Enhanced Random Forest Algorithm for Faster and Accurate Fraud Detection

Mrs. G. Jaculine Priya¹ Ph.D. Research Scholar, Department of Computer Science, VISTAS, Pallavaram, Chennai, Tamil Nadu, India, Dr. S. Saradha² Assistant Professor, MEASI Institute of Information Technology, Chennai, Tamil Nadu, India,

Abstract: Paper In this Digital Era security is the important concern. The more using digital communication and digital fraud also becomes more. When the technology grows in a faster pace, another side fraudsters also becoming strong and fast. Customers have to be very careful using their user's name, password, pin number, OTP etc. Without Customers knowledge the fraudsters are stealing their secured details. Industry also spending more time, effort and money to stop/reduce these illegal attempts keeping their client's security. But unfortunately, still there's a loop hole we are facing and fraudsters are well ahead in their job. AI&ML is the endowment in this current digital technological world. This niche technology can be used Finance, Medical, Insurance, Ecommerce almost across all the domains. Each company/industry has its own security approaches for providing a safe environment to their clients. Here, suggested solution aims to create a global centralized framework for sharing their fraud patterns by getting into a Digital handshake across the organization from various domains. Here AI&ML algorithms helps to detect and prevent fraud attacks intelligently. These algorithms help to find out all the hidden patterns. Choosing the right algorithm definitely it improves the detecting and preventing the fraud patterns and also gives secured environment to the customers.

Keywords: Artificial Intelligence, Machine Learning, Fraud Detection, Prevention, Credit Card

I. INTRODUCTION

Cyber Security is the major concern when an information is conveyed across a medium in our technological era, the applications we use keep on changing in a daily basis. The digital transition is greatly assisting all sectors. At the same time, we leave our digital footprints wherever and whenever we go. It makes users extremely cautious about their credentials. Users want us to be proactive in preventing fraud. In this paper, we will examine how to improve security features. However, scammers do not lag behind when it comes to embracing new technologies. Corporates must make sure that their security measures stay up with their operations with the innovated solutions.

A large number of fraudulent transactions occur in the banking sector on a daily basis. As an example, if we take only credit/debit card fraudulent transactions, the fraud will be notified to the appropriate bank management. As a temporary solution, the bank has blocked the specific customer's card from future transactions and advised the customers to reset his password and other sensitive details.

Only after that did the specific bank/corporate begin investigating how this fraudulent activity occurred, and what

strategy the thieves employed to breach the client account. Once the pattern is detected, the banking/corporate sectors will endeavor to close loopholes and implement the required security measures to effectively stop the same pattern-related malicious activity. Meanwhile, criminals would have committed several fraudulent transactions following the same pattern, resulting in massive losses.

Fraudsters are constantly inventing new ways of accessing banks and their customers in the context of a rapidly changing global financial market, where demand for face-to-face banks is falling, volumes of online wallets are expanding, and transactions are made in seconds. Banks must be adaptable in responding to emerging dangers and embracing new strategies and technologies in order to foresee and prevent fraud.

Financial firms' greatest threat is cyber-related theft. To limit fraud risks in the future, financial institutions will need to make a massive change in their approach. Fundamentally, financial institutions must comprehend the rapid digital revolution occurring all around us, recognize the resulting evolving fraud threats, and create fraud risk management practices capable of mitigating these fraud risks in a consistent, efficient, and intelligent manner. Existing financial institution solutions, while expensive to operate, are not capable of coping with rising fraud concerns because they are too scattered and unsophisticated. Computer fraud risk management in the future should be able to adapt to the rapidly evolving digital transformation, discover previously unidentified fraud threats, take advantage of technology, and reduce compliance costs. Fixed deposits, loan disbursement or granting credit facilities for bribery, hacking and other web ATM-based scams are some of the latest fraud instances in India covered by the media. These high-profile incidents in recent years have demonstrated that scams may harm an organization's brand as well as its earnings, operating efficiencies, and customer satisfaction. It can have a severe influence on staff morale and confidence of investors, in addition to potential regulatory sanctions. Although no business can completely eliminate the risk of fraud, it is critical to have procedures in place that can prevent and detect fraud.

The following is the structure of this paper:

The Solving the Global Issue and Data Preprocessing are discussed in section 2. Various Machine Learning Algorithms and Comparison of Algorithms are discussed in section 3. Proposed Enhanced ML Algorithm steps are detailed in section 4. The Results and Discussions are stated in Section 5.

II. SOLVING THE GLOBAL ISSUE

After fraudulent activity reported, fraudulent transactions are identified by the employer of the organization. This information/pattern can be saved in a centralized database. Based on this experience security of the organization can be improved. Each company has its own security approach for providing a safe environment for its clients. In my previous papers suggested a solution aims to create a centralized global framework for sharing fraud patterns, allowing any business to retain and see fresh fraudulent transactions carried out by criminals.

Geotagging, IP Address, Date, Time, as well as other fraudrelated information may all be maintained in a record, just like any other fraud transaction. We might reach a digital agreement between organizations across the domains. These organizations can share the knowledge in a standard format by the industry who signed the agreement. Now that all fraud transactions are visible to everyone, the surviving organizations may take preventative security measures to avoid large losses. A fraudster's relationship with a victim may be compared to a cat-and-mouse fight, in which each side is constantly learning and adapting, employing inventiveness and understanding of the other's objectives to develop new fraud prevent and detect techniques. In this paper discussed how to get to fast and accurate fraud detection solution by using the proposed enhanced machine learning algorithms in real time.

A. Importance of AI-ML Algorithms in Solving the Global Issue

Machine Learning is a subset of AI. It has 3 types of algorithms such as supervised, unsupervised and reinforcement learning algorithms. The Supervised Algorithms address regression and classification/anomaly findings problems. The unsupervised algorithms help to cluster unlabeled components. Reinforcement Learning helps industry to address action-reward-punishment the environment-based issues. It helps the model to learn and correct the error on its own. The Credit Card Fraud Detection system uses supervised algorithms such as Logistic Regression, Naïve Bayes and Support Vector Machine (SVM) to identify the fraud transactions. The Anomaly detection is a kind of classification model, where the sensitivity is key to find the anomaly. In general, if the model threshold is greater than 0.5, then it will classify it as +positive case but here the sensitivity/recall/true positive rate is very important to classify +ve/-ve transactions.

Machine Learning Techniques in Ai Technology now give intelligent answers to the majority of issues of human and traditional approaches in FP & FD. In the past, industries relied on their employees and their traditional methodology to FP & FD. The government, as well as all the private sector organizations, is investing more money on fraud prevention. The algorithms are created by fraud domain specialists in the traditional method. The algorithms and procedures are strictly based on rules. Traditional methods are no longer sufficient to address this issue. This is the current state, and it's clear existing FP & FD operate in a soiled approach. Its more over reactive rather than proactive mode. Recovery time of fraud issues are more and not real time. We don't learn from our history. We don't learn from our/another business verticals. Industry has to work well ahead of fraudsters. Because of the popularity and precision of Artificial Intelligence, every sector is currently transitioning from the traditional approach to ML-based solutions for FP & FD.

Fraud detection and prevention algorithms surely find out the hidden and new patterns based on the experience. It is a continuous learning procedure. It needs full life cycle which includes monitoring, learning, identifying, preventing, and decision-making in real time. Fraudulent transactions may be prevented and detected using the correct Supervised and Unsupervised learning algorithms. The benefit is that we may use either of the models individually or in combination to find abnormalities.

B. Proposed Global Fraud Prevention High Level Architecture

Transaction information is published and stored in a centralized database. The Global Fraud Prevention High-Level Architecture is depicted in the above Figure. The source organization with a NO or YES for the "fraud help indication." A "NO" would indicate that the publishing company has previously detected a transaction as fraudulent and shared this information in the centralized database for information purposes. This will help other companies involved in this experiment avoid similar suspicious transactions using the same patterns in the future. A "YES" would indicate that the publishing organization's is looking for an indicator from the Global model to allow or prevent the transaction from being completed in real-time.



Fig. 1.Proposed Global Fraud Prevention High Level Architecture- Level I

Similarly, remaining participating organizations in this digital handshake can communicate information across sectors. It aids in the prevention and detection of fraudulent transactions worldwide. The accuracy of a model is determined by the number of data sets utilized for training and testing. The model, once trained, can detect and prevent fraud patterns. Based on the fraud indicator signal supplied by the model, parent apps hosted by the company where the fraud originated can prevent the transaction from successfully completing or execute extra security validations before allowing the transaction to pass through. Using the suggested global architecture model, newly discovered fraud tendencies may then be shared across enterprises, allowing for proactive fraud prevention. Thereby its extremely important for organizations to join hands together in proactively preventing fraud Globally.



Fig. 2.Proposed Global Fraud Prevention High Level Architecture- Level II

C. Step wise execution for implementing a model

1. Load Dataset into the system.

2. Class Count function to check if the dataset is balanced or imbalanced

3. Scale Non- Anonymised features such as Time and Amount

4. Concatenate both Scaled Amount and Time with Actual Dataset/frame.

- 5. Drop old Amount and Time features
- 6. Random under-sampling function
- 6.1 Split the scaled dataset into Train and Test

6.2.Reset index of scaled train and test dataset

6.3.Find no of fraud transactions in the (random) Training dataset

6.4.Segregate normal and fraud transactions from the training data

6.5.Randomly selecting the same no of normal transactions as fraud transactions

6.6.Reset Index of both selected normal and fraud transactions

6.7.Concatenate both selected normal and fraud transactions. 6.8.Shuffle the subsample/final data frame.

7.Remove extreme outliners and create final dataset

8.Split final dataset into Train and Test

9.Spot-check a couple of Classification Algorithms (using cross validation)

10. Make Predictions on Test data

11. End

D. Data Preprocessing

Steps in Data Pre-processing in Machine Learning

- 1. Obtain the dataset
- 2. Add all necessary libraries
- 3. Load the dataset
- 4. Recognizing and dealing with missing values
- 5. Categorical data encoding
- 6. Dividing the dataset
- 7. Scaling of features

Because of their varied origin, the bulk of real-world datasets for machine learning are very sensitive to absent, irregular, and incomplete data. As a result, data preprocessing is critical for improving overall data quality. Duplicate or missing data may provide an inaccurate picture of the overall statistics of the data. Deviations and inconsistent data can disrupt the model's overall learning, resulting in incorrect predictions. Quality data must be used to make quality judgments. To obtain this high-quality data, data pre-processing is required.

Implementation of Dataset

Data set represents real-world data. This data set contains all credit card transactions. It's an imbalanced data set.

- It has 30 features and 1 target. Of 30 features, 28 features are labelled V1 to V28.
- The remaining 2 features are Time and Amount.
- The 28 features are in the form of PCA complaint.

Both Amount and Time are not normalised, that is, not in line with other variables in terms of scale.

1	A		8	С	D	E	F	G	н	1	J	ĸ	L	м	N	0	P	Q	R	s	E
1	Time		V1	V2	V3	V4	V5	V6	V7	V8	V9	V10	V11	V12	V13	V14	V15	V16	V17	V18	٧1
2		0	-1.35981	-0.07278	2.536347	1.378155	-0.33832	0.462388	0.239599	0.098698	0.363787	0.090794	-0.5516	-0.6178	-0.99139	-0.31117	1.468177	-0.4704	0.207971	0.025791	0.
3		0	1.191857	0.266151	0.16648	0.448154	0.060018	-0.08236	-0.0788	0.085102	-0.25543	-0.16697	1.612727	1.065235	0.489095	-0.14377	0.635558	0.463917	-0.1148	-0.18336	-(
4		1	-1.35835	-1.34016	1.773209	0.37978	-0.5032	1.800499	0.791461	0.247676	-1.51465	0.207643	0.624501	0.066084	0.717293	-0.16595	2.345865	-2.89008	1.109969	-0.12136	4
5		1	-0.96627	-0.18523	1.792993	-0.86329	-0.01031	1.247203	0.237609	0.377436	-1.38702	-0.05495	-0.22649	0.178228	0.507757	-0.28792	-0.63142	-1.05965	-0.68409	1.965775	-
6		2	-1.15823	0.877737	1.548718	0.403034	-0.40719	0.095921	0.592941	-0.27053	0.817739	0.753074	-0.82284	0.538196	1.345852	-1.11967	0.175121	-0.45145	-0.23703	-0.03819	0.
7		2	-0.42597	0.960523	1.141109	-0.16825	0.420987	-0.02973	0.476201	0.260314	-0.56867	-0.37141	1.341262	0.359894	-0.35809	-0.13713	0.517617	0.401726	-0.05813	0.068653	-(
8		4	1.229658	0.141004	0.045371	1.202613	0.191881	0.272708	-0.00516	0.081213	0.45495	-0.09925	-1.41691	-0.15383	-0.75106	0.167372	0.050144	-0.44359	0.002821	-0.61199	-0
9		7	-0.64427	1.417964	1.07438	-0.4922	0.948934	0.428118	1.120631	-3.80786	0.615375	1.249376	-0.61947	0.291474	1.757964	-1.32387	0.686133	-0.07613	-1.22213	-0.35822	0.
10		7	-0.89429	0.286157	-0.11319	-0.27153	2.669599	3.721818	0.370145	0.851084	-0.39205	-0.41043	-0.70512	-0.11045	-0.28625	0.074355	-0.32878	-0.21008	-0.49977	0.118765	0.
11		9	-0.33826	1.119593	1.044367	-0.22219	0.499361	-0.24676	0.651583	0.069539	-0.73673	-0.36685	1.017614	0.83539	1.005844	-0.44352	0.150219	0.739453	-0.54098	0.476577	0.
12		10	1.449044	-1.17634	0.91386	-1.37567	-1.97138	-0.62915	-1.42324	0.048456	-1.72041	1.626659	1.199644	-0.67144	-0.51395	-0.09505	0.23093	0.031967	0.253415	0.854344	-(
13		10	0.384978	0.616109	-0.8743	-0.09402	2.924584	3.317027	0.470455	0.538247	-0.55889	0.309755	-0.25912	-0.32614	-0.09005	0.362832	0.928904	-0.12949	-0.80998	0.359985	0.
14		10	1.249999	-1.22164	0.38393	-1.2349	-1.48542	-0.75323	-0.6894	-0.22749	-2.09401	1.323729	0.227666	-0.24268	1.205417	-0.31763	0.725675	-0.81561	0.873936	-0.84779	-4
15		11	1.069374	0.287722	0.828613	2.71252	-0.1784	0.337544	-0.09672	0.115982	-0.22108	0.46023	-0.77366	0.323387	-0.01108	-0.17849	-0.65556	-0.19993	0.124005	-0.9805	-6
16		12	-2.79185	-0.32777	1.64175	1.767473	-0.13659	0.807596	-0.42291	-1.90711	0.755713	1.151087	0.844555	0.792944	0.370448	-0.73498	0.406796	-0.30305	-0.15587	0.778265	2.
17		12	-0.75242	0.345485	2.057323	-1.46864	-1.15839	-0.07785	-0.60858	0.003503	-0.43617	0.747731	-0.79398	-0.77041	1.047627	-1.0566	1.106953	1.660114	-0.27927	-0.41999	0.
18		12	1.103215	-0.0403	1.267332	1.289091	-0.735	0.288069	-0.58606	0.18938	0.782333	-0.26798	-0.45031	0.936708	0.70838	-0.46865	0.354574	-0.24663	-0.00921	-0.59591	-0
19		13	-0.43691	0.918966	0.924591	-0.72722	0.915679	-0.12787	0.707642	0.087962	-0.66527	-0.73798	0.324098	0.277192	0.252624	-0.2919	-0.18452	1.143174	-0.92871	0.68047	0. ,
		c	reditcard	۲																	

Fig. 3. Data Preprocessing Dataset Representation

Identifying Fraud Transactions

n	[11]: N	df[df.	Class==1] # Fro	ud trans	actions											
	Out[11]:	V5	V6	V7	V8	V9		V21	V22	V23	V24	V25	V26	V27	V28	Amount	Class
		.522188	-1.426545	-2.537387	1.391657	-2.770089		0.517232	-0.035049	-0.465211	0.320198	0.044519	0.177840	0.261145	-0.143276	0.00	1
		.359805	-1.064823	0.325574	-0.067794	-0.270953		0.661696	0.435477	1.375966	-0.293803	0.279798	-0.145362	-0.252773	0.035764	529.00	1
		.821628	-0.075788	0.562320	-0.399147	-0.238253	-	-0.294166	-0.932391	0.172726	-0.087330	-0.156114	-0.542628	0.039566	-0.153029	239.93	1
		.128131	-1.706536	-3.496197	-0.248778	-0.247768		0.573574	0.176968	-0.436207	-0.053502	0.252405	-0.657488	-0.827136	0.849573	59.00	1
		.624201	-1.357746	1.713445	-0.496358	-1.282858		-0.379068	-0.704181	-0.656805	-1.632653	1.488901	0.566797	-0.010016	0.146793	1.00	1
			1.00		144					144				244			
		.566487	-2.010494	-0.882650	0.697211	-2.084945		0.778584	-0.319189	0.639419	-0.294885	0.537503	0.788395	0.292680	0.147968	390.00	1
		.442581	-1.326536	-1.413170	0.248525	-1.127396	-	0.370612	0.028234	-0.145640	-0.081049	0.521875	0.739467	0.389152	0.186637	0.76	1
		.120541	-0.003346	-2.234739	1.210158	-0.652250		0.751826	0.834108	0.190944	0.032070	-0.739695	0.471111	0.385107	0.194361	77.89	1
		.840618	-2.943548	-2.208002	1.058733	-1.632333		0.583276	-0.269209	-0.456108	-0.183659	-0.328168	0.606116	0.884876	-0.253700	245.00	1
		.151147	-0.096695	0.223050	-0.068384	0.577829		-0.164350	-0.295135	-0.072173	-0.450261	0.313267	-0.289617	0.002988	-0.015309	42.53	1

Fig. 4 . Data Preprocessing Fraud Transactions Class 1

Identifying Good Transactions

0.01101.																
our[10].	V5	V6	V7	V8	V9	•••	V21	V22	V23	V24	V25	V26	V27	V28	Amount	Class
	.338321	0.462388	0.239599	0.098698	0.363787		-0.018307	0.277838	-0.110474	0.066928	0.128539	-0.189115	0.133558	-0.021053	149.62	0
	.060018	-0.082361	-0.078803	0.085102	-0.255425		-0.225775	-0.638672	0.101288	-0.339846	0.167170	0.125895	-0.008983	0.014724	2.69	0
	.503198	1.800499	0.791461	0.247676	-1.514654		0.247998	0.771679	0.909412	-0.689281	-0.327642	-0.139097	-0.055353	-0.059752	378.66	0
	.010309	1.247203	0.237609	0.377436	-1.387024	***	-0.108300	0.005274	-0.190321	-1.175575	0.647376	-0.221929	0.062723	0.061458	123.50	0
	.407193	0.095921	0.592941	-0.270533	0.817739		-0.009431	0.798278	-0.137458	0.141267	-0.206010	0.502292	0.219422	0.215153	69.99	0
	-	144	-					-						200	-	347
	.364473	-2.606837	-4.918215	7.305334	1.914428		0.213454	0.111864	1.014480	-0.509348	1.436807	0.250034	0.943651	0.823731	0.77	0
	.868229	1.058415	0.024330	0.294869	0.584800	***	0.214205	0.924384	0.012463	-1.016226	-0.606624	-0.395255	0.068472	-0.053527	24.79	0
	.630515	3.031260	-0.296827	0.708417	0.432454		0.232045	0.578229	-0.037501	0.640134	0.265745	-0.087371	0.004455	-0.026561	67.88	0
	.377961	0.623708	-0.686180	0.679145	0.392087		0.265245	0.800049	-0.163298	0.123205	-0.569159	0.546668	0.108821	0.104533	10.00	0
	.012546	-0.649617	1.577005	-0.414650	0.486180		0.261057	0.643078	0.376777	0.008797	.0.473549	.0.818267	.0.002415	0.013649	217.00	0

Fig. 5 . Data Preprocessing Fraud Transactions Class 0 Fraud transactions are very less comparing to normal transactions. Feature class is the respective variable it takes 1 in case of fraud transaction and 0 in case of good transaction. We have to handle this imbalanced data set first.

Removing Duplicate Records

				0															
Fil	е	Edit	Vier	w In	sert Ce	l Kerne	I Widge	ets Help	2								Not Truste	d P	ython 3
ß	+	3<	0	6 ^	↓ ► F	Run 🔳 🕻	C 🗰 Ca	de	~ 🖂										
				dtype	: object														
			D	upli	cate F	Record	ds												
	Ir	[12]	I: H	df[df	.duplicat	ed()]													
		Out	[12]:		Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	-	V21	V22	V23	
				3	3 26.0	-0.529912	0.873892	1.347247	0.145457	0.414209	0.100223	0.711206	0.176066	-0.286717		0.046949	0.208105	-0.185548	0.00
				3	5 26.0	-0.535388	0.865268	1.351076	0.147575	0.433680	0.086983	0.693039	0.179742	-0.285642		0.049526	0.206537	-0.187108	0.00
				11	3 74.0	1.038370	0.127486	0.184456	1.109950	0.441699	0.945283	-0.036715	0.350995	0.118950		0.102520	0.605089	0.023092	-0.62
				11	4 74.0	1.038370	0.127486	0.184456	1.109950	0.441699	0.945283	-0.036715	0.350995	0.118950		0.102520	0.605089	0.023092	-0.62
				11	5 74.0	1.038370	0.127486	0.184456	1.109950	0.441699	0.945283	-0.036715	0.350995	0.118950		0.102520	0.605089	0.023092	-0.62
				28295	7 171288.0	1.912550	-0.455240	-1.750654	0.454324	2.089130	4.160019	-0.881302	1.081750	1.022928		-0.524067	-1.337510	0.473943	0.61
				28348	3 171627.0	-1.464380	1.368119	0.815992	-0.601282	-0.689115	-0.487154	-0.303778	0.884953	0.054065	-	0.287217	0.947825	-0.218773	0.08
				28348	171627.0	-1.457978	1.378203	0.811515	-0.603760	-0.711883	-0.471672	-0.282535	0.880654	0.052808		0.284205	0.949659	-0.216949	0.08
				28419	1 172233.0	-2.667936	3.160505	-3.355984	1.007845	-0.377397	-0.109730	-0.667233	2.309700	-1.639306		0.391483	0.266536	-0.079853	-0.09
				1000													the second of		

Fig. 6 . Data Preprocessing Removal of Duplicate Record

	host.8888/notebooks/Credit%20Card%20Data%20Preprocessing.pynb	6 G	E1	'⊕
💭 jupyter C	edit Card Data Preprocessing Last Checkcoint: 4 hours ago (autosaved)		2	Logout
File Edil Vie	v Insert Coll Kennel Widgels Help	Not Truste	01	Python 3
0 + × 0	Δ ★ Ψ ▶ Run ■ C ≫ Code ~ 100			
	284191 172233.0 -2.667396 3.160505 -3.355584 1.007845 -0.377307 -0.109730 -0.667233 2.300700 -1.635306 0.301483	0.266535	-0.07985	13 -0.096
	284193 172203.0 -2.801642 3.123188 -3.339407 1.017018 -0.200095 -0.167054 -0.745886 2.325916 -1.834651 0.402639	0.250746	-0.08660	6 -0.090
	1061 rows × 31 columns			
	C			
In [19]: 🕨	# Renew DaylCotte Accord: point("he of records inform: ", lm(df)) df = df.drog.dplicats() print("he of records after (", lm(df))			
	No of records before: 284807 No of records after: 283726			

Fig. 7 . Data Preprocessing After Removal of Duplicate Records

In data preprocessing we will often need to find out the duplicate data rows and how to handle them. Finding the duplicate data , counting the duplicate and non-duplicate datas, extracting the duplicted rows with correct location , determining to keep the remaining rows and drop the duplicated rows. Similarly we can consider and do the columns as well by dropping duplicates.

Describing Dataset Descriptive Statistics

Dataset description helps to understand how datas have been collected in a defined structure. Each row and column describes a particular variable like mean, count, standard deviation, percentiles and minimum-maximum value ranges are included. We can get a descriptive statistics summary of given data frame.

File Ed	R Vew	Inse	rt Cell	Kernel	Widgets H	icip					Not Trusted	- 1
+ 1	< 2 B	•	↓ R	n 🔳 C 🕨	Code	× 18						
	Da	ntase	t Des	criptive	e Statis	tics						
In [14]: M	# Dotos df.iloc	et Descr [:, :-1]	<pre>iptive Stati .describe().</pre>	istics (exc .T	Luded clas:	s variable)					
	Dut[14]:		count	mean	std	min	25%	50%	75%	max		
		Time	283726.0	94811.077600	47481.047891	0.000000	54204.750000	84692.500000	139298.000000	172792.000000		
		V1	283726.0	0.005917	1.048026	-56.407510	-0.015951	0.020384	1.316058	2.454930		
		V2	283726.0	-0.001135	1.646703	.72.715728	-0.600321	0.063949	0.800283	22.057729		
		V3	283728.0	0.001613	1.506882	-48.325589	-0.889682	0.179963	1.026980	9.382558		
		V4	283728.0	-0.002968	1.414184	-5.883171	-0.850134	-0.022248	0.739647	16.875344		
		¥5	283726.0	0.001828	1.377008	-113.743307	-0.689830	-0.053468	0.612218	34.801665		
		V6	283726.0	-0.001139	1.331931	-26.160505	-0.759031	-0.275168	0.396792	73.301626		
		¥7	283726.0	0.001801	1.227664	-43.557242	-0.552509	0.040859	0.570474	120.589494		
		V8	283726.0	-0.000854	1.179054	-73.216718	-0.208828	0.021898	0.325704	20.007208		
		V9	283728.0	-0.001598	1.095492	-13,434068	-0.644221	-0.052598	0.595977	15.591995		
		V10	283728.0	-0.001441	1.076407	-24.588262	-0.535578	-0.083237	0.453819	23.745138		
		V11	283726.0	0.000202	1.018720	-4.797473	-0.761649	-0.032306	0.739579	12.018913		
		140		0.000747	0.0000774		A 404.400	P. 47002373	A ####03#	10.00.000		

Fig. 8 . Data Preprocessing- Dataset Descriptive Statistics

Heat Map

Heat map can infer features are not correlated mostly that means the variation of one feature minutely affects the other feature that either has a positive or a negative correlation with each other.



Fig. 9 . Data Preprocessing - Heat Map of the Dataset

III. MACHINE LEARNING ALGORITHMS

A. Support Vector Machine Algorithm:

The SVM method is best understood by concentrating on its fundamental kind, the SVM classifier. The SVM classifier is designed to create a hyper-lane in an N-dimensional environment that splits pieces of data into multiple groups. This hyperplane, however, is chosen based on margin, as the hyperplane with the greatest margin between the two classes is evaluated. These margins are determined using Support Vectors, which are data points. Support Vectors are data points that are close to the hyperplane and aid in its orientation.

If the operation of an SVM classifier needs to be understood theoretically, it may be divided into the following steps -

Step 1: The SVM algorithm predicts the classes. One of the classes is labelled as 1, while the other is labelled as -1. Load the necessary libraries.

Step 2: Import the dataset and extract the X and Y variables independently.

Step 3: Separate the dataset into train and test subsets.

Step 4: Set up the SVM classifier model

Step 5: SVM classifier model fitting

Step 6: Making predictions

Step 7: Assessing the model's performance

When there is no mistake in the classification, the gradients are just updated using the regularization parameter, whereas the loss function is additionally employed when misclassification occurs.

B. KNN Algorithm:

K Nearest Neighbors is a simple method that uses a similarity metric to predict the categorization of unlabeled data. We calculate the distance between the points when two parameters are shown in a 2D Cartesian system to get an idea of how similar they are. The KNN algorithm is based on the idea that similar objects are close to one another.

The K Nearest Neighbors method is one of the most basic and straightforward classification techniques. Based on how it operates, it is also classified as a "Lazy Learning Algorithm." The K-value that everyone achieves while training the model is usually an odd integer, however this is not required. However, there are a few drawbacks to employing KNN. A few of them are as follows:

- It is incompatible with categorical data because we cannot determine the distance between two category characteristics.
- It also does not perform well with high-dimensional data since the algorithm will struggle to calculate the distance in each dimension.

The KNN Algorithm in Python: A Step-by-Step Guide

Step 1: Import Libraries. Importing the libraries required to execute KNN is demonstrated below.

Step 2: Import the Dataset We can see the dataset being imported here....

Step 3: Split the dataset...

Step 4: Design a Training Model.

- Step 5: Make Running Predictions.
- Step 6: Validation Check.

C. Random Forest Algorithm:

Random Forest operates in two stages: First Step: Generating the random forest by combining N decision trees. Second Step: Making predictions for each tree generated in the first step. Random Forest is a classifier that uses a several decision trees on various subsets of a given dataset and averages them to enhance the predicted accuracy of that dataset. "Rather than depending on a single decision tree, the random forest calculates the forecast from each tree and choosing the final output result based on the majority vote of predictions." The greater number of trees in the forest, the stronger the precision and smaller the chance of errors.

The steps to illustrate the working process:

Step 1: Select M data points at random from the training set. Step 2: Generate decision trees for the specified data points (Subsets).

Step 3: Construct the N number of decision trees

Step 4: Reverse steps 1 and 2.

Step 5: Find each decision tree predictions for new data points and allocate the new data points to the category that has received the most votes.

Random Forest Algorithm Implementation

- 1. Pre-Processing of Data
- 2. RFA fitting to the training set
- 3. Estimating the Test Set outcome
- 4. Creation of the Confusion Matrix
- 5. Visualizing the Training Set's Outcome
- 6. Projecting the results of the test set

D. Comparison of Machine Learning Algorithms

There are two distinctions in the efficiency of random forest and gradient boosting. The random forest can create each tree individually, but gradient boosting can only build one tree at a time, hence the random forest performs worse than gradient boosting. Random Forest is designed for multi-class issues, whereas SVM is designed for two-class problems. Multiclass problem must be divided into numerous binary classification tasks. Random Forest performs effectively when characteristics are both categorical and numerical. We all know that KNN is a lazy learner; it memorizes the information, resulting in a 0-training time. Due to the fact that it does not train for parameters or weights.

While predicting, it really does all of the work. It has a complexity on the order of n * m *d, where n is the volume of the training data, m is the amount of the test data, and d is the number of operations that can be performed for every test. So, after approximately 10 minutes, it stops making

predictions. As seen above, Decision Tree completed instantaneously with 85 percent accuracy, KNN with 92 percent accuracy but significant running time and consuming resources all along and Random Forest with 93 percent accuracy and very little running time.

Table 1.Comparsion of various ML Algorithms for Data Set-1

	S	NO	AI	GO	С	PR	RI	EC	AC	CCUR	F1
			RI	TH	L	EC	A	I.	AC	Y	SCO
S		ALG	Q	CLA	Ā	PŘĚC	IS	RE	Ĉ	ACCU	ŔĔĬ
Γ	I	RITI	H	SS	S	IØN		AL	L	RACY	SC
()	Μ			S						OR
	1		SV	M	0	1.00	1.(00	1.0	0	0.85
1		SVM	-	0	1	1,009	0.1	751.0	00.8	51.00	0.7
	2		Lo	gisti	0	0.960	1.(₀₀ 0.6	$3_{1.0}$	0 ^{0.76}	0.90
2		Logi	stį F	eg	1	1094	0.8	$88^{1.0}$	00.9	11.00	0.7
	3	c Reg	3 KN	JA.	0	0. <mark>8ã</mark> 0	1 (₀₀ 0.6	3_{10}	₀ 0.72	0 92
3		KNN	[0	1	1,009	0.8	81.0	009	<u>3</u> 1.00	0.8
	4		Ra	ndo	0	0.940	1.(ы ^{р.7}	61.0	₀ 0.83	0.93
4		Rand	om	0	1	1097	0.8	81.0	00.9	<u>2</u> 1.00	0.8
		m		1		0.94		0.7	7	0.85	9

Table 2 .Comparsion of various ML Algorithms for Data Set-2

E. Confusion Matrix



Precision:

Precision is one of the indicators of ML Algorithms model's performance i.e., the Quality of the positive prediction made by model. The number of true positives divided by the total number of positive predictions.

Recall:

The number of true positives divided by the number of positive values in the test data. Low Recall indicates a high number of false negatives.

F1-Score:

The weighted average of precision and recall.

Confusion Matrix:

A table showing correct predictions and types of incorrect predictions

		Prediction							
Actual	0(Negative)	1(Positive)						
0(-ve)	TN	(True	FP	(False					
	Negative)		Positive)						
1(+ve)	FN	(False	TP	(True					
	Negative)		Positive)						

$$Accuracy = \frac{TN+TP}{TN+TP+FP+FN}$$
(1)
TP

$$Precision = \overline{FP+TP}$$
(2)
$$\underline{TP}$$

$$Recall = \overline{TP+FN}$$
(3)

$$F1 Score = 2* \frac{Precision*Recall}{Precision+Recall}$$
(4)

Precision: 0.746	Recall: 0.837	Accuracy: 0.999	F1 Score:
0.740	0.057	0.999	Score.
			0.789

Table 3 . Performance Results

IV. PROPOSED ENHANCED RANDOM FOREST ALGORITHM

Above steps mentioned in II-D more over standard way of how machine learning algorithms works. Since proposing global model solution, it is important that algorithms used here are efficient and speed in real time. Need of the hour is enhanced machine learning algorithm for best performance. Performance has been improved in Random Forest Algorithm by selecting the wise trees alone to predict accurately in terms of identifying fraud attacks.

Revised Steps wise execution for Proposed Model

1. Load Dataset into the system.

2. Class Count function to check if the dataset is balanced or imbalanced

3. Scale Non- Anonymised features such as Time and Amount

4. Concatenate both Scaled Amount and Time with Actual Dataset/frame.

- 5. Drop old Amount and Time features
- 6. Random under-sampling function
- 6.1 Split the scaled dataset into Train and Test
- 6.2.Reset index of scaled train and test dataset

6.3.Find no of fraud transactions in the (random) Training dataset

6.4. Segregate normal and traud transactions from the training data

6.5.Randomly selecting the same no of normal transactions as fraud transactions6.6.Reset Index of both selected normal and fraud

transactions

6.7.Concatenate both selected normal and fraud transactions. 6.8.Shuffle the subsample/final data frame.

7.Remove extreme outliners and create final dataset

8.Split final dataset into Train and Test

9.Spot-check a couple of Classification Algorithms (using cross validation)

10. Make Predictions on Test data

11. Tune the Global model with Enhanced Random

Forest Algorithm 12.End

u

V. RESULTS AND DISCUSSION

Early detection, achieved by improving algorithms to detect both developing risks and the actions of fraudster, may be a critical step toward minimizing and reducing losses. On an enterprise-wide scale, incident detection that integrates complicated, adaptive, signaling, and reporting systems may automate the correlation and analysis of enormous volumes of data across the sectors, as well as numerous danger indicators. Monitoring systems in organizations should be operational 24/7, seven days a week, with enough support for fast incident response and remediation processes. A thorough awareness of recognized risks and controls, as well as industry norms and laws, may help financial services businesses protect their systems through the design and implementation of proactive, risk-informed controls. Based on best practices, banks can implement a "defense-in-depth" strategy to combat known and developing threats. This entails sharing common security layers, both to offer stability and to potentially impede, if not prevent, the advancement of ongoing threats.

To summarize, firms cannot manage fraudulent transactions in silos any more. Operating in silos will cost enterprises throughout the world trillions of dollars in lost revenue. Increased preventative measures must be considered and implemented, and the idea of centrally keeping and managing fraud data would meet this demand. Simply allowing organizational systems to communicate and learn from one other fast will dramatically lessen the impact of fraud, therefore actively limiting fraudsters from attacking them.

This proposed framework accomplishes this by allowing organizations to safely access the central database via exposed API (Application Program Interfaces) for every payment in their system, inspecting for fraud while also sharing any new transactions to centralized repository for other organizations to benefit in real time. In this regard, three distinct models – Support Vector Machine, KNN, and Random Forest have been tested by using historical data thus far. Random Forest appears to be producing good results, with the highest success rate among the algorithms tested. With the correct technology in place, firms can work together to avoid fraud by sharing their fraud experiences and employing the latest and best strategies.

REFERENCES

^{1.} G. Jaculine Priya and Dr.S.Saradha, "Fraud Detection and Prevention Using Machine Learning Algorithms: A Review," 7th International

Conference on Electrical Energy Systems (ICEES), 2021, pp. 304-308, DOI: 10.1109/ICEES51510.2021.9383631.

- G. Jaculine Priya and Dr.S.Saradha., "Real Time Global Fraud Detection and Prevention"., International E-Conference on Advances in Information Technology., June-2020 BIHER., ISBN NO.978-93-5407-796-8.
- E. M. S. W. Balagolla, W. P. C. Fernando, R. M. N. S. Rathnayake, M. J. M. R. P. Wijesekera, A. N. Senarathne and K. Y. Abeywardhana, "Credit Card Fraud Prevention Using Blockchain," 2021 6th International Conference for Convergence in Technology (I2CT), 2021, pp. 1-8, doi: 10.1109/I2CT51068.2021.9418192.
- B. A. Smadi, A. A. S. AlQahtani and H. Alamleh, "Secure and Fraud Proof Online Payment System for Credit Cards," 2021 IEEE 12th Annual Ubiquitous Computing, Electronics & Mobile Communication Conference (UEMCON), 2021, pp. 0264-0268, doi: 10.1109/UEMCON53757.2021.9666549.
- N. Chumuang, S. Hiranchan, M. Ketcham, W. Yimyam, P. Pramkeaw and S. Tangwannawit, "Developed Credit Card Fraud Detection Alert Systems via Notification of LINE Application," 2020 15th International Joint Symposium on Artificial Intelligence and Natural Language Processing (iSAI-NLP), 2020, pp. 1-6, doi: 10.1109/iSAI-NLP51646.2020.9376829.
- A. M. Zinjurde and V. B. Kamble, "Credit Card Fraud Detection and Prevention by Face Recognition," 2020 International Conference on Smart Innovations in Design, Environment, Management, Planning and Computing (ICSIDEMPC), 2020, pp. 86-90, doi: 10.1109/ICSIDEMPC49020.2020.9299587.
- H. Dar, A. Abbasi and A. Naveed, "Credit Card Fraud Prevention Planning using Fuzzy Cognitive Maps and Simulation," 2020 8th International Conference on Reliability, Infocom Technologies and Optimization (Trends and Future Directions) (ICRITO), 2020, pp. 289-294, doi: 10.1109/ICRITO48877.2020.9198002.
- N. Boutaher, A. Elomri, N. Abghour, K. Moussaid and M. Rida, "A Review of Credit Card Fraud Detection Using Machine Learning Techniques," 2020 5th International Conference on Cloud Computing and Artificial Intelligence: Technologies and Applications (CloudTech), 2020, pp. 1-5, doi: 10.1109/CloudTech49835.2020.9365916.
- V. Shah, P. Shah, H. Shetty and K. Mistry, "Review of Credit Card Fraud Detection Techniques," 2019 IEEE International Conference on System, Computation, Automation and Networking (ICSCAN), 2019, pp. 1-7, doi: 10.1109/ICSCAN.2019.8878853.
- Mehak Mahajan and Sandeep Sharma., "Detect Fraud in Credit Card using Data Mining Techniques". International Journal of Innovative and Exploring Engineering, 9, pp.2278-3075. December 2019.
- 11. Mandeep Singh and Sunny Kumar and Tushant Garg., "Credit Card Fraud Detection Using Hidden Markov Model". International Journal of Engineering and Computer Science, 8, pp.24878-24882. November 2019.
- Shiv Shankar Singh., "Electronic Credit Card Fraud Detection System by Collaboration of Machine Learning Models". International Journal of Innovative and Exploring Engineering, 8(12S), pp.92-95. October 2019.
- Maniraj, S.P., Aditya Saini., Swarna Deep Sarkar and Shadap Ahmed., "Credit Card Fraud Detection Using Machine Learning and Data Science". International Journal of Engineering Research & Technology,8(9), pp.110-115. September 2019.
- Yashvi Jain., Namrata Tiwari., Shripriya Dupey., and Sarika Jain., "A Comparative Analysis of Various Credit Card Fraud Detection Techniques". International Journal of Recent Technology and Engineering, 7(5S2), pp.402-407. January 2019.
- Olawale Adepoju., Julius Wosowei., Shiwani lawte and Hemaint Jaiman., "Comparative Evaluation of Credit Card Fraud Detection Using Machine Learning Techniques"., Global Conference for Advancement in Technology, IEEE, 2019, 978-1-7281-3694-3
- Thennakoon, Anuruddha, et al. "Real-time credit card fraud detection using machine learning." 2019 9th International Conference on Cloud Computing, Data Science & Engineering (Confluence). IEEE, 2019.
- Jhangiani, Resham, Doina Bein, and Abhishek Verma. "Machine learning pipeline for fraud detection and prevention in e-commerce transactions." 2019 IEEE 10th Annual Ubiquitous Computing, Electronics & Mobile Communication Conference (UEMCON). IEEE, 2019.
- Debachudamani Prusti., and Santanu Kumar Rath., "Fraudulent Transaction Detection in Credit Card by Applying Ensemble Machine Learning techniques"., 10th ICCCNT July 2019., IEEE
- Anish Halimaa A and Dr. K.Sundarakantham "Machine Learning Based Intrusion Detection System", ISBN: 978-1-5386-9439-8, IEEE 2019.

- Masoumen Zareapoor and Pourya Shamsolmoali., "Application of Credit Card Fraud Detection: Based on Bagging Ensemble Classifier". Procedia Computer Science, 48, pp. 679-685. 2015
- Hunt,W.(2020).Artificial Intelligence's Role in Finance and How Financial Companies are Leveraging the Technology to Their Advantage. Available at SSRN 3707908.
- 22. Brynjolfsson, E., & Mcafee, A.N.D.R.E.W. (2017). The Business of artificial intelligence. Harvard Business Review,7,3-11.
- 23. Kumarapandian, S. (2018). Melanoma classification using multiwavelet transform and support vector machine. International Journal of MC Square Scientific Research, 10(3),01-07.
- Yang, Q., Liu, Y., Chen, T., & Tong, Y. (2019). Federated machine learning: Concept and applications. ACM Transactions on Intelligent Systems and Technology (TIST), 10(2), 1-19.
- 25. Bao, Y., Hilary, G., & Ke, B. (2020). Artificial intelligence and fraud detection. Available at SSRN 3738618
- V. Mareeswari and G. Gunasekaran, "Prevention of credit card fraud detection based on HSVM," 2016 International Conference on Information Communication and Embedded Systems (ICICES), 2016, pp. 1-4, doi: 10.1109/ICICES.2016.7518889.
- K. T. Hafiz, S. Aghili and P. Zavarsky, "The use of predictive analytics technology to detect credit card fraud in Canada," 2016 11th Iberian Conference on Information Systems and Technologies (CISTI), 2016, pp. 1-6, doi: 10.1109/CISTI.2016.7521522.
- F. Ghobadi and M. Rohani, "Cost sensitive modeling of credit card fraud using neural network strategy," 2016 2nd International Conference of Signal Processing and Intelligent Systems (ICSPIS), 2016, pp. 1-5, doi: 10.1109/ICSPIS.2016.7869880.
- E. Caldeira, G. Brandao and A. C. M. Pereira, "Fraud Analysis and Prevention in e-Commerce Transactions," 2014 9th Latin American Web Congress, 2014, pp. 42-49, doi: 10.1109/LAWeb.2014.23.
- Seeja, K.R. and Masoumeh Zareapoor., "A Novel Credit Card Fraud Detection Model Based on Frequent Itemset Mining". The Scientific World Journal, Volume 2014, Article ID 252797,10 pages. September 2014
- M. D. H. Mahdi, K. M. Rezaul and M. A. Rahman, "Credit Fraud Detection in the Banking Sector in UK: A Focus on E-Business," 2010 Fourth International Conference on Digital Society, 2010, pp. 232-237, doi: 10.1109/ICDS.2010.45

AUTHORS PROFILE



Mrs. G. Jaculine Priya is currently pursuing PhD in Computer Science at VISTAS, Pallavaram. Her Area of interest and research includes Artificial Intelligence, Machine Learning, Computer Security and Fraud. She has been actively taken part and presented and published various papers in International and National Conferences in his research area.



Dr.S. Saradha, working as Assistant Professor, MEASI Institute of Information Technology Chennai. She has more than 10 years of experience in Educational Institute. Her area of research includes Data Mining, Artificial Neural Networks, Machine Learning, Big Data Analytics. She has published more than 17 Research Papers in

National and International journals. She is interested in writing textbooks for the students to make them understand any concept in an easy way. She is a recognised supervisor, guiding M.Phil. and PhD scholars.