

```

1 import numpy as np
2 import pandas as pd
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 from scipy import stats
6 from sklearn.svm import SVC
7 from sklearn.ensemble import RandomForestClassifier
8 from sklearn.model_selection import train_test_split
9 from imblearn.over_sampling import SMOTE
10 from sklearn.linear_model import LogisticRegression
11 from sklearn.metrics import roc_curve, roc_auc_score
12 from sklearn.tree import DecisionTreeClassifier
13 from sklearn.preprocessing import StandardScaler
14 from xgboost import XGBClassifier
15 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, confusion_matrix, roc
16 from imblearn.under_sampling import NearMiss
17 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
18

```

```

1 from google.colab import drive
2 drive.mount('/content/drive')

```

Mounted at /content/drive

```

1 #read data from dataset
2 bank_data = pd.read_csv("/content/drive/MyDrive/FYP2025/res_data.csv")
3
4 bank_data.head()

```

Mounted at /content/drive

	ORGNR	Year	ZipCode	Municipality	Bankrupt	Business_Category	Size_Category	OperatingProfit	F
0	5560024712	1998	23163	1287.0	0	30	2	-558.0	
1	5560024712	1999	23163	1287.0	0	30	2	-634.0	
2	5560024712	2000	23163	1287.0	0	30	2	3809.0	
3	5560024712	2001	23163	1287.0	0	30	2	-1147.0	
4	5560024712	2002	23163	1287.0	0	30	2	-1250.0	

5 rows × 102 columns

```

1 #0 bankcrpt no
2 #1 bancrpt yes
3 print(bank_data['Bankrupt'].value_counts())
4

```

Mounted at /content/drive

```

Bankrupt
0    90966
1     335
Name: count, dtype: int64

```

```

1 #Visual representation of dataset
2 bank_data.info()
3

```

Mounted at /content/drive

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 91301 entries, 0 to 91300
Columns: 102 entries, ORGNR to Profit_Margin
dtypes: float64(68), int64(34)
memory usage: 71.1 MB

```

```

1 #calculate average inventory
2 average_inventory_series =bank_data.groupby('ORGNR')['Inventory'].transform(lambda x: (x + x.shift(-1)) /

```

```

1
2
3 #add average inventory as column in dataframe
4 bank_data['Average_Inventory'] = average_inventory_series
5

```

```

1 # Replace NaN values in 'Average_Inventory' with the previous values
2 bank_data['Average_Inventory'] = bank_data['Average_Inventory'].ffill()
3

```

```

1
2
3
4 #calculating total liabilities
5 bank_data['Total_Liabilities'] = bank_data['Total_Curr_Liabilities'] + bank_data['Total_Non_curr_lia']
6

```

```

1
2 bank_data

```



	ORGNR	Year	ZipCode	Municipality	Bankrupt	Business_Category	Size_Category	OperatingProf
0	5560024712	1998	23163	1287.0	0	30	2	-558
1	5560024712	1999	23163	1287.0	0	30	2	-634
2	5560024712	2000	23163	1287.0	0	30	2	3809
3	5560024712	2001	23163	1287.0	0	30	2	-1147
4	5560024712	2002	23163	1287.0	0	30	2	-1250
...
91296	5567417315	2009	43430	1384.0	0	30	2	134
91297	5567417315	2010	43430	1384.0	0	30	2	70
91298	5567417315	2011	43430	1384.0	0	30	2	125
91299	5567417315	2012	43430	1384.0	0	30	2	76
91300	5567417315	2013	40013	1480.0	0	30	2	26

91301 rows × 104 columns

```

1 #eleminating assests ROE
2 # Check if columns exist before dropping
3 if 'Equity_to_Assets' in bank_data.columns and 'ROE' in bank_data.columns:
4     bank_data = bank_data.drop(['Equity_to_Assets', 'ROE'], axis=1)
5 else:
6     # Print a message indicating which columns are not found
7     missing_columns = [col for col in ['Equity_to_Assets', 'ROE'] if col not in bank_data.columns]
8     print(f"Columns not found in DataFrame: {missing_columns}")
9 bank_data.head()

```



	ORGNR	Year	ZipCode	Municipality	Bankrupt	Business_Category	Size_Category	OperatingProfit	F
0	5560024712	1998	23163	1287.0	0	30	2	-558.0	
1	5560024712	1999	23163	1287.0	0	30	2	-634.0	
2	5560024712	2000	23163	1287.0	0	30	2	3809.0	
3	5560024712	2001	23163	1287.0	0	30	2	-1147.0	
4	5560024712	2002	23163	1287.0	0	30	2	-1250.0	

5 rows × 102 columns

```
1
2 #opreating cost
3 bank_data['Operating_CF'] = (bank_data['NetProfit'] + bank_data['Depreciation']) - (bank_data['Total_CurrA'] - bank_data['Total_CurrLiabilities'])
4
```

```
1 #calculating liablites to asset
2 bank_data['Liabilities_to_Assets'] = bank_data['Total_Liabilities'] / bank_data['Total_Assets']
3
```

```
1
2
3 #calculating return of assets
4 bank_data['Return_on_Assets'] = bank_data['NetProfit'] / bank_data['Total_Assets']
5
```

```
1 #deb to equity
2 bank_data['Debt_to_Equity'] = bank_data['Total_Liabilities'] / bank_data['Total_Equity']
3
```

```
1 #profit margin and gross margin
2 bank_data['Profit_Margin'] = bank_data['NetProfit'] / bank_data['Net_Sales']
3 bank_data['Gross_Margin'] = (bank_data['Net_Sales'] - bank_data['Production_Costs']) / bank_data['Net_Sales']
4
```

```
1 #networth to assets
2 bank_data['NetWC_to_Assets'] = (bank_data['Total_CurrA'] - bank_data['Total_CurrLiabilities']) / bank_data['Total_Assets']
3
```

```
1 #ebit
2 bank_data['EBIT_to_Assets'] = bank_data['OperatingProfit'] / bank_data['Total_Assets']
3 bank_data['EBIT_to_Sales'] = bank_data['OperatingProfit'] / bank_data['Net_Sales']
4
```

```
1 #finanical ratio
2 bank_data['Current_Ratio'] = bank_data['Total_CurrA'] / bank_data['Total_CurrLiabilities']
3 bank_data['Cash_Ratio'] = (bank_data['CashnBank'] + bank_data['Securities_Investments']) / bank_data['Total_receivables']
4 bank_data['Working_Capital_Ratio'] = bank_data['Total_CurrA'] / bank_data['Total_CurrLiabilities']
5
```

```
1 bank_data['Equity_to_Liabilities'] = bank_data['Total_Equity'] / bank_data['Total_Liabilities']
2 bank_data['Net_Sales_to_Assets'] = bank_data['Net_Sales'] / bank_data['Total_Assets']
3 bank_data['Net_Worth_to_Debt'] = (bank_data['Total_Assets'] - bank_data['Intangible_FxdA']) / bank_data['Total_Liabilities']
4
```

```
1 bank_data["CF_to_debt"] = bank_data["Operating_CF"] / bank_data["Total_Liabilities"]
2 bank_data["CF_to_sales"] = bank_data["Operating_CF"] / bank_data["Net_Sales"]
```

```

3 bank_data["CF_coverage"] = bank_data["Operating_CF"] / (bank_data["Other_oper_exp"] + bank_data["Financia
4 bank_data["CF_to_equity"] = bank_data["Operating_CF"] / bank_data["Total_Equity"]
5 bank_data["CF_to_Assets"] = bank_data["Operating_CF"] / bank_data["Total_Assets"]
6

```

```

1 bank_data["RE_to_Sales"] = bank_data["Accumulated_profit"] / bank_data["Net_Sales"]
2 bank_data["RE_to_Assets"] = bank_data["Accumulated_profit"] / bank_data["Total_Assets"]
3 bank_data["RE_to_NetIncome"] = bank_data["Accumulated_profit"] / bank_data["NetProfit"]
4

```

```

1
2 bank_data["Inventory_Turnover"] = bank_data["Production_Costs"] / bank_data["Average_Inventory"]
3 bank_data["Inventory_to_Sales"] = bank_data["Inventory"] / bank_data["Net_Sales"]
4

```

```

1 bank_data.head(7)
2

```

	ORGNR	Year	ZipCode	Municipality	Bankrupt	Business_Category	Size_Category	OperatingProfit	F
0	5560024712	1998	23163	1287.0	0	30	2	-558.0	
1	5560024712	1999	23163	1287.0	0	30	2	-634.0	
2	5560024712	2000	23163	1287.0	0	30	2	3809.0	
3	5560024712	2001	23163	1287.0	0	30	2	-1147.0	
4	5560024712	2002	23163	1287.0	0	30	2	-1250.0	
5	5560024712	2003	23163	1287.0	0	30	2	-1492.0	
6	5560024712	2004	23163	1287.0	0	30	2	-1024.0	

7 rows × 125 columns

```

1 predictor_columns = [
2     'Asset_Turnover_ratio', 'ROA', 'Interest_on_Debt', 'Debt_to_Equity',
3     'WC_to_Sales', 'Quick_ratio', 'Operating_Margin',
4     'Net_Margin', 'Profit_Margin', 'Operating_CF',
5     'Liabilities_to_Assets', 'Return_on_Assets', 'NetWC_to_Assets',
6     'EBIT_to_Assets', 'EBIT_to_Sales', 'Current_Ratio', 'Cash_Ratio',
7     'Working_Capital_Ratio', 'Equity_to_Liabilities', 'Net_Sales_to_Assets',
8     'Net_Worth_to_Debt', 'CF_to_debt', 'CF_to_sales', 'CF_coverage',
9     'CF_to_equity', 'CF_to_Assets', 'Gross_Margin', 'RE_to_Sales',
10    'RE_to_Assets', 'RE_to_NetIncome', 'Inventory_Turnover', 'Inventory_to_Sales'
11 ]
12

```

```

1
2
3
4
5
6
7 # Create a DataFrame indicating missing and inf values
8 missing_inf_df = pd.DataFrame(index=predictor_columns)
9 # Fill the DataFrame with zeros (no missing or inf values)
10 missing_inf_df['Missing'] = 0
11 missing_inf_df['Inf'] = 0
12 # Find missing (NaN) values and inf values in the DataFrame
13 for column in predictor_columns:
14     missing_inf_df.at[column, 'Missing'] = bank_data[column].isnull().sum()
15     missing_inf_df.at[column, 'Inf'] = np.isinf(bank_data[column]).sum()
16

```

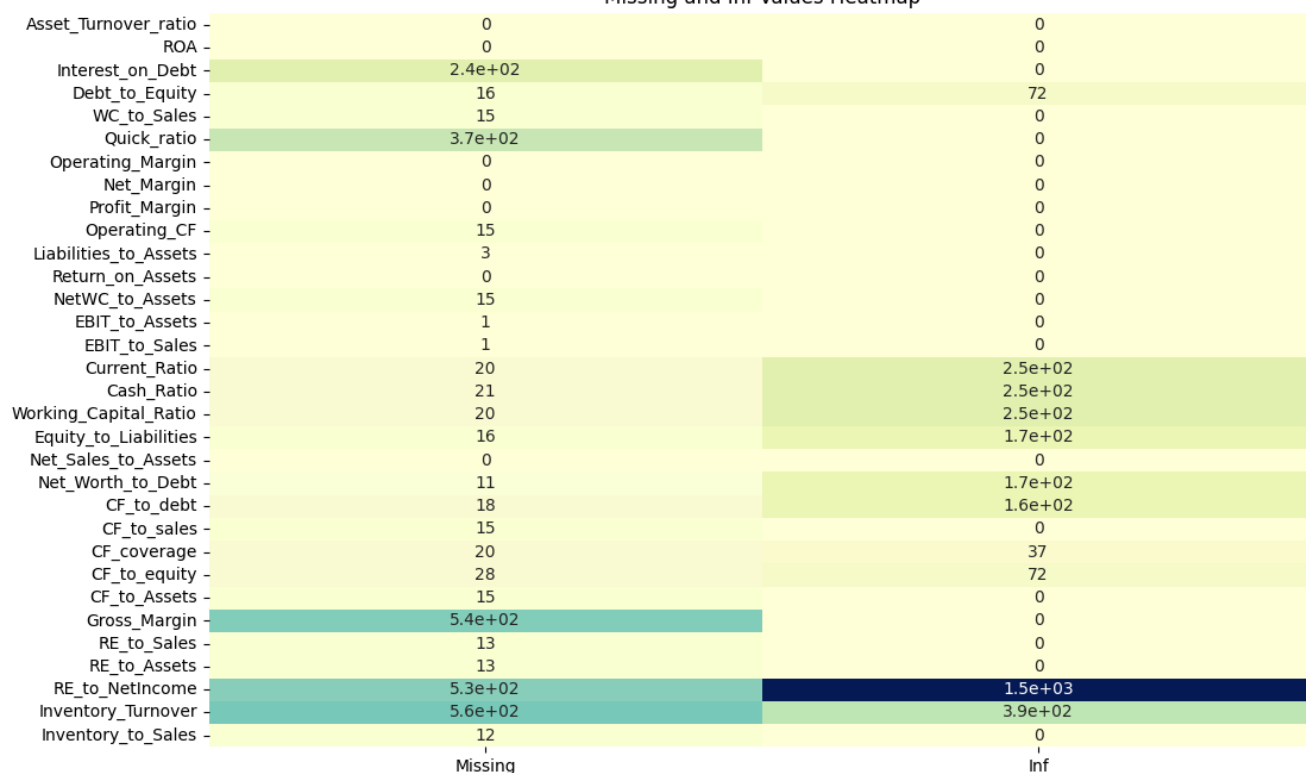
```

1 # Create a heatmap to visualize missing and inf values
2 plt.figure(figsize=(12, 8))
3 sns.heatmap(missing_inf_df, annot=True, cmap='YlGnBu', cbar=False)
4 plt.title('Missing and Inf Values Heatmap')
5 plt.show()
6

```



Missing and Inf Values Heatmap



```

1 # Count missing and inf values in each predictor column
2 missing_counts = {}
3 inf_counts = {}
4 for column in predictor_columns:
5     missing_counts[column] = bank_data[column].isnull().sum()
6     inf_counts[column] = np.isinf(bank_data[column]).sum()
7
8 missing_counts_df = pd.DataFrame(missing_counts, index=['Missing Count'])
9 inf_counts_df = pd.DataFrame(inf_counts, index=['Inf Count'])
10
11 # Create a mask to identify rows with at least one missing or inf value in predictor columns
12 rows_with_missing_or_inf = bank_data[predictor_columns].isnull().any(axis=1) | bank_data[predictor_columns].isinf().any(axis=1)
13
14 missing_or_inf_df = bank_data[rows_with_missing_or_inf]
15
16 # Print the missing and inf value counts and the DataFrame with affected observations
17 print("Missing Value Counts:")
18 print(missing_counts_df)
19 print("\nInf Value Counts:")
20 print(inf_counts_df)
21 print("\nDataFrame with Observations Containing Missing or Inf Values:")
22 missing_or_inf_df.head(10)
23

```



Missing Value Counts:

	Asset_Turnover_ratio	ROA	Interest_on_Debt	Debt_to_Equity	\	
Missing Count	0	0	238	16		
	WC_to_Sales	Quick_ratio	Operating_Margin	Net_Margin	\	
Missing Count	15	368	0	0		
	Profit_Margin	Operating_CF	...	CF_to_sales	CF_coverage	\
Missing Count	0	15	...	15	20	
	CF_to_equity	CF_to_Assets	Gross_Margin	RE_to_Sales	\	
Missing Count	28	15	543	13		
	RE_to_Assets	RE_to_NetIncome	Inventory_Turnover	\		
Missing Count	13	531	562			
	Inventory_to_Sales					
Missing Count	12					

[1 rows x 32 columns]

Inf Value Counts:

	Asset_Turnover_ratio	ROA	Interest_on_Debt	Debt_to_Equity	\	
Inf Count	0	0	0	72		
	WC_to_Sales	Quick_ratio	Operating_Margin	Net_Margin	\	
Inf Count	0	0	0	0		
	Profit_Margin	Operating_CF	...	CF_to_sales	CF_coverage	\
Inf Count	0	0	...	0	37	
	CF_to_equity	CF_to_Assets	Gross_Margin	RE_to_Sales	\	
Inf Count	72	0	0	0		
	RE_to_Assets	RE_to_NetIncome	Inventory_Turnover	\		
Inf Count	0	1480	387			
	Inventory_to_Sales					
Inf Count	0					

[1 rows x 32 columns]

DataFrame with Observations Containing Missing or Inf Values:

	ORGNR	Year	ZipCode	Municipality	Bankrupt	Business_Category	Size_Category	OperatingProfit
7	5560064403	1998	24131	1285.0	0	30	0	40.0
13	5560064403	2004	24131	1285.0	0	30	0	32.0
16	5560064403	2007	24131	1285.0	0	30	0	37.0
130	5560234238	1998	18205	162.0	0	30	3	-519.0
131	5560234238	1999	18205	162.0	0	30	3	-499.0
153	5560284761	2001	11152	180.0	0	30	0	0.0
157	5560284761	2007	11122	180.0	0	30	0	0.0
169	5560312000	2006	18205	162.0	0	30	0	10.0
173	5560312000	2010	18205	162.0	0	30	0	216.0
174	5560312000	2012	18205	162.0	0	30	0	58.0

10 rows x 125 columns

```

1 bank_data.replace([np.inf, -np.inf], 0, inplace=True)
2 bank_data = bank_data.dropna()
3 print(bank_data['Bankrupt'].value_counts())
4 print(bank_data.shape)
5

```

```

Bankrupt
0    73069
1     298
Name: count, dtype: int64
(73367, 125)

```

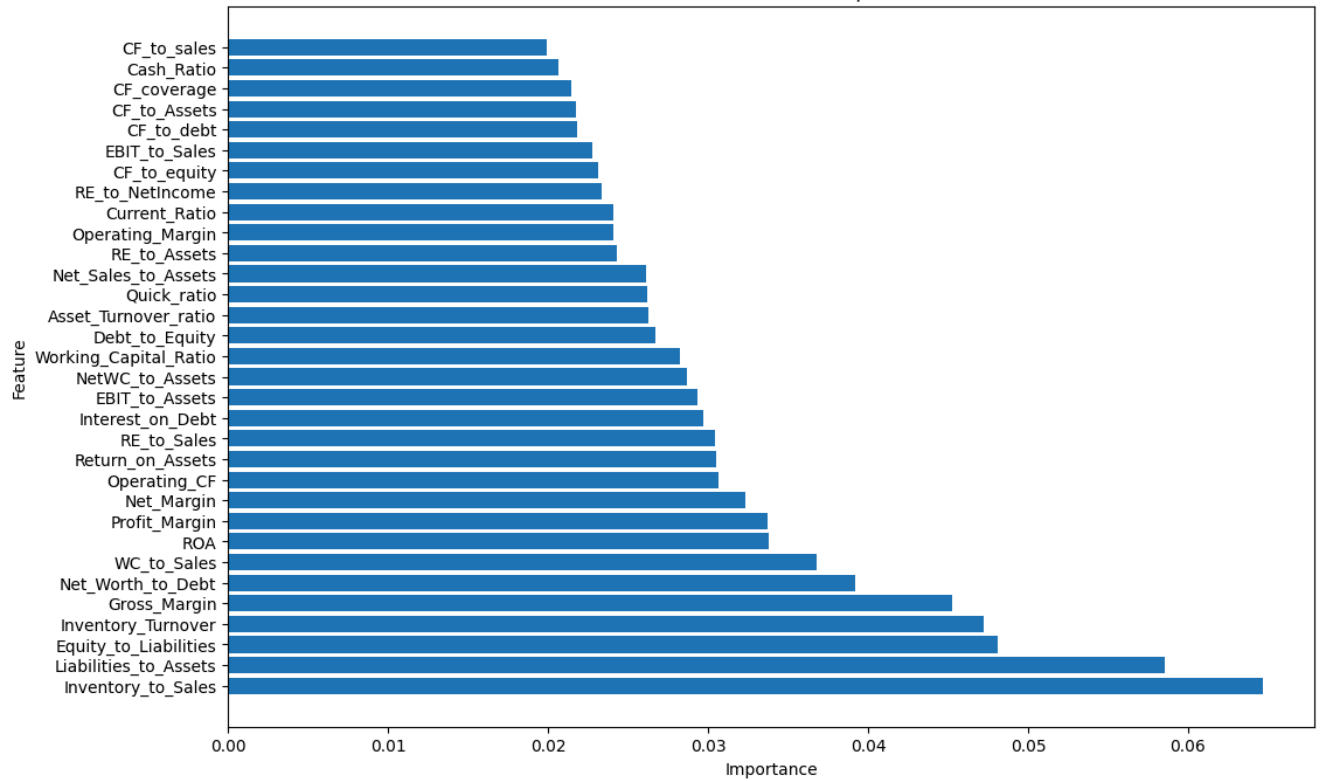
```

1 # Select predictors and target column
2 predictors = [
3     'Asset_Turnover_ratio', 'ROA', 'Interest_on_Debt', 'Debt_to_Equity',
4     'WC_to_Sales', 'Quick_ratio', 'Operating_Margin',
5     'Net_Margin', 'Profit_Margin', 'Operating_CF',
6     'Liabilities_to_Assets', 'Return_on_Assets', 'NetWC_to_Assets',
7     'EBIT_to_Assets', 'EBIT_to_Sales', 'Current_Ratio', 'Cash_Ratio',
8     'Working_Capital_Ratio', 'Equity_to_Liabilities', 'Net_Sales_to_Assets',
9     'Net_Worth_to_Debt', 'CF_to_debt', 'CF_to_sales', 'CF_coverage',
10    'CF_to_equity', 'CF_to_Assets', 'Gross_Margin', 'RE_to_Sales',
11    'RE_to_Assets', 'RE_to_NetIncome', 'Inventory_Turnover', 'Inventory_to_Sales'
12 ]
13
14 target = 'Bankrupt'
15
16 # Separate predictors and target
17 X = bank_data[predictors]
18 y = bank_data[target]
19
20 # Fit a Random Forest model
21 rf_model = RandomForestClassifier(random_state=42, class_weight='balanced')
22 rf_model.fit(X, y)
23
24 # Get feature importances
25 feature_importances = rf_model.feature_importances_
26
27 # Create a DataFrame to store feature names and their importance scores
28 feature_importance_df = pd.DataFrame({'Feature': predictors, 'Importance': feature_importances})
29
30 # Sort the DataFrame by importance in descending order
31 feature_importance_df = feature_importance_df.sort_values(by='Importance', ascending=False)
32
33 # Plot the feature importance
34 plt.figure(figsize=(12, 8))
35 plt.barh(feature_importance_df['Feature'], feature_importance_df['Importance'])
36 plt.xlabel('Importance')
37 plt.ylabel('Feature')
38 plt.title('Random Forest Feature Importance')
39 plt.show()
40
41 # Print the sorted feature importance
42 print("Feature Importance:")
43 print(feature_importance_df)
44

```



Random Forest Feature Importance



Feature Importance:

	Feature	Importance
31	Inventory_to_Sales	0.064670
10	Liabilities_to_Assets	0.058577
18	Equity_to_Liabilities	0.048071
30	Inventory_Turnover	0.047256
26	Gross_Margin	0.045235
20	Net_Worth_to_Debt	0.039217
4	WC_to_Sales	0.036796
1	ROA	0.033806
8	Profit_Margin	0.033724
7	Net_Margin	0.032320
9	Operating_CF	0.030627
11	Return_on_Assets	0.030475
27	RE_to_Sales	0.030462
2	Interest_on_Debt	0.029720
13	EBIT_to_Assets	0.029335
12	NetWC_to_Assets	0.028678
17	Working_Capital_Ratio	0.028277
3	Debt_to_Equity	0.026739
0	Asset_Turnover_ratio	0.026295
5	Quick_ratio	0.026226
19	Net_Sales_to_Assets	0.026163
28	RE_to_Assets	0.024271
6	Operating_Margin	0.024085
15	Current_Ratio	0.024064
29	RE_to_NetIncome	0.023322
24	CF_to_equity	0.023149
14	EBIT_to_Sales	0.022773
21	CF_to_debt	0.021807
25	CF_to_Assets	0.021771
23	CF_coverage	0.021480
16	Cash_Ratio	0.020682
22	CF_to_sales	0.019927

```
1 feature_scores = pd.Series(rf_model.feature_importances_, index=X.columns).sort_values(ascending=False)
2
3 feature_scores
4
```




0

Inventory_to_Sales	0.064670
Liabilities_to_Assets	0.058577
Equity_to_Liabilities	0.048071
Inventory_Turnover	0.047256
Gross_Margin	0.045235
Net_Worth_to_Debt	0.039217
WC_to_Sales	0.036796
ROA	0.033806
Profit_Margin	0.033724
Net_Margin	0.032320
Operating_CF	0.030627
Return_on_Assets	0.030475
RE_to_Sales	0.030462
Interest_on_Debt	0.029720
EBIT_to_Assets	0.029335
NetWC_to_Assets	0.028678
Working_Capital_Ratio	0.028277
Debt_to_Equity	0.026739
Asset_Turnover_ratio	0.026295
Quick_ratio	0.026226
Net_Sales_to_Assets	0.026163
RE_to_Assets	0.024271
Operating_Margin	0.024085
Current_Ratio	0.024064
RE_to_NetIncome	0.023322
CF_to_equity	0.023149
EBIT_to_Sales	0.022773
CF_to_debt	0.021807
CF_to_Assets	0.021771
CF_coverage	0.021480
Cash_Ratio	0.020682
CF_to_sales	0.019927

dtype: float64

```
1 # Split the data into train and test sets
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.4, random_state=42)
3
4 # Use SMOTE to oversample the minority class
5 smote = SMOTE(random_state=42)
6 X_train_resampled, y_train_resampled = smote.fit_resample(X_train, y_train)
7
8
```

```

9
10
11 print("X_train shape: ", X_train.shape)
12 print("X_test shape: ", X_test.shape)
13 print("y_train shape: ", y_train.shape)
14 print("y_test shape: ", y_test.shape)
15

```

```

⇒ X_train shape: (44020, 32)
   X_test shape: (29347, 32)
   y_train shape: (44020,)
   y_test shape: (29347,)

```

```

1 # Count the number of bankrupt and non-bankrupt companies
2 bankrupt_count_train = y_train.sum()
3 non_bankrupt_count_train = len(y_train) - bankrupt_count_train
4 bankrupt_count_test = y_test.sum()
5 non_bankrupt_count_test = len(y_test) - bankrupt_count_test
6
7 print("Train Set:")
8 print("Bankrupt Companies:", bankrupt_count_train)
9 print("Non-Bankrupt Companies:", non_bankrupt_count_train)
10 print("\nTest Set:")
11 print("Bankrupt Companies:", bankrupt_count_test)
12 print("Non-Bankrupt Companies:", non_bankrupt_count_test)
13

```

```

⇒ Train Set:
   Bankrupt Companies: 182
   Non-Bankrupt Companies: 43838

   Test Set:
   Bankrupt Companies: 116
   Non-Bankrupt Companies: 29231

```

```

1 # ✅ Import the classifier first
2 from sklearn.ensemble import RandomForestClassifier
3
4 # 1. Define the Random Forest model
5 rf_model = RandomForestClassifier(
6     n_estimators=50,
7     max_depth=10,
8     class_weight='balanced',
9     random_state=42
10 )#hyper parameter tuning
11
12 # 2. Fit on resampled data
13 rf_model.fit(X_train_resampled, y_train_resampled)
14
15 # 3. Predict on test data
16 y_pred = rf_model.predict(X_test)

```

```

1 import pickle
2
3 # Save the trained Random Forest model
4 with open("rf_model.pkl", "wb") as f:
5     pickle.dump(rf_model, f)
6

```

```

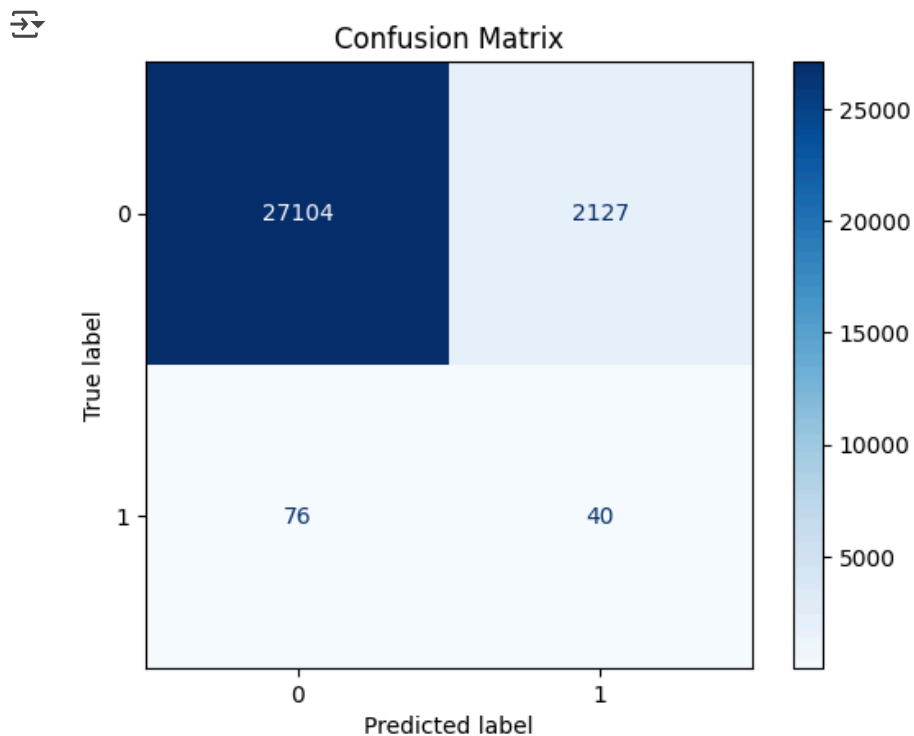
1 from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
2
3 cm = confusion_matrix(y_test, y_pred)
4
5 print('Confusion matrix\n\n', cm)
6

```

⇒ Confusion matrix

```
[[27104 2127]
 [   76   40]]
```

```
1
2 # Plot the confusion matrix
3 disp = ConfusionMatrixDisplay(confusion_matrix=cm)
4 disp.plot(cmap=plt.cm.Blues)
5 plt.title('Confusion Matrix')
6 plt.show()
7
```



```
1 from sklearn.metrics import classification_report
2
3 print(classification_report(y_test, y_pred))
4
```

⇒

	precision	recall	f1-score	support
0	1.00	0.93	0.96	29231
1	0.02	0.34	0.04	116
accuracy			0.92	29347
macro avg	0.51	0.64	0.50	29347
weighted avg	0.99	0.92	0.96	29347

```
1 from sklearn.linear_model import LogisticRegression
2 from imblearn.over_sampling import SMOTE
3 from sklearn.metrics import roc_curve, roc_auc_score
4 import joblib
5 import matplotlib.pyplot as plt
6
7 # Assuming you already have your X_train, y_train, X_test, y_test data loaded
8 # If using resampling techniques like SMOTE for imbalanced data
9 smote = SMOTE(random_state=42)
10 X_train_resampled, y_train_resampled = smote.fit_resample(X_train, y_train)
11
12 # Initialize Logistic Regression model
13 logreg_model = LogisticRegression(random_state=426)
```

```
14
15 # Fit the model to the resampled data
16 logreg_model.fit(X_train_resampled, y_train_resampled)
17
18 # Save the logistic regression model
19 joblib.dump(logreg_model, 'logistic_regression_model.pkl')
20
21 # Make predictions on the test set using the logistic regression model
22 y_pred = logreg_model.predict(X_test)
23
24 # Predict probabilities for the positive class
25 y_prob = logreg_model.predict_proba(X_test)[:, 1]
26
27 # Calculate the ROC curve
28 fpr, tpr, thresholds = roc_curve(y_test, y_prob)
29
30 # Calculate the AUC (Area Under the Curve)
31 auc = roc_auc_score(y_test, y_prob)
32
33 # Plot the ROC curve
34 plt.figure(figsize=(8, 6))
35 plt.plot(fpr, tpr, color='darkorange', lw=2, label='ROC curve (AUC = {:.2f})'.format(auc))
36 plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
37 plt.xlabel('False Positive Rate (FPR)')
38 plt.ylabel('True Positive Rate (TPR)')
39 plt.title('Logistic Regression ROC Curve')
40 plt.legend(loc='lower right')
41 plt.show()
42
43 # Print the AUC score
44 print("AUC Score:", auc)
45
46 # Now your logistic regression model is ready for predictions.
47
```

➔ /usr/local/lib/python3.11/dist-packages/sklearn/linear_model/_logistic.py:465: ConvergenceWarning: lbfgs ·
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

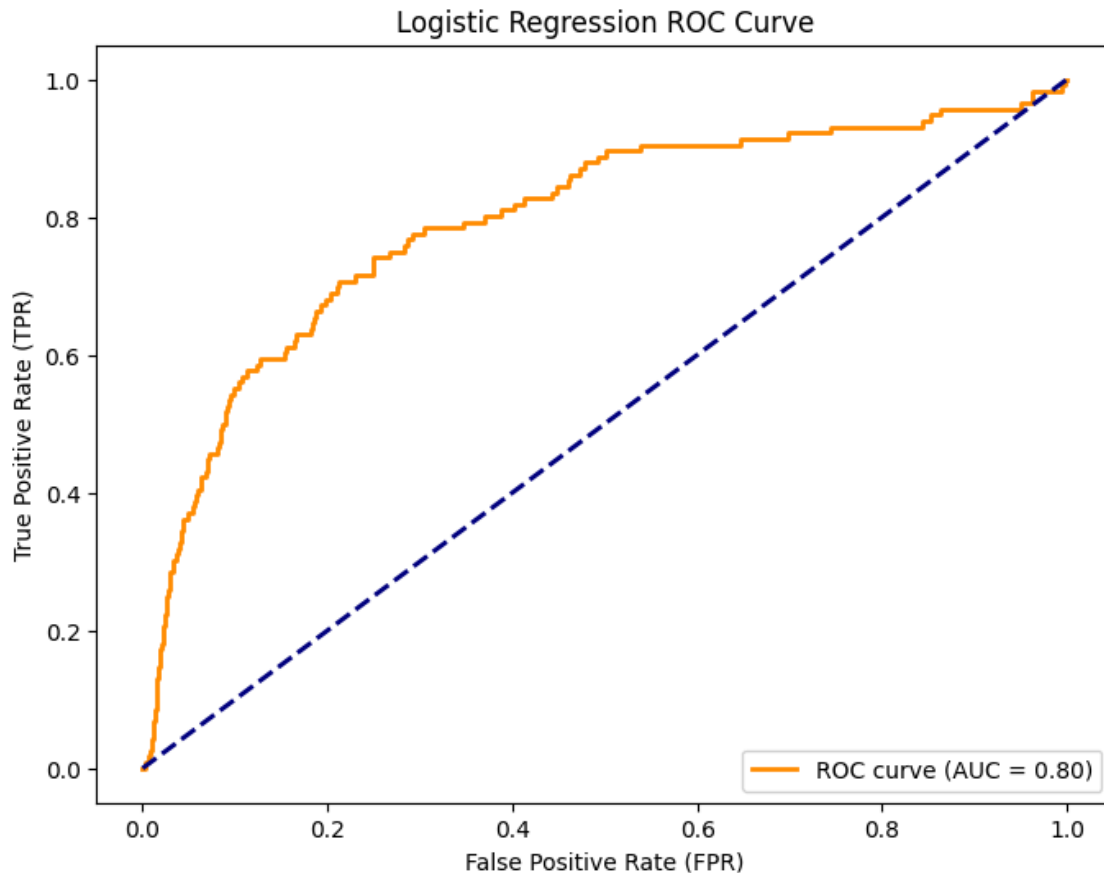
Increase the number of iterations (max_iter) or scale the data as shown in:

<https://scikit-learn.org/stable/modules/preprocessing.html>

Please also refer to the documentation for alternative solver options:

https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression

```
n_iter_i = _check_optimize_result(
```



AUC Score: 0.7956131834530888

```
1 pip install shap
2
```

➔ Requirement already satisfied: shap in /usr/local/lib/python3.11/dist-packages (0.47.1)
Requirement already satisfied: numpy in /usr/local/lib/python3.11/dist-packages (from shap) (2.0.2)
Requirement already satisfied: scipy in /usr/local/lib/python3.11/dist-packages (from shap) (1.14.1)
Requirement already satisfied: scikit-learn in /usr/local/lib/python3.11/dist-packages (from shap) (1.6.1)
Requirement already satisfied: pandas in /usr/local/lib/python3.11/dist-packages (from shap) (2.2.2)
Requirement already satisfied: tqdm>=4.27.0 in /usr/local/lib/python3.11/dist-packages (from shap) (4.67.
Requirement already satisfied: packaging>20.9 in /usr/local/lib/python3.11/dist-packages (from shap) (24.
Requirement already satisfied: slicer==0.0.8 in /usr/local/lib/python3.11/dist-packages (from shap) (0.0.
Requirement already satisfied: numba>=0.54 in /usr/local/lib/python3.11/dist-packages (from shap) (0.60.0
Requirement already satisfied: cloudpickle in /usr/local/lib/python3.11/dist-packages (from shap) (3.1.1)
Requirement already satisfied: typing-extensions in /usr/local/lib/python3.11/dist-packages (from shap) (4.
Requirement already satisfied: llvmlite<0.44,>=0.43.0dev0 in /usr/local/lib/python3.11/dist-packages (from
Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.11/dist-packages (from pa
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.11/dist-packages (from pandas->shap
Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.11/dist-packages (from pandas->sh
Requirement already satisfied: joblib>=1.2.0 in /usr/local/lib/python3.11/dist-packages (from scikit-lear
Requirement already satisfied: threadpoolctl>=3.1.0 in /usr/local/lib/python3.11/dist-packages (from scik
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.11/dist-packages (from python-dateutil>

```
1 # Only use a sample (e.g., 500 rows) for explanation
2 import shap
3 sample_data = X_test.sample(500, random_state=42)
4
5 explainer = shap.TreeExplainer(rf_model)
```

```


6 shap_values = explainer.shap_values(sample_data)
7

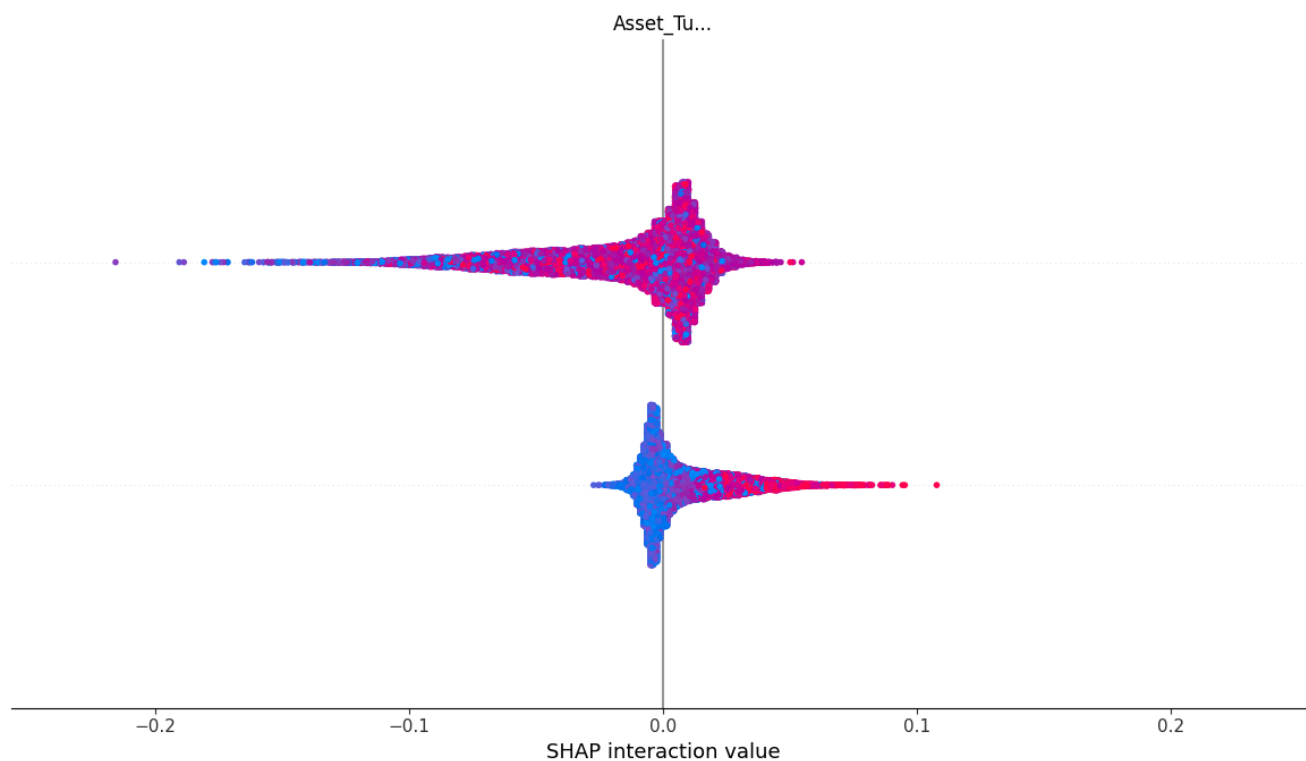
```

```

1 import shap
2
3
4 # Create a SHAP explainer using the base logistic regression model
5 explainer = shap.Explainer(rf_model, X_train_resampled)
6
7 # Compute SHAP values for the test set (or whichever data you want to explain)
8 shap_values = explainer.shap_values(X_test)
9
10 # Plot SHAP summary plot
11 shap.summary_plot(shap_values, X_test)
12

```

 100%|=====| 58666/58694 [15:44<00:00]
 <Figure size 640x480 with 0 Axes>



```

1 import pandas as pd
2 import shap
3
4
5 # If X_train is a NumPy array, make sure to convert it to a DataFrame with column names
6 X_train_df = pd.DataFrame(X_train_resampled, columns=predictor_columns)
7 X_test_df = pd.DataFrame(X_test, columns=predictor_columns)
8
9 # Initialize the SHAP explainer
10 explainer = shap.Explainer(logreg_model, X_train_df)
11
12 # Calculate SHAP values
13 shap_values = explainer(X_test_df)
14
15 # Calculate mean absolute SHAP value for each feature (global importance)
16 shap_summary = pd.DataFrame({
17     'Feature': X_test_df.columns,
18     'Mean_Abs_SHAP_Value': np.abs(shap_values.values).mean(axis=0)
19 })
20

```

```

21 # Sort by importance
22 shap_summary_sorted = shap_summary.sort_values(by='Mean_Abs_SHAP_Value', ascending=False)
23
24 # Save top 10 (or any number) to CSV
25 top_n = 10
26 shap_summary_sorted.head(top_n).to_csv('top_shap_features.csv', index=False)
27
28 print("Top SHAP features saved to 'top_shap_features.csv'")
29

```

➡ Top SHAP features saved to 'top_shap_features.csv'

```

1 # Step 1: Import libraries
2 from sklearn.linear_model import LogisticRegression
3 import joblib
4
5
6 # Step 3: Save the model
7 joblib.dump(logreg_model, 'logistic_regression_model.pkl')
8
9 print("✅ Model saved as 'logistic_regression_model.pkl'")
10

```

➡ ✅ Model saved as 'logistic_regression_model.pkl'

```

1 ! pip install streamlit -q
2

```

➡

_____	44.3/44.3 kB	2.8 MB/s	eta 0:00:00
_____	9.8/9.8 MB	41.9 MB/s	eta 0:00:00
_____	6.9/6.9 MB	49.1 MB/s	eta 0:00:00
_____	79.1/79.1 kB	5.2 MB/s	eta 0:00:00

```

1 !wget -q -O - ipv4.icanhazip.com

```

➡ 35.233.148.243

```

1 !streamlit run /content/drive/MyDrive/FYP2025/app_py.py & npx localtunnel --port 8501

```

...
Collecting usage statistics. To deactivate, set browser.gatherUsageStats to false.

👋👋👋

You can now view your Streamlit app in your browser.

Local URL: <http://localhost:8501>

Network URL: <http://172.28.0.12:8501>

External URL: <http://35.233.148.243:8501>

🌐:::your url is: <https://violet-comics-bathe.loca.lt>

```

1 !ls /content/drive/MyDrive
2

```

➡

```

'1st YEAR.zip'
'2023_01_02 3_53 pm Office Lens.pdf'
'311521104033 (1).jpg'
311521104033arw.pdf

```