

# **FRA Project - Coded A**

**Finance and Risk Analytics**

Isha Shukla  
26 July 2024

## Contents

SL. No.	Title	Page No.
1	Exploratory Data Analysis	4
2	Data Pre-processing	9
3	Model Building	12
4	Model Performance Improvement	15
5	Model Performance Comparison and Final Model Selection	25
6	Actionable Insights & Recommendations	28

## Plots

SL. No.	Plots
1	Count of Default
2	Boxplot
3	Heatmap
4	ROC Curve
5	Feature Importances

## Part A

### Context

In the realm of modern finance, businesses encounter the perpetual challenge of managing debt obligations effectively to maintain a favorable credit standing and foster sustainable growth. Investors keenly scrutinize companies capable of navigating financial complexities while ensuring stability and profitability. A pivotal instrument in this evaluation process is the balance sheet, which provides a comprehensive overview of a company's assets, liabilities, and shareholder equity, offering insights into its financial health and operational efficiency. In this context, leveraging available financial data, particularly from preceding fiscal periods, becomes imperative for informed decision-making and strategic planning.

### Objective

A group of venture capitalists want to develop a Financial Health Assessment Tool. With the help of the tool, it endeavors to empower businesses and investors with a robust mechanism for evaluating the financial well-being and creditworthiness of companies. By harnessing machine learning techniques, they aim to analyze historical financial statements and extract pertinent insights to facilitate informed decision-making via the tool. Specifically, they foresee facilitating the following with the help of the tool:

1. Debt Management Analysis: Identify patterns and trends in debt management practices to assess the ability of businesses to fulfill financial obligations promptly and efficiently, and identify potential cases of default.
2. Credit Risk Evaluation: Evaluate credit risk exposure by analyzing liquidity ratios, debt-to-equity ratios, and other key financial indicators to ascertain the likelihood of default and inform investment decisions.

They have hired you as a data scientist and provided you with the financial metrics of different companies. The task is to analyze the data provided and develop a predictive model leveraging machine learning techniques to identify whether a given company will be tagged as a defaulter in terms of net worth next year. The predictive model will help the organization anticipate potential challenges with the financial performance of the companies and enable proactive risk mitigation strategies.

### Data Dictionary

The data consists of financial metrics from the balance sheets of different companies.

# Exploratory Data Analysis

## I. Shape of the data

(4256, 51)

## II. Data Information

- There are 50 features of float data type, 1 features of integer data type and 1 feature of object data types which I created.

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4256 entries, 0 to 4255
Data columns (total 51 columns):
#   Column                                     Non-Null Count  Dtype
---  -
0    Num                                         4256 non-null   int64
1    Networth Next Year                         4256 non-null   float64
2    Total assets                             4256 non-null   float64
3    Net worth                                4256 non-null   float64
4    Total income                             4025 non-null   float64
5    Change in stock                          3706 non-null   float64
6    Total expenses                          4091 non-null   float64
7    Profit after tax                         4102 non-null   float64
8    PBDITA                                   4102 non-null   float64
9    PBT                                       4102 non-null   float64
10   Cash profit                             4102 non-null   float64
11   PBDITA as % of total income              4177 non-null   float64
12   PBT as % of total income                 4177 non-null   float64
13   PAT as % of total income                 4177 non-null   float64
14   Cash profit as % of total income         4177 non-null   float64
15   PAT as % of net worth                    4256 non-null   float64
16   Sales                                    3951 non-null   float64
17   Income from fancial services             3145 non-null   float64
18   Other income                            2700 non-null   float64
19   Total capital                            4251 non-null   float64
20   Reserves and funds                      4158 non-null   float64
21   Borrowings                              3825 non-null   float64
22   Current liabilities & provisions          4146 non-null   float64
23   Deferred tax liability                  2887 non-null   float64
24   Shareholders funds                      4256 non-null   float64
25   Cumulative retained profits              4211 non-null   float64
26   Capital employed                        4256 non-null   float64
27   TOL/TNW                                4256 non-null   float64
28   Total term liabilities / tangible net worth 4256 non-null   float64
29   Contingent liabilities / Net worth (%)    4256 non-null   float64
30   Contingent liabilities                   2854 non-null   float64
31   Net fixed assets                        4124 non-null   float64
32   Investments                             2541 non-null   float64
33   Current assets                          4176 non-null   float64
34   Net working capital                     4219 non-null   float64
35   Quick ratio (times)                    4151 non-null   float64
36   Current ratio (times)                   4151 non-null   float64
37   Debt to equity ratio (times)            4256 non-null   float64
38   Cash to current liabilities (times)      4151 non-null   float64
39   Cash to average cost of sales per day    4156 non-null   float64
40   Creditors turnover                      3865 non-null   float64
41   Debtors turnover                        3871 non-null   float64
42   Finished goods turnover                 3382 non-null   float64
43   WIP turnover                           3492 non-null   float64
44   Raw material turnover                   3828 non-null   float64
45   Shares outstanding                      3446 non-null   float64
46   Equity face value                       3446 non-null   float64
47   EPS                                     4256 non-null   float64
48   Adjusted EPS                           4256 non-null   float64
49   Total liabilities                       4256 non-null   float64
50   PE on BSE                              1629 non-null   float64
dtypes: float64(50), int64(1)
memory usage: 1.7 MB
```

Data Info

- There were no duplicate values.

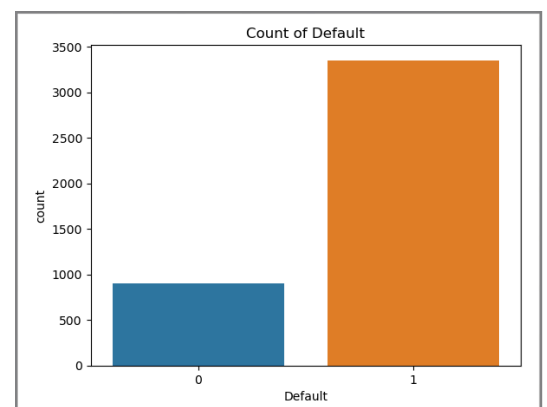
- Unique entries in the dataset.

Num	4256
Networth Next Year	2574
Total assets	2961
Net worth	2376
Total income	2870
Change in stock	1164
Total expenses	2898
Profit after tax	1467
PBDITA	1826
PBT	1568
Cash profit	1655
PBDITA as % of total income	2032
PBT as % of total income	1878
PAT as % of total income	1692
Cash profit as % of total income	1867
PAT as % of net worth	2385
Sales	2847
Income from fincial services	561
Other income	406
Total capital	1525
Reserves and funds	2361
Borrowings	2135
Current liabilities & provisions	2095
Deferred tax liability	950
Shareholders funds	2413
Cumulative retained profits	2265
Capital employed	2783
TOL/TNW	841
Total term liabilities / tangible net worth	508
Contingent liabilities / Net worth (%)	1926
Contingent liabilities	1351
Net fixed assets	2234
Investments	894
Current assets	2488
Net working capital	2065
Quick ratio (times)	409
Current ratio (times)	517
Debt to equity ratio (times)	642
Cash to current liabilities (times)	249
Cash to average cost of sales per day	2051
Creditors turnover	1608
Debtors turnover	1640
Finished goods turnover	2201
WIP turnover	1941
Raw material turnover	1601
Shares outstanding	2370
Equity face value	18
EPS	1815
Adjusted EPS	1730
Total liabilities	2961
PE on BSE	1142
dtype: int64	

Unique entries count in each column.

Created a default column.

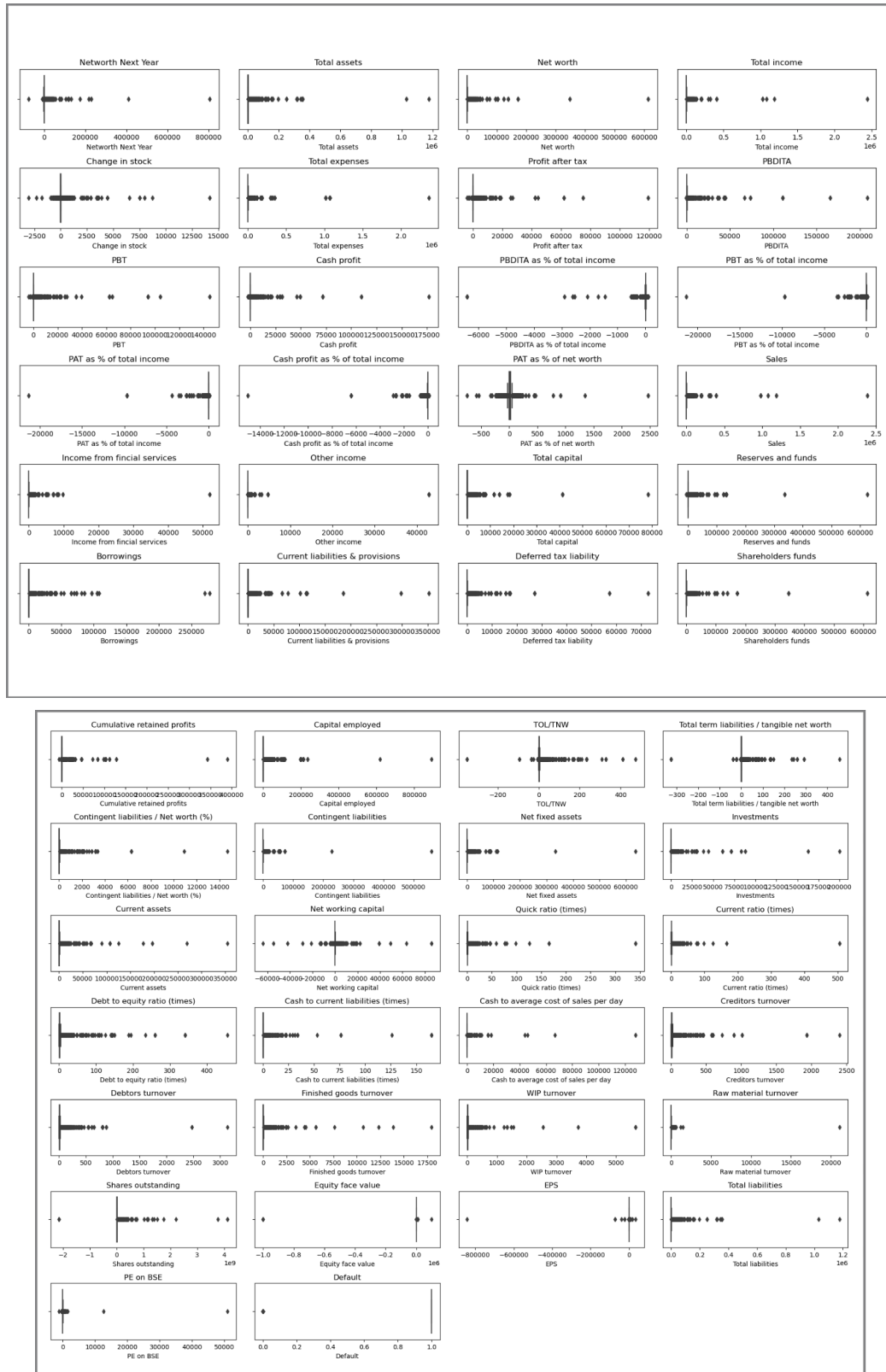
**III. Univariate analysis - Count of Default** - The number of instances where default did not occur (category '1') is significantly higher than the number of instances where default did occur (category '0'). This indicates that non-default cases are more prevalent than default cases within this dataset. A company will not be tagged as a defaulter if its net worth next year is positive, or else, it'll be tagged as a defaulter.



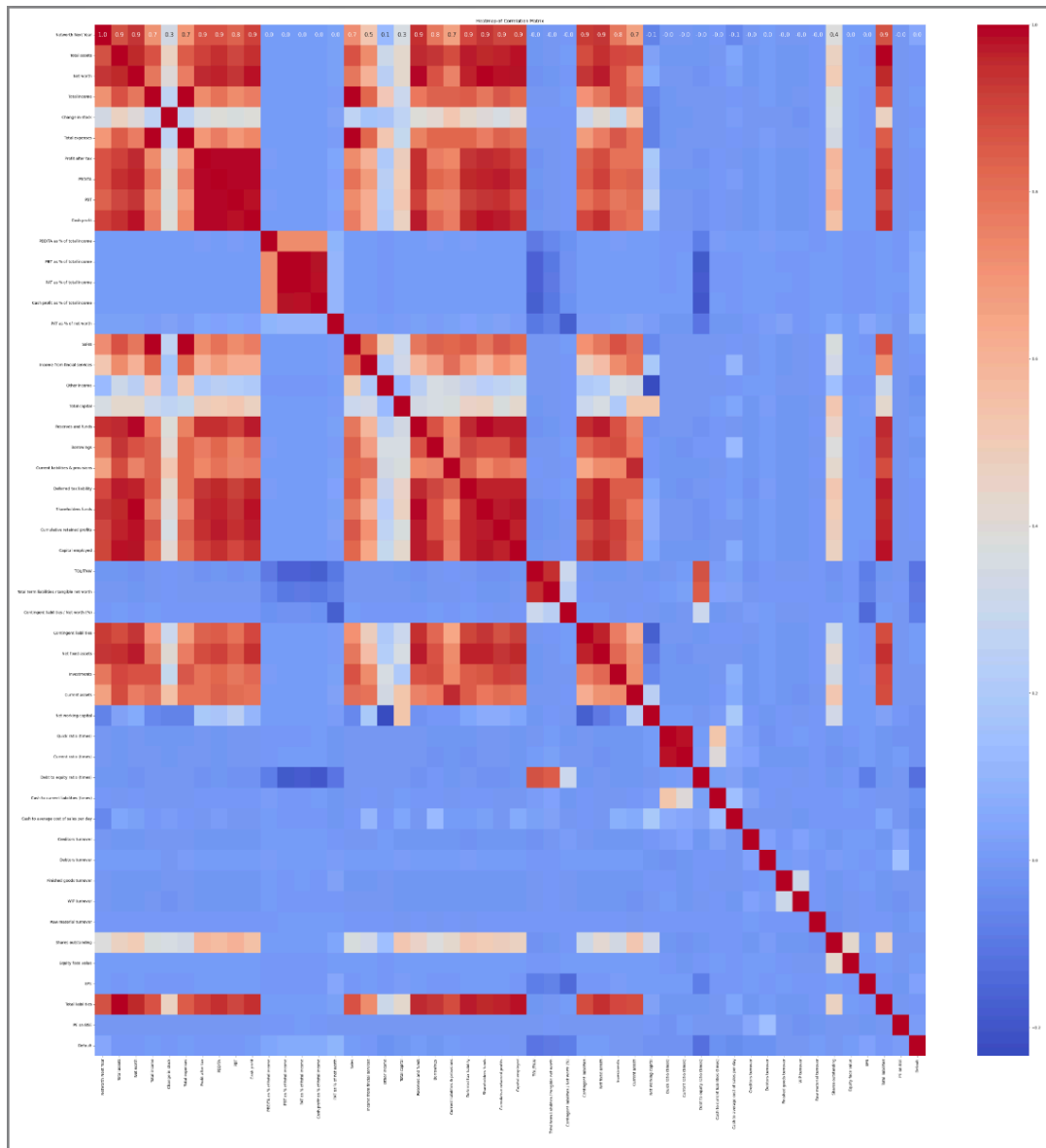
- Percentage of Default

Percentage of defaulters 78.76 %

- Boxplot



## IV. Bivariate Analysis



Heatmap

### 1. **Networth Next Year:**

- Highly positively correlated with Net Worth (0.930), Total Assets (0.878), and Profit After Tax (0.868).
- Moderate positive correlation with Total Income (0.711) and Total Expenses (0.691).
- Weak positive correlation with Change in Stock (0.345).

### 2. **Total Assets:**

- Highly positively correlated with Net Worth (0.959), PBDITA (0.943), and Cash Profit (0.940).

- Strong positive correlation with Profit After Tax (0.908), PBT (0.895), and Total Income (0.869).
3. **Net Worth:**
- Very high positive correlation with PBDITA (0.963), Cash Profit (0.978), and Profit After Tax (0.954).
  - Strong positive correlation with Total Assets (0.959), PBT (0.932), and Total Income (0.784).
4. **Total Income:**
- Extremely high positive correlation with Total Expenses (0.999).
  - High positive correlation with PBDITA (0.793), Cash Profit (0.763), and Net Worth (0.784).
  - Moderate positive correlation with Change in Stock (0.276).
5. **Change in Stock:**
- Weak positive correlation with Total Income (0.276), Total Assets (0.471), and Net Worth (0.394).
  - Moderate positive correlation with Total Expenses (0.274) and Profit After Tax (0.367).
6. **Total Expenses:**
- Extremely high positive correlation with Total Income (0.999).
  - High positive correlation with PBDITA (0.769), Profit After Tax (0.700), and Cash Profit (0.737).
7. **Profit After Tax:**
- Very high positive correlation with Cash Profit (0.990), PBDITA (0.990), and PBT (0.995).
  - Strong positive correlation with Net Worth (0.954), Total Assets (0.908), and Networth Next Year (0.868).
8. **PBDITA:**
- Very high positive correlation with Profit After Tax (0.990), PBT (0.989), and Cash Profit (0.992).
  - Strong positive correlation with Net Worth (0.963), Total Assets (0.943), and Total Income (0.793).
9. **PBT:**
- Very high positive correlation with Profit After Tax (0.995), PBDITA (0.989), and Cash Profit (0.978).
  - Strong positive correlation with Net Worth (0.932), Total Assets (0.895), and Networth Next Year (0.834).



#### 10. **Cash Profit:**

- Extremely high positive correlation with PBDITA (0.992), Profit After Tax (0.990), and PBT (0.978).
- Strong positive correlation with Net Worth (0.978), Total Assets (0.940), and Networth Next Year (0.907).

#### **Key Insights:**

- **Net Worth, Total Assets, and Profit After Tax** have the most significant positive impact on the **Networth Next Year**.
- **Total Income** and **Total Expenses** are almost perfectly correlated, indicating they move in tandem.
- **PBDITA, PBT, and Cash Profit** are highly correlated with each other and with **Net Worth** and **Total Assets**.
- **Change in Stock** has the weakest correlations among the metrics analyzed, indicating it has less impact on the other financial metrics.

### **V. Data Preprocessing**

#### **1. Drop Columns with Few Unique Values:**

- We can drop the columns Equity face value and Cash to current liabilities (times) as they have very few unique values

#### **2. Outliers Check:**

- Examine and address outliers in the dataset.

	Column	No. of outliers
0	Networth Next Year	624
1	Total assets	585
2	Net worth	595
3	Total income	508
4	Change in stock	750
5	Total expenses	518
6	Profit after tax	712
7	PBDITA	584
8	PBT	704
9	Cash profit	627

10	PBDITA as % of total income	346
11	PBT as % of total income	546
12	PAT as % of total income	610
13	Cash profit as % of total income	426
14	PAT as % of net worth	427
15	Sales	500
16	Income from fincial services	517
17	Other income	389
18	Total capital	551
19	Reserves and funds	643
20	Borrowings	532
21	Current liabilities & provisions	581
22	Deferred tax liability	406
23	Shareholders funds	588
24	Cumulative retained profits	699
25	Capital employed	572
26	TOL/TNW	414
27	Total term liabilities / tangible net worth	406
28	Contingent liabilities / Net worth (%)	478
29	Contingent liabilities	393

29	Contingent liabilities	393
30	Net fixed assets	569
31	Investments	451
32	Current assets	532
33	Net working capital	806
34	Quick ratio (times)	371
35	Current ratio (times)	397
36	Debt to equity ratio (times)	381
37	Cash to average cost of sales per day	583
38	Creditors turnover	442
39	Debtors turnover	408
40	Finished goods turnover	399
41	WIP turnover	378
42	Raw material turnover	296
43	Shares outstanding	476
44	EPS	638
45	Total liabilities	585
46	PE on BSE	237
47	Default	904

Outliers in each column

### 3. Data Preparation for Modeling:

- Separate the target variable (`default` column) from the rest of the data.

### 4. Split Data

- Divide the data into training and testing sets in the ratio 75:25.

## 5. Missing Values Detection and Treatment:

- Identify and handle missing values in the dataset.
- Missing value in **train** dataset
- Missing value in **test** dataset

Networth Next Year	0
Total assets	0
Net worth	0
Total income	177
Change in stock	424
Total expenses	125
Profit after tax	114
PBDITA	114
PBT	114
Cash profit	114
PBDITA as % of total income	65
PBT as % of total income	65
PAT as % of total income	65
Cash profit as % of total income	65
PAT as % of net worth	0
Sales	237
Income from fincial services	848
Other income	1180
Total capital	5
Reserves and funds	75
Borrowings	328
Current liabilities & provisions	88
Deferred tax liability	1043
Shareholders funds	0
Cumulative retained profits	34
Capital employed	0
TOL/TNW	0
Total term liabilities / tangible net worth	0
Contingent liabilities / Net worth (%)	0
Contingent liabilities	1061
Net fixed assets	100
Investments	1290
Current assets	65
Net working capital	33
Quick ratio (times)	84
Current ratio (times)	84
Debt to equity ratio (times)	0
Cash to average cost of sales per day	75
Creditors turnover	286
Debtors turnover	287
Finished goods turnover	662
WIP turnover	578
Raw material turnover	315
Shares outstanding	616
EPS	0
Total liabilities	0
PE on BSE	1992
dtype: int64	

Missing value in train dataset

Networth Next Year	0
Total assets	0
Net worth	0
Total income	54
Change in stock	126
Total expenses	40
Profit after tax	40
PBDITA	40
PBT	40
Cash profit	40
PBDITA as % of total income	14
PBT as % of total income	14
PAT as % of total income	14
Cash profit as % of total income	14
PAT as % of net worth	0
Sales	68
Income from fincial services	263
Other income	376
Total capital	0
Reserves and funds	23
Borrowings	103
Current liabilities & provisions	22
Deferred tax liability	326
Shareholders funds	0
Cumulative retained profits	11
Capital employed	0
TOL/TNW	0
Total term liabilities / tangible net worth	0
Contingent liabilities / Net worth (%)	0
Contingent liabilities	341
Net fixed assets	32
Investments	425
Current assets	15
Net working capital	4
Quick ratio (times)	21
Current ratio (times)	21
Debt to equity ratio (times)	0
Cash to average cost of sales per day	25
Creditors turnover	105
Debtors turnover	98
Finished goods turnover	212
WIP turnover	186
Raw material turnover	113
Shares outstanding	194
EPS	0
Total liabilities	0
PE on BSE	635
dtype: int64	

Missing value in test dataset

- Use KNN Imputer to replace missing values.

## 6. Scaling the Data:

- Apply `StandardScaler()` to standardize the dataset.

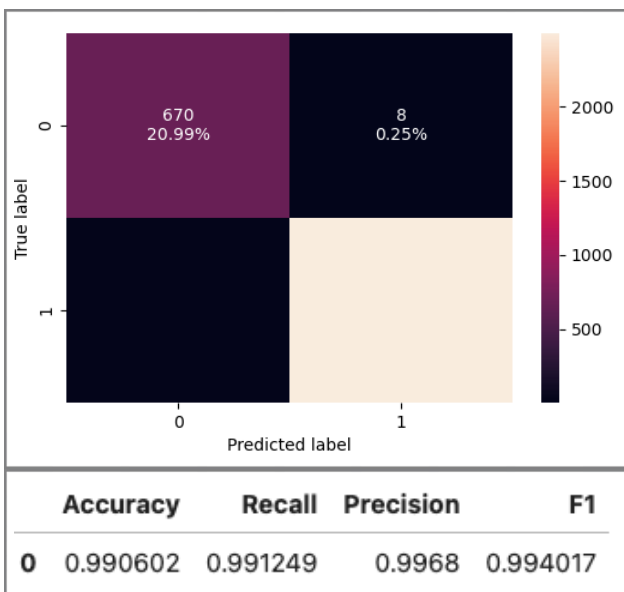
	Networth Next Year	Total assets	Net worth	Total income	Change in stock	Total expenses	Profit after tax	PBDITA	PBT	Cash profit	...	average cost of sales per day	Creditors turnover	Debt turnover
0	-0.069510	-0.093515	-0.090218	-0.063618	-0.082629	-0.062063	-0.080287	-0.085284	-0.080341	-0.080944	...	-0.059548	-0.165366	-0.135366
1	-0.071415	-0.094434	-0.094256	-0.066741	-0.134937	-0.065198	-0.087894	-0.092591	-0.089593	-0.088819	...	-0.053807	-0.070216	-0.142807
2	-0.057220	-0.088881	-0.076723	-0.074186	-0.061308	-0.073015	-0.078668	-0.080500	-0.077617	-0.071464	...	-0.058170	-0.063630	-0.129517
3	-0.075341	-0.081726	-0.096148	-0.071645	-0.105088	-0.070028	-0.089451	-0.086834	-0.088940	-0.079249	...	-0.051081	-0.158920	-0.031101

PBDITA	PBT	Cash profit	...	Cash to average cost of sales per day	Creditors turnover	Debtors turnover	Finished goods turnover	WIP turnover	Raw material turnover	Shares outstanding	EPS	Total liabilities	PE on BSE
-0.087066	-0.091867	-0.071213	...	-0.001312	-0.135427	-0.095051	-0.155083	-0.139530	-0.395188	-0.099742	0.024322	-0.102046	-0.235560
-0.140266	-0.117507	-0.134049	...	-0.127513	-0.098252	-0.189336	-0.098467	-0.122520	1.939744	-0.074977	0.015168	-0.164107	1.507136
-0.124323	-0.101013	-0.120684	...	-0.130336	-0.151305	-0.169148	-0.010361	-0.116255	0.033344	-0.075071	0.018001	-0.155785	0.232957
-0.117828	-0.099445	-0.111102	...	-0.105261	-0.097218	-0.145823	-0.078031	-0.042809	0.626194	-0.094305	0.018693	-0.113262	-0.094247

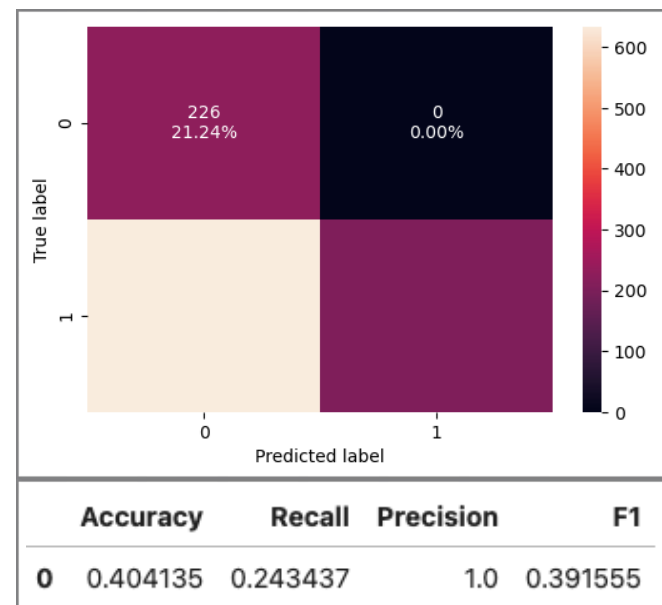
Scaled dataset for test

## VI. Model Building

### 1. Logistic Regression



Training Performance



Testing Performance

Confusion Matrix:  
[[ 670 8]  
[ 22 2492]]  
True Negatives (TN): 670 (20.99%)  
False Positives (FP): 8 (0.25%)  
False Negatives (FN): 22 (0.69%)  
True Positives (TP): 2492 (78.07%)

Confusion Matrix:  
[[226 0]  
[634 204]]  
True Negatives (TN): 226 (21.24%)  
False Positives (FP): 0 (0.00%)  
False Negatives (FN): 634 (59.59%)  
True Positives (TP): 204 (19.17%)

## Summary

Optimization terminated successfully.						
Current function value: 0.034979						
Iterations: 749						
Function evaluations: 850						
Gradient evaluations: 850						
Logit Regression Results						
=====						
Dep. Variable:	Default	No. Observations:	3192			
Model:	Logit	Df Residuals:	3145			
Method:	MLE	Df Model:	46			
Date:	Sat, 27 Jul 2024	Pseudo R-squ.:	0.9324			
Time:	20:00:24	Log-Likelihood:	-111.65			
converged:	True	LL-Null:	-1650.7			
Covariance Type:	nonrobust	LLR p-value:	0.000			
=====						
	coef	std err	z	P> z	[0.025	0.975]
-----						
const	1887.9543	221.952	8.506	0.000	1452.936	2322.972
Networth Next Year	2.609e+04	2982.784	8.746	0.000	2.02e+04	3.19e+04
Total assets	167.6406	6.36e+06	2.64e-05	1.000	-1.25e+07	1.25e+07
Net worth	-696.3138	1128.393	-0.617	0.537	-2907.924	1515.296
Total income	148.2940	4587.390	0.032	0.974	-8842.826	9139.414
Change in stock	-3.1482	28.423	-0.111	0.912	-58.856	52.559
Total expenses	59.8403	4513.136	0.013	0.989	-8785.744	8905.425
Profit after tax	484.6741	346.252	1.400	0.162	-193.966	1163.315
PBDITA	-1093.5041	313.608	-3.487	0.000	-1708.164	-478.844
PBT	-194.8066	337.367	-0.577	0.564	-856.033	466.420
Cash profit	-265.0958	299.712	-0.885	0.376	-852.520	322.328
PBDITA as % of total income	0.2883	0.432	0.667	0.505	-0.559	1.135
PBT as % of total income	0.1212	4.023	0.030	0.976	-7.764	8.007
PAT as % of total income	1.0255	4.082	0.251	0.802	-6.974	9.025
Cash profit as % of total income	-1.2495	0.646	-1.935	0.053	-2.515	0.016
PAT as % of net worth	0.1374	0.230	0.598	0.550	-0.313	0.588
Sales	152.5703	1135.485	0.134	0.893	-2072.939	2378.079
Income from fincial services	-287.7507	157.728	-1.824	0.068	-596.893	21.391
Other income	-37.3978	109.176	-0.343	0.732	-251.379	176.583
Total capital	68.0306	58.146	1.170	0.242	-45.934	181.995
Reserves and funds	766.2684	629.738	1.217	0.224	-467.995	2000.532
Borrowings	227.7502	540.589	0.421	0.674	-831.784	1287.285
Current liabilities & provisions	21.8777	410.301	0.053	0.957	-782.297	826.052
Deferred tax liability	-41.0763	111.638	-0.368	0.713	-259.882	177.730
Shareholders funds	-686.5134	1640.700	-0.418	0.676	-3902.227	2529.200
Cumulative retained profits	-103.2423	192.535	-0.536	0.592	-480.605	274.120
Capital employed	-332.2072	1741.046	-0.191	0.849	-3744.594	3080.180
TOL/TNW	-0.4120	0.608	-0.677	0.498	-1.604	0.780
Total term liabilities / tangible net worth	-0.1837	0.769	-0.239	0.811	-1.691	1.323
Contingent liabilities / Net worth (%)	-0.3469	0.418	-0.829	0.407	-1.167	0.473
Contingent liabilities	-33.1874	119.670	-0.277	0.782	-267.737	201.362
Net fixed assets	38.4789	99.475	0.387	0.699	-156.488	233.445
Investments	-212.9387	175.492	-1.213	0.225	-556.896	131.019
Current assets	-167.3960	170.397	-0.982	0.326	-501.369	166.576
Net working capital	21.5820	58.889	0.366	0.714	-93.837	137.002
Quick ratio (times)	0.0189	1.137	0.017	0.987	-2.210	2.248
Current ratio (times)	-0.2745	1.230	-0.223	0.823	-2.685	2.136
Debt to equity ratio (times)	0.5928	0.742	0.799	0.424	-0.861	2.047
Cash to average cost of sales per day	0.1092	0.463	0.236	0.813	-0.797	1.016
Creditors turnover	-0.3688	0.413	-0.892	0.372	-1.179	0.442
Debtors turnover	-0.5449	0.920	-0.592	0.554	-2.349	1.259
Finished goods turnover	0.7461	1.805	0.413	0.679	-2.793	4.285
WIP turnover	-2.0273	1.235	-1.641	0.101	-4.448	0.394
Raw material turnover	-3.4472	0.682	-5.053	0.000	-4.784	-2.110
Shares outstanding	0.1149	0.824	0.139	0.889	-1.501	1.731
EPS	-0.0265	4.096	-0.006	0.995	-8.054	8.001
Total liabilities	167.6406	6.36e+06	2.64e-05	1.000	-1.25e+07	1.25e+07
PE on BSE	-0.0685	0.200	-0.342	0.732	-0.461	0.324
=====						

### Logit Regression Results

## Observation

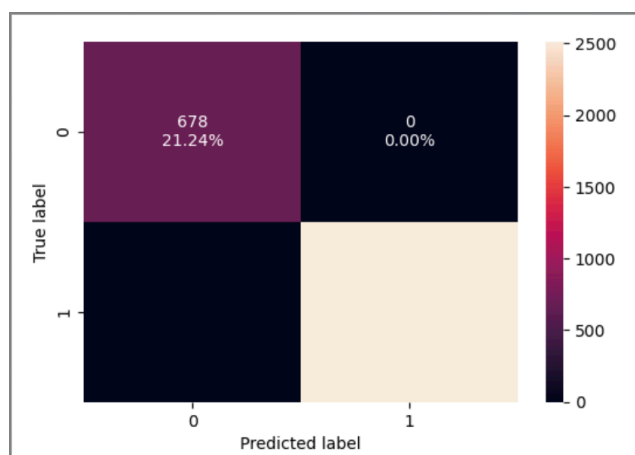
### Logistic Regression Training Performance:

- **Accuracy:** 0.990602
- **Recall:** 0.991249
- **Precision:** 0.9968
- **F1 Score:** 0.994017

## Logistic Regression Testing Performance:

- **Accuracy:** 0.404135
- **Recall:** 0.243437
- **Precision:** 1.0
- **F1 Score:** 0.391555

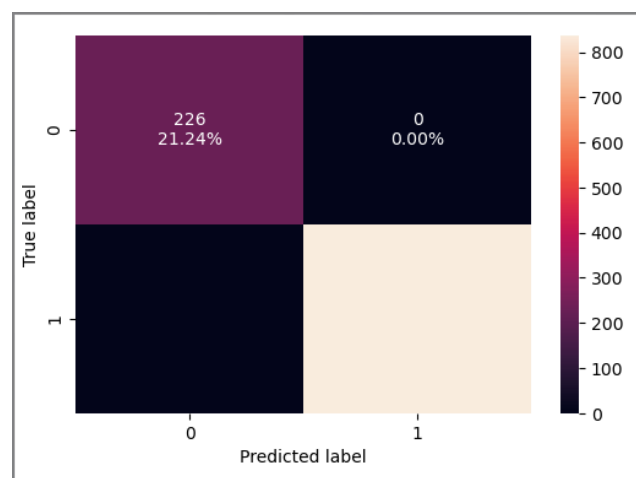
## 2. Random Forest



	Accuracy	Recall	Precision	F1
0	1.0	1.0	1.0	1.0

Training Performance

```
Confusion Matrix:
[[ 678   0]
 [   0 2514]]
True Negatives (TN): 678
False Positives (FP): 0
False Negatives (FN): 0
True Positives (TP): 2514
True Negatives (TN): 678 (21.24%)
False Positives (FP): 0 (0.00%)
False Negatives (FN): 0 (0.00%)
True Positives (TP): 2514 (78.76%)
```



	Accuracy	Recall	Precision	F1
0	1.0	1.0	1.0	1.0

Testing Performance

```
Confusion Matrix:
[[226   0]
 [   0 838]]
True Negatives (TN): 226
False Positives (FP): 0
False Negatives (FN): 0
True Positives (TP): 838
True Negatives (TN): 226 (21.24%)
False Positives (FP): 0 (0.00%)
False Negatives (FN): 0 (0.00%)
True Positives (TP): 838 (78.76%)
```

## VII. Model Performance Improvement

Variance Inflation Factor (VIF) is a measure of multicollinearity in a set of multiple regression variables. A high VIF indicates that the associated independent variable is highly collinear with the other variables in the model. Here are the steps to calculate VIF using statsmodels and pandas:

1. Fit the OLS model: We fit an Ordinary Least Squares (OLS) model to each independent variable against all the other independent variables.
2. Calculate VIF: The VIF for each variable is calculated using the formula:

$$VIF = 1 / (1 - R^2)$$

where  $R^2$  is the coefficient of determination of the regression of that variable against all the other variables.

The output `high_vif_columns` contains a list of variables that have a Variance Inflation Factor (VIF) greater than or equal to 5. This indicates that these variables are highly collinear with other independent variables in the dataset.

### Dropping the columns which have VIF > 5:

```
['Networth Next Year',
 'Total assets',
 'Net worth',
 'Total income',
 'Change in stock',
 'Total expenses',
 'Profit after tax',
 'PBDITA',
 'PBT',
 'Cash profit',
 'PBDITA as % of total income',
 'PBT as % of total income',
 'PAT as % of total income',
 'Cash profit as % of total income',
 'Sales',
 'Income from financial services',
 'Other income',
 'Total capital',
 'Reserves and funds',
 'Borrowings',
 'Current liabilities & provisions',
 'Deferred tax liability',
 'Shareholders funds',
 'Cumulative retained profits',
 'Capital employed',
 'TOL/TNW',
 'Total term liabilities / tangible net worth',
 'Contingent liabilities',
 'Net fixed assets',
 'Investments',
 'Current assets',
 'Net working capital',
 'Quick ratio (times)',
 'Current ratio (times)',
 'Debt to equity ratio (times)',
 'Total liabilities']
```

Dropping these columns which have VIF > 5

Based on the Variance Inflation Factor (VIF) analysis, columns with VIF values greater than 5 have been identified and subsequently dropped to address multicollinearity concerns. The following columns have been removed from the dataset: Networth Next Year, Total assets, Net worth, Total income, Change in stock, Total expenses, Profit after tax, PBDITA, PBT, Cash profit, PBDITA as % of total income, PBT as % of total income, PAT as % of total income, Cash profit as % of total income, Sales, Income from financial services, Other income, Total capital, Reserves and funds, Borrowings, Current liabilities & provisions, Deferred tax liability, Shareholders funds, Cumulative retained profits, Capital employed, TOL/TNW, Total term liabilities / tangible net worth, Contingent liabilities, Net fixed assets, Investments, Current assets, Net working capital, Quick ratio (times), Current ratio (times), Debt to equity ratio (times), and Total liabilities. By dropping these columns, the dataset is now optimized for further analysis, ensuring that multicollinearity does not adversely affect the results. This refined dataset will allow for more accurate and reliable statistical modeling and interpretation.

Shape of the data after dropping columns which had vif > 5:

(3192, 11)

Shape of scaled train  
dataset

(1064, 11)

Shape of scaled test  
dataset

Optimal Threshold value for Improved logistic regression **0.779**.



Optimization terminated successfully.  
Current function value: 0.505812  
Iterations: 42  
Function evaluations: 43  
Gradient evaluations: 43

Logit Regression Results

Dep. Variable:	Default	No. Observations:	3192
Model:	Logit	Df Residuals:	3180
Method:	MLE	Df Model:	11
Date:	Sat, 27 Jul 2024	Pseudo R-squ.:	0.02188
Time:	20:05:09	Log-Likelihood:	-1614.6
converged:	True	LL-Null:	-1650.7
Covariance Type:	nonrobust	LLR p-value:	4.578e-11

	coef	std err	z	P> z	[0.025	0.975]
const	1.3928	0.057	24.373	0.000	1.281	1.505
PAT as % of net worth	0.4012	0.077	5.193	0.000	0.250	0.553
Contingent liabilities / Net worth (%)	-0.1327	0.048	-2.764	0.006	-0.227	-0.039
Cash to average cost of sales per day	-0.1351	0.068	-1.980	0.048	-0.269	-0.001
Creditors turnover	-0.0439	0.038	-1.158	0.247	-0.118	0.030
Debtors turnover	-0.0514	0.037	-1.390	0.164	-0.124	0.021
Finished goods turnover	-0.0093	0.043	-0.216	0.829	-0.094	0.075
WIP turnover	-0.0216	0.040	-0.543	0.587	-0.100	0.056
Raw material turnover	1.5052	0.897	1.678	0.093	-0.253	3.263
Shares outstanding	-0.0286	0.041	-0.694	0.488	-0.109	0.052
EPS	0.0673	0.165	0.407	0.684	-0.257	0.391
PE on BSE	-0.0382	0.044	-0.872	0.383	-0.124	0.048

#### Logit Regression Improved Summary

#### Model Summary:

The optimization of the logistic regression model has successfully terminated, with the model converging after 42 iterations. The current function value stands at 0.505812, indicating the log-likelihood of the final model. Here's the summary of the logistic regression results:

- **Dependent Variable:** Default
- **Number of Observations:** 3192
- **Method:** Maximum Likelihood Estimation (MLE)
- **Log-Likelihood:** -1614.6
- **Pseudo R-squared:** 0.021882.

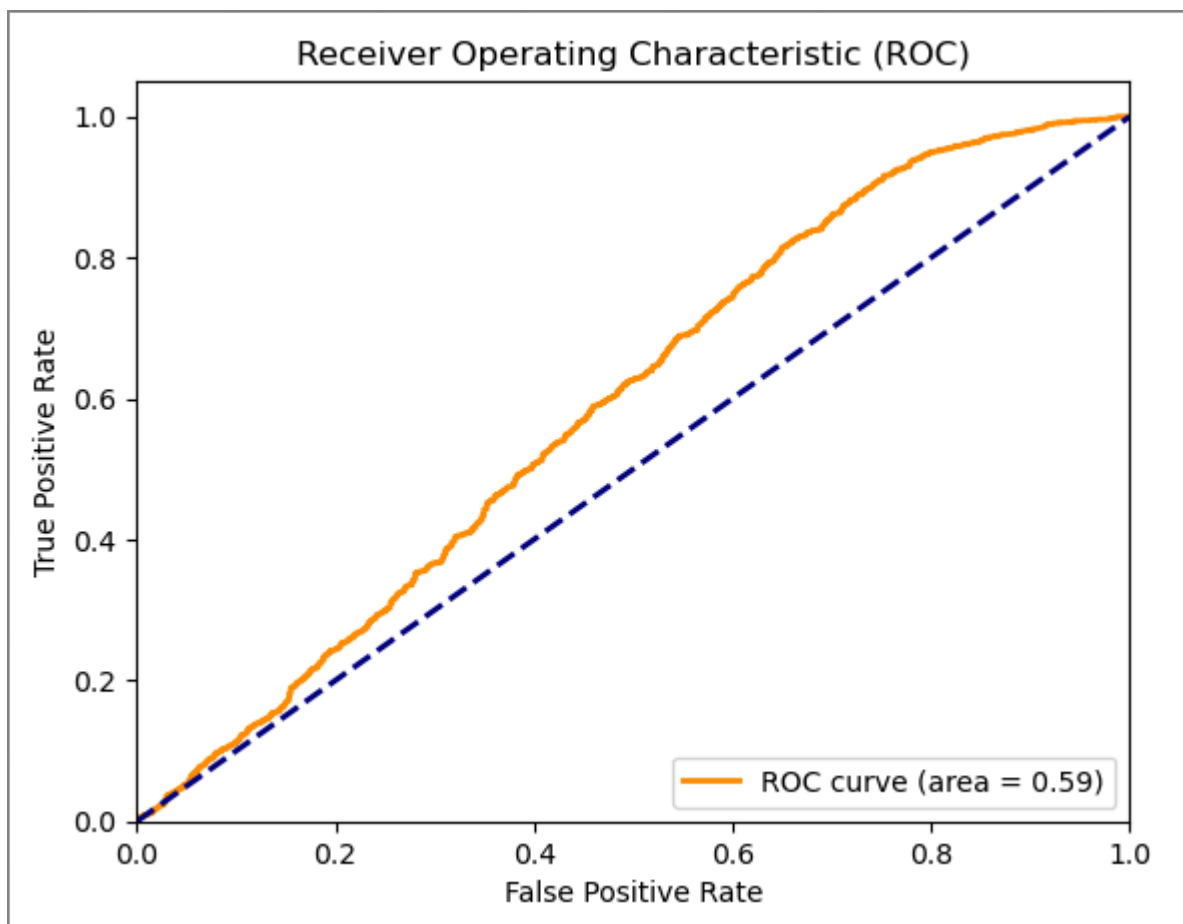
#### Significant Variables:

1. **PAT as % of net worth:** Positive coefficient (0.4012) and highly significant ( $p < 0.001$ ), indicating that higher profitability as a percentage of net worth increases the likelihood of default.

2. **Contingent liabilities / Net worth (%):** Negative coefficient (-0.1327) and significant ( $p = 0.006$ ), suggesting that higher contingent liabilities relative to net worth decrease the likelihood of default.

3. **Cash to average cost of sales per day:** Negative coefficient (-0.1351) and marginally significant ( $p = 0.048$ ), implying that better liquidity reduces the probability of default.

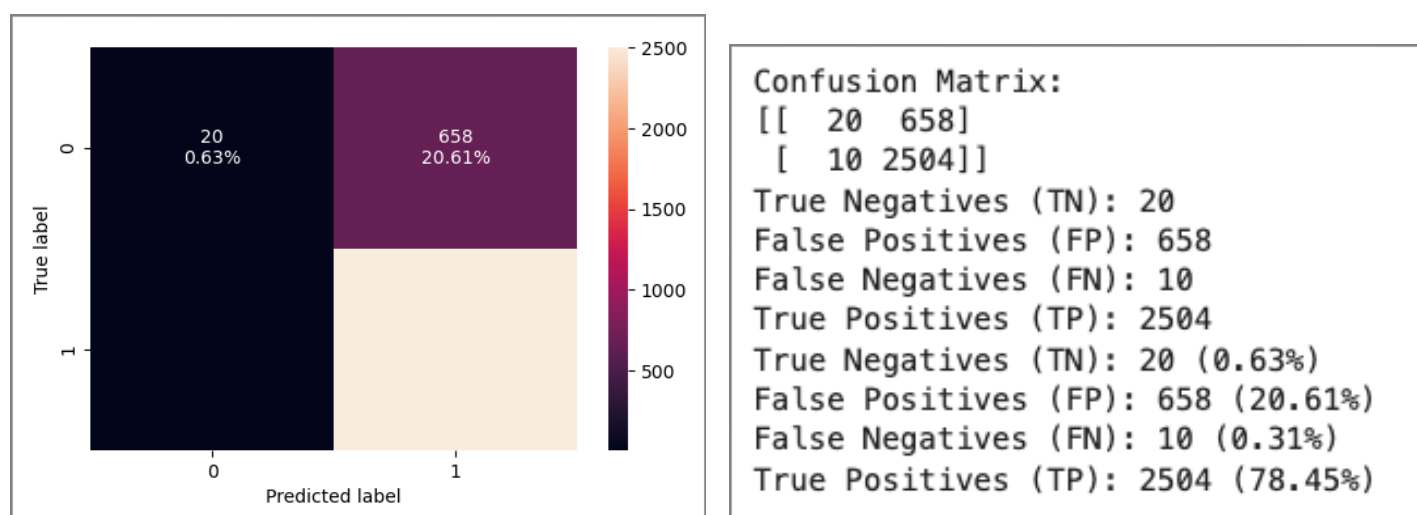
4. **Creditors turnover, Debtors turnover, Finished goods turnover, WIP turnover, Shares outstanding, EPS, and PE on BSE:** Not statistically significant, indicating these variables do not have a strong effect on default likelihood in this model.
5. **Raw material turnover:** Positive coefficient (1.5052) with marginal significance ( $p = 0.093$ ), which might suggest a relationship where higher raw material turnover could increase the risk of default, though this is less conclusive.
6. **Other Observations:**
  - **Pseudo R-squared (0.02188):** This value indicates the proportion of variance in the dependent variable that is explained by the independent variables. While this value is relatively low, it is not uncommon in logistic regression models dealing with complex, real-world data.



ROC Curve

The ROC (Receiver Operating Characteristic) curve is a graphical representation of a binary classification model's performance, plotting the True Positive Rate (TPR) against the False Positive Rate (FPR) at various threshold settings. The Area Under the ROC Curve (AUC) is a single value that ranges from 0 to 1 and summarizes the model's ability to distinguish between the positive and negative classes. An AUC of 1 indicates a perfect model, while an AUC of 0.5 signifies no discriminative power, equivalent to random guessing. An AUC of 0.59, as seen in our model, suggests that the model has a modest ability to differentiate between the classes, performing better than random guessing but still leaving room for significant improvement.

## Logistic Regression Performance - Training Set



- **True Negatives (TN): 20 (0.63%):** The model correctly predicted 20 instances as negative out of the total predictions.
- **False Positives (FP): 658 (20.61%):** The model incorrectly predicted 658 instances as positive when they were actually negative.
- **False Negatives (FN): 10 (0.31%):** The model incorrectly predicted 10 instances as negative when they were actually positive.
- **True Positives (TP): 2504 (78.45%):** The model correctly predicted 2504 instances as positive out of the total predictions.

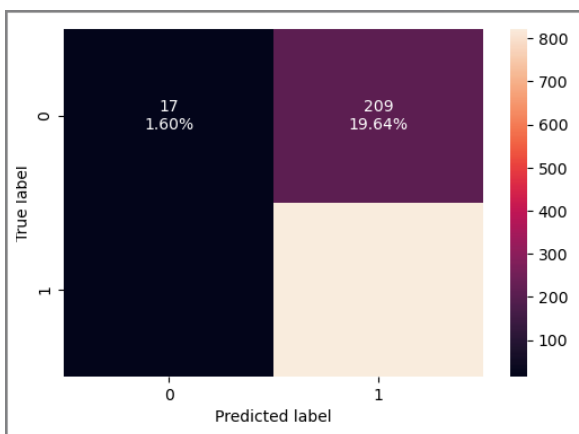
	Accuracy	Recall	Precision	F1
0	0.790727	0.996022	0.791904	0.882311

Model Performance

The accuracy of the model is approximately 79%, indicating the overall correctness of the model's predictions. The recall is very high at 99.60%, showing that the model is excellent at identifying positive instances. Precision is about 79.19%, meaning that when the model predicts a positive instance, it is correct 79.19% of the time. The F1 score, which is the harmonic mean of precision and recall, is 88.23%, suggesting a good balance between precision and recall.

In summary, while the model demonstrates strong recall and a reasonable F1 score, the precision and accuracy indicate there is room for improvement, especially considering the relatively high rate of false positives (20.61%).

## Logistic Regression Performance - Test Set



Confusion Matrix:  
[[ 17 209]  
[ 16 822]]  
True Negatives (TN): 17  
False Positives (FP): 209  
False Negatives (FN): 16  
True Positives (TP): 822  
True Negatives (TN): 17 (1.60%)  
False Positives (FP): 209 (19.64%)  
False Negatives (FN): 16 (1.50%)  
True Positives (TP): 822 (77.26%)

	Accuracy	Recall	Precision	F1
0	0.788534	0.980907	0.797284	0.879615

Test Set

The performance of the logistic regression model on the training set is summarized by key metrics: Accuracy, Recall, Precision, and F1 Score.

- **Accuracy:** 0.788534 (78.85%)
  - This metric represents the proportion of total correct predictions (both true positives and true negatives) made by the model. An accuracy of 78.85% indicates that the model correctly predicted the outcome for approximately 79% of the instances in the training set.
- **Recall (Sensitivity):** 0.980907 (98.09%)
  - Recall measures the model's ability to identify true positive cases. With a recall of 98.09%, the model successfully detected almost all actual positive cases in the training set.
- **Precision:** 0.797284 (79.73%)
  - Precision indicates the proportion of positive predictions that are actually correct. A precision of 79.73% means that when the model predicts a positive outcome, it is correct approximately 80% of the time.
- **F1 Score:** 0.879615 (87.96%)
  - The F1 score is the harmonic mean of precision and recall, providing a single metric that balances the two. An F1 score of 87.96% suggests that the model has a good balance between identifying positive cases and ensuring that positive predictions are correct.

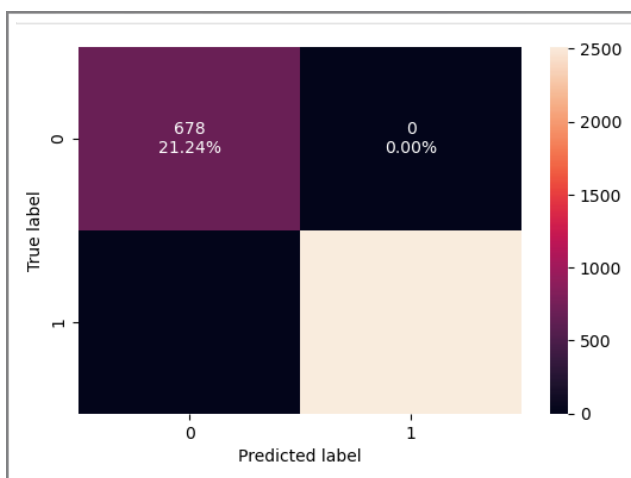
In summary, the logistic regression model shows strong recall and a good F1 score, indicating it effectively identifies positive instances while maintaining a reasonable balance with precision. However, the accuracy and precision suggest there is room for improvement, particularly in reducing false positives and enhancing overall prediction correctness.

## Model Performance Improvement - Random Forest

```
Parameters used in the Random Forest Classifier:  
bootstrap: True  
ccp_alpha: 0.0  
class_weight: balanced  
criterion: gini  
max_depth: 7  
max_features: sqrt  
max_leaf_nodes: None  
max_samples: None  
min_impurity_decrease: 0.0  
min_samples_leaf: 5  
min_samples_split: 2  
min_weight_fraction_leaf: 0.0  
n_estimators: 200  
n_jobs: None  
oob_score: False  
random_state: 42  
verbose: 0  
warm_start: False
```

Parameters used in the Random Forest Classifier

## Random Forest Performance - Training Set



Confusion Matrix:  
[[ 678 0]  
 [ 0 2514]]  
True Negatives (TN): 678  
False Positives (FP): 0  
False Negatives (FN): 0  
True Positives (TP): 2514  
True Negatives (TN): 678 (21.24%)  
False Positives (FP): 0 (0.00%)  
False Negatives (FN): 0 (0.00%)  
True Positives (TP): 2514 (78.76%)

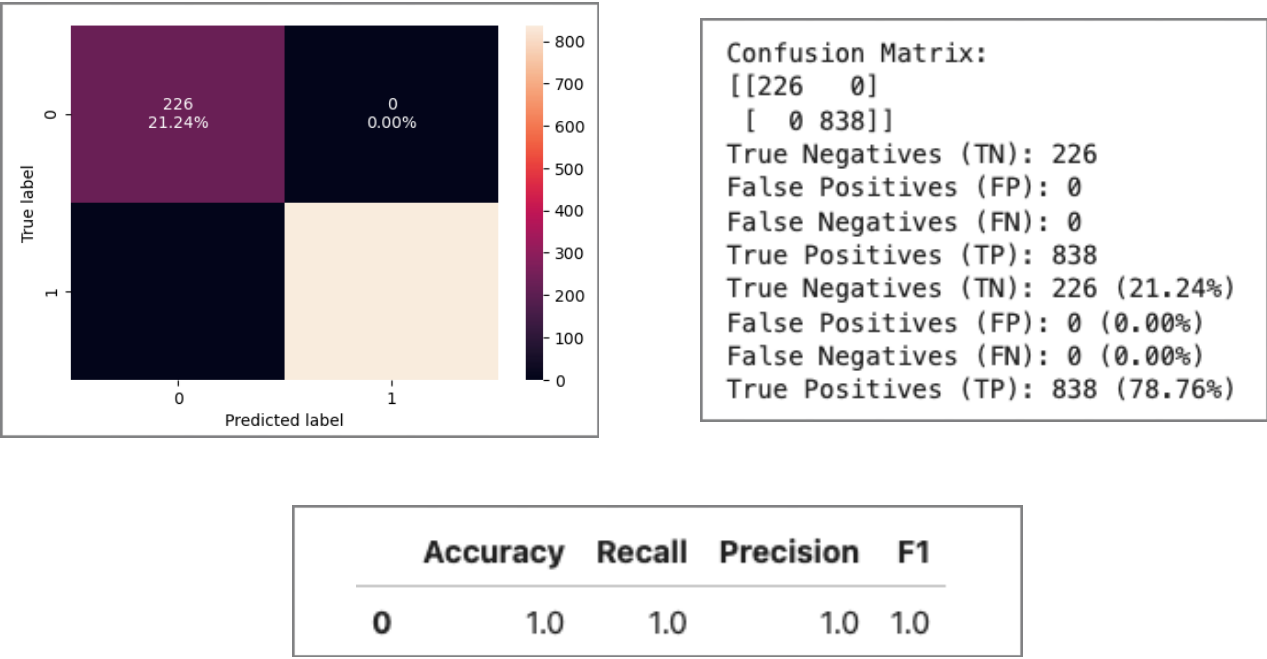
	Accuracy	Recall	Precision	F1
0	1.0	1.0	1.0	1.0

Training set Random Forest

The Random Forest model demonstrates perfect performance on the training set, achieving 100% in all key metrics. This includes accuracy, recall, precision, and F1 score. The confusion matrix shows that the model correctly classified all

instances without any false positives or false negatives. However, this level of performance on the training set may indicate overfitting, where the model performs exceptionally well on training data but may not generalize as effectively to unseen test data.

### Random Forest Performance - Testing Set



Random Forest Performance - Test Set

The Random Forest model's performance on the testing set shows perfect results, similar to its training set performance. Key metrics are as follows:

- **Accuracy:** 1.0 (100%)
  - This indicates that the model correctly predicted the outcome for all instances in the testing set.
- **Recall (Sensitivity):** 1.0 (100%)
  - The model successfully detected all actual positive cases in the testing set.
- **Precision:** 1.0 (100%)
  - All positive predictions made by the model are correct.
- **F1 Score:** 1.0 (100%)
  - The harmonic mean of precision and recall, indicating perfect balance between identifying positive cases and ensuring positive predictions are correct.

These metrics suggest that the Random Forest model maintains its perfect performance when applied to unseen data, with 100% in accuracy, recall, precision, and F1 score. This consistent performance across both training and testing sets indicates a highly effective model, although it may also suggest potential overfitting, warranting further validation on additional datasets.

## VIII.Model Comparison

Training performance comparison:				
	Logistic Regression	Tuned Logistic Regression	Random Forest	Tuned Random Forest
<b>Accuracy</b>	0.990602	0.790727	1.0	1.0
<b>Recall</b>	0.991249	0.996022	1.0	1.0
<b>Precision</b>	0.996800	0.791904	1.0	1.0
<b>F1</b>	0.994017	0.882311	1.0	1.0

Training performance comparison

Testing performance comparison:				
	Logistic Regression	Tuned Logistic Regression	Random Forest	Tuned Random Forest
<b>Accuracy</b>	0.404135	0.788534	1.0	1.0
<b>Recall</b>	0.243437	0.980907	1.0	1.0
<b>Precision</b>	1.000000	0.797284	1.0	1.0
<b>F1</b>	0.391555	0.879615	1.0	1.0

Testing performance comparison

The performance of the models on the training set shows some significant differences. The non-tuned Logistic Regression model achieved high accuracy (0.990602), recall (0.991249), precision (0.996800), and F1 score (0.994017), which might indicate potential overfitting due to its near-perfect metrics. The tuned Logistic Regression model, on the other hand, exhibited a lower accuracy (0.790727) and precision (0.791904) but maintained a high recall (0.996022), suggesting it was particularly good at identifying true positives even if it occasionally misclassified some negatives.

In contrast, both the Random Forest and the Tuned Random Forest models achieved perfect scores across all metrics (accuracy, recall, precision, F1 score), indicating that they fit the training data exceptionally well.



The testing set results reinforce these observations. The non-tuned Logistic Regression model showed poor performance on the testing set, with low accuracy (0.404135) and recall (0.243437) despite a perfect precision score (1.000000), resulting in a low F1 score (0.391555). The tuned Logistic Regression model performed significantly better, with high accuracy (0.788534), recall (0.980907), precision (0.797284), and F1 score (0.879615), indicating it generalized well to new data.

Both the Random Forest and Tuned Random Forest models maintained perfect performance on the testing set, with accuracy, recall, precision, and F1 scores all at 1.0, demonstrating exceptional generalization from training to testing data.

Based on these comparisons, the Random Forest models, both tuned and non-tuned, exhibit superior performance and generalization capability compared to Logistic Regression models. Given their perfect metrics on both training and testing sets, the Random Forest models are likely the best choice for this classification problem.

## **Final Model Selection**

After comparing the performance of the models on both training and testing sets, it is evident that the Random Forest models, both tuned and non-tuned, significantly outperform the Logistic Regression models. The non-tuned Logistic Regression model showed signs of overfitting with near-perfect metrics on the training set but poor generalization to the testing set, with low accuracy (0.404135) and recall (0.243437). The tuned Logistic Regression model improved generalization, achieving higher accuracy (0.788534) and recall (0.980907) on the testing set, but still lagged behind the Random Forest models.

The Random Forest and Tuned Random Forest models both achieved perfect scores across all metrics (accuracy, recall, precision, F1 score) on both the training and testing sets, demonstrating their superior capability to generalize and accurately classify new data. This consistent performance indicates that these models are highly effective for the given classification problem.

Given the perfect performance metrics and robust generalization ability, the Random Forest models are selected as the final models for this classification

task. Their ability to handle complex data structures and interactions, coupled with their high accuracy and reliability, make them the optimal choice for predicting outcomes accurately.

## Feature Importance

Feature importance helps identify which features contribute the most to the predictions of a model. In the context of Random Forest, feature importance is typically derived from the average of the decrease in impurity (e.g., Gini impurity or entropy) brought by each feature across all trees in the forest.

### To Determine Feature Importance in Random Forest:

1. *Train the Model*
2. *Extract Feature Importance*
3. *Interpret Feature Importance*
  - *Higher values indicate that the feature is more important in making predictions.*



In the Random Forest model, the feature importance values provide insights into which variables significantly influence the prediction of defaults. Here are some key observations based on the feature importance values:

1. **Networth Next Year:** This feature stands out with the highest importance score of 0.701159, indicating it plays a crucial role in predicting defaults. This suggests that the net worth projection is a strong indicator of a company's financial health and stability.
2. **TOL/TNW (Total Outside Liabilities to Tangible Net Worth):** With an importance score of 0.021158, this ratio is also a significant predictor. It measures a company's leverage and financial risk, where higher values may indicate higher default risk.
3. **Cash Profit:** This feature has an importance score of 0.018662, highlighting the importance of liquidity and cash flow in assessing financial stability. Cash profit reflects the company's ability to generate cash from operations.
4. **PAT as % of Net Worth:** With an importance score of 0.017156, this ratio indicates the profitability relative to the net worth, showing the company's efficiency in generating profits from its equity base.
5. **Debt to Equity Ratio:** This feature, with an importance score of 0.015290, underscores the significance of the company's debt levels relative to its equity. High debt levels can increase the risk of default.
6. **Other Significant Features:**
  - **Cash Profit as % of Total Income** (0.012409)
  - **Cumulative Retained Profits** (0.011428)
  - **PBT as % of Total Income** (0.010218)
  - **Reserves and Funds** (0.009870)
  - **Total Term Liabilities / Tangible Net Worth** (0.008768)

These features collectively highlight the importance of profitability, liquidity, and leverage in predicting financial defaults. The importance scores decrease gradually, but each feature still contributes to the model's predictive power.

## Conclusion

The Random Forest model identifies a combination of profitability, liquidity, leverage, and projected net worth as the most critical factors in predicting

defaults. This information can be used by financial analysts and decision-makers to focus on the most influential variables when assessing the risk of defaults and making informed decisions.

	Feature	Importance
0	Networth Next Year	0.701159
26	TOL/TNW	0.021158
9	Cash profit	0.018662
14	PAT as % of net worth	0.017156
36	Debt to equity ratio (times)	0.015290
13	Cash profit as % of total income	0.012409
24	Cumulative retained profits	0.011428
11	PBT as % of total income	0.010218
19	Reserves and funds	0.009870
27	Total term liabilities / tangible net worth	0.008768
6	Profit after tax	0.007633
23	Shareholders funds	0.007561
8	PBT	0.007308
35	Current ratio (times)	0.007292
2	Net worth	0.007227
44	EPS	0.006579
12	PAT as % of total income	0.006476
34	Quick ratio (times)	0.006445
46	PE on BSE	0.006193
7	PBDITA	0.005558
38	Creditors turnover	0.005531
18	Total capital	0.005321
25	Capital employed	0.004742
42	Raw material turnover	0.004710
10	PBDITA as % of total income	0.004585
33	Net working capital	0.004581
43	Shares outstanding	0.004562
29	Contingent liabilities	0.004379
3	Total income	0.004272
41	WIP turnover	0.004119
37	Cash to average cost of sales per day	0.004079
15	Sales	0.004039
40	Finished goods turnover	0.003945
1	Total assets	0.003893
28	Contingent liabilities / Net worth (%)	0.003798
30	Net fixed assets	0.003669
22	Deferred tax liability	0.003538
5	Total expenses	0.003446
31	Investments	0.003345
21	Current liabilities & provisions	0.003321
32	Current assets	0.003293
17	Other income	0.003282
45	Total liabilities	0.003237
39	Debtors turnover	0.003114
20	Borrowings	0.002944
4	Change in stock	0.002937
16	Income from fincial services	0.002929

Feature name and their importances

## IX. Actionable Insights & Recommendations

### 1. Networth Next Year:

- **Insight:** “Networth Next Year” is the most significant predictor of defaults.
- **Recommendation:** Regularly project and analyze future net worth to ensure financial stability. Focus on strategies that enhance net worth, such as reinvesting profits, reducing liabilities, and increasing assets.

2. **TOL/TNW (Total liabilities of the customer divided by Total net worth):**

- **Insight:** A higher ratio indicates higher financial leverage and risk.
- **Recommendation:** Aim to keep this ratio within industry benchmarks. Consider debt restructuring and avoid taking on excessive liabilities to maintain a balanced financial structure.

3. **Cash Profit:**

- **Insight:** Cash profit is an important measure of operational efficiency.
- **Recommendation:** Focus on increasing cash profits by optimizing operational processes, reducing waste, and enhancing revenue streams. Regularly review and improve cost management practices.

4. **PAT as % of net worth:**

- **Insight:** This ratio measures profitability relative to net worth.
- **Recommendation:** Enhance profitability by implementing cost-saving measures, optimizing pricing strategies, and exploring new business opportunities. Continuously monitor this ratio to ensure sustainable growth.

5. **Debt to equity ratio (times):**

- **Insight:** Indicates the company's financial leverage.
- **Recommendation:** Maintain a balanced debt-to-equity ratio by managing debt levels and considering equity financing options. Regularly review the ratio to ensure it aligns with industry standards and financial goals.

6. **Cash to average cost of sales per day:**

- **Insight:** This ratio measures liquidity and the ability to cover sales costs with available cash.
- **Recommendation:** Improve liquidity management by maintaining adequate cash reserves. Implement efficient cash flow management practices, such as speeding up receivables and managing payables effectively.

7. **Reserves and funds:**

- **Insight:** Strong reserves and funds indicate financial stability.
- **Recommendation:** Build and maintain robust reserves to cushion against economic downturns and unforeseen expenses. Allocate a portion of profits to reserves regularly to ensure long-term financial health.

8. **Current ratio (times) and Quick ratio (times):**

- **Insight:** These ratios measure short-term liquidity and financial health.
- **Recommendation:** Regularly monitor these ratios to ensure sufficient liquidity. Optimize working capital by managing inventory levels, accelerating receivables, and extending payables where possible.

9. **Contingent Liabilities / Net Worth :**

- **Insight:** Higher contingent liabilities relative to net worth increase financial risk.
- **Recommendation:** Minimize contingent liabilities by carefully assessing and managing potential risks. Maintain comprehensive insurance coverage and regularly review contingent liabilities to mitigate their impact on financial health.

#### 10. **EPS (Earnings Per Share) and PE Ratio:**

- **Insight:** These metrics provide insights into profitability and market valuation.
- **Recommendation:** Focus on improving earnings per share by enhancing operational efficiency and revenue growth. Monitor the PE ratio to ensure the company is valued appropriately in the market.

Implementing these recommendations based on the significant columns can help enhance financial stability, improve profitability, and mitigate risks, leading to better overall performance and reduced likelihood of defaults.

### **Summary**

By focusing on these actionable insights and implementing the recommended strategies based on specific column analysis, companies can better manage their financial health and reduce the risk of defaults. This proactive approach can lead to more sustainable growth and financial stability, benefiting all stakeholders involved.

To effectively mitigate the risk of default, the company should consider the following refined strategies:

#### 1. **Enhance Equity Position:**

- **Strengthen Capital Base:** Improve key metrics such as *Net worth* and *Shareholders funds*. Pursue new equity financing options or reinvest retained earnings to improve the equity-to-liability ratio. Converting existing debt to equity can also help reduce financial leverage and strengthen the company's balance sheet.

#### 2. **Optimize Debt Management:**

- **Restructure Debt:** Focus on reducing the *Debt to equity ratio (times)* and the *TOL/TNW* (Total Outside Liabilities/Tangible Net Worth). Engage in

proactive negotiations with creditors to restructure existing debt. Explore options such as extending repayment periods, reducing interest rates, or converting debt to equity to alleviate immediate financial pressures and improve cash flow.

**3. Implement Rigorous Cost Control:**

- **Streamline Expenses:** Analyze and optimize *Total expenses* and *Operating expenses*. Conduct a comprehensive review of operating expenses to identify cost-saving opportunities. Focus on optimizing essential expenditures and eliminating inefficiencies to enhance overall cost management.

**4. Drive Revenue Growth:**

- **Expand Market Reach:** Focus on increasing *Total income* and *Sales*. Develop and execute strategies to increase sales through targeted marketing campaigns, diversification of product offerings, or entering new markets. Expanding revenue streams can provide a more stable financial foundation.

**5. Strengthen Liquidity Management:**

- **Optimize Cash Flow:** Improve metrics such as *Cash profit* and *Quick ratio (times)*. Enhance cash flow management practices, defer non-essential expenditures, and optimize working capital. Ensure the company maintains sufficient liquidity to meet short-term obligations.

**6. Invest in Strategic Innovation:**

- **Foster Growth:** Invest in innovation to drive long-term growth and competitiveness. Focus on improving *EPS* (Earnings Per Share) and leveraging *Investments*. Explore cost-effective innovation strategies, such as partnerships, joint ventures, or accessing grants, to support research and development efforts.

**7. Establish Robust Risk Monitoring:**

- **Continuous Assessment:** Implement a comprehensive risk monitoring system to regularly assess financial health and identify early warning signs of potential issues. Continuously review key financial ratios such as *PAT as % of net worth*, *Contingent liabilities / Net worth (%)*, and *Creditors turnover* to proactively address risks and ensure timely interventions.

By adopting these strategies, the company can strengthen its financial stability, improve its ability to meet obligations, and position itself for sustained growth and resilience against default risk.