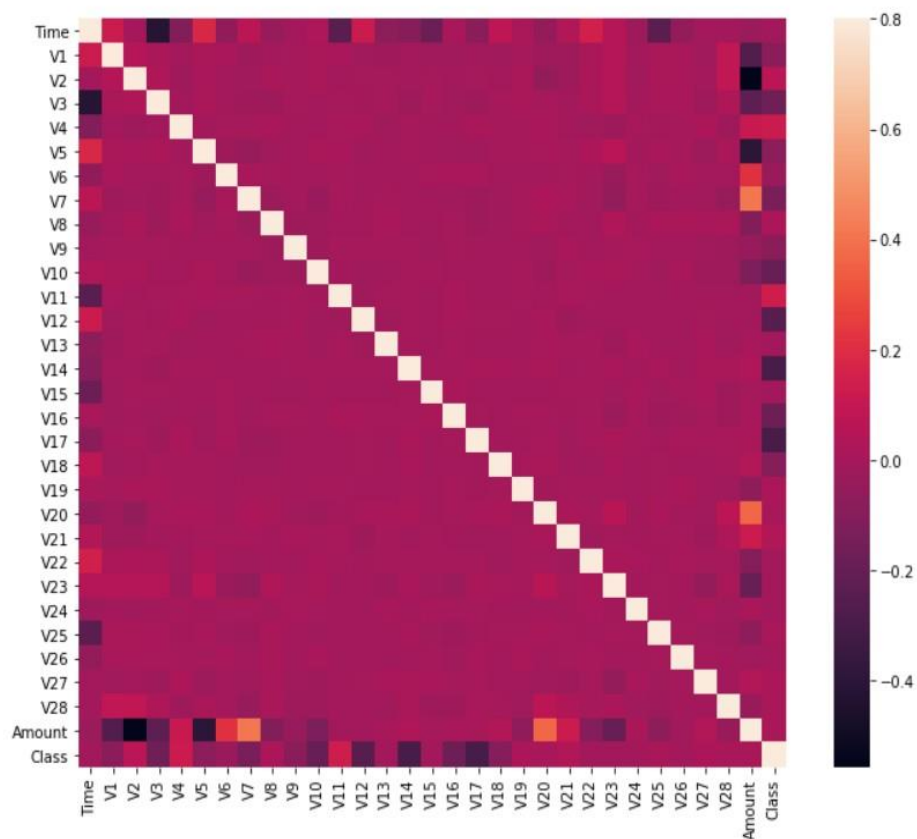## ANALYSIS OF INDEPENDENT VARIABLES:

**Generating HeatMap to determine correlation between different columns in the dataset:**

## CORE – IMPLEMENTATION:

```
In [13]: from sklearn.metrics import classification_report, accuracy_score
         from sklearn.ensemble import IsolationForest
         from sklearn.neighbors import LocalOutlierFactor

         # define random states
         state = 1

         # define outlier detection tools to be compared
         classifiers = {
             "Isolation Forest": IsolationForest(max_samples=len(X),
                                                 contamination=outlier_fraction,
                                                 random_state=state),
             "Local Outlier Factor": LocalOutlierFactor(
                 n_neighbors=20,
                 contamination=outlier_fraction)}
```

```
In [14]: # Fit the model
         plt.figure(figsize=(9, 7))
         n_outliers = len(Fraud)


         for i, (clf_name, clf) in enumerate(classifiers.items()):

             # fit the data and tag outliers
             if clf_name == "Local Outlier Factor":
                 y_pred = clf.fit_predict(X)
                 scores_pred = clf.negative_outlier_factor_
             else:
                 clf.fit(X)
                 scores_pred = clf.decision_function(X)
                 y_pred = clf.predict(X)
```

## EVALUATION METRICS:

```
Isolation Forest: 71
0.99750711000316
              precision    recall  f1-score   support

           0       1.00      1.00      1.00     28432
           1       0.28      0.29      0.28        49

    accuracy                           1.00     28481
   macro avg       0.64      0.64      0.64     28481
weighted avg       1.00      1.00      1.00     28481

Local Outlier Factor: 97
0.9965942207085425
              precision    recall  f1-score   support

           0       1.00      1.00      1.00     28432
           1       0.02      0.02      0.02        49

    accuracy                           1.00     28481
   macro avg       0.51      0.51      0.51     28481
weighted avg       1.00      1.00      1.00     28481
```

## PROJECT OUTCOME:

- The model is simple and fast enough to detect the anomaly and classify it as a fraudulent transaction as quickly as possible.
- Imbalance can be dealt with by properly using some Isolation forest algorithm and Local Outlier factor (LOF).
- For protecting the privacy of the user the dimensionality of the data will be reduced. A more trustworthy source will be taken which double-checks the data, at least for training the model.
- The model is simple and interpretable so that when the scammer adapts to it with just some tweaks, we can have a new model up and running to deploy.