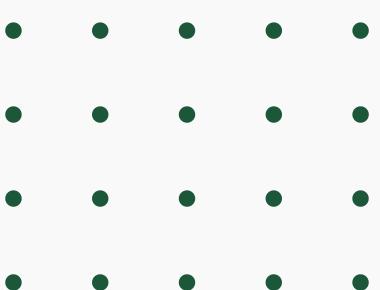




Project Exhibition - 1

Hybrid Machine Learning Algorithm for Credit Card Fraud Detection

Under the Guidance of Dr. Sandip Mal



Our Team



**ANJALI
KOLHATKAR**
21BCE11171



**SHEIKH MOHD
SHOAIB**
21BCE11346



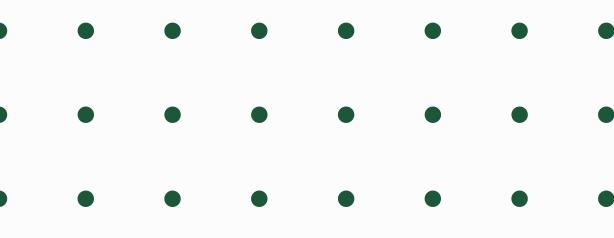
**SHUBHAM
SINGH**
21BCE10456



**HARSHITA
NIMJE**
21BCE10524



**AJINKIYA
FULLPATIL**
21BCE11160



Content

01
02
03
04
05
06
07
08

Credit Card Fraud Detection

Why It is Necessary?

Some ML Algorithms for fraud detection

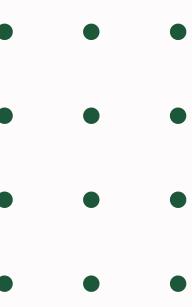
Hybrid Model For Fraud Detection

Architecture

Implementataon Of Hybrid Model

Evaluation and Results

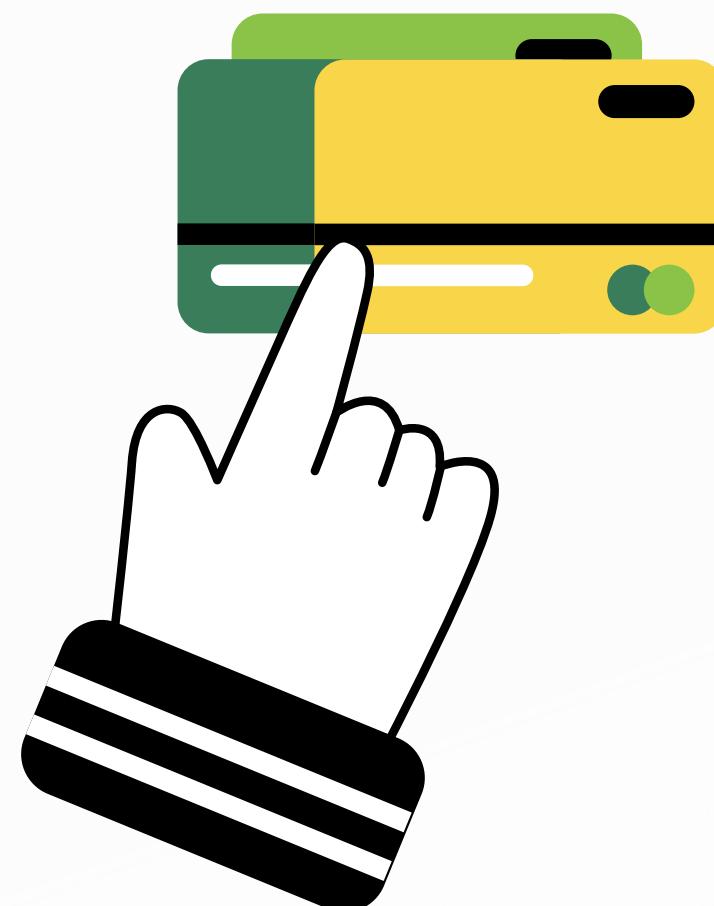
Conclusion



What is Credit Card Fraud?



- Credit card fraud refers to the **unauthorized or deceptive use** of someone else's credit card information for financial gain.
- It encompasses various illegal activities in which a person's credit card details, such as the **card number, expiration date, and security code**, are used by fraudsters to make unauthorized transactions or gain access to funds. This can include making unauthorized purchases, withdrawing cash, or transferring funds, all without the cardholder's consent.
- Credit card fraud is a widespread problem that can result in financial losses for individuals and damage to the reputation of financial institutions.



Why Credit Card Detection is Necessary ?

- Losses related to credit card fraud will grow to **\$43 billion** within five years and climb to **\$408.5 billion** globally within the next decade, according to a recent Nilson Report – meaning that credit card fraud detection has become more important than ever.
- In 2021, credit card fraud ranked as the **second most common type of identity theft in the U.S.**



ML Algorithms For Fraud Detection

Credit card fraud detection depends on various machine learning algorithms to identify fraud transactions. Here are some commonly used ML algorithms for credit card fraud detection:



Logistic Regression:

A classic binary classification algorithm that's used to predict whether a transaction is fraudulent or not based on historical data.

Random Forest:

An ensemble learning method that combines multiple decision trees to improve accuracy. It's effective at handling imbalanced datasets common in fraud detection.

Support Vector Machine (SVM):

SVM is used to classify transactions into fraudulent and non-fraudulent categories. It can handle high-dimensional data effectively.

Convolution Neural Networks:

Deep learning models like Convolutional Neural Networks (CNN) are used for pattern recognition and fraud detection.



Naive Bayes:

Naive Bayes classifiers are probabilistic models that can estimate the likelihood of a transaction being fraudulent.



Gradient Boosting:

Gradient boosting is a machine learning technique that combines multiple weak models, like decision trees, to create a strong predictive model. It's an iterative process where each model corrects the mistakes made by the previous ones.



Accuracy of Different Models

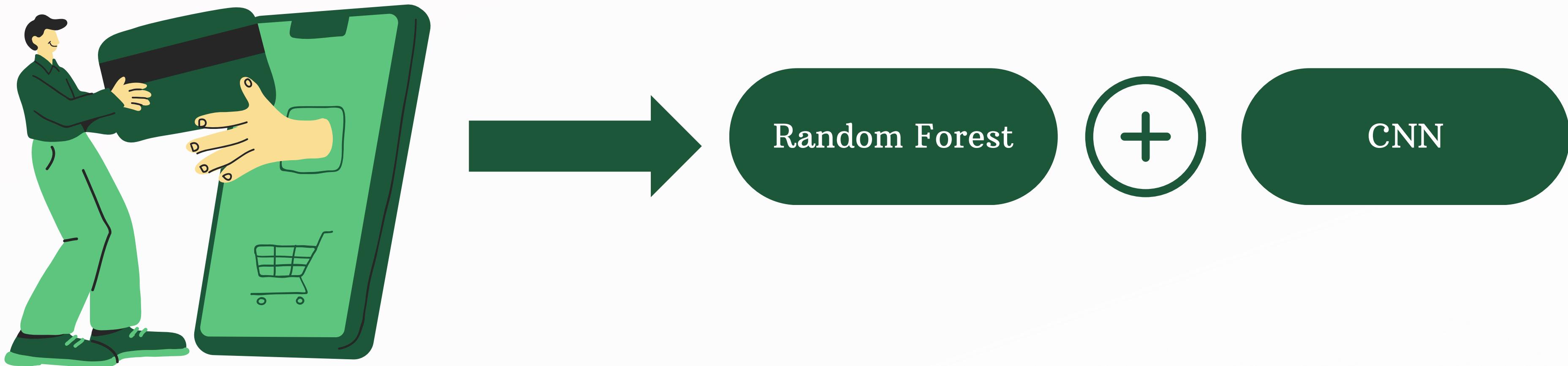
Algorithm	Accuracy	
Decision Tree	94%	
SVM Model	95.99%	
Logistic Regression	97.2%	
Gradient Boosting	99.801%	
CNN	99.901%	
Random Forest	99.90%	

The background image shows a modern office space with a high ceiling featuring exposed pipes and ductwork. The walls are covered in lush green plants and vines. Large windows provide natural light. The floor is made of light-colored wood. In the foreground, there are several wooden conference tables with black office chairs around them. A large sofa and armchairs are arranged in a seating area. A sign on the wall reads "752 Digital P".

Hybrid Model For fraud Detection

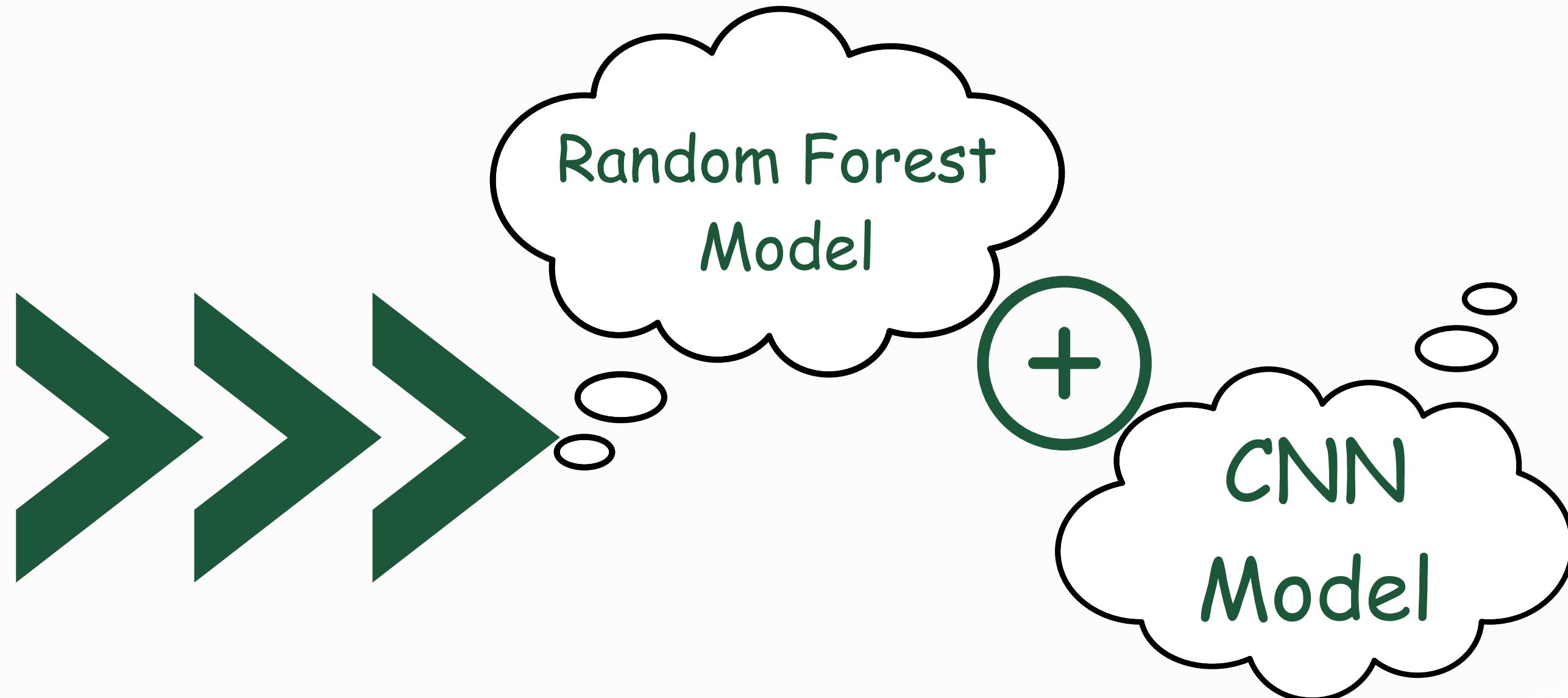
Hybrid Model

Hybrid machine learning combines different approaches such as supervised and unsupervised learning, reinforcement learning, or deep learning to address complex problems. It can help to overcome the limitations of individual machine learning algorithms.



Hybrid Model

Hybrid Model



Hybrid architectures can be more interpretable. This is because normal machine learning algorithms are often easier to interpret than deep learning algorithms. By combining a deep learning model with a traditional machine learning algorithm, it is possible to create a hybrid architecture that is both accurate and interpretable.

Architecture



Process Flow

01



Data Source And Analytics

https://docs.google.com/document/d/1DVAPJuDUlO2clF4iwRvesunIJM3mPhbUZBOE_vYiHK0/edit

It's important to check or analyse data so that you can remove any null value or unexpected values and errors. This is an important a step because the data we acquire could have these errors and it is important to remove it else the model will give errors. In this analyzation:

- Checking the data set to know the form of data entries and the number of row and column:

```
1 # I have called the file from its path
2 df = pd.read_csv("creditcard.csv")
3 df.head()

Time V1 V2 V3 V4 V5 V6 V7 V8 V9 ... V21 V22 V23 V24 V25 V26 V27 V28 Amount Class
0 -1.359607 -0.072781 2.536347 1.378155 -0.338321 0.462368 0.239599 0.098698 0.363787 ... -0.018307 0.277638 -0.110474 0.066928 0.128539 -0.189115 0.133568 -0.021053 149.82 0.0
1 0.1191867 0.266151 0.166480 0.448154 0.060018 -0.082361 -0.078803 0.085102 -0.265426 ... -0.225775 -0.638672 0.101288 -0.339846 0.167170 0.125896 -0.008983 0.014724 2.69 0.0
2 -1.358364 -1.340163 1.773209 0.379780 -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 378.86 0.0
3 -0.966272 -0.195226 1.792993 -0.063291 -0.010309 1.247203 0.237609 0.377436 -1.387024 ... -0.108300 0.005274 -0.190321 -1.176575 0.647376 -0.221929 0.062723 0.061468 123.50 0.0
4 -1.158233 0.877737 1.548718 0.403034 -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 0.0
5 rows × 31 columns

T 1 df.tail()

Time V1 V2 V3 V4 V5 V6 V7 V8 V9 ... V21 V22 V23 V24 V25 V26 V27 V28 Amount Class
21873 31906 1.285161 -0.701466 0.221870 -0.744753 -0.674967 0.034267 -0.631373 0.097860 -0.836749 ... 0.083025 0.190069 -0.074942 -0.278466 0.463616 -0.223877 0.011735 0.000997 39.75 0.0
21874 31906 -0.345979 1.094568 1.282208 0.068241 -0.016966 -0.989567 0.693968 -0.060856 -0.338656 ... -0.266852 -0.723026 -0.003364 0.322646 -0.173496 0.073145 0.241997 0.097610 5.99 0.0
21875 31907 1.171052 0.718042 -0.227521 1.373722 0.152349 -0.949040 0.271427 -0.153000 -0.184297 ... -0.013584 0.088769 -0.059149 0.278992 0.569996 -0.319374 0.056064 0.058347 1.00 0.0
21876 31907 1.209330 -1.130242 1.614261 0.000265 -1.974416 0.262320 -1.369564 0.173769 0.639865 ... -0.620157 -0.673057 0.037646 0.488680 0.161446 1.067798 0.026745 0.028549 40.00 0.0
21877 31907 -0.184653 0.376893 0.632988 -2.282365 0.403553 -0.371666 0.695359 -0.064835 0.610998 ... -0.130997 -0.312833 -0.156567 -0.899683 0.093186 -0.130187 0.029255 0.009226 NaN NaN
5 rows × 31 columns

[ ] 1 df.shape
(2878, 31)
```

Process Flow

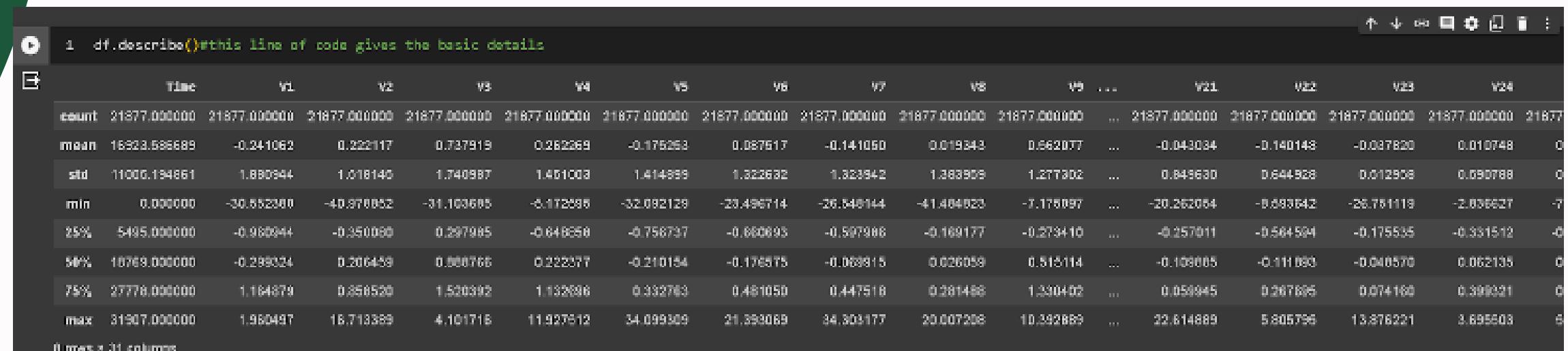
Data Analysis Processing

```
1 # cleaning the dataset
2 df.isnull().sum()# this line of code tell how many null values are there in a column
```

	0
0	0
1	0
2	0
3	0
4	0
5	0
6	0
7	0
8	0
9	0
10	0
11	0
12	0
13	0
14	0
15	0
16	0
17	0
18	0
19	0
20	0
21	0
22	0
23	0
24	0
25	0
26	0
27	0
28	0
29	0
30	0
31	0
32	0
33	0
34	0
35	0
36	0
37	0
38	0
39	0
40	0
41	0
42	0
43	0
44	0
45	0
46	0
47	0
48	0
49	0
50	0
51	0
52	0
53	0
54	0
55	0
56	0
57	0
58	0
59	0
60	0
61	0
62	0
63	0
64	0
65	0
66	0
67	0
68	0
69	0
70	0
71	0
72	0
73	0
74	0
75	0
76	0
77	0
78	0
79	0
80	0
81	0
82	0
83	0
84	0
85	0
86	0
87	0
88	0
89	0
90	0
91	0
92	0
93	0
94	0
95	0
96	0
97	0
98	0
99	0
100	0
101	0
102	0
103	0
104	0
105	0
106	0
107	0
108	0
109	0
110	0
111	0
112	0
113	0
114	0
115	0
116	0
117	0
118	0
119	0
120	0
121	0
122	0
123	0
124	0
125	0
126	0
127	0
128	0
129	0
130	0
131	0
132	0
133	0
134	0
135	0
136	0
137	0
138	0
139	0
140	0
141	0
142	0
143	0
144	0
145	0
146	0
147	0
148	0
149	0
150	0
151	0
152	0
153	0
154	0
155	0
156	0
157	0
158	0
159	0
160	0
161	0
162	0
163	0
164	0
165	0
166	0
167	0
168	0
169	0
170	0
171	0
172	0
173	0
174	0
175	0
176	0
177	0
178	0
179	0
180	0
181	0
182	0
183	0
184	0
185	0
186	0
187	0
188	0
189	0
190	0
191	0
192	0
193	0
194	0
195	0
196	0
197	0
198	0
199	0
200	0
201	0
202	0
203	0
204	0
205	0
206	0
207	0
208	0
209	0
210	0
211	0
212	0
213	0
214	0
215	0
216	0
217	0
218	0
219	0
220	0
221	0
222	0
223	0
224	0
225	0
226	0
227	0
228	0
229	0
230	0
231	0
232	0
233	0
234	0
235	0
236	0
237	0
238	0
239	0
240	0
241	0
242	0
243	0
244	0
245	0
246	0
247	0
248	0
249	0
250	0
251	0
252	0
253	0
254	0
255	0
256	0
257	0
258	0
259	0
260	0
261	0
262	0
263	0
264	0
265	0
266	0
267	0
268	0
269	0
270	0
271	0
272	0
273	0
274	0
275	0
276	0
277	0
278	0
279	0
280	0
281	0
282	0
283	0
284	0
285	0
286	0
287	0
288	0
289	0
290	0
291	0
292	0
293	0
294	0
295	0
296	0
297	0
298	0
299	0
300	0
301	0
302	0
303	0
304	0
305	0
306	0
307	0
308	0
309	0
310	0
311	0
312	0
313	0
314	0
315	0
316	0
317	0
318	0
319	0
320	0
321	0
322	0
323	0
324	0
325	0
326	0
327	0
328	0
329	0
330	0
331	0
332	0
333	0
334	0
335	0
336	0
337	0
338	0
339	0
340	0
341	0
342	0
343	0
344	0
345	0
346	0
347	0
348	0
349	0
350	0
351	0
352	0
353	0
354	0
355	0
356	0
357	0
358	0
359	0
360	0
361	0
362	0
363	0
364	0
365	0
366	0
367	0
368	0
369	0
370	0
371	0
372	0
373	0
374	0
375	0
376	0
377	0
378	0
379	0
380	0
381	0
382	0
383	0
384	0
385	0
386	0
387	0
388	0
389	0
390	0
391	0
392	0
393	0
394	0
395	0
396	0
397	0
398	0
399	0
400	0
401	0
402	0
403	0
404	0
405	0
406	0
407	0
408	0
409	0
410	0
411	0
412	0
413	0
414	0
415	0
416	0
417	0
418	0
419	0
420	0
421	0
422	0
423	0
424	0
425	0
426	0
427	0
428	

Process Flow

- Using describe function we will obtain the basic description about the data:



	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	...	V21	V22	V23	V24
count	21877.000000	21877.000000	21877.000000	21877.000000	21877.000000	21877.000000	21877.000000	21877.000000	21877.000000	21877.000000	...	21877.000000	21877.000000	21877.000000	21877.000000
mean	16323.586689	-0.241062	0.222117	0.737819	0.262369	-0.176263	0.087617	-0.141050	0.019343	0.562077	...	-0.043034	-0.140148	-0.037820	0.010748
std	11000.194861	1.880344	1.018140	1.740987	1.461003	1.414399	1.322632	1.323942	1.283959	1.277302	...	0.843630	0.644328	0.012956	0.090788
min	0.000000	-30.552300	-40.978052	-31.103605	-6.172585	-32.082129	-23.496714	-26.546144	-41.404823	-7.175097	...	-20.262054	-8.583642	-26.751119	-2.836627
25%	5495.000000	-0.960944	-0.350090	0.297985	-0.646550	-0.756737	-0.660693	-0.597986	-0.169177	-0.273410	...	-0.257011	-0.564594	-0.175535	-0.331512
50%	16769.000000	-0.298024	0.206459	0.888766	0.222377	-0.210154	-0.176575	-0.063915	0.026059	0.515114	...	-0.108605	-0.111080	-0.048570	0.062135
75%	27778.000000	1.164879	0.858520	1.520392	1.132896	0.332783	0.461050	0.447518	0.281488	1.330402	...	0.050945	0.267696	0.074160	0.390921
max	31907.000000	1.960497	16.713389	4.101716	11.937612	34.099369	21.393089	34.303177	29.007208	10.392889	...	22.614889	5.805796	13.378221	3.895803

- Getting the count of fraud count and non fraud count in the Class column:

```
[1]: 1 # getting the total of all the non-fraud entries in the dtas
2 non_fraud = len(df[df.Class == 0])
3 print(non_fraud)
4
52791
[2]: 1 # getting the total fraud form the complete entries
2 fraud = len(df[df.Class == 1])
3 print(fraud)
45
```

Process Flow

Data Source

- Normalizing the values which are varied from the range other columns are present.

```
[ ] 1 # here we will be using RFM which dont require standardise the data
[ ] 2 scaler = StandardScaler()
[ ]

[ ] 1
[2] df['Normalized_amount'] = scaler.fit_transform(df['Amount'].values.reshape(-1,1))

[ ] 1
[2] df.drop(['Amount', 'Time'], axis = 1, inplace = True)

[ ] 1 df.head()

      V2   V5   V6   V5   V6   V7   V8   V9   V38 ... V21   V22   V23   V24   V25   V26   V27   V28 class Normalized_amount
0 -0.072781 2.636347 1.378166 -0.338321 0.462368 0.239599 0.090698 0.363787 0.096794 ... -0.018107 0.277838 -0.110474 0.066928 0.126639 -0.189116 0.133568 -0.021053 0.0 0.379919
1 0.266151 0.166480 0.446154 0.060018 -0.082361 -0.079803 0.085102 -0.256425 -0.166974 ... -0.225775 -0.638672 0.101298 -0.339845 0.167170 0.125895 -0.008983 0.014724 0.0 -0.336789
2 -1.340163 1.773209 0.379780 -0.563199 1.800499 0.791461 0.247676 -1.514654 0.207643 ... 0.247998 0.771679 0.509412 -0.689281 -0.327642 -0.138097 -0.068363 -0.059752 0.0 1.497025
3 -0.105226 1.791293 -0.863291 -0.010309 1.241203 0.237609 0.377436 -1.387824 -0.054952 ... -0.106309 0.086274 -0.196321 -1.175675 0.647376 -0.221929 0.062723 0.061450 0.0 0.252623
4 0.877737 1.548718 0.403034 -0.407193 0.095921 0.592941 -0.270533 0.817739 0.751074 ... -0.029431 0.798278 -0.137456 0.141267 -0.206010 0.502292 0.219422 0.215153 0.0 -0.039484

[ ] 1
[2] df.drop(['Class'], axis = 1)
[3] y = df['Class']
```

Process Flow

02



- Data will be divided in dependable and independent parts namely “y” and “x” respectively. After this we will use these “x”(independent data) and “y”(dependent data) data.

```
] 1 # making all the independent variable
2 x = df.drop(["Class"], axis = 1)
3 y = df["Class"]

D 1 from sklearn.model_selection import train_test_split

] 1 # making the trian and test data by dividing the data with the testsize of 0.3
2 x_train, x_test, y_train, y_test = train_test_split(x,y, test_size = 0.3, random_state = 1)
```

Process Flow

Hybrid Model Making



03

- Here we will be making a hybrid model using the random forest method and CNN model, as these have the best accuracy when performed individually.

```
i cnn_model.summary()

Model: "sequential_1"
_________________________________________________________________
Layer (type)        Output Shape         Params
---                
conv2d_1 (Conv2D)  (None, 27, 32)      128
max_pooling2d_1 (MaxPooling2D) (None, 13, 4)    0
flatten_1 (Flatten) (None, 436)        0
dense_2 (Dense)   (None, 128)          53376
dense_3 (Dense)   (None, 1)            129
_________________________________________________________________
Total params: 53633 (209.50 KB)
Trainable params: 53633 (209.50 KB)
Non-trainable params: 0 (0.00 Byte)
```

Accuracy

99.969%

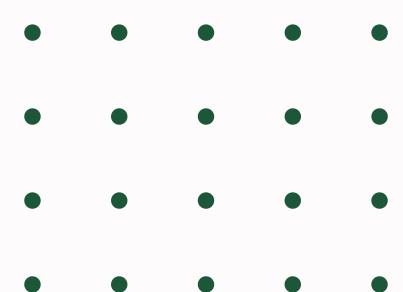
```
1 # Create the CNN model
2 cnn_model = Sequential()
3 # cnn_model = Sequential()
4 cnn_model.add(Conv1D(32, kernel_size=3, activation='relu', input_shape=(x_train.shape[0], x_train.shape[1], 1)))
5 cnn_model.add(MaxPooling1D(pool_size=2))
6 cnn_model.add(Flatten())
7 cnn_model.add(Dense(128, activation='relu'))
8 cnn_model.add(Dense(1, activation='sigmoid'))
9 # Create the random forest model
10 rf_model = RandomForestClassifier(n_estimators=100)
11
12 # Train the CNN model
13 cnn_model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
14 cnn_model.fit(x_train, y_train, epochs=10)
15
16 # Extract the features from the CNN model
17 cnn_features = cnn_model.predict(x_test)
18
19 # Train the random forest model on the CNN features
20 rf_model.fit(cnn_features, y_test)
21
22 # Make predictions on the test set
23 y_pred = rf_model.predict(cnn_model.predict(x_test))
24 accuracy = np.sum(y_pred == y_test) / len(y_test)
25 print('Accuracy:', accuracy*100)
```



```
Epoch 1/10
479/479 [=====] - 2s 3ms/step - loss: 0.0258 - accuracy: 0.9957
Epoch 2/10
479/479 [=====] - 2s 3ms/step - loss: 0.0076 - accuracy: 0.9978
Epoch 3/10
479/479 [=====] - 2s 5ms/step - loss: 0.0049 - accuracy: 0.9982
Epoch 4/10
479/479 [=====] - 2s 4ms/step - loss: 0.0050 - accuracy: 0.9981
Epoch 5/10
479/479 [=====] - 2s 4ms/step - loss: 0.0039 - accuracy: 0.9987
Epoch 6/10
479/479 [=====] - 2s 4ms/step - loss: 0.0036 - accuracy: 0.9999
Epoch 7/10
479/479 [=====] - 2s 3ms/step - loss: 0.0033 - accuracy: 0.9991
Epoch 8/10
479/479 [=====] - 2s 3ms/step - loss: 0.0030 - accuracy: 0.9993
Epoch 9/10
479/479 [=====] - 2s 3ms/step - loss: 0.0030 - accuracy: 0.9999
Epoch 10/10
479/479 [=====] - 2s 3ms/step - loss: 0.0029 - accuracy: 0.9999
206/206 [=====] - 1s 2ms/step
206/206 [=====] - 0s 2ms/step
Accuracy: 99.9693077391835
```

Accuracy of Hybrid Model

By utilizing random forest and gradient boosting algorithms for credit card fraud detection, we achieved an **accuracy rate of 99.969%**.



Conclusion

In conclusion we can say that, the **hybrid model, with an accuracy of 99.97%**, outperforms other algorithms that had an accuracy of up to 99.908% at detecting credit card fraud.

We have combined two powerful algorithms, Random forest and Gradient Boosting Algorithm, which resulted as **hybrid model being the best solution for detecting credit card fraud**.

Moreover ,We demonstrated that this hybrid model isn't just suitable for one dataset; it's an excellent solution for this type of problem since it can manage the variety of data types commonly present in fraud detection.



Thank you

