

---

# FRA CODED PROJECT – PART A

## Business Report

---

DSBA

Submitted By: Maheep Singh

Batch : PGP-DSBA (PGPDSBA.O.AUG24.A)

## Table of Contents

List of Figures .....	4
Business Context & Data Dictionary .....	6
Context.....	6
Objective .....	6
Data Description .....	6
Rubric Question 1: Define the Problem and Perform Exploratory Data Analysis .....	8
Problem Definition .....	8
Data Overview.....	8
Univariate & Bivariate/Multivariate (w.r.t. Default) Analysis.....	13
Rubric Question 2: Data Preprocessing .....	26
Feature Engineering.....	26
Outlier Treatment .....	26
Data Preparation for Modelling .....	27
Duplicate/Missing Treatment & Data Scaling .....	27
Check & Treat Imbalance in Target Variable (Default) .....	29
Rubric Question 3: Model building .....	30
Model Evaluation Criteria .....	30
Logistic Regression Model .....	31
Build Model.....	31
Checking Model Performance .....	32
Random Forest Model .....	33
Build Model.....	33
Checking Model Performance .....	33
Rubric Question 4: Model Performance Improvement .....	34
Logistic Regression Model – Tuning .....	34
Build Tuned Model.....	34
Check Model Performance (Tuned Model) .....	37
Random Forest – Tuning .....	39
Build Tuned Model.....	39
Check Model Performance (Tuned Model) .....	40
Rubric Question 5: Model Performance Comparison and Final Model Selection .....	41
Training Dataset Performance Comparison .....	41
Test Dataset Performance Comparison .....	41
Final Model Selection .....	41

Most Important Features.....	42
Rubric Question 6: Actionable Insights & Recommendations .....	44
Actionable Insights.....	44
Business Recommendations .....	44

## List of Figures

Figure 1: Top 5 rows of the dataset.....	8
Figure 2: Datatypes in the Dataset.....	8
Figure 3: Missing/Duplicate Value-check.....	9
Figure 4: Unique Value Check .....	10
Figure 5: Statistical Summary of the Dataset.....	11
Figure 6: Univariate Analysis – Boxplot (Part A) .....	13
Figure 7: Univariate Analysis – Boxplot (Part B).....	14
Figure 8: Univariate Analysis – Boxplot (Part C).....	15
Figure 9: Univariate Analysis – Boxplot (Part C).....	16
Figure 10: Univariate Analysis – Histogram (Part A) .....	16
Figure 11: Univariate Analysis – Histogram (Part B) .....	17
Figure 12: Univariate Analysis – Histogram (Part C) .....	18
Figure 13: Univariate Analysis – Histogram (Part D).....	19
Figure 14: Bivariate/Multivariate (w.r.t. 'Default') Analysis – Boxplot (Part A).....	20
Figure 15: Bivariate/Multivariate (w.r.t. 'Default') Analysis – Boxplot (Part B).....	21
Figure 16: Bivariate/Multivariate (w.r.t. 'Default') Analysis – Boxplot (Part C).....	22
Figure 17: Bivariate/Multivariate (w.r.t. 'Default') Analysis – Boxplot (Part D).....	23
Figure 18: Bivarariate Analysis – Heatmap .....	25
Figure 19: Outlier Inspection .....	26
Figure 20: Missing-value check on Train/Test datasets.....	27
Figure 21: Subset of Train dataset post Data Scaling .....	28
Figure 22: Missing-value check on Train/Test datasets post data-treatment (KNNImputer) .....	28
Figure 23: Check Imbalance in Target Variable .....	29
Figure 24: Train dataset summary - pre/post SMOTE.....	29
Figure 25: Subset of Train dataset with Intercept.....	31
Figure 26: Logistic Regression Summary.....	31
Figure 27: Logistic Regression: Confusion Matrix & Metric Performance on Training dataset .....	32
Figure 28: Logistic Regression: Confusion Matrix & Metric Performance on Test dataset.....	32
Figure 29: Random Forest: Confusion Matrix & Metric Performance on Training dataset .....	33
Figure 30: Random Forest: Confusion Matrix & Metric Performance on Test dataset.....	33
Figure 31: VIF for Independent Variables .....	35
Figure 32: VIF for Independent Variables (post-Multicollinearity fix).....	36
Figure 33: Logistic Regression (Tuned) Output Summary .....	36
Figure 34: ROC Curve & Optimal Threshold Value .....	37
Figure 35: Logistic Regression (Tuned): Confusion Matrix & Metric Performance on Training Dataset.....	37
Figure 36: Logistic Regression (Tuned): Confusion Matrix & Metric Performance on Test Dataset .....	38
Figure 37: Best Parameters for Tuning Random Forest Classifier .....	39
Figure 38: Random Forest (Tuned): Confusion Matrix & Metric Performance on Training Dataset.....	40
Figure 39: Random Forest (Tuned): Confusion Matrix & Metric Performance on Test Dataset .....	40
Figure 40: Training Dataset Performance Comparison .....	41
Figure 41: Test Dataset Performance Comparison .....	41
Figure 42: Feature Importance .....	42

## List of Tables

Table 1: Model Comparison .....	41
Table 2: Feature Importance Inference (Top 7) .....	43

# Business Context & Data Dictionary

## Context

In the realm of modern finance, businesses encounter the perpetual challenge of managing debt obligations effectively to maintain a favourable 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 endeavours 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 analyse 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 fulfil financial obligations promptly and efficiently, and identify potential cases of default.
2. Credit Risk Evaluation: Evaluate credit risk exposure by analysing 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 analyse 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 Description

The data consists of financial metrics from the balance sheets of different companies. The detailed data dictionary is given below: -

- Networth Next Year: Net worth of the customer in the next year
- Total assets: Total assets of customer
- Net worth: Net worth of the customer of the present year
- Total income: Total income of the customer
- Change in stock: Difference between the current value of the stock and the value of stock in the last trading day
- Total expenses: Total expenses done by the customer
- Profit after tax: Profit after tax deduction
- PBDITA: Profit before depreciation, income tax, and amortization
- PBT: Profit before tax deduction
- Cash profit: Total Cash profit
- PBDITA as % of total income:  $PBDITA / Total\ income$
- PBT as % of total income:  $PBT / Total\ income$
- PAT as % of total income:  $PAT / Total\ income$
- Cash profit as % of total income:  $Cash\ Profit / Total\ income$
- PAT as % of net worth:  $PAT / Net\ worth$
- Sales: Sales done by the customer
- Income from financial services: Income from financial services
- Other income: Income from other sources
- Total capital: Total capital of the customer
- Reserves and funds: Total reserves and funds of the customer
- Borrowings: Total amount borrowed by the customer
- Current liabilities & provisions: current liabilities of the customer
- Deferred tax liability: Future income tax customer will pay because of the current transaction
- Shareholders funds: Amount of equity in a company which belongs to shareholders
- Cumulative retained profits: Total cumulative profit retained by customer
- Capital employed: Current asset minus current liabilities
- TOL/TNW: Total liabilities of the customer divided by Total net worth

- Total term liabilities / tangible net worth: Short + long term liabilities divided by tangible net worth
- Contingent liabilities / Net worth (%): Contingent liabilities / Net worth
- Contingent liabilities: Liabilities because of uncertain events
- Net fixed assets: The purchase price of all fixed assets
- Investments: Total invested amount
- Current assets: Assets that are expected to be converted to cash within a year
- Net working capital: Difference between the current liabilities and current assets
- Quick ratio (times): Total cash divided by current liabilities
- Current ratio (times): Current assets divided by current liabilities
- Debt to equity ratio (times): Total liabilities divided by its shareholder equity
- Cash to current liabilities (times): Total liquid cash divided by current liabilities
- Cash to average cost of sales per day: Total cash divided by the average cost of the sales
- Creditors turnover: Net credit purchase divided by average trade creditors
- Debtors turnover: Net credit sales divided by average accounts receivable
- Finished goods turnover: Annual sales divided by average inventory
- WIP turnover: The cost of goods sold for a period divided by the average inventory for that period
- Raw material turnover: Cost of goods sold is divided by the average inventory for the same period
- Shares outstanding: Number of issued shares minus the number of shares held in the company
- Equity face value: cost of the equity at the time of issuing
- EPS: Net income divided by the total number of outstanding shares
- Adjusted EPS: Adjusted net earnings divided by the weighted average number of common shares outstanding on a diluted basis during the plan year
- Total liabilities: Sum of all types of liabilities
- PE on BSE: Company's current stock price divided by its earnings per share

**Note:** 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.

## Rubric Question 1: Define the Problem and Perform Exploratory Data Analysis

### Problem Definition

- Objective is to develop a Financial Health Assessment Tool by analysing financial metrics of different companies.
- Develop 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.
- This tool is expected to help the organization anticipate potential challenges with the financial performance of the companies and enable proactive risk mitigation strategies, using which, share recommendations to the business.

### Data Overview

- Load dataset & display top 5 rows (Truncated view due to high no. of columns): -**

	Num	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
0	1	395.30000	827.60000	336.50000	534.10000	13.50000	508.70000	38.90000	124.40000	64.60000	95.20000	23.29000	12.10000	7.28000	17.82000
1	2	36.20000	67.70000	24.30000	137.90000	-3.70000	131.00000	3.20000	5.50000	1.00000	3.80000	3.99000	0.73000	2.32000	2.76000
2	3	84.00000	238.40000	78.90000	331.20000	-18.10000	309.20000	3.90000	25.80000	10.50000	9.40000	7.79000	3.17000	1.18000	2.84000
3	4	2041.40000	6883.50000	1443.30000	8448.50000	212.20000	8482.40000	178.30000	418.40000	185.10000	178.00000	4.95000	2.19000	2.11000	2.11000
4	5	41.80000	90.90000	47.00000	388.60000	3.40000	392.70000	-0.70000	7.20000	-0.60000	3.90000	1.85000	-0.15000	-0.18000	1.00000

Figure 1: Top 5 rows of the dataset

- There are **4256 rows & 52 columns** in the dataset
- Creating 'Default' feature (Target Variable): -**
  - ✓ New column 'Default' is created using the below logic (as instructed in the problem):
    - Value = 1 | 'Networth Next Year' < 0
    - Value = 0 | 'Networth Next Year' >= 0
- Checking datatypes: -**

Data #	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 financial 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
51	Default	4256 non-null	int64

dtypes: float64(50), int64(2)

Figure 2: Datatypes in the Dataset



✓ There are 52 columns - 50 float64 (numeric) & 2 Int64 (numeric) datatypes in the dataset.

▪ **Check Missing & Duplicate Values: -**

- ✓ No Duplicated values found.
- ✓ However, we have a lot of missing values that would require treatment (covered in subsequent section).

**Duplicated Values: 0**

Missing Values:-

Num	0
Networth Next Year	0
Total assets	0
Net worth	0
Total income	231
Change in stock	550
Total expenses	165
Profit after tax	154
PBDITA	154
PBT	154
Cash profit	154
PBDITA as % of total income	79
PBT as % of total income	79
PAT as % of total income	79
Cash profit as % of total income	79
PAT as % of net worth	0
Sales	305
Income from fincial services	1111
Other income	1556
Total capital	5
Reserves and funds	98
Borrowings	431
Current liabilities & provisions	110
Deferred tax liability	1369
Shareholders funds	0
Cumulative retained profits	45
Capital employed	0
TOL/TNW	0
Total term liabilities / tangible net worth	0
Contingent liabilities / Net worth (%)	0
Contingent liabilities	1402
Net fixed assets	132
Investments	1715
Current assets	80
Net working capital	37
Quick ratio (times)	105
Current ratio (times)	105
Debt to equity ratio (times)	0
Cash to current liabilities (times)	105
Cash to average cost of sales per day	100
Creditors turnover	391
Debtors turnover	385
Finished goods turnover	874
WIP turnover	764
Raw material turnover	428
Shares outstanding	810
Equity face value	810
EPS	0
Adjusted EPS	0
Total liabilities	0
PE on BSE	2627
Default	0

Figure 3: Missing/Duplicate Value-check

✓ Missing value treatment is covered in the subsequent section.

- Checking for **Unique Values** for each column: -

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 financial 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
Default	2

Figure 4: Unique Value Check

- **Dropping Redundant Features:** -
  - ✓ Dropping '**Equity face value**' as it has only 18 unique values in comparison to others
  - ✓ Dropping '**Num**' as it is a unique identifier & serves no purpose in the analysis
  - ✓ Dropping '**Networth Next Year**' as it has already been translated into 'Default' variable. Retaining this would lead to errors while building Logistic Regression Model (Singular Matrix Error).

Statistical Summary of the dataset: -

	count	mean	std	min	25%	50%	75%	max
Total assets	4256.00000	3573.61715	30074.44344	0.10000	91.30000	315.50000	1120.80000	1176509.20000
Net worth	4256.00000	1351.94960	12961.31165	0.00000	31.47500	104.80000	389.85000	613151.60000
Total income	4025.00000	4688.18979	53918.94661	0.00000	107.10000	455.10000	1485.00000	2442828.20000
Change in stock	3706.00000	43.70248	436.91505	-3029.40000	-1.80000	1.60000	18.40000	14185.50000
Total expenses	4091.00000	4356.30110	51398.08712	-0.10000	96.80000	426.80000	1395.70000	2366035.30000
Profit after tax	4102.00000	295.05059	3079.90207	-3908.30000	0.50000	9.00000	53.30000	119439.10000
PBDITA	4102.00000	605.94064	5646.23063	-440.70000	6.92500	36.90000	158.70000	208576.50000
PBT	4102.00000	410.25904	4217.41531	-3894.80000	0.80000	12.60000	74.17500	145292.60000
Cash profit	4102.00000	408.26748	4143.92639	-2245.70000	2.90000	19.40000	96.25000	176911.80000
PBDITA as % of total income	4177.00000	3.17989	172.25656	-6400.00000	4.97000	9.68000	16.47000	100.00000
PBT as % of total income	4177.00000	-18.19683	419.91109	-21340.00000	0.56000	3.34000	8.94000	100.00000
PAT as % of total income	4177.00000	-20.03367	423.57619	-21340.00000	0.35000	2.37000	6.42000	150.00000
Cash profit as % of total income	4177.00000	-9.02128	299.95743	-15020.00000	2.00000	5.66000	10.73000	100.00000
PAT as % of net worth	4256.00000	10.16786	61.53240	-748.72000	0.00000	8.04000	20.20250	2466.67000
Sales	3951.00000	4645.68454	53080.90330	0.10000	113.35000	468.60000	1481.20000	2384984.40000
Income from fincial services	3145.00000	81.36006	1042.75868	0.00000	0.50000	1.90000	9.80000	51938.20000
Other income	2700.00000	55.95289	1178.41526	0.00000	0.40000	1.50000	6.20000	42856.70000
Total capital	4251.00000	224.55766	1684.95129	0.10000	13.20000	42.60000	103.15000	78273.20000
Reserves and funds	4158.00000	1210.56193	12816.22922	-6525.90000	5.30000	55.15000	282.52500	625137.80000
Borrowings	3825.00000	1176.24808	8581.24892	0.10000	24.40000	99.80000	358.30000	278257.30000
Current liabilities & provisions	4146.00000	960.63143	9140.53613	0.10000	17.50000	70.30000	265.92500	352240.30000
Deferred tax liability	2887.00000	234.49512	2106.25316	0.10000	3.20000	13.50000	51.30000	72796.60000
Shareholders funds	4256.00000	1376.48672	13010.69116	0.00000	32.30000	107.60000	408.90000	613151.60000
Cumulative retained profits	4211.00000	937.18198	9853.09609	-6534.30000	1.10000	37.40000	206.20000	390133.80000
Capital employed	4256.00000	2433.61758	20496.40388	0.00000	61.30000	221.20000	790.30000	891408.90000
TOL/TNW	4256.00000	4.02534	20.87909	-350.48000	0.60000	1.42000	2.83000	473.00000
Total term liabilities / tangible net worth	4256.00000	1.85429	15.87506	-325.60000	0.05000	0.34500	1.00000	456.00000
Contingent liabilities / Net worth (%)	4256.00000	55.70750	369.16567	0.00000	0.00000	5.36000	31.01250	14704.27000
Contingent liabilities	2854.00000	948.55224	12056.73758	0.10000	6.00000	37.85000	195.32500	559506.80000
Net fixed assets	4124.00000	1209.48652	12502.39664	0.00000	26.20000	93.85000	352.82500	636604.60000
Investments	2541.00000	721.86588	6793.85987	0.00000	1.00000	8.20000	63.80000	199978.60000
Current assets	4176.00000	1350.36001	10155.57275	0.10000	36.60000	148.35000	515.00000	354815.20000
Net working capital	4219.00000	162.87424	3182.02996	-63839.00000	-1.10000	16.70000	86.50000	85782.80000
Quick ratio (times)	4151.00000	1.49735	9.32752	0.00000	0.41000	0.67000	1.03000	341.00000
Current ratio (times)	4151.00000	2.25740	12.47829	0.00000	0.93000	1.23000	1.72000	505.00000
Debt to equity ratio (times)	4256.00000	2.87156	15.59997	0.00000	0.22000	0.79000	1.75000	456.00000
Cash to current liabilities (times)	4151.00000	0.52842	4.79634	0.00000	0.02000	0.07000	0.19000	165.00000
Cash to average cost of sales per day	4156.00000	145.15793	2521.99181	0.00000	2.88000	8.04000	21.97000	128040.76000
Creditors turnover	3865.00000	16.81226	75.67492	0.00000	3.72000	6.17000	11.69000	2401.00000
Debtors turnover	3871.00000	17.92903	90.16443	0.00000	3.81000	6.47000	11.85000	3135.20000
Finished goods turnover	3382.00000	84.36999	562.63736	-0.09000	8.19000	17.32000	40.01250	17947.60000
WIP turnover	3492.00000	28.68451	169.65092	-0.18000	5.10000	9.86000	20.24000	5651.40000
Raw material turnover	3828.00000	17.73393	343.12586	-2.00000	3.02000	6.41000	11.82250	21092.00000
Shares outstanding	3446.00000	23764909.55543	170979041.32987	-2147483647.00000	1308382.50000	4750000.00000	10906020.00000	4130400545.00000
EPS	4256.00000	-196.21747	13061.95342	-843181.82000	0.00000	1.49000	10.00000	34522.53000
Adjusted EPS	4256.00000	-197.52761	13061.92951	-843181.82000	0.00000	1.24000	7.61500	34522.53000
Total liabilities	4256.00000	3573.61715	30074.44344	0.10000	91.30000	315.50000	1120.80000	1176509.20000
PE on BSE	1629.00000	55.46229	1304.44530	-1116.64000	2.97000	8.69000	17.00000	51002.74000
Default	4256.00000	0.21241	0.40906	0.00000	0.00000	0.00000	0.00000	1.00000

Figure 5: Statistical Summary of the Dataset

✓ Below is the **concise statistical summary** of the dataset: -

- **High Variability in Financial Health:** Core metrics like **Net worth**, **Total income**, and **Cash profit** have extreme standard deviations and ranges, indicating significant disparities in company financial performance.
- **Presence of Negative Indicators:** Many companies have negative values for **Profit after tax**, **Net working capital**, and **EPS**, signalling financial stress and potential risk of default.
- **Skewed Liquidity Ratios:** **Quick ratio** and **Current ratio** have low medians (0.93 and 1.23 respectively) but high max values, suggesting most companies operate with tight liquidity while a few have substantial buffers.
- **Debt Levels Show Risk Exposure:** **Debt to equity ratio** has a high mean (2.82) and extreme outliers (max: 456), indicating heavy leveraging by some companies—critical for credit risk assessment.
- **Outliers in Turnover Ratios:** Metrics like **Creditors turnover** and **Debtors turnover** show extremely large max values, suggesting irregular financial activity or data anomalies that could skew modelling.
- **Low Mean Profitability Ratios:** Ratios like **PBDITA as % of total income** (mean  $\approx 3.17\%$ ) and **PAT as % of net worth** ( $\approx 10.17\%$ ) are low, reinforcing that many companies operate on thin margins.
- **Default Rate is Imbalanced:** Only  $\sim 21.2\%$  of companies are labelled as **defaulters (Default = 1)**, highlighting the need to handle class imbalance carefully during model training.
- **Large Negative EPS & Adjusted EPS:** Both metrics show heavy negative skew and extremely low minima (e.g., Adjusted EPS min: -843181.22), likely from companies with major losses—strong signals for default.
- **Contingent Liabilities Show Red Flags:** **Contingent liabilities / Net worth (%)** has a huge spread (0 to 14,704%)—some companies are exposed to massive off-balance-sheet risks.
- **Wide Asset Distribution:** **Total assets** and **Net fixed assets** range from near zero to over 117 million and 66 million respectively, confirming a highly diverse company size distribution—important for segmentation in modelling.
- **Significant Missingness in Financial Services Income:** Only 3145 out of 4256 records have Income from financial services, suggesting either missing values or that many companies do not operate in this segment—this feature may have limited predictive value or need imputation.
- **Extreme Outliers in Valuation Ratios:** **PE on BSE** has a minimum of -1116 and a maximum over 15000, which are unrealistic and likely due to negative or near-zero EPS values—requires transformation or capping.
- **High Variance in Shareholders' Funds:** **Shareholders funds** show wide dispersion (mean: 1376, std: 13010), indicating some companies are massively equity-backed while others operate with minimal owner funding—important for understanding leverage and sustainability.
- **Capital Employed is Skewed:** Median **capital employed** is 221K vs. a max over 8.9 million, pointing to a few large companies dominating the asset base—segmentation by size could enhance model performance.
- **Potential Data Quality Issues in Turnover Ratios:** Negative or near-zero values in **WIP turnover** and **Raw material turnover** are unexpected and may indicate data entry issues or require careful preprocessing to avoid misleading patterns.

## Univariate & Bivariate/Multivariate (w.r.t. Default) Analysis

- Perform Univariate Analysis – Use Boxplots & Histograms to analyse each numerical variable: -

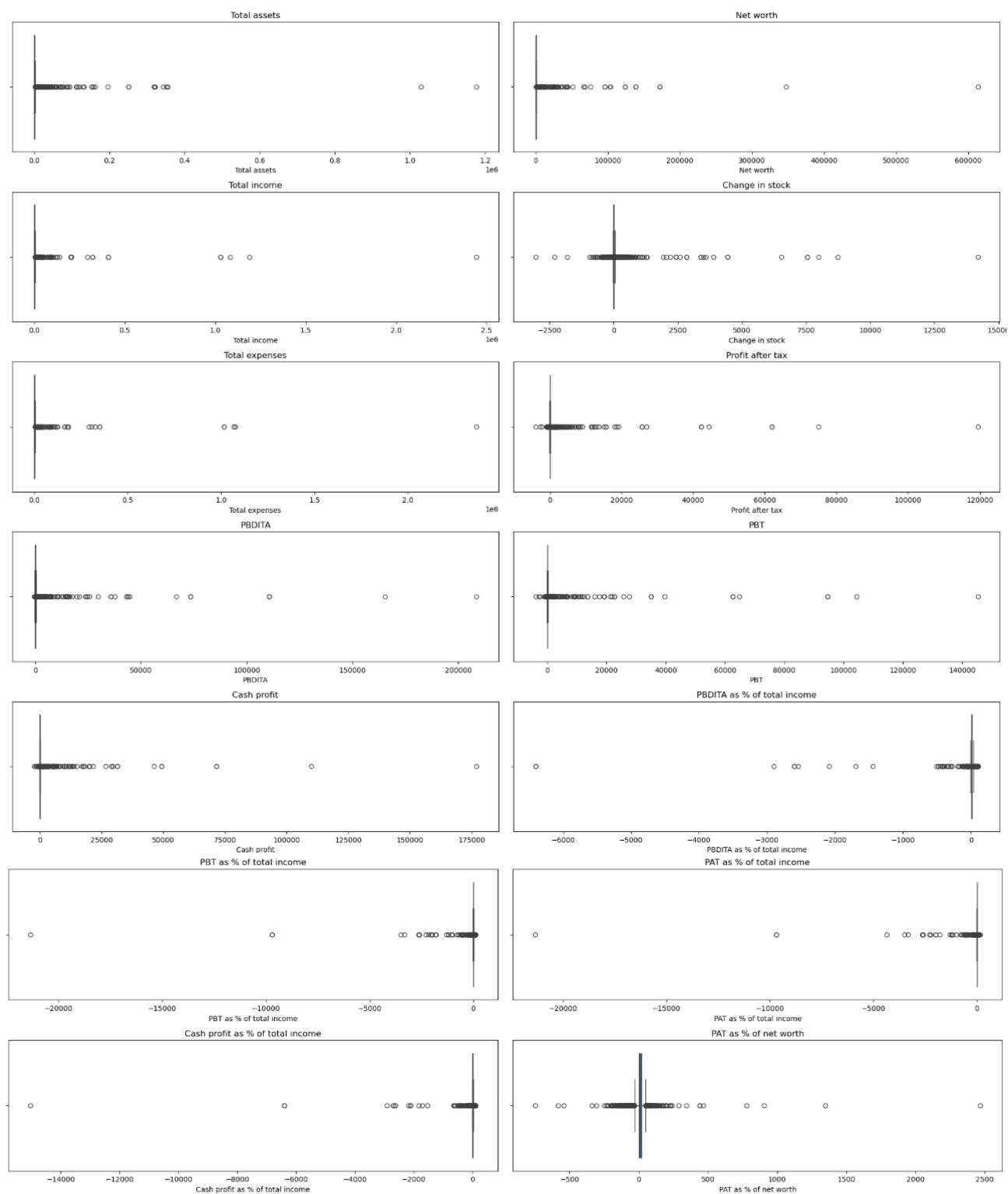


Figure 6: Univariate Analysis – Boxplot (Part A)

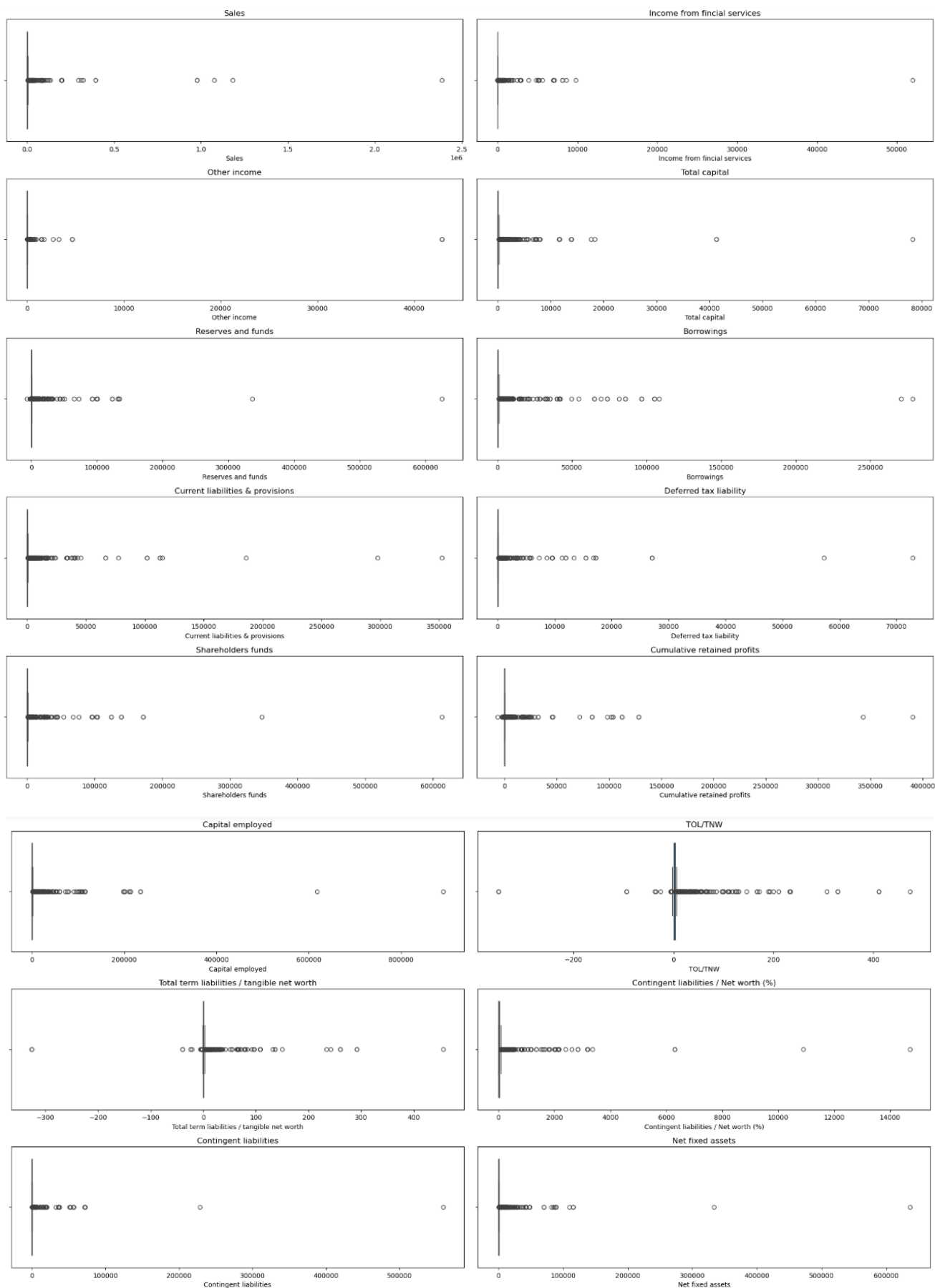


Figure 7: Univariate Analysis – Boxplot (Part B)

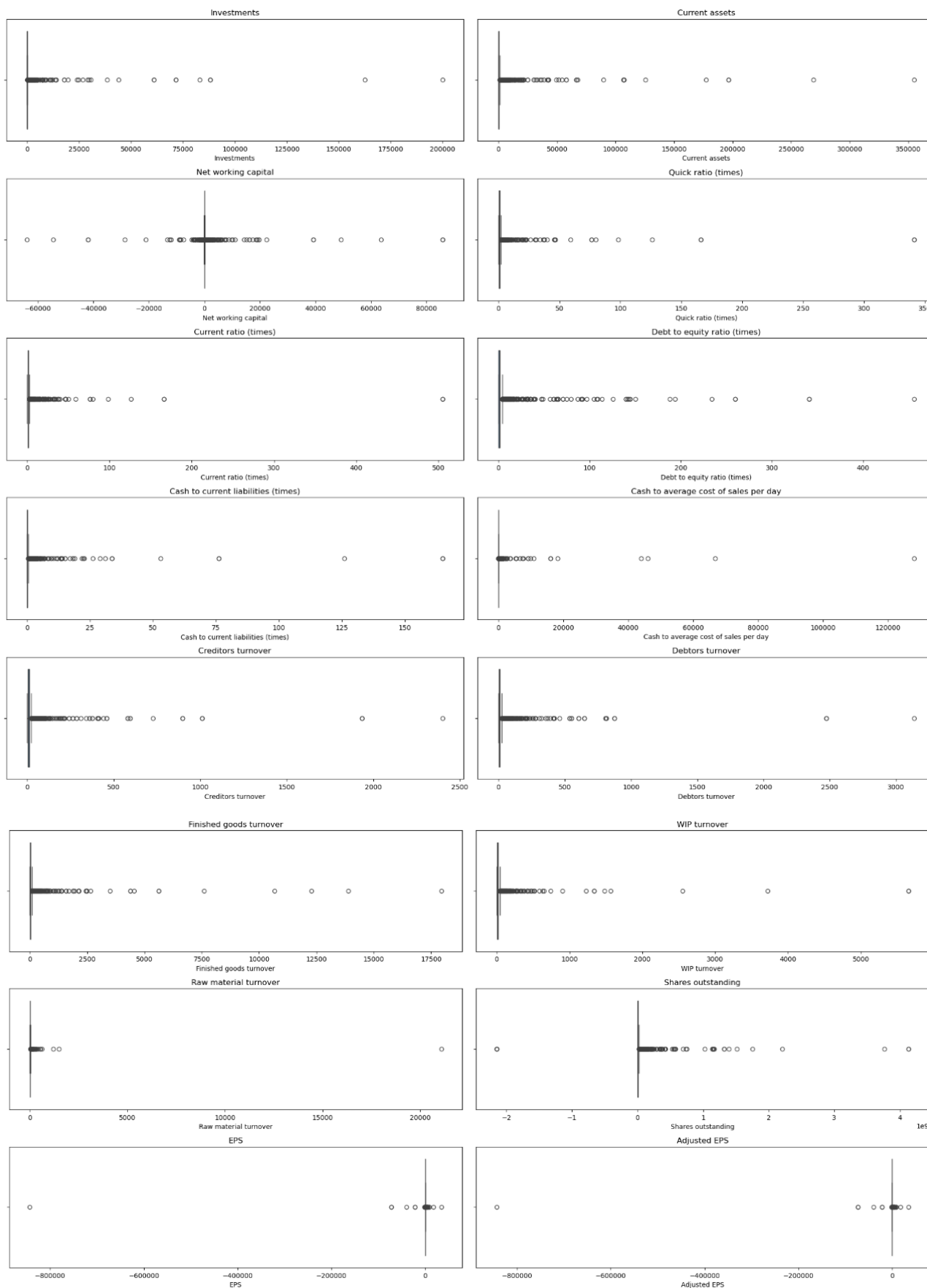


Figure 8: Univariate Analysis – Boxplot (Part C)

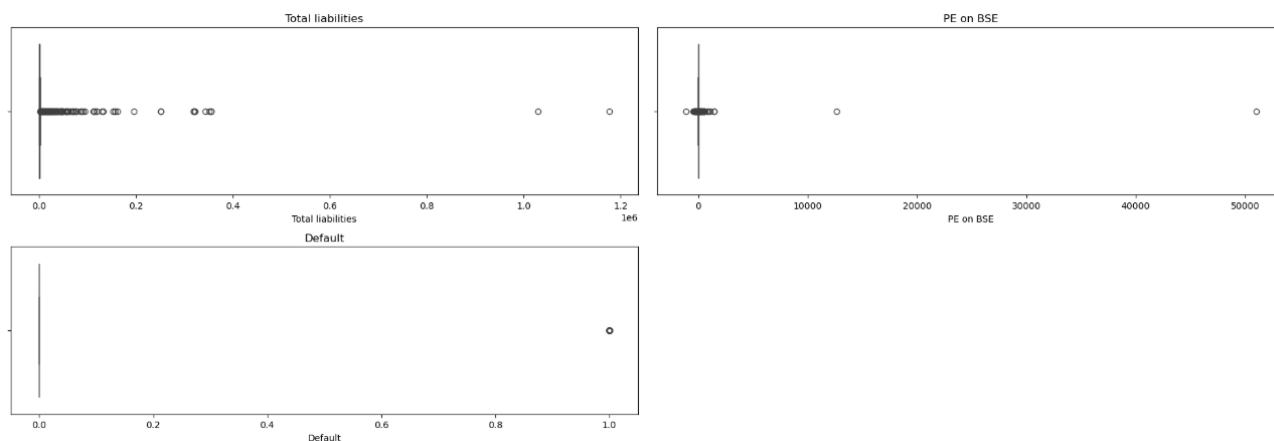


Figure 9: Univariate Analysis – Boxplot (Part C)

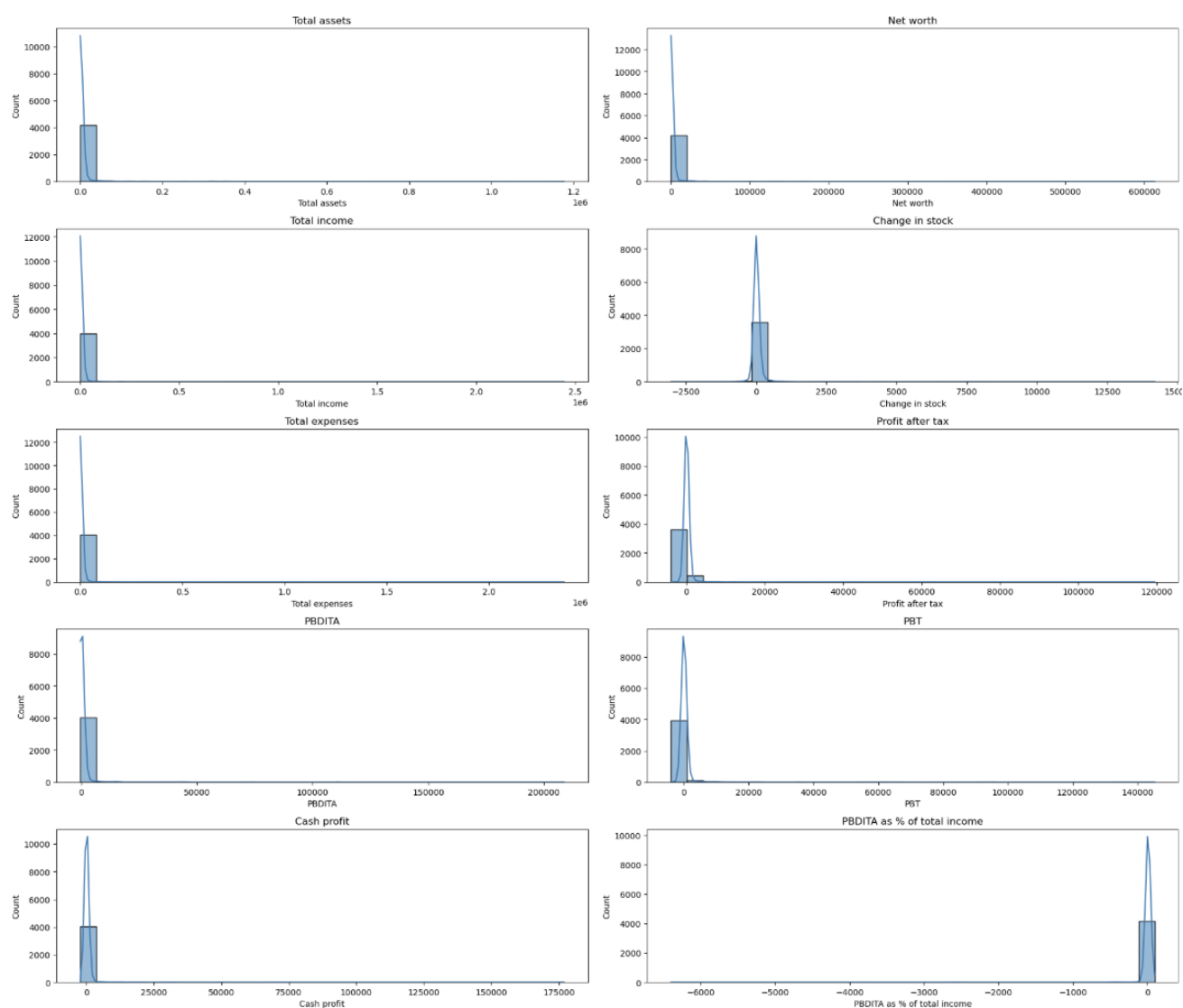


Figure 10: Univariate Analysis – Histogram (Part A)



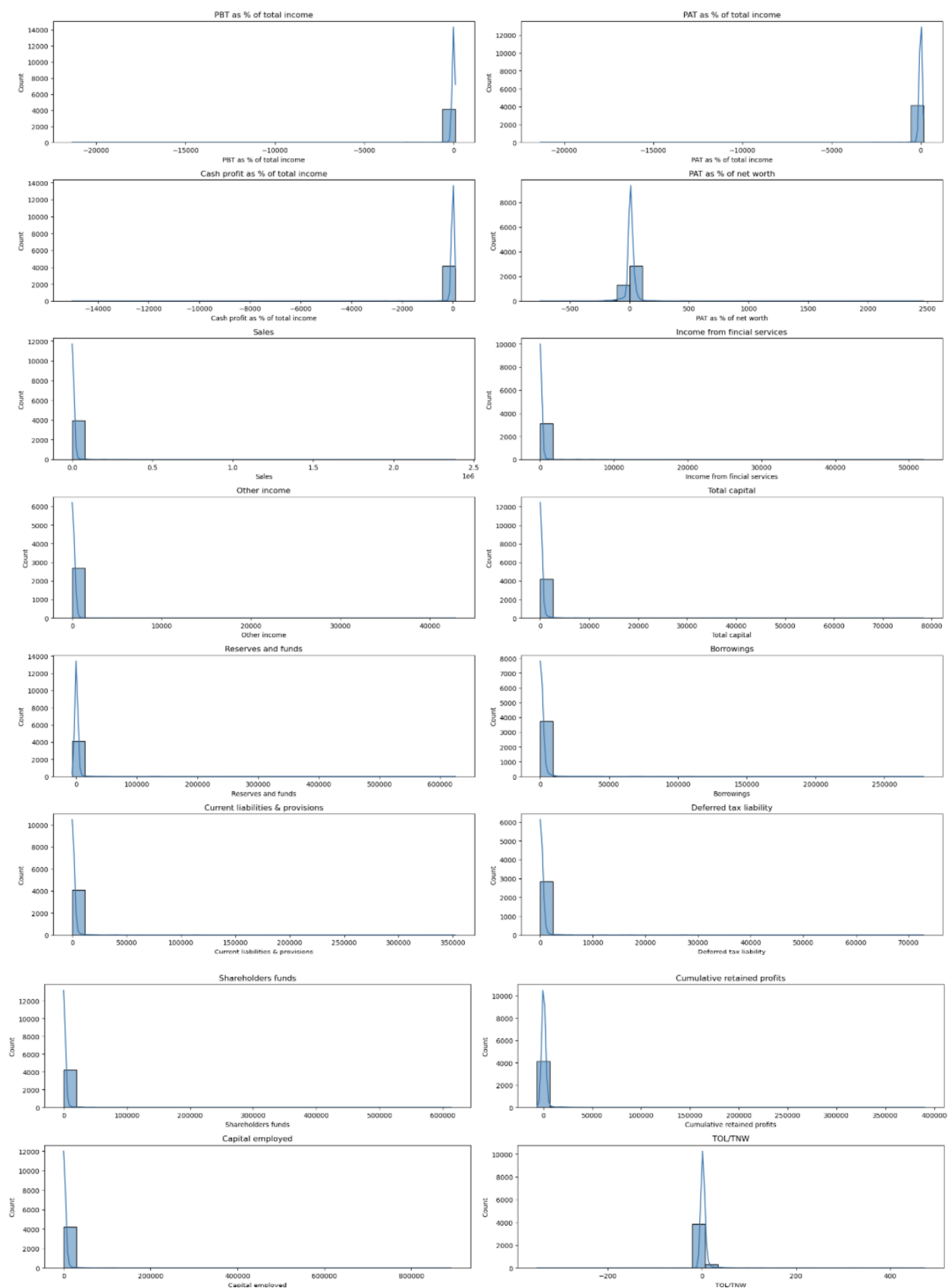


Figure 11: Univariate Analysis – Histogram (Part B)

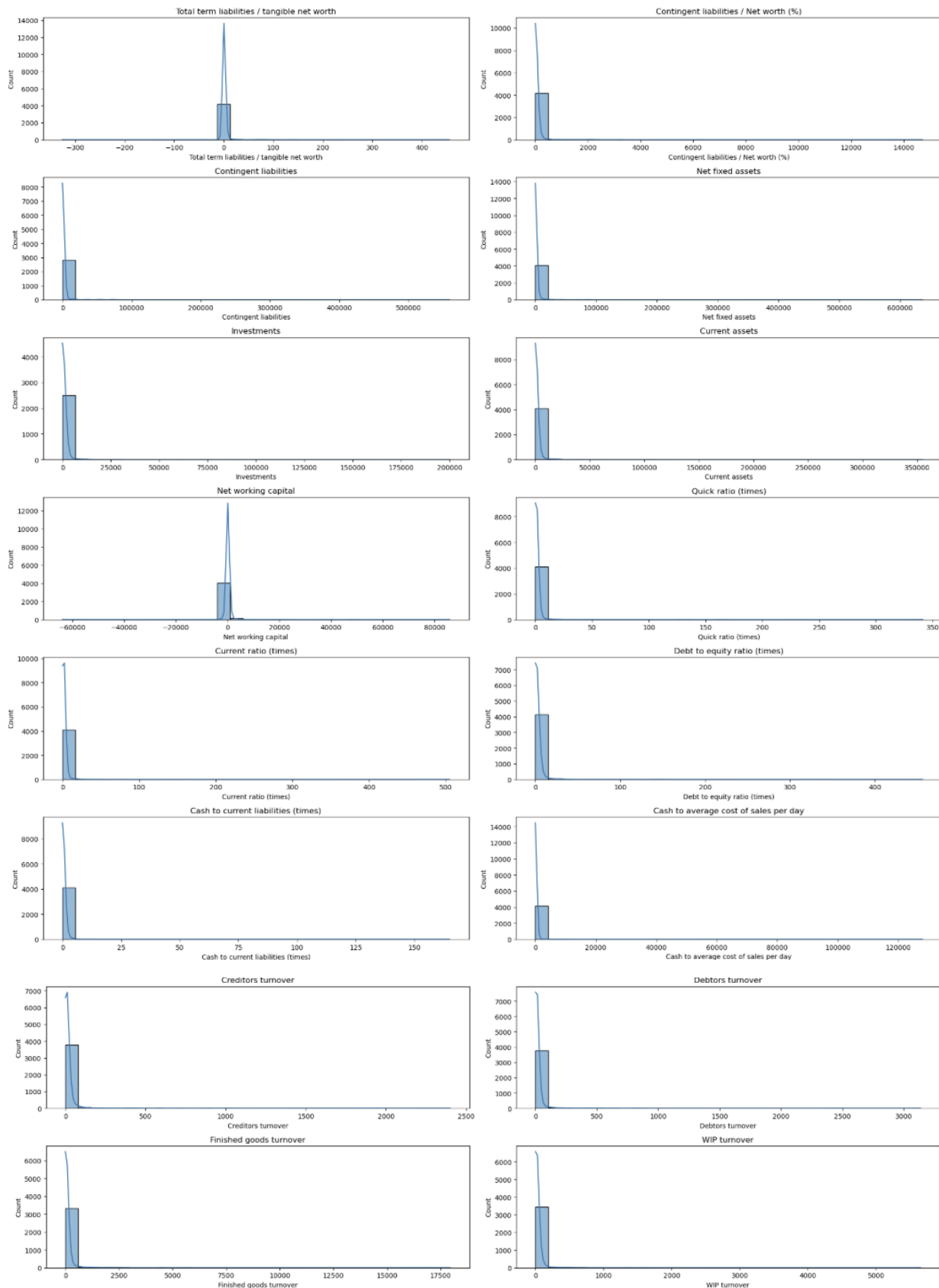


Figure 12: Univariate Analysis – Histogram (Part C)

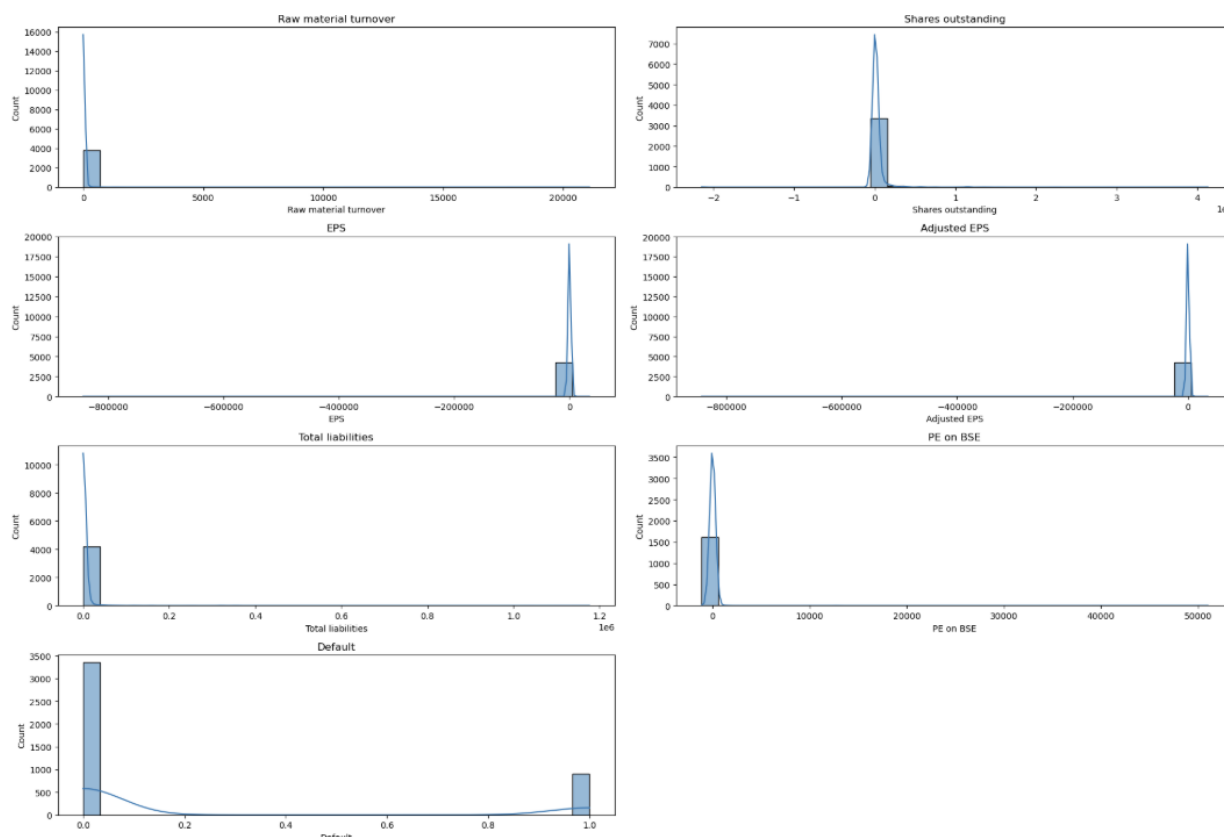


Figure 13: Univariate Analysis – Histogram (Part D)

✓ **Observations: Combined Histogram & Boxplot Insights: -**

- **Adjusted EPS:** Highly right-skewed with sharp concentration near zero and outliers from rare high-earning firms.
- **Cash profit:** Right-skewed with several extreme outliers indicating a few firms have strong cash generation.
- **Cash profit as % of total income:** Most firms operate on tight margins with occasional highly efficient firms standing out.
- **Depreciation:** Mostly low values with long right tail and heavy outliers from capital-heavy firms.
- **Interest:** Majority incur minimal interest, with isolated firms having very high debt servicing costs.
- **Operating profit margin:** Narrow spread and clustering near zero, but boxplot reveals outliers with strong margins.
- **PBDITA:** Modest operational profits are common, but significant positive outliers inflate the scale.
- **PBT (Profit before tax):** Strong right skew and presence of extreme outliers in highly profitable firms.
- **Profit after tax:** Matches PBT profile with low central values and a few post-tax success outliers.
- **Total assets:** Most firms are asset-light, with extreme right tail caused by large corporations.
- **Net fixed assets:** Heavily right-skewed with a handful of capital-intensive firms driving outliers.
- **Current assets:** Fairly balanced with some liquidity-rich outliers seen in the boxplot.
- **Investments:** Mostly low investments, but with wide histogram tail and sharp boxplot outliers.
- **Reserves and funds:** Commonly small reserves with very high spikes among financially strong firms.
- **Share capital:** Values are clustered due to standard capital structure, but a few firms raise notably more equity.
- **Shareholders' funds:** Right-skewed with many low-cap firms and a few significantly capitalized ones.
- **Borrowings:** Histogram shows concentration at zero with major outliers on the high end in boxplot.
- **Total income:** Most firms earn modestly, but histogram and boxplot confirm a few dominate income.
- **Net sales:** Reflects income pattern—right-skewed with rare yet extreme top-line outliers.
- **Total expenses:** Expense patterns mirror income with long tails and high-value outliers.
- **Value of production:** Wide spread and outliers suggest varied production capacities across firms.
- **Quick ratio (times):** Mostly cantered between 1–2, but outliers suggest both over- and under-liquidity.
- **TOL/TNW:** Low leverage for most, but few firms with dangerously high ratios act as boxplot outliers.
- **Current liabilities & provisions:** Most firms maintain low short-term debt with long tails in the boxplot.
- **Deferred tax liability:** Near-zero for most with some large, infrequent outliers visible.
- **Equity dividend:** Histogram is zero-inflated and boxplot confirms sparse dividend distribution.
- **Cumulative retained profits:** Heavily concentrated at lower end; boxplot exposes mature firms with high retention.

- **Retained earnings:** Mostly near zero with a few firms consistently reinvesting profits.
- **Change in stock:** Centered around zero with symmetric histogram but sharp outliers in both directions.
- **Stock:** Low stock values dominate, but some firms hold disproportionately large inventories.
- **Capital employed:** Capital usage is minimal in most, while outliers reveal capital-intensive operations.
- **Net worth:** Histogram skews right with strong boxplot outliers in well-capitalized firms.
- **Sales:** Modest revenue for most firms, sharply skewed with massive sales outliers.
- **Total liabilities:** Most companies carry limited liabilities; some outliers carry a high burden.
- **Equity face value:** Uniform with no spread; histogram shows fixed nominal values across firms.
- **Cash to current liabilities (times):** Reasonable liquidity for many; outliers suggest poor cash planning or hoarding.
- **Cash to average cost of sales per day:** Varies widely—histogram reveals operational gaps; boxplot confirms instability.
- **Creditors' turnover:** Histogram shows wide variance, with some firms paying suppliers exceptionally quickly.
- **Debtors' turnover:** Similar wide spread; some firms collect receivables far more efficiently than others.
- **Finished goods turnover:** Turnover is low in most; outliers suggest efficient inventory management in few.
- **Total expenses:** Expense distribution is heavily right-skewed with large company outliers.
- **Total income:** Repeats earlier income pattern; histogram confirms sharp revenue imbalance.
- **Borrowings:** Histogram confirms dominance of debt-free firms; boxplot reaffirms outliers with high leverage.
- **Net fixed assets:** Matches earlier asset insight—few firms dominate capital ownership.
- **Current assets:** Central concentration with wide tail; firms differ in short-term asset strategy.
- **Retained earnings:** Reiterates that most firms retain little profit; outliers are financially mature firms.
- **Stock:** Inventory distribution is positively skewed; boxplot confirms high outliers.
- **Reserves and funds:** Most firms keep small reserves; skew and outliers highlight few highly prudent firms.
- **Cash profit:** Wide variance in cash-based profitability, sharply skewed with strong boxplot outliers.

- **Perform Bivariate/Multivariate Analysis** – Use Boxplots with respect to Default variable, followed by Heatmap for all numerical variables

✓ **Display Boxplots against 'Default' variable: -**

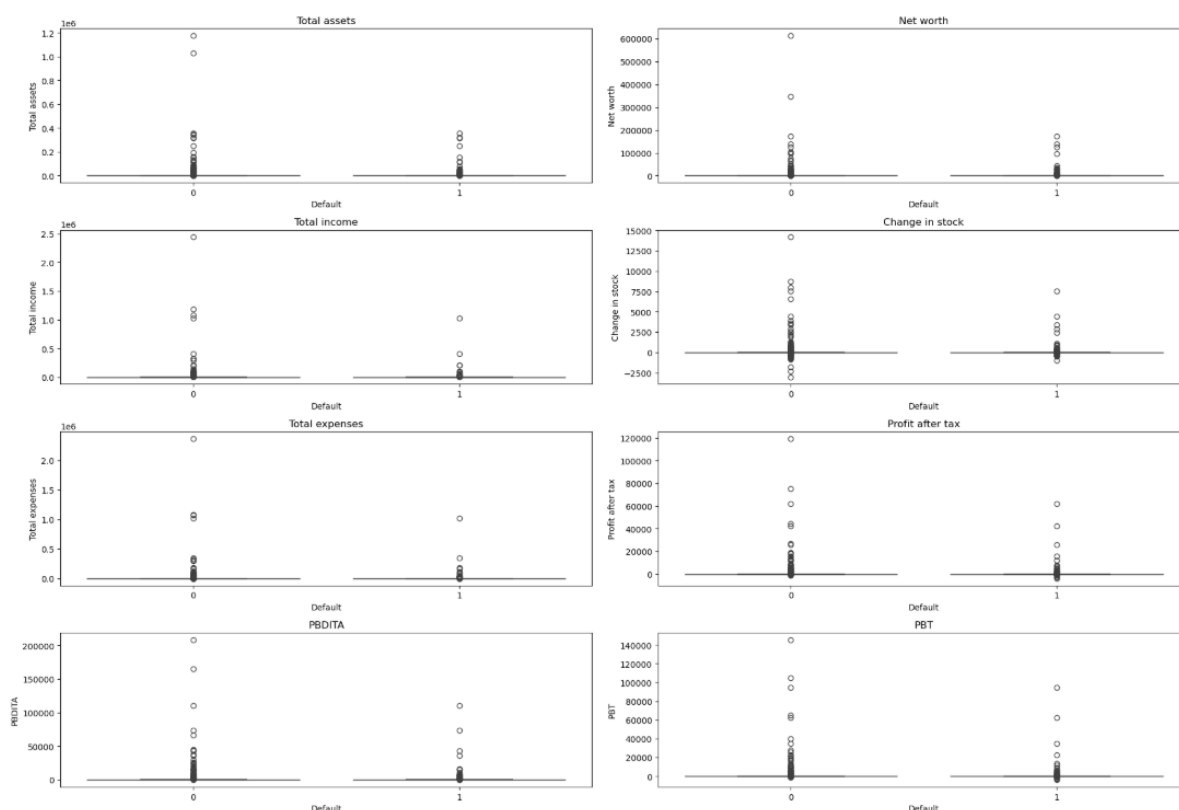


Figure 14: Bivariate/Multivariate (w.r.t. 'Default') Analysis – Boxplot (Part A)

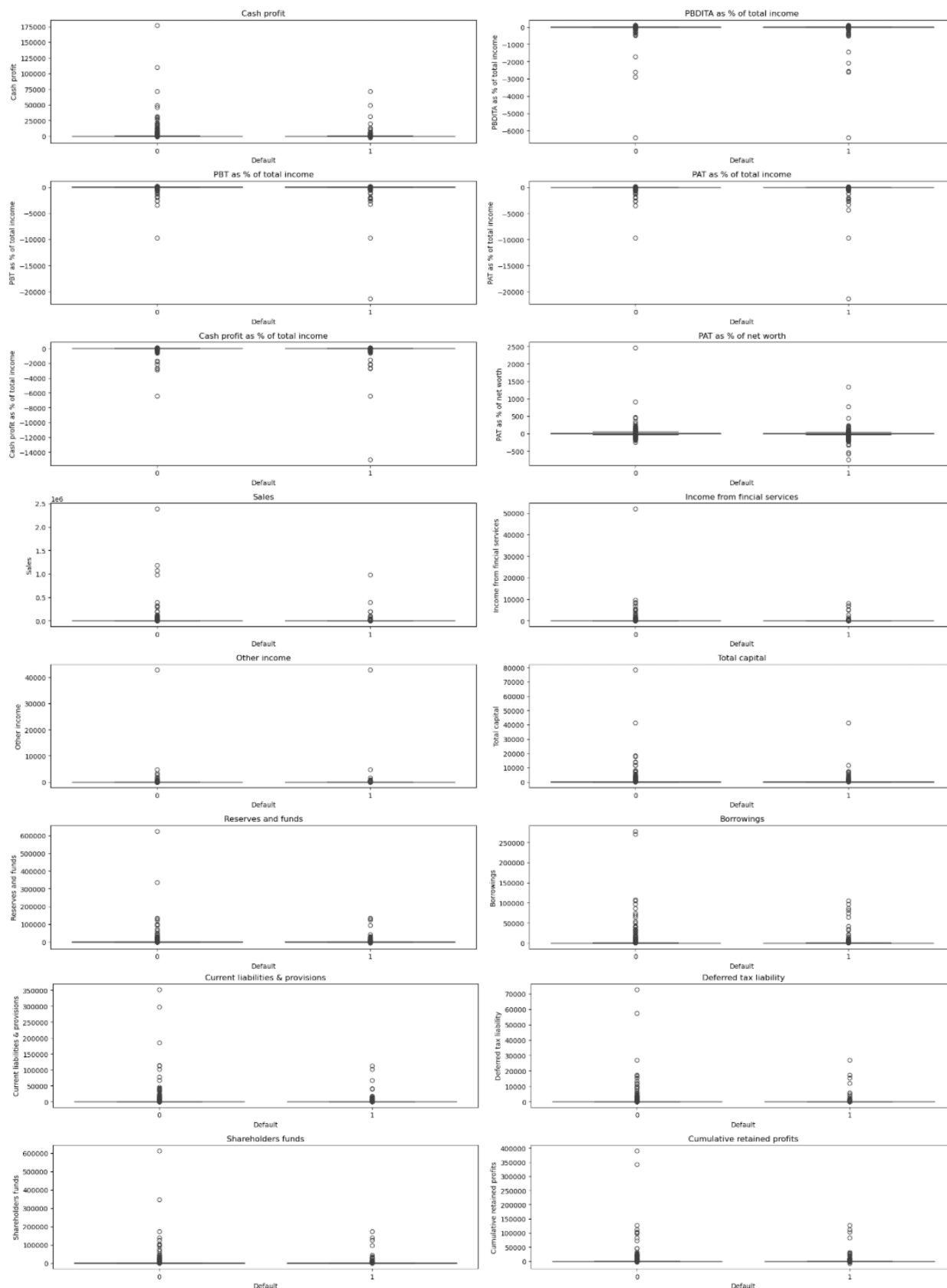


Figure 15: Bivariate/Multivariate (w.r.t. 'Default') Analysis – Boxplot (Part B)

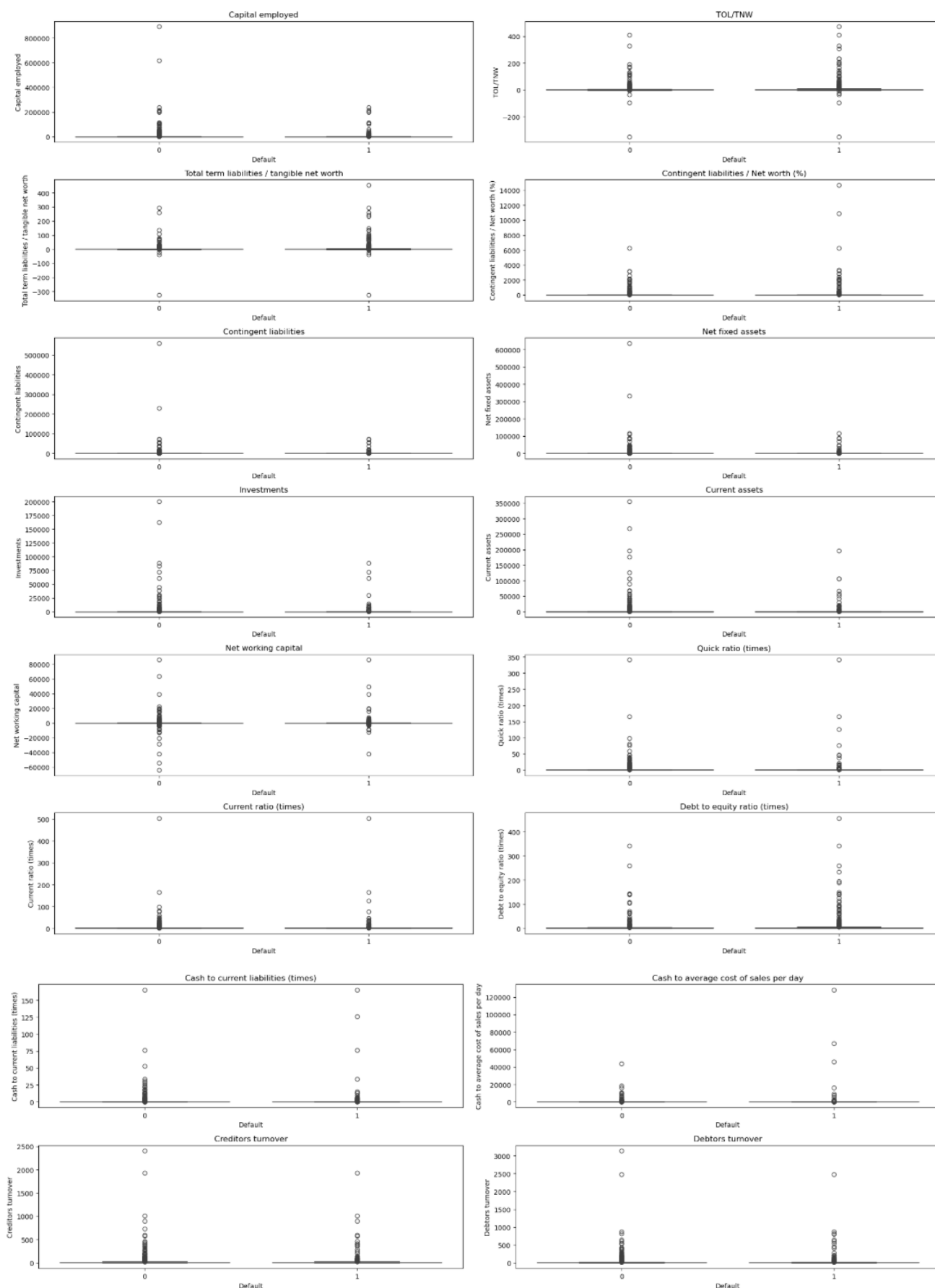


Figure 16: Bivariate/Multivariate (w.r.t. 'Default') Analysis – Boxplot (Part C)

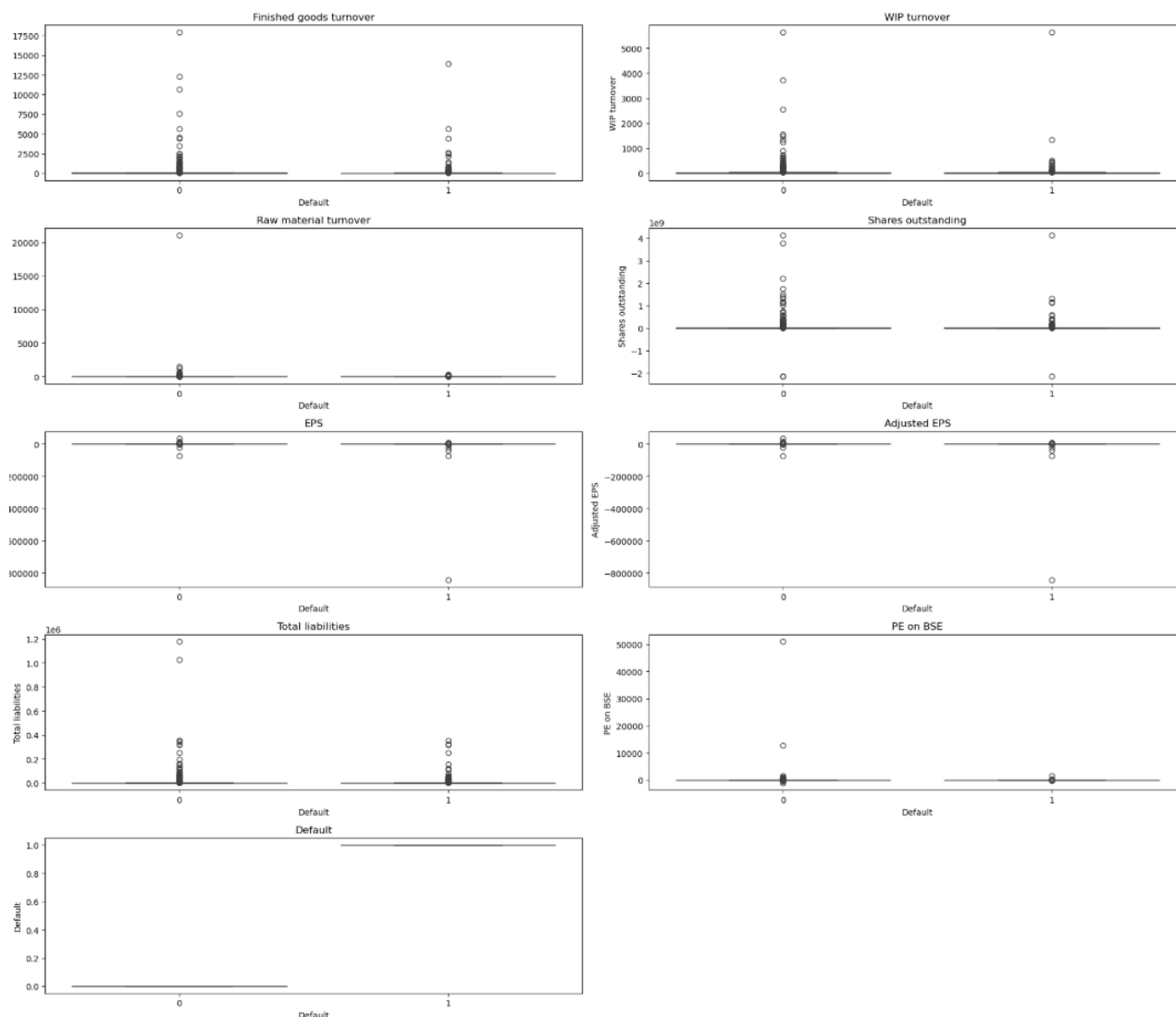


Figure 17: Bivariate/Multivariate (w.r.t. 'Default') Analysis – Boxplot (Part D)

- **Observations: -**
  - **Adjusted EPS:** Defaulters tend to have significantly lower earnings per share.
  - **Cash profit:** Firms that default typically show much lower cash profits.
  - **Cash profit as % of total income:** Profit margins are notably lower in defaulted firms.
  - **Depreciation:** No clear distinction; both groups show overlapping depreciation values.
  - **Interest:** Higher interest expenses are more common among defaulters.
  - **Operating profit margin:** Defaulting firms show weaker operational margins.
  - **PBDITA:** Lower operational profits observed in defaulters.
  - **PBT (Profit before tax):** Defaulted firms often report lower or negative pre-tax profits.
  - **Profit after tax:** Post-tax profitability is considerably lower in defaulted companies.
  - **Equity dividend:** Firms that default generally pay no dividends.
  - **Retained earnings:** Non-defaulters tend to retain significantly more earnings.
  - **Cumulative retained profits:** Strongly lower in default-prone firms.
  - **Total assets:** Defaulters operate with much smaller asset bases.
  - **Net fixed assets:** Lower fixed asset levels are more common among defaulters.
  - **Current assets:** Defaulters generally possess fewer current assets.
  - **Investments:** Very limited investment seen among defaulters.
  - **Reserves and funds:** Defaulters have notably smaller reserves.
  - **Share capital:** Little distinction; most firms show similar equity issuance.
  - **Shareholders' funds:** Defaulting firms are typically undercapitalized.
  - **Net worth:** Clearly lower for defaulters, validating the default condition.

- **Capital employed:** Much lower in default-prone companies.
- **Total income:** Defaulters report significantly less income.
- **Net sales:** Sales volumes are far lower in defaulted firms.
- **Total expenses:** Lower expenses correlate with defaulters, likely due to smaller scale.
- **Value of production:** Firms that default often operate at smaller production scales.
- **Total liabilities:** Higher liabilities are often present in defaulted firms.
- **Total expenses:** Reinforces smaller expense base in defaulters.
- **Borrowings:** Defaulters carry visibly higher borrowings.
- **TOL/TNW:** Highly leveraged firms are more likely to default.
- **Current liabilities & provisions:** Defaulters carry higher short-term obligations.
- **Deferred tax liability:** Slightly higher values observed in defaulters.
- **Quick ratio (times):** Lower liquidity ratios are common among defaulters.
- **Change in stock:** No significant difference; both groups vary similarly.
- **Stock:** Slightly higher stock levels in defaulters, possibly due to inefficiency.
- **Current assets:** Reconfirms reduced liquidity in defaulters.
- **Cash to current liabilities (times):** Lower cash buffers in defaulted firms.
- **Cash to avg. cost of sales/day:** Defaulters have shorter operational cash runways.
- **Creditors turnover:** Slightly slower payment cycles in defaulters.
- **Debtors' turnover:** Defaulters often take longer to collect payments.
- **Finished goods turnover:** Inventory efficiency is weaker in defaulting firms.
- **Sales:** Lower sales levels are consistently tied to defaulters.
- **Retained earnings:** Confirms lower retention in distressed firms.
- **Total income:** Reaffirms income disparity favouring non-defaulters.
- **Borrowings:** Confirms elevated debt levels in defaulters.
- **Net fixed assets:** Matches original—lower in defaulting companies.
- **Stock:** Reinforces inefficiency in defaulted firms.
- **Reserves and funds:** Again, highlights poor financial cushions in defaulters.
- **Cash profit:** Strong disparity—non-defaulters outperform significantly.



✓ **Display Heatmap for all numerical variables: -**

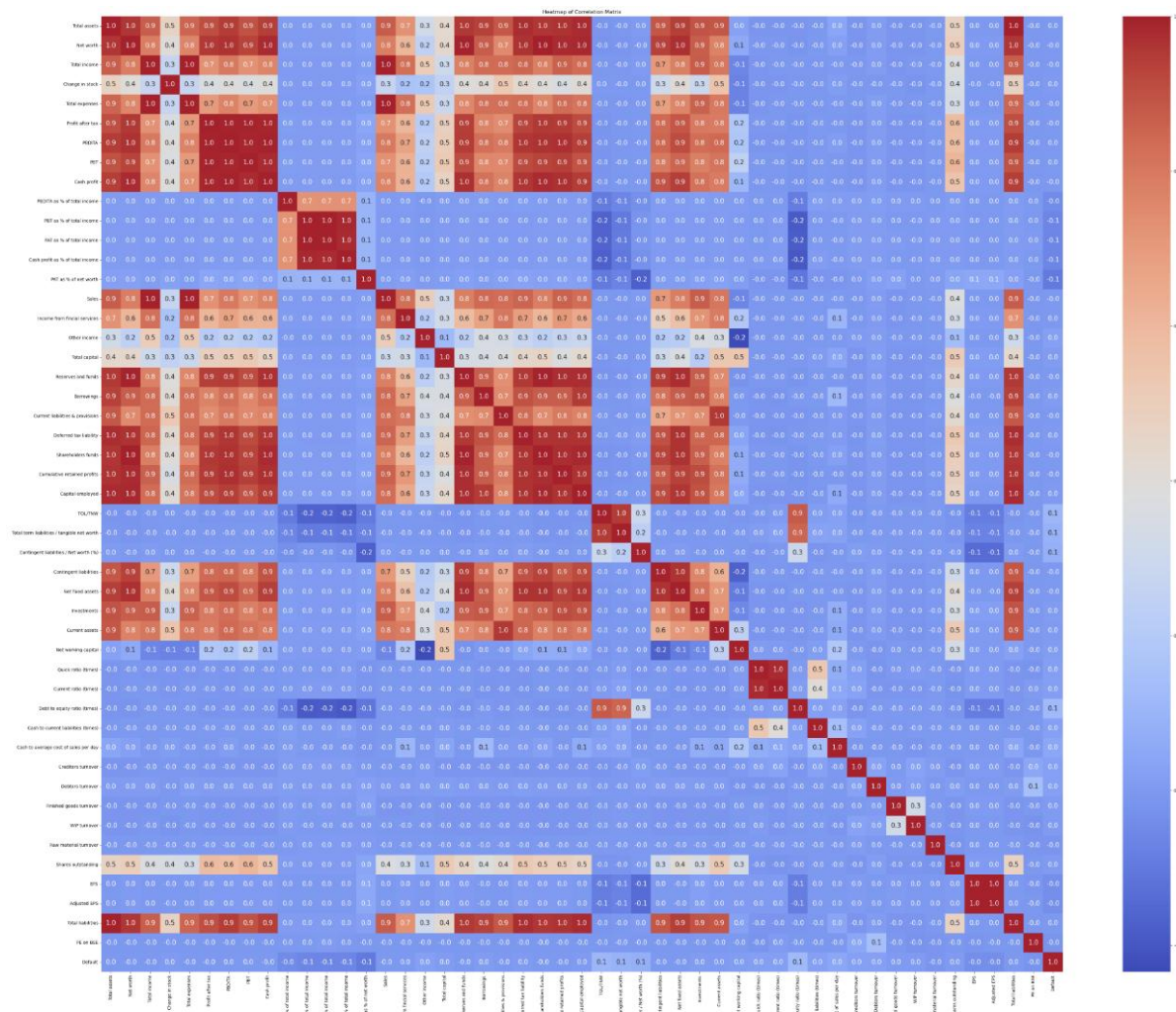


Figure 18: Bivariate Analysis – Heatmap

- **Observations: -**
  - **Strong Positive Correlations:**
    - **Total income, Net sales, Total expenses, and Value of production** are highly interrelated, reflecting firm size and operational scale.
    - **Net worth, Shareholders funds, Capital employed, and Reserves and funds** also exhibit strong mutual correlation, highlighting core equity structure.
  - **Highly Redundant Variables:**
    - **Many variables (e.g., Cash profit, PBDITA, and PBT)** show high overlap and may be redundant in modelling.
    - **Duplicates and derived variables** such as Retained earnings and Cumulative retained profits are tightly correlated.
  - **Negative Correlation with Default:**
    - **Variables like Net worth, Reserves and funds, and Cash profit** show negative correlation with Default, indicating firms with higher values are less likely to default.
  - **Low or No Correlation:**
    - **Equity face value, Deferred tax liability, and Change in stock** show minimal correlation with other variables, suggesting low predictive value.

## Rubric Question 2: Data Preprocessing

### Feature Engineering

- Target Variable Creation – ‘Default’ variable already created from ‘Networth Next Year’ → refer section [Creating ‘Default’ feature](#)
- Drop Redundant Variables → refer section [Drop Redundant features](#)

### Outlier Treatment

- Refer the Univariate section for the boxplots of all variables → refer section [Perform Univariate Analysis](#)
- Clearly there are lot of outliers in the dataset. Below is the outlier percentage for each variable: -

	Column	Outliers Percentage
0	Total assets	13.74530
1	Net worth	13.98026
2	Total income	11.93609
3	Change in stock	17.62218
4	Total expenses	12.17105
5	Profit after tax	16.72932
6	PBDITA	13.72180
7	PBT	16.54135
8	Cash profit	14.73214
9	PBDITA as % of total income	8.12970
10	PBT as % of total income	12.82895
11	PAT as % of total income	14.33271
12	Cash profit as % of total income	10.00940
13	PAT as % of net worth	10.03289
14	Sales	11.74812
15	Income from financial services	12.14756
16	Other income	9.14004
17	Total capital	12.94643
18	Reserves and funds	15.10808
19	Borrowings	12.50000
20	Current liabilities & provisions	13.65132
21	Deferred tax liability	9.53947
22	Shareholders funds	13.81579
23	Cumulative retained profits	16.42387
24	Capital employed	13.43985
25	TOL/TNW	9.72744
26	Total term liabilities / tangible net worth	9.53947
27	Contingent liabilities / Net worth (%)	11.23120
28	Contingent liabilities	9.23402
29	Net fixed assets	13.36936
30	Investments	10.59680
31	Current assets	12.50000
32	Net working capital	18.93797
33	Quick ratio (times)	8.71711
34	Current ratio (times)	9.32801
35	Debt to equity ratio (times)	8.95207
36	Cash to current liabilities (times)	12.66447
37	Cash to average cost of sales per day	13.69831
38	Creditors turnover	10.38534
39	Debtors turnover	9.58647
40	Finished goods turnover	9.37500
41	WIP turnover	8.88158
42	Raw material turnover	6.95489
43	Shares outstanding	11.18421
44	EPS	14.99060
45	Adjusted EPS	16.30639
46	Total liabilities	13.74530
47	PE on BSE	5.56861
48	Default	21.24060

Figure 19: Outlier Inspection

- Since the outlier count is significant & considering the business case where different companies of different sizes can have varied financial information, treating or removing these outliers don't make sense. Hence, **we choose not to treat outliers.**
- Please note, the 'Default' variable is the Target variable & treating its outliers makes no sense.

## Data Preparation for Modelling

- As a data scientist, we want to analyse the data provided to find which factors have a high influence on companies turning Default, build a predictive model that can predict which companies have a higher likelihood of defaulting.
  - We split the data into Dependent (Default) & Independent variable (remaining others) Data Frames.
  - Before we proceed to build a model, we must encode categorical features. Since, we don't have any categorical variables, we skip this step for this problem.
  - We split the data into Train and Test datasets to be able to evaluate the model that we build on the Train data. We build a model using the Train dataset and then check its performance on Test dataset.
- We use 'random\_state' value 42 to split the data into Train and Test Datasets in the ratio of 75:25.

## Duplicate/Missing Treatment & Data Scaling

- Please refer [Check Missing & Duplicate Values](#) section for inspection.
- No treatment required for Duplicated values.
- However, missing value treatment is definitely required.
- Now that we have split the data into train & test datasets, let's inspect for missing values again on these datasets: -

Checking missing values in X_train		Checking missing values in X_test	
Total assets	0	Total assets	0
Net worth	0	Net worth	0
Total income	183	Total income	48
Change in stock	427	Change in stock	123
Total expenses	133	Total expenses	32
Profit after tax	125	Profit after tax	29
PBDITA	125	PBDITA	29
PBT	125	PBT	29
Cash profit	125	Cash profit	29
PBDITA as % of total income	59	PBDITA as % of total income	20
PBT as % of total income	59	PBT as % of total income	20
PAT as % of total income	59	PAT as % of total income	20
Cash profit as % of total income	59	Cash profit as % of total income	20
PAT as % of net worth	0	PAT as % of net worth	0
Sales	237	Sales	68
Income from financial services	843	Income from financial services	268
Other income	1179	Other income	377
Total capital	4	Total capital	1
Reserves and funds	78	Reserves and funds	20
Borrowings	330	Borrowings	101
Current liabilities & provisions	85	Current liabilities & provisions	25
Deferred tax liability	1036	Deferred tax liability	333
Shareholders funds	0	Shareholders funds	0
Cumulative retained profits	35	Cumulative retained profits	10
Capital employed	0	Capital employed	0
TOL/TNW	0	TOL/TNW	0
Total term liabilities / tangible net worth	0	Total term liabilities / tangible net worth	0
Contingent liabilities / Net worth (%)	0	Contingent liabilities / Net worth (%)	0
Contingent liabilities	1052	Contingent liabilities	350
Net fixed assets	101	Net fixed assets	31
Investments	1290	Investments	425
Current assets	61	Current assets	19
Net working capital	29	Net working capital	8
Quick ratio (times)	82	Quick ratio (times)	23
Current ratio (times)	82	Current ratio (times)	23
Debt to equity ratio (times)	0	Debt to equity ratio (times)	0
Cash to current liabilities (times)	82	Cash to current liabilities (times)	23
Cash to average cost of sales per day	78	Cash to average cost of sales per day	22
Creditors turnover	297	Creditors turnover	94
Debtors turnover	296	Debtors turnover	89
Finished goods turnover	675	Finished goods turnover	199
WIP turnover	588	WIP turnover	176
Raw material turnover	327	Raw material turnover	101
Shares outstanding	609	Shares outstanding	201
EPS	0	EPS	0
Adjusted EPS	0	Adjusted EPS	0
Total liabilities	0	Total liabilities	0
PE on BSE	1972	PE on BSE	655

Checking missing values in y\_train

0

Checking missing values in y\_test

0

Figure 20: Missing-value check on Train/Test datasets

- Technique used to impute missing values – **KNNImputer**
- KNNImputer** is an unsupervised imputer that replaces missing values in a dataset with the distance-weighted average of the samples' k nearest neighbours' values
- Pre-requisite for KNNImputer is that **dataset should be scaled** for accurate calculations. We use **StandardScaler()** to carry out **scaling of train & test datasets**. Below is the subset of trained dataset post scaling (truncated view): -

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	PAT as % of net worth	Sales
-0.14457	-0.10217	-0.09561	-0.34679	-0.09400	-0.10925	-0.12604	-0.10489	-0.11953	0.03142	0.05371	0.05322	0.05353	-0.10615	-0.09682
-0.14401	-0.13169	-0.04922	-0.10685	-0.04365	-0.10763	-0.12247	-0.10193	-0.11691	0.00519	0.04820	0.04962	0.04024	-0.15907	-0.04919
-0.16773	-0.15526	-0.11017	-0.11074	-0.10640	-0.11504	-0.13656	-0.11150	-0.13270	0.25968	0.15482	0.14771	0.17395	-0.00215	-0.11151
-0.13771	-0.14204	-0.09002	0.05817	-0.08419	-0.10877	-0.11259	-0.10498	-0.10124	0.05540	0.05131	0.05207	0.06649	0.06462	-0.09099
-0.12843	-0.14670	-0.09146	-0.51040	-0.09116	-0.11238	-0.11773	-0.10972	-0.13551	0.04593	0.04736	0.05008	0.03236	0.02893	-0.09371

Figure 21: Subset of Train dataset post Data Scaling

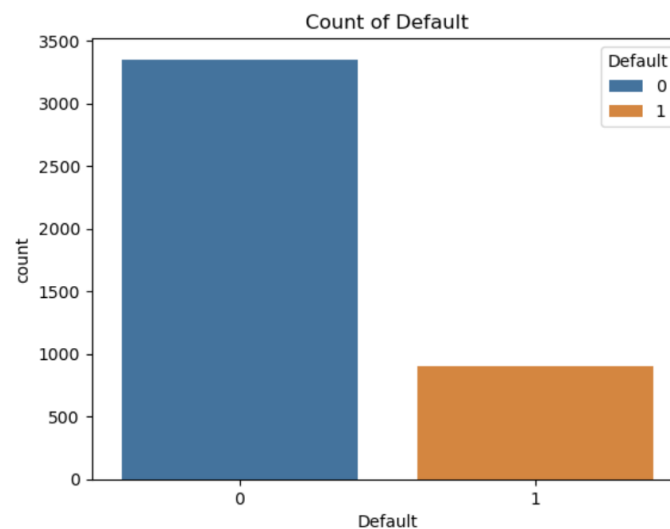
- Post data-scaling, both the train & test datasets are imputed with missing values using **KNNImputer(n\_neighbors=5)**. Below are the missing value results post missing-value treatment: -

Missing Values in X_train post Missing Value Treatment:		Missing Values in X_test post Missing Value Treatment:	
Total assets	0	Total assets	0
Net worth	0	Net worth	0
Total income	0	Total income	0
Change in stock	0	Change in stock	0
Total expenses	0	Total expenses	0
Profit after tax	0	Profit after tax	0
PBDITA	0	PBDITA	0
PBT	0	PBT	0
Cash profit	0	Cash profit	0
PBDITA as % of total income	0	PBDITA as % of total income	0
PBT as % of total income	0	PBT as % of total income	0
PAT as % of total income	0	PAT as % of total income	0
Cash profit as % of total income	0	Cash profit as % of total income	0
PAT as % of net worth	0	PAT as % of net worth	0
Sales	0	Sales	0
Income from financial services	0	Income from financial services	0
Other income	0	Other income	0
Total capital	0	Total capital	0
Reserves and funds	0	Reserves and funds	0
Borrowings	0	Borrowings	0
Current liabilities & provisions	0	Current liabilities & provisions	0
Deferred tax liability	0	Deferred tax liability	0
Shareholders funds	0	Shareholders funds	0
Cumulative retained profits	0	Cumulative retained profits	0
Capital employed	0	Capital employed	0
TOL/TNW	0	TOL/TNW	0
Total term liabilities / tangible net worth	0	Total term liabilities / tangible net worth	0
Contingent liabilities / Net worth (%)	0	Contingent liabilities / Net worth (%)	0
Contingent liabilities	0	Contingent liabilities	0
Net fixed assets	0	Net fixed assets	0
Investments	0	Investments	0
Current assets	0	Current assets	0
Net working capital	0	Net working capital	0
Quick ratio (times)	0	Quick ratio (times)	0
Current ratio (times)	0	Current ratio (times)	0
Debt to equity ratio (times)	0	Debt to equity ratio (times)	0
Cash to current liabilities (times)	0	Cash to current liabilities (times)	0
Cash to average cost of sales per day	0	Cash to average cost of sales per day	0
Creditors turnover	0	Creditors turnover	0
Debtors turnover	0	Debtors turnover	0
Finished goods turnover	0	Finished goods turnover	0
WIP turnover	0	WIP turnover	0
Raw material turnover	0	Raw material turnover	0
Shares outstanding	0	Shares outstanding	0
EPS	0	EPS	0
Adjusted EPS	0	Adjusted EPS	0
Total liabilities	0	Total liabilities	0
PE on BSE	0	PE on BSE	0

Figure 22: Missing-value check on Train/Test datasets post data-treatment (KNNImputer)

## Check & Treat Imbalance in Target Variable (Default)

- Let's look at the distribution of 2 values (1 & 0) of 'Default' in the dataset: -



Percentage of Defaulters in the Dataset = 21.24 %

Figure 23: Check Imbalance in Target Variable

- Clearly, with 21% defaulters in the dataset will have a big impact during model building. Untreated imbalance will lead to **under-detection of defaulters**, reducing Recall for the minority class.
- To **address imbalance** in dataset, we shall use one of the oversampling techniques, **SMOTE**. Please note that we shall only **treat the training dataset only** & skip any treatment in the test dataset to avoid any noise-creep in the dataset kept for final testing.
- SMOTE (Synthetic Minority Oversampling Technique)** is a **machine learning resampling technique** used to address class imbalance in datasets. It works by generating synthetic minority class samples, effectively increasing the number of minority examples in the dataset. This helps to prevent biased models that might overfit to the majority class. SMOTE generates synthetic minority samples by finding the k-nearest neighbours of each minority instance.
- Below is the **Train dataset summary**, before & after running SMOTE: -

Before SMOTE, counts of label 'Yes': 678  
Before SMOTE, counts of label 'No': 2514

After SMOTE, counts of label 'Yes': 2514  
After SMOTE, counts of label 'No': 2514

After SMOTE, the shape of train\_X: (5028, 48)  
After SMOTE, the shape of train\_y: (5028,)

Figure 24: Train dataset summary - pre/post SMOTE

## Rubric Question 3: Model building

### Model Evaluation Criteria

- **Accuracy**
  - Definition: Measures the proportion of correct predictions (both True Positives and True Negatives).
  - Business Implications: Accuracy is useful only if the dataset is balanced (equal distribution of cancelled and not-cancelled bookings), which is not in our case.
  - Verdict: **Accuracy is not the best metric for this problem**, especially if the dataset is imbalanced.
- **Precision**
  - Definition: Measures the proportion of Correctly predicted "No Defaulters" out of all predicted as "No Defaulters".
  - Precision =  $TP / (TP + FP)$
  - Business Implications: High precision implies that there are **few false positives** (non-defaulters wrongly flagged as defaulters).
  - Precision is more important when the cost of False Positives (FP) (e.g., wasted resources) is critical.
  - Verdict: **Precision is useful but less critical than Recall in this problem** as focus on high precision may lower recall—i.e., the model may miss some actual defaulters to avoid false alarms.
- **Recall – Metric of Choice for this Problem**
  - Definition: Measures the proportion of correctly predicted "Not Cancelled" bookings out of all actual "Not Cancelled" bookings.
  - Recall =  $TP / (TP + FN)$
  - Business Implications: High recall ensures the hotel accurately identifies most of the Defaulters. A model with high recall will result in fewer False Negatives (FN) (fewer missed defaulters), ensuring **risky entities are not overlooked**, thereby, **reducing the risk of Non-Performing Assets (NPAs)**.
  - Recall is critical in this problem because minimizing False Negatives (overbooking) is vital for business success.
  - Verdict: **Recall is the most important metric for the hotel booking problem**.
- ✓ **F1-Score**
  - Definition: The harmonic means of Precision and Recall, balancing the two.
  - $F1 = 2 \times (Precision \times Recall) / (Precision + Recall)$
  - Business Implications: F1-Score balances Precision and Recall. It's useful when you need to ensure both (a) Avoiding wasted resources (Precision) (b) Avoiding defaulters (Recall).
  - Verdict: **F1-Score** may be a good secondary metric but **Recall remains the primary metric to prioritize** given the context of the problem where the objective is to predict defaulters.
- ✓ **Final Metric Verdict**
  - **Primary Metric: Recall:** Missing a defaulter is more damaging than false alarms



## Logistic Regression Model

### Build Model

- Before building the model, we add a constant to the train & test dataset. Below are the top 5 rows of train dataset post intercept addition (truncated view): -

	const	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	PAT as % of net worth	Sales	Income from financial services
0	1.00000	-0.14457	-0.10217	-0.09561	-0.34679	-0.09400	-0.10925	-0.12604	-0.10489	-0.11953	0.03142	0.05371	0.05322	0.05353	-0.10615	-0.09682	-0.13051
1	1.00000	-0.14401	-0.13169	-0.04922	-0.10685	-0.04365	-0.10763	-0.12247	-0.10193	-0.11691	0.00519	0.04820	0.04962	0.04024	-0.15907	-0.04919	-0.13542
2	1.00000	-0.16773	-0.15526	-0.11017	-0.11074	-0.10640	-0.11504	-0.13656	-0.11150	-0.13270	0.25968	0.15482	0.14771	0.17395	-0.00215	-0.11151	-0.12850
3	1.00000	-0.13771	-0.14204	-0.09002	0.05817	-0.08419	-0.10877	-0.11259	-0.10498	-0.10124	0.05540	0.05131	0.05207	0.06649	0.06462	-0.09099	-0.13475
4	1.00000	-0.12843	-0.14670	-0.09146	-0.51040	-0.09116	-0.11238	-0.11773	-0.10972	-0.13551	0.04593	0.04736	0.05008	0.03236	0.02893	-0.09371	-0.13542

Figure 25: Subset of Train dataset with Intercept

- Train Dataset (with Intercept) is passed through the Logit Function.
- Below is the summary of the model output: -

```
Optimization terminated successfully.
Current function value: 0.659264
Iterations 9
```

Logit Regression Results						
Dep. Variable:	Default	No. Observations:	5028			
Model:	Logit	Df Residuals:	4980			
Method:	MLE	Df Model:	47			
Date:	Wed, 28 May 2025	Pseudo R-squ.:	0.04888			
Time:	22:37:20	Log-Likelihood:	-3314.8			
converged:	True	LL-Null:	-3485.1			
Covariance Type:	nonrobust	LLR p-value:	3.540e-46			
	coef	std err	z	P> z	[0.025	0.975]
const	-0.0646	0.080	-0.805	0.421	-0.222	0.093
Total assets	0.2243	nan	nan	nan	nan	nan
Net worth	1.8278	1.685	1.085	0.278	-1.474	5.129
Total income	-62.9139	29.330	-2.145	0.032	-120.399	-5.429
Change in stock	-0.4173	0.253	-1.652	0.099	-0.912	0.078
Total expenses	38.7625	24.925	1.555	0.120	-10.089	87.614
Profit after tax	5.1744	2.071	2.499	0.012	1.116	9.233
PBDITA	1.5004	1.192	1.259	0.208	-0.836	3.837
PBT	-1.0786	1.251	-0.862	0.389	-3.531	1.374
Cash profit	-3.6320	0.866	-4.195	0.000	-5.329	-1.935
PBDITA as % of total income	0.2134	0.129	1.657	0.098	-0.039	0.466
PBT as % of total income	1.6512	1.418	1.165	0.244	-1.127	4.430
PAT as % of total income	-1.9006	1.387	-1.370	0.171	-4.619	0.818
Cash profit as % of total income	-0.1657	0.291	-0.571	0.568	-0.735	0.404
PAT as % of net worth	-0.4490	0.070	-6.391	0.000	-0.587	-0.311
Sales	22.1156	12.605	1.755	0.079	-2.589	46.821
Income from financial services	0.0831	0.196	0.425	0.671	-0.301	0.467
Other income	0.8360	0.560	1.492	0.136	-0.262	1.934
Total capital	-0.3634	0.491	-0.741	0.459	-1.325	0.598
Reserves and funds	-3.3437	3.145	-1.063	0.288	-9.508	2.821
Borrowings	0.6086	1.119	0.544	0.586	-1.584	2.801
Current liabilities & provisions	0.9801	2.266	0.433	0.665	-3.461	5.422
Deferred tax liability	-0.1340	0.662	-0.203	0.839	-1.431	1.163
Shareholders funds	2.4088	3.746	0.643	0.520	-4.934	9.752
Cumulative retained profits	-0.0182	0.378	-0.048	0.962	-0.759	0.722
Capital employed	-1.8214	5.649	-0.322	0.747	-12.893	9.250
TOL/TNW	0.1646	0.157	1.045	0.296	-0.144	0.473
Total term liabilities / tangible net worth	-0.2515	0.170	-1.483	0.138	-0.584	0.081
Contingent liabilities / Net worth (%)	0.0934	0.056	1.665	0.096	-0.017	0.203
Contingent liabilities	0.3598	0.217	1.660	0.097	-0.065	0.785
Net fixed assets	0.4634	0.418	1.108	0.268	-0.356	1.283
Investments	0.1357	0.265	0.512	0.609	-0.384	0.655
Current assets	-1.0324	0.616	-1.675	0.094	-2.240	0.176
Net working capital	0.3188	0.243	1.311	0.190	-0.158	0.795
Quick ratio (times)	-0.0633	0.263	-0.241	0.810	-0.578	0.452
Current ratio (times)	0.0618	0.248	0.249	0.803	-0.425	0.548
Debt to equity ratio (times)	0.3994	0.106	3.780	0.000	0.192	0.607
Cash to current liabilities (times)	0.0191	0.051	0.372	0.710	-0.082	0.120
Cash to average cost of sales per day	0.1286	0.081	1.587	0.112	-0.030	0.287
Creditors turnover	0.0642	0.034	1.875	0.061	-0.003	0.131
Debtors turnover	0.0699	0.038	1.849	0.065	-0.004	0.144
Finished goods turnover	-0.0277	0.042	-0.663	0.507	-0.110	0.054
WIP turnover	-0.0923	0.056	-1.655	0.098	-0.202	0.017
Raw material turnover	-0.8754	0.477	-1.834	0.067	-1.811	0.060
Shares outstanding	0.0402	0.064	0.628	0.530	-0.085	0.166
EPS	210.7929	4.959	42.507	0.000	201.073	220.512
Adjusted EPS	-211.0712	4.764	-44.309	0.000	-220.408	-201.735
Total liabilities	0.2243	nan	nan	nan	nan	nan
PE on BSE	-0.0147	0.049	-0.300	0.764	-0.111	0.082

Figure 26: Logistic Regression Summary

▪ **Concise summary of the model output: -**

- ✓ The model is statistically significant (LLR p-value < 0.001) but explains limited variance (Pseudo  $R^2 = 0.0489$ ).
- ✓ Key predictors reducing default risk: Higher total income and Adjusted EPS (strongest negative coefficient).
- ✓ Some profit-related variables (e.g., PAT, Cash profit) have unexpected positive signs, indicating possible multicollinearity.
- ✓ Liquidity and solvency ratios (e.g., current ratio, debt to equity) are not significant in isolation.
- ✓ Overall, the model needs refinement for better predictive power.

## Checking Model Performance

▪ Below is the Confusion Matrix & Model Performance Metrics on Training dataset: -

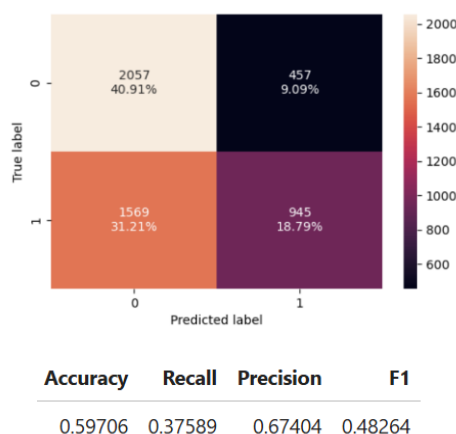


Figure 27: Logistic Regression: Confusion Matrix & Metric Performance on Training dataset

▪ Below is the Confusion Matrix & Model Performance Metrics on Test dataset: -

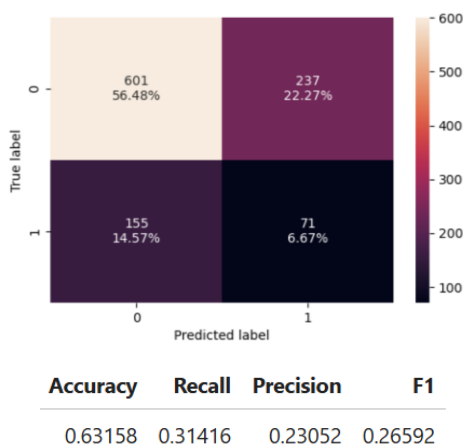


Figure 28: Logistic Regression: Confusion Matrix & Metric Performance on Test dataset

▪ **Summary of Performance: -**

- **Model Generalization:** The Recall metrics for both the training and testing datasets are very similar, indicating that the logistic regression model generalizes well and is not overfitting.
- **Strengths:**
  - ✓ **Decent Accuracy** (~63%) suggesting Indicates the model performs better than random guessing.
- **Weaknesses:**
  - ✓ **Very Low Recall** (~31%): The model misses ~69% of actual defaulters (only 71 out of 226 detected). This model cannot be trusted fully as we have not dealt with multicollinearity yet & removed non-significant p-values. This is a critical weakness, as the main business goal is to flag as many defaulters as possible.
  - ✓ **Very Low Precision** (~23%): This suggests poor reliability of the model's predictions and will lead to many false alarms. When the model predicts a company will default, it has only 23% chance of being accurate.
  - ✓ **Low F1 Score** (~0.26): Reflects poor balance between precision and recall.
- On the test set, the **model fails to serve its primary purpose** of identifying defaulters due to low recall and F1 score. Despite acceptable accuracy, it is not reliable for real-world risk detection.



## Random Forest Model

### Build Model

- A Random Forest is a supervised ensemble machine learning algorithm that combines multiple decision trees to make predictions.
- We use `RandomForestClassifier(random_state=42)` to build the Random Forest Model using the training dataset.

### Checking Model Performance

- Below is the Confusion Matrix & Model Performance Metrics on Training dataset: -

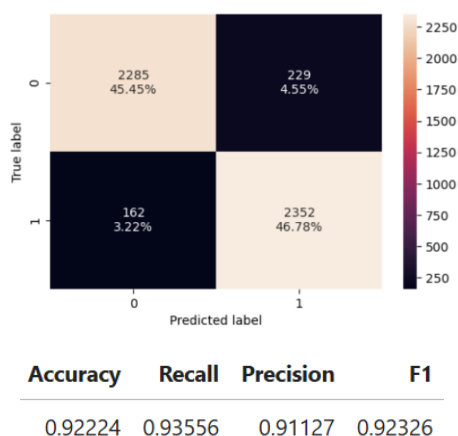


Figure 29: Random Forest: Confusion Matrix & Metric Performance on Training dataset

- Below is the Confusion Matrix & Model Performance Metrics on Test dataset: -

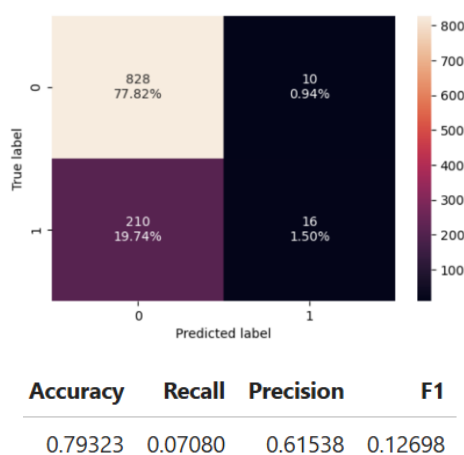


Figure 30: Random Forest: Confusion Matrix & Metric Performance on Test dataset

- **Summary of Performance: -**
  - **Model Generalization:** The difference in the Recall metrics for both the training and testing datasets is very high, indicating that the Random Forest model **suffers from overfitting** & doesn't generalize well.
  - **Strengths:**
    - ✓ **High Overall Accuracy (~79%):** Performs well for majority class (non-defaulters).
    - ✓ **Moderate Precision (~61):** Model predicts a company will default; it has 61% chance of being accurate
  - **Weaknesses:**
    - ✓ **Extremely Low Recall (~8%):** Model misses 92% of actual defaulters, defeating the core business goal of default detection.
    - ✓ **Very Low F1 Score (~0.12):** Indicates poor balance between precision and recall.
    - ✓ **Overfitting to Training Data:** Huge performance gap between training and test recall (93.5% → 7.08%). Model learns training data too well but fails to generalize.
  - The **Random Forest model shows severe overfitting**. While it performs excellently on training data, it fails to generalize; capturing only **7% of defaulters** on the test set, making it **unreliable for real-world risk detection**.

## Rubric Question 4: Model Performance Improvement

### Logistic Regression Model – Tuning

#### Build Tuned Model

- **Dealing with Multicollinearity: -**
  - Multicollinearity occurs when predictor variables in a regression model are correlated. This correlation is a problem because predictor variables should be independent. If the correlation between variables is high, it can cause problems when we fit the model and interpret the results. When we have multicollinearity in the model, the coefficients that the model suggests are unreliable.
  - **Variance Inflation factor:** Variance inflation factors measure the inflation in the variances of the regression coefficients estimates due to collinearities that exist among the predictors. It is a measure of how much the variance of the estimated regression coefficient  $\beta_k$  is 'inflated' by the existence of correlation among the predictor variables in the model.
  - General Rule of Thumb while interpreting VIF: -
    - ✓ If VIF is 1, then there is no correlation among the  $k$ th predictor and the remaining predictor variables, and hence, the variance of  $\beta_k$  is not inflated at all.
    - ✓ If VIF exceeds 5, we say there is moderate VIF, and if it is 10 or exceeding 10, it shows signs of high multi-collinearity.
  - Below is the VIF for all independent variables: -

Variance Inflation Factors:		
	Variable	VIF
0	Total assets	inf
46	Total liabilities	inf
45	Adjusted EPS	2089147.54510
44	EPS	2089144.50303
2	Total income	385487.18437
4	Total expenses	192959.62521
14	Sales	67242.32902
24	Capital employed	9608.67052
22	Shareholders funds	4851.65224
1	Net worth	2367.61412
5	Profit after tax	2206.24072
10	PBT as % of total income	2129.57506
11	PAT as % of total income	2008.59364
7	PBT	1235.05823
20	Current liabilities & provisions	960.51923
19	Borrowings	866.84859
6	PBDITA	745.96346
18	Reserves and funds	743.67904
8	Cash profit	516.34866
23	Cumulative retained profits	128.79489
31	Current assets	118.32839
33	Quick ratio (times)	81.27745
29	Net fixed assets	78.97122
16	Other income	78.26072
34	Current ratio (times)	69.35457
12	Cash profit as % of total income	57.29069
21	Deferred tax liability	49.77574
17	Total capital	33.26338
3	Change in stock	20.19297
32	Net working capital	20.05024
30	Investments	16.55961
25	TOL/TNW	14.93245
28	Contingent liabilities	13.51993
26	Total term liabilities / tangible net worth	12.24888
15	Income from financial services	12.05318
35	Debt to equity ratio (times)	5.60751
36	Cash to current liabilities (times)	3.76160
43	Shares outstanding	2.97301
37	Cash to average cost of sales per day	2.62771
9	PBDITA as % of total income	2.54937
13	PAT as % of net worth	1.16188
27	Contingent liabilities / Net worth (%)	1.15495
38	Creditors turnover	1.06640
41	WIP turnover	1.06042
40	Finished goods turnover	1.05487
39	Debtors turnover	1.02353
47	PE on BSE	1.00490
42	Raw material turnover	1.00049

Figure 31: VIF for Independent Variables

- Following steps to be taken to fix Multicollinearity (if any): -
  - ✓ In case any of the VIF was greater than 5, we follow the below steps iteratively: -
    - Drop every column one by one that has a VIF score greater than 5.
    - Check the VIF scores again.
    - Continue till we have all variables with VIF scores under 5.
- Below is the final VIF Table after running the above algorithm, that ran for 24 iterations, removing 24 variables. Below is the list of variables with VIF < 5: -

Variable	VIF
Debt to equity ratio (times)	4.63060
Total term liabilities / tangible net worth	4.27768
Total capital	3.25930
Income from fincial services	2.95900
PAT as % of total income	2.69848
Investments	2.67654
Contingent liabilities	2.58505
PBDITA as % of total income	2.50820
Shares outstanding	2.25843
Net working capital	1.96956
Other income	1.86568
Change in stock	1.57035
Cash to average cost of sales per day	1.34563
Cash to current liabilities (times)	1.32256
Current ratio (times)	1.19539
Contingent liabilities / Net worth (%)	1.09640
PAT as % of net worth	1.09295
WIP turnover	1.05588
Finished goods turnover	1.05279
EPS	1.02250
Creditors turnover	1.01111
Debtors turnover	1.00890
PE on BSE	1.00295
Raw material turnover	1.00036

Figure 32: VIF for Independent Variables (post-Multicollinearity fix)

- Building the Logistic Regression Model again & displaying the summary: -

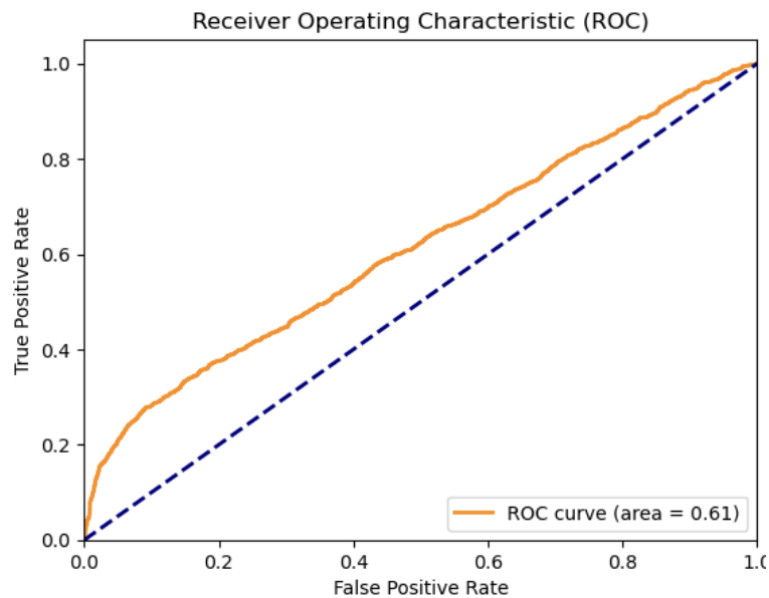
Optimization terminated successfully.  
Current function value: 0.664444  
Iterations 9

Logit Regression Results						
Dep. Variable:	Default	No. Observations:	5028			
Model:	Logit	Df Residuals:	5003			
Method:	MLE	Df Model:	24			
Date:	Wed, 28 May 2025	Pseudo R-squ.:	0.04141			
Time:	22:45:27	Log-Likelihood:	-3340.8			
converged:	True	LL-Null:	-3485.1			
Covariance Type:	nonrobust	LLR p-value:	3.225e-47			
	coef	std err	z	P> z	[0.025	0.975]
const	-0.0546	0.032	-1.698	0.089	-0.118	0.008
Change in stock	-0.0926	0.053	-1.746	0.081	-0.197	0.011
PBDITA as % of total income	0.1693	0.105	1.618	0.106	-0.036	0.374
PAT as % of total income	-0.3543	0.168	-2.106	0.035	-0.684	-0.025
PAT as % of net worth	-0.4381	0.059	-7.484	0.000	-0.553	-0.323
Income from fincial services	-0.1457	0.092	-1.583	0.113	-0.326	0.035
Other income	0.0206	0.066	0.314	0.754	-0.108	0.149
Total capital	0.0632	0.056	1.122	0.262	-0.047	0.174
Total term liabilities / tangible net worth	-0.1114	0.099	-1.128	0.259	-0.305	0.082
Contingent liabilities / Net worth (%)	0.1030	0.053	1.926	0.054	-0.002	0.208
Contingent liabilities	0.1780	0.070	2.540	0.011	0.041	0.315
Investments	-0.0720	0.091	-0.790	0.430	-0.251	0.107
Net working capital	-0.0141	0.051	-0.274	0.784	-0.115	0.087
Current ratio (times)	0.0018	0.031	0.058	0.954	-0.059	0.063
Debt to equity ratio (times)	0.4342	0.099	4.369	0.000	0.239	0.629
Cash to current liabilities (times)	0.0057	0.034	0.168	0.867	-0.061	0.073
Cash to average cost of sales per day	0.1555	0.071	2.178	0.029	0.016	0.295
Creditors turnover	0.0597	0.033	1.806	0.071	-0.005	0.124
Debtors turnover	0.0765	0.038	2.038	0.042	0.003	0.150
Finished goods turnover	-0.0240	0.041	-0.580	0.562	-0.105	0.057
WIP turnover	-0.0919	0.055	-1.678	0.093	-0.199	0.015
Raw material turnover	-0.8756	0.474	-1.847	0.065	-1.805	0.054
Shares outstanding	0.0172	0.054	0.321	0.748	-0.088	0.122
EPS	-0.2519	0.268	-0.940	0.347	-0.777	0.274
PE on BSE	-0.0126	0.049	-0.257	0.797	-0.108	0.083

Figure 33: Logistic Regression (Tuned) Output Summary

▪ Determining the **Optimal Threshold using ROC Curve:**

- The ROC curve illustrates the trade-off between the True Positive Rate (TPR) (Recall) and the False Positive Rate (FPR) for a binary classifier across various threshold values. The diagonal blue line represents a random classifier, with an area under the curve (AUC) of 0.5. A perfect classifier would achieve an AUC of 1.0.



**Optimal Threshold: 0.525**

Figure 34: ROC Curve & Optimal Threshold Value

## Check Model Performance (Tuned Model)

▪ Checking Model Performance post running the model with Optimal Threshold Value: -

- Below is the Confusion Matrix & Model Performance Metrics on Training dataset: -

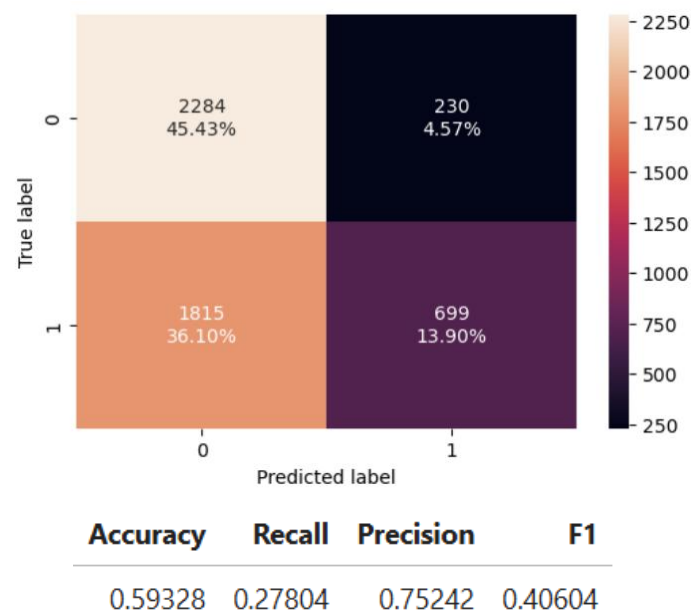


Figure 35: Logistic Regression (Tuned): Confusion Matrix & Metric Performance on Training Dataset

- Below is the Confusion Matrix & Model Performance Metrics on Test dataset: -

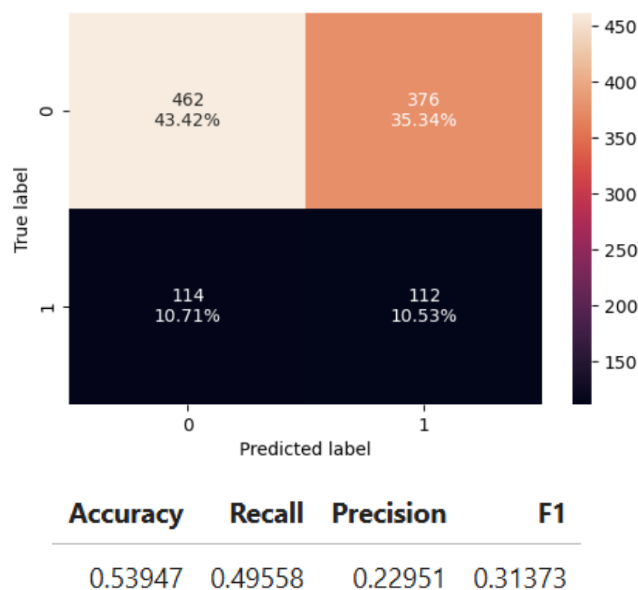


Figure 36: Logistic Regression (Tuned): Confusion Matrix & Metric Performance on Test Dataset

▪ **Summary of Performance: -**

- **Model Generalization:** The gap between Recall metrics for both the training and testing datasets have reduced, indicating that the model is tad more reliable.
- **Strengths:**
  - ✓ **Much improved Recall (~50%):** Indicates that the model is tad better at catching defaulters. The model certainly performs better than random guessing.
  - ✓ **Model generalizes better:** Test Recall is higher than training recall & the gap has also reduced, suggesting threshold tuning helped.
- **Weaknesses:**
  - ✓ **Lower precision (~23%):** The model generates more false alarms (many non-defaulters flagged)
  - ✓ **Low F1 Score (~0.31):** Reflects poor balance between precision and recall.
- **Tuned Logistic Regression improves recall significantly (to ~50%)** and offers a **better balance between catching defaulters and false alarms**, making it a more practical baseline model than the earlier versions.

## Random Forest – Tuning

### Build Tuned Model

- In order to tune the performance of Random Forest, we shall be passing **hyperparameters space** in the **GridSearchCV** function.
- The following **hyperparameter space** is used: -
  - 'n\_estimators': [50, 100] | Number of trees in the random forest | More trees generally improve performance and stability, but increase training time
  - 'max\_depth': [2, 3] | Maximum depth each tree can grow to | Limits model complexity to prevent overfitting
  - 'min\_samples\_split': [500, 1000] | Minimum number of samples required to split an internal node | Prevents splits when a node has fewer than the specified samples. Larger values mean fewer branches, preventing overfitting (Very high values are effective for large datasets and when overfitting is a concern)
  - 'min\_samples\_leaf': [500, 1000] | Minimum number of samples that must be at a leaf node | Larger values mean smoother model with fewer, more meaningful decisions. Works well to reduce noise and overfitting, especially for imbalanced classes
  - 'max\_features': ['sqrt'] | Number of features to consider when looking for the best split | Encourages diversity among trees (reduces correlation), improving ensemble performance
  - 'class\_weight': ['balanced'] | Adjusts weights inversely proportional to class frequencies | Assigns higher weight to minority class (defaulters) so the model doesn't ignore them. Crucial in imbalanced datasets to improve recall of minority class
- **GridSearchCV** stands for Grid Search with Cross-Validation, and it is a technique used in machine learning to tune hyperparameters by systematically searching through a grid of predefined combinations to find the best-performing model. Its purpose is to find the **optimal combination of hyperparameters** that give the **best performance** on the chosen evaluation metric **without overfitting**. It is executed in following steps: -
  - Define a set (grid) of possible values for each hyperparameter.
  - For each combination, the model is Trained on a subset of the data & then Validated on a different subset (via k-fold cross-validation).
  - The combination that yields the best average performance across folds is selected.
- Upon running the GridSearchCV, Tuned Random Forest Model is built again using the following best parameters (output from GridSearchCV): -

```
Best parameters: {'class_weight': 'balanced', 'max_depth': 2, 'max_features': 'sqrt', 'min_samples_leaf': 1000, 'min_samples_split': 500, 'n_estimators': 100}
```

```
Parameters used in the Random Forest Classifier:
bootstrap: True
ccp_alpha: 0.0
class_weight: balanced
criterion: gini
max_depth: 2
max_features: sqrt
max_leaf_nodes: None
max_samples: None
min_impurity_decrease: 0.0
min_samples_leaf: 1000
min_samples_split: 500
min_weight_fraction_leaf: 0.0
monotonic_cst: None
n_estimators: 100
n_jobs: None
oob_score: False
random_state: 42
verbose: 0
warm_start: False
```

Figure 37: Best Parameters for Tuning Random Forest Classifier

## Check Model Performance (Tuned Model)

- Checking Model Performance post tuning Random Forest: -

✓ Below is the Confusion Matrix & Model Performance Metrics on Training dataset: -

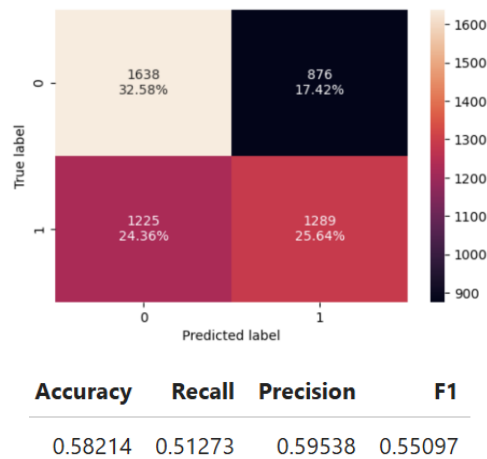


Figure 38: Random Forest (Tuned): Confusion Matrix & Metric Performance on Training Dataset

✓ Below is the Confusion Matrix & Model Performance Metrics on Test dataset: -

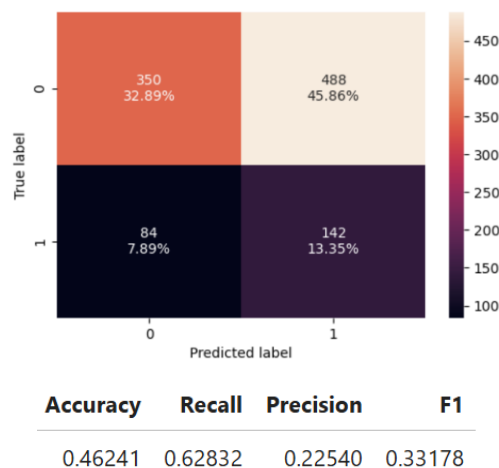


Figure 39: Random Forest (Tuned): Confusion Matrix & Metric Performance on Test Dataset

- Summary of Performance: -

- Model Generalization:** The gap between Recall metrics for both the training and testing datasets have significantly reduced, indicating that the model is pretty reliable & generalized well.
- Strengths:**
  - ✓ **High Recall (~63%):** The model is now strong at catching defaulters, which is crucial for risk mitigation.
  - ✓ **Model generalizes better:** Performance remains stable between train and test, especially for Recall.
- Weaknesses:**
  - ✓ **Low Precision (~23%):** Many false positives. In real-world use, this might strain investigative resources or raise false alarms. Expected as the model performance has been done in favour of Recall which has a trade-off with Precision.
  - ✓ **Low Accuracy (~46%):** Reflects the cost of improving Recall in an imbalanced setting; though, not a priority metric here but still notable.
  - ✓ **Low F1 Score (~0.33):** Reflects poor balance between precision and recall. Since, the focus is on Recall, Precision scores have been compromised.
- Tuned Random Forest Model significantly improves Recall (to ~63%), making it well-suited for default prediction.** While it generates more false positives, it **effectively captures high-risk companies**, fulfilling the tool's risk mitigation objective.



## Rubric Question 5: Model Performance Comparison and Final Model Selection

### Training Dataset Performance Comparison

- Below is the summary of performance comparison of Training dataset

Training performance comparison:

	Logistic Regression	Tuned Logistic Regression	Random Forest	Tuned Random Forest
<b>Accuracy</b>	0.59706	0.59328	0.92224	0.58214
<b>Recall</b>	0.37589	0.27804	0.93556	0.51273
<b>Precision</b>	0.67404	0.75242	0.91127	0.59538
<b>F1</b>	0.48264	0.40604	0.92326	0.55097

Figure 40: Training Dataset Performance Comparison

### Test Dataset Performance Comparison

- Below is the summary of performance comparison of Test dataset

Testing performance comparison:

	Logistic Regression	Tuned Logistic Regression	Random Forest	Tuned Random Forest
<b>Accuracy</b>	0.63158	0.53947	0.79323	0.46241
<b>Recall</b>	0.31416	0.49558	0.07080	0.62832
<b>Precision</b>	0.23052	0.22951	0.61538	0.22540
<b>F1</b>	0.26592	0.31373	0.12698	0.33178

Figure 41: Test Dataset Performance Comparison

### Final Model Selection

- Below is the comparison of all models with their Strengths & Weakness: -

Model	Strengths	Weaknesses
Logistic Regression	<ul style="list-style-type: none"> <li>✓ Decent precision (23%)</li> <li>✓ Simple and interpretable</li> </ul>	<ul style="list-style-type: none"> <li>✓ Low recall (31%)</li> <li>✓ Low F1 (0.27)</li> <li>✓ Struggles with class imbalance</li> </ul>
Tuned Logistic Regression	<ul style="list-style-type: none"> <li>✓ Improved recall (49%) and F1 (0.31)</li> <li>✓ Balanced generalization</li> </ul>	<ul style="list-style-type: none"> <li>✓ Precision drops (22%)</li> <li>✓ Accuracy lowest amongst all (54%)</li> </ul>
Random Forest	<ul style="list-style-type: none"> <li>✓ Excellent training performance (93% recall &amp; F1)</li> <li>✓ High precision (61%)</li> </ul>	<ul style="list-style-type: none"> <li>✓ Overfits: recall plummets to 7%</li> <li>✓ Worst F1 (0.13)</li> </ul>
Tuned Random Forest	<ul style="list-style-type: none"> <li>✓ Best recall (63%)</li> <li>✓ Best F1 (0.33)</li> <li>✓ Captures defaulters well</li> <li>✓ High model generalization for recall</li> </ul>	<ul style="list-style-type: none"> <li>✓ Low precision (22.5%)</li> <li>✓ Lowest accuracy (46%)</li> <li>✓ Many false positives</li> </ul>

Table 1: Model Comparison

### Best Model: Tuned Random Forest

- Tuned Random Forest is the **best model for credit default prediction** because it **maximizes the ability to detect true defaulters** (crucial for financial risk management), even if it sacrifices some precision. It **aligns perfectly with the business need to be proactive** rather than reactive in handling potential defaults.
- Key Strengths Aligned with Business Goals:** -
  - ✓ **High Defaulter Detection Rate (Recall):** Captures over 60% of risky companies & helps avoid missed defaulters, the most dangerous type of error in credit risk.
  - ✓ **Risk Mitigation Focus:** Prefers catching risky cases even if some healthy firms are flagged, which aligns with conservative lending and investment strategy.
  - ✓ **Balanced Performance:** While accuracy is low (46%), that's acceptable in imbalanced classification. Best trade-off among all models in terms of real-world impact.
  - ✓ **High Model Generalization:** Recall gap between Train & Test dataset performance is minimal, meaning it performs consistently well on both the training and unseen test data. This means there is no significant overfitting or underfitting & The model can adapt to new, unseen data with confidence.
- Trade-Offs (Acceptable Given the Context):** -
  - ✓ **Low Precision (22.5%):** More false positives (non-defaulters flagged as risky), but this is acceptable in early-stage screening & can be mitigated by human review or second-tier models.

■ **PLEASE NOTE:** *The Random Forest Model could have been further tuned for better performance metrics, but, was restrained due to limited computing capacity. There was surely more room for hyperparameter-space tuning, but was deliberately restricted since the running time (& its respective troubleshooting!) was going as high as 1-2 hours per code execution.*

## Most Important Features

- Below are the features contributing to the determination of whether a company would default: -

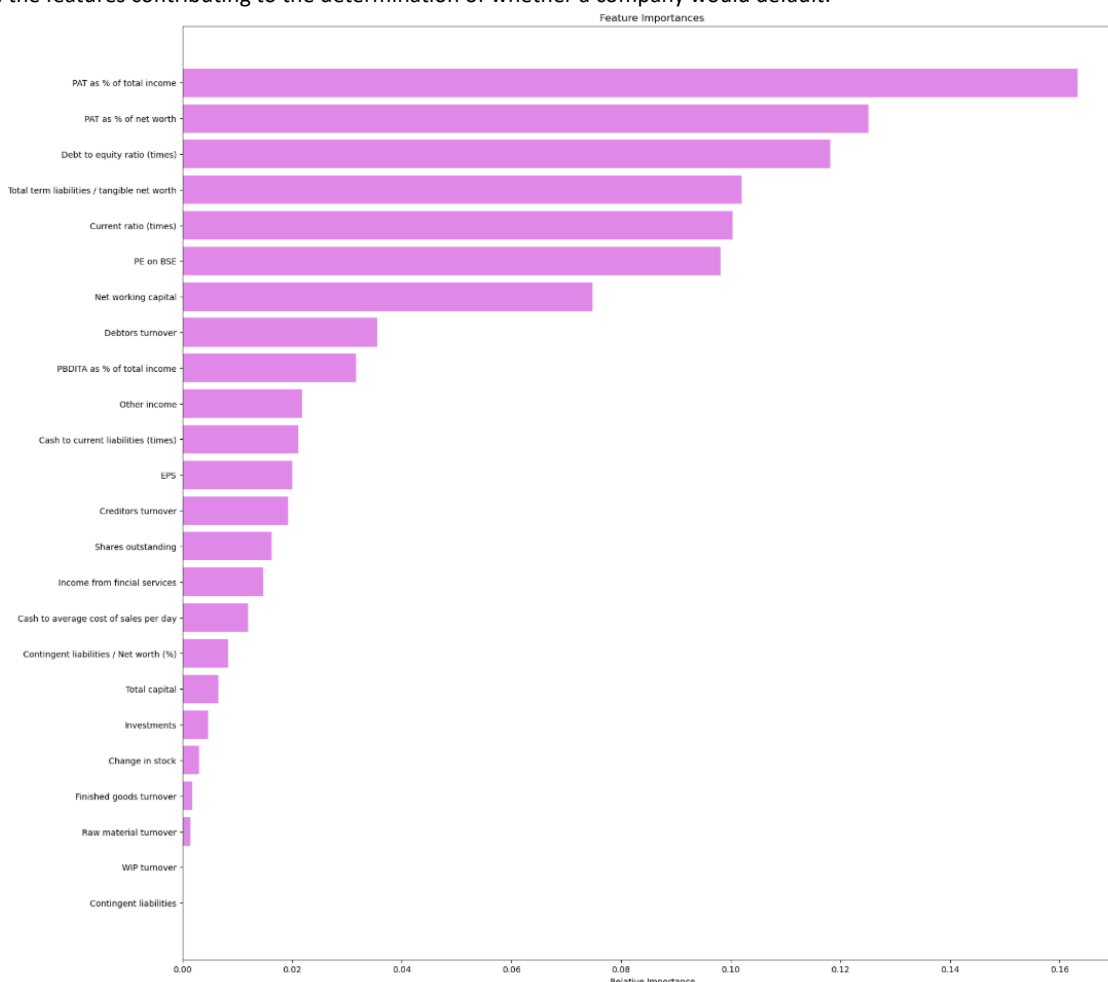


Figure 42: Feature Importance

- **Feature Importance (Top 7) inferences** can be summarized in the below table: -

Feature	Business Inference
<b>PAT % of Total Income</b>	Indicates overall profitability. Low PAT margins may signal poor earnings efficiency → a strong indicator of financial stress.
<b>PAT % of Net Worth</b>	Measures return on equity. Low values may mean the company is not effectively utilizing its capital → potential default risk.
<b>Debt to Equity Ratio (times)</b>	High values imply over-leveraging. Companies heavily funded by debt are more vulnerable to repayment failure.
<b>Total Term Liabilities / Tangible Net Worth</b>	A measure of long-term solvency. High values suggest a company may be structurally risky and reliant on long-term borrowing.
<b>Current Ratio (times)</b>	Indicates short-term liquidity. A low current ratio shows inability to meet short-term obligations → common cause of default.
<b>PE on BSE</b>	Reflects market perception of earnings potential. Very low or negative P/E may imply investor distrust or declining earnings outlook.
<b>Net Working Capital</b>	Shows operational liquidity. Negative or weak working capital limits day-to-day operations and signals potential cash flow crises.

*Table 2: Feature Importance Inference (Top 7)*

## Rubric Question 6: Actionable Insights & Recommendations

### Actionable Insights

1. **Closely Monitor PAT Margins**
  - ✓ Low **PAT % of total income** signals inefficient earnings and weak profitability.
  - ✓ **Recommended Action:** Flag companies with declining PAT margins for financial stress audits.
2. **Evaluate Return on Equity Regularly**
  - ✓ Low **PAT % of net worth** implies poor utilization of shareholder funds.
  - ✓ **Recommended Action:** Require turnaround plans or restrict additional funding to underperforming firms.
3. **Limit Exposure to Over-Leveraged Firms**
  - ✓ High **Debt to Equity Ratio** is a strong predictor of default.
  - ✓ **Recommended Action:** Set stricter borrowing caps or require collateral for high D/E borrowers.
4. **Assess Long-Term Solvency Health**
  - ✓ High **Total Term Liabilities / Tangible Net Worth** reflects poor capital structure.
  - ✓ **Recommended Action:** Use this ratio as a filter in lending or investment eligibility criteria.
5. **Use Current Ratio as a Short-Term Risk Flag**
  - ✓ A low **Current Ratio** suggests difficulty in meeting short-term obligations.
  - ✓ **Recommended Action:** Apply liquidity thresholds for credit renewal or loan disbursement.
6. **Incorporate Market Signals via PE on BSE**
  - ✓ Unusual or negative **PE ratios** may indicate declining investor confidence or unsustainable earnings.
  - ✓ **Recommended Action:** Use as a soft trigger for early warning and enhanced financial scrutiny.
7. **Focus on Working Capital Efficiency**
  - ✓ Weak **Net Working Capital** indicates potential cash flow mismatches.
  - ✓ **Recommended Action:** Offer credit with conditions tied to working capital improvement plans.
8. **Track Debtors Turnover to Gauge Credit Discipline**
  - ✓ Low **Debtors Turnover** suggests poor collection cycles and strained cash flow.
  - ✓ **Recommended Action:** Flag companies with declining turnover for collection risk monitoring.
9. **Review PBDITA Margins for Operating Strength**
  - ✓ Low **PBDITA as % of total income** means weak operational profitability before depreciation and taxes.
  - ✓ **Recommended Action:** Prioritize operational restructuring support or advisory for such firms.
10. **Monitor "Other Income" Reliance**
  - ✓ Heavy reliance on **Other Income** may indicate non-core or unstable revenue sources.
  - ✓ **Recommended Action:** Evaluate sustainability of income and push for diversification if it's non-operational.

### Business Recommendations

1. **Implement Financial Health Scoring**
  - ✓ Design a **risk scoring framework** that gives higher weight to PAT % of total income, Debt to equity ratio, Current ratio.
  - ✓ Use this score as a **pre-screening tool** for lending or investing.
2. **Set Profitability Thresholds for Approvals**
  - ✓ Require a **minimum PAT margin or return on equity** for loan approvals or investment consideration.
  - ✓ Flag companies with consistently low profitability for **review or intervention**.
3. **Enforce Leverage Limits in Credit Policies**
  - ✓ Define **industry-specific maximum debt-to-equity ratios**.
  - ✓ Automatically trigger **credit limit reductions** or **additional collateral requests** for companies that breach this limit.
4. **Prioritize Liquidity in Credit Evaluation**
  - ✓ Mandate a **minimum current ratio** (e.g.,  $\geq 1.2$ ) for short-term loans or working capital facilities.
  - ✓ Regularly monitor liquidity ratios and **issue alerts** for deterioration.
5. **Incorporate Market Sentiment Metrics**
  - ✓ Use metrics like **PE ratio on BSE** as a proxy for investor confidence.
  - ✓ Apply **soft rejections or delayed disbursements** when market sentiment around a company deteriorates sharply.
6. **Monitor Working Capital Health Proactively**
  - ✓ Track **Net Working Capital trends** monthly or quarterly.
  - ✓ Offer **working capital optimization advisory** to at-risk firms as a preventive measure.

7. **Evaluate Collection Efficiency Metrics**
  - ✓ Use **Debtors Turnover** as a proxy for receivables risk.
  - ✓ Impose **stricter payment terms or collection targets** before approving new lines of credit.
8. **Assess Earnings Quality, Not Just Size**
  - ✓ Be cautious of companies with high **other income** or non-operating earnings.
  - ✓ Require financial disclosures to **distinguish recurring from one-off income sources**.
9. **Build Early Warning Systems**
  - ✓ Set thresholds on key features (e.g., debt to equity > 2.5, current ratio < 1.0, PAT margin < 5%) that **automatically flag risky companies** for manual intervention.
10. **Conduct Periodic Re-Risk Reviews**
  - ✓ Re-evaluate financials of existing borrowers/investees every 6–12 months using the top risk indicators.
  - ✓ This enables **proactive mitigation** before default risk escalates.