

Ahmed Awadallah

Binary Classification using RandomForest and XGBoost

Introduction

Companies rise and fall. Knowing when they rise and fall can be incredibly lucrative for some interested parties and thus predicting such occurrences has great value. This case study focuses on predicting whether a company will or will not go bankrupt. The dataset contains data on companies across 5 years. The total feature count is 64 with 43405 observations total.

Methods

The dataset is split into 5 different years and thus there are 5 different files to work with. The filetype associated with this data is arff. The data was loaded in, transformed into a dataframe, and then coalesced into one large dataset.

Next the data was explored for general understanding as well as presence of NAs. All the features were found to be continuous so imputation was far easier. There was NAs found all of type np.nan which were handled by feature. If a feature had less than 50% NA then it was imputed using mean imputation otherwise the feature was removed. None of the features had excessive NAs so all were kept and imputed.

Scaling was not performed on the dataset as the algorithms we are using do not suffer from features with different scales.

The target variable was stored as an object datatype so that was converted into integer type as this is expected for some of the models and analysis tools.

The dataset was shuffled to remove any ordering present in the dataset.

Random Forest hyperparameters tuning was accomplished using sklearn random cross validation. Splits of 10 were chosen as this was the number chosen by the paper. Below are the ranges of values searched across.

```
distributions = {
    "criterion": ["gini", "entropy"],
    "max_depth": np.linspace(4, 500, 8).astype(int),
    "min_samples_split": np.linspace(2, 10, 9).astype(int),
    "min_samples_leaf": np.linspace(1, 10, 10).astype(int),
    "min_weight_fraction_leaf": np.logspace(-6, -1, 20),
    "min_impurity_decrease": np.logspace(-6, -1, 20)
}
```

XGBoost hyperparameter tuning was accomplished through a custom function which can be viewed in appendix A. Cross validation splits of 10 were chosen to maintain the same choice as the paper. Below are the range of values searched across.

```
distributions = {
    'eta': np.linspace(0, 1, 16),
    'max_depth': np.linspace(0, 1000, 16).astype(int),
    'min_child_weight': np.logspace(-6, 6, 16),
    'subsample': np.logspace(-6, 0, 8),
    'colsample_bytree': np.logspace(-6, 0, 8),
    'colsample_bylevel': np.logspace(-6, 0, 8),
    'colsample_bynode': np.logspace(-6, 0, 8),
    'lambda': np.logspace(-6, 6, 16),
    'alpha': np.logspace(-6, 6, 16)
}
```

The optimizing metric was auc for both models as this was the metric used in the paper.

Results

AUC Scores RF

```
RF Cross Val AUC
~~~~~
auc=0.950493, std=0.008503
```

Figure 1: Random Forest 10 Fold Cross Val
AUC Score and Standard Deviation

Random Forest Achieved an AUC score of 0.95. We can additionally say that we are 95% confident that the AUC lies between 0.944 and 0.956.

Classification Report RF

Random Forest Classification Report							
			precision	recall	f1-score	support	
		0	0.97	1.00	0.98	4136	
		1	1.00	0.31	0.48	204	
	accuracy				0.97	4340	
	macro avg		0.98	0.66	0.73	4340	
	weighted avg		0.97	0.97	0.96	4340	

Figure 2: Random Forest Classification Report for a Single Fold

A value of 0 indicates the company did not go bankrupt while a value of 1 indicates a company did go bankrupt.

Overall the model performs decently. It is important to remember that this dataset is imbalanced as evidenced by the support values. The overall accuracy looks good but we see that specifically predicting if a company goes bankrupt is hard. The precision score is high which shows the model rarely predicts a company bankrupt when it is in fact not bankrupt. The recall score is low which shows the model will often classify bankrupt companies as not bankrupt when they are.

Feature Importance RF

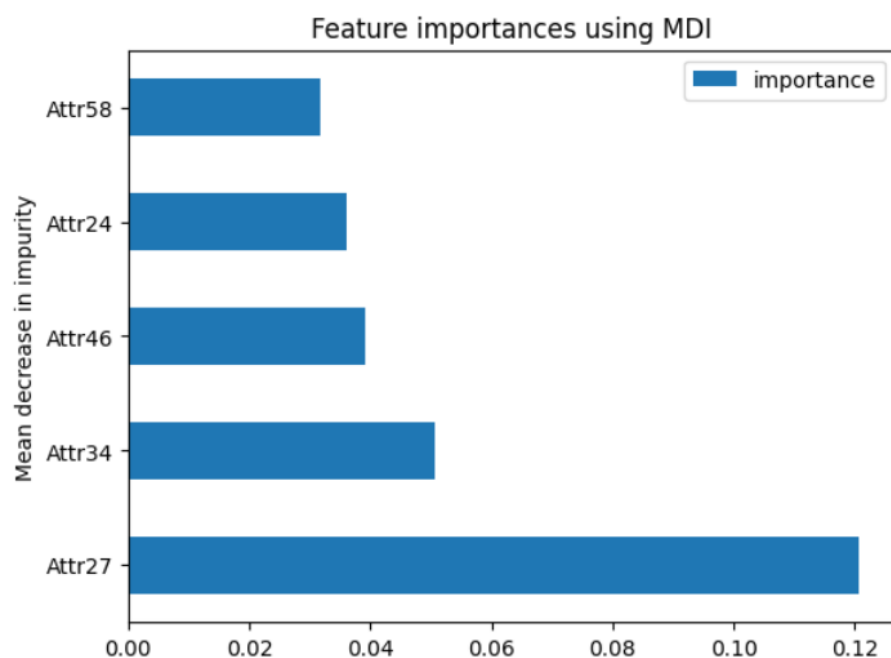


Figure 3: Random Forest Feature Importance

Below are the references for the codified attribute numbers:

Attr58: total costs /total sales

Attr24: gross profit (in 3 years) / total assets

Attr46: (current assets - inventory) / short-term liabilities

Attr34: operating expenses / total liabilities

Attr27: profit on operating activities / financial expenses

The codified names were chosen to remain as the actual attributes themselves are all largely complex mathematical combinations that are equally as vague as just an attribute number.

For determining feature importance in a Random Forest we decided to use the MDI metric which stands for Mean Decrease in Impurity. The higher this value the better.

For our model it appears the most important attribute was Attr27 which is profit on operating activities / financial expenses. This attribute is a bit interesting as profit by itself already takes into account expenses as it is a measure of net gain. This attribute goes one step further by taking the net gain and producing a ratio of net gain to expenses. The main takeaway is the differences between gain of money and loss of money was useful to the model.

Attr34 and Attr46 both involve a metric with liabilities so that may be worth looking into further.

Attr 24 is a relationship between profit and assets. Profit was a component in the Attr27 so further evidence profit is important.

Attr58 is the ratio between costs and sales.

AUC Scores XG

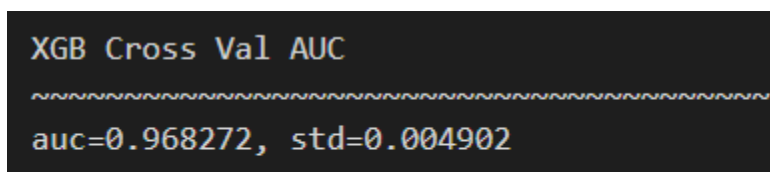


Figure 4: XGB 10 Fold Cross Validation AUC score and Standard Deviation

XGB achieved an auc score of 0.968. Additionally we can say with 95% confidence that the true auc score lies between 0.964 and 0.972. Since the confidence interval of random forest and auc don't overlap we can say statistically they have different values and that the XGB model is better. (RF: 0.944 and 0.956, XGB: 0.964 and 0.972)

Classification Report XG

XGB Classification Report						
			precision	recall	f1-score	support
		0	0.98	1.00	0.99	4136
		1	0.92	0.58	0.71	204
	accuracy				0.98	4340
	macro avg		0.95	0.79	0.85	4340
	weighted avg		0.98	0.98	0.98	4340

Figure 5: XGB Classification Report

A value of 0 indicates the company did not go bankrupt while a value of 1 indicates a company did go bankrupt.

Overall the model performs well. Previously the random first struggled with the bankruptcy class but here we see the XGB model improved in this regard. It was able to achieve an f1 score of 0.71 compared to the random forest which got 0.48. This increase in performance was gained by improving recall but it came at a small cost of precision.

With both classification reports reviewed a small aside on metrics is worth discussing. The choice of auc as the optimizing metric was done because the paper we are comparing to used this metric. This appears to prefer optimizing precision over recall but in practice preferring recall is likely preferable here. For example if we were investing in a company and someone told us that the company will go bankrupt and it turned out down the road they never did that sucks but it's not a huge deal. On the other hand if they told us it wont go bankrupt and it ended up going bankrupt we may be frustrated with that result. The asymmetry in these failures is mirrored within the precision and recall metrics respectively.

Feature Importance XG

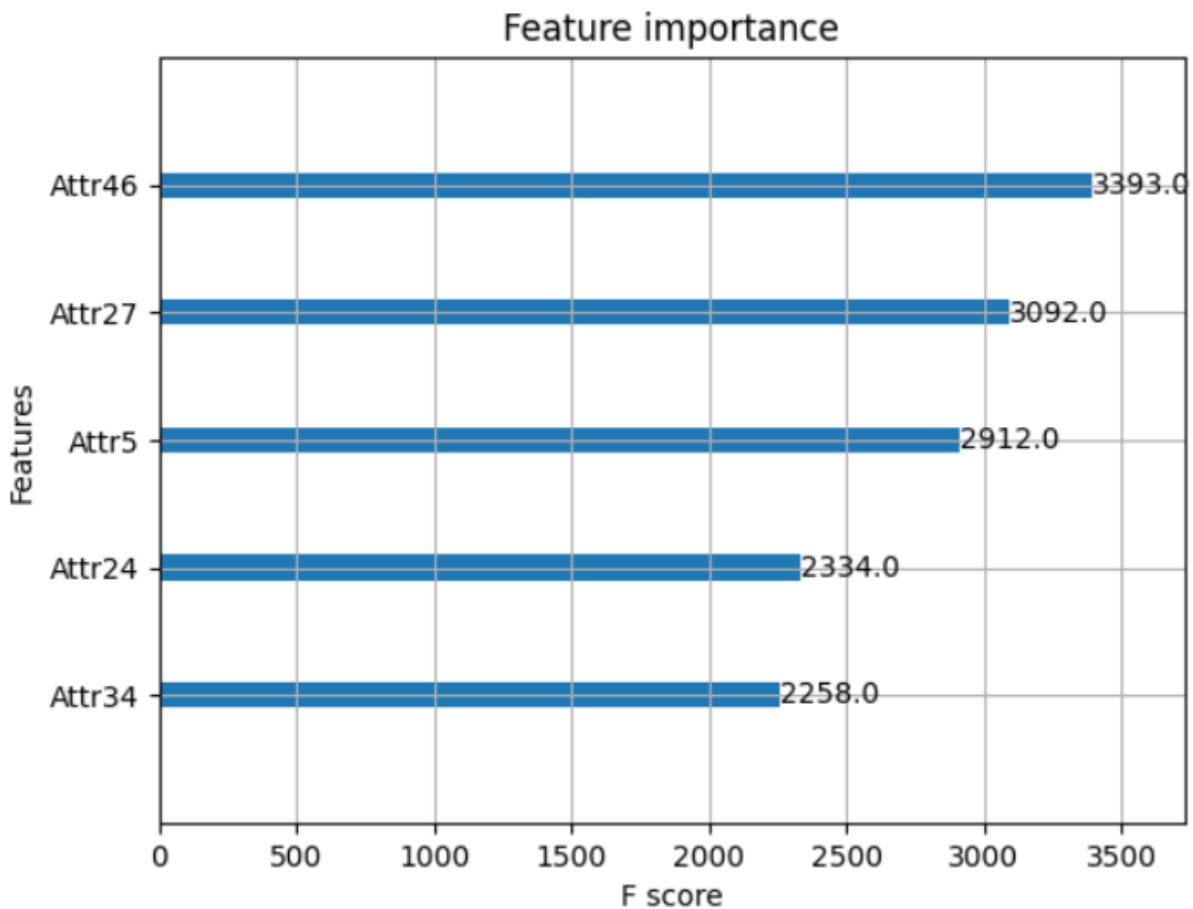


Figure 6: XGB Feature Importance

Below are the references for the codified attribute numbers:

Attr46: $(\text{current assets} - \text{inventory}) / \text{short-term liabilities}$

Attr27: $\text{profit on operating activities} / \text{financial expenses}$

Attr5: $[(\text{cash} + \text{short-term securities} + \text{receivables} - \text{short-term liabilities}) / (\text{operating expenses} - \text{depreciation})] * 365$

Attr24: $\text{gross profit (in 3 years)} / \text{total assets}$

Attr34: $\text{operating expenses} / \text{total liabilities}$

First off it's important to observe that 4/5 of these features were also present in the top 5 features of the random forest albeit in a different order. This gives extra confidence that these features are important in bankruptcy prediction. A deeper exploration into the inherent behavior of these complex mathematical combinations is warranted.

For the XGB model the most important feature was Attr46 which is $(\text{current assets} - \text{inventory}) / \text{short-term liabilities}$. In the random forest model this was the third most important feature. This

indicates that there is something important about this ratio but like stated earlier this is a complex relationship. A more thorough examination of this metric should be done to be able to derive why the model found this so important.

Since the other features were present in random forest the last attribute we will discuss is Attr5 which was not present in random forest and is the third most important feature for XGB. Attr5 is $[(\text{cash} + \text{short-term securities} + \text{receivables} - \text{short-term liabilities}) / (\text{operating expenses} - \text{depreciation})] * 365$. This is an incredibly complex feature combination and is unlike any of the other features we've discussed. This is not easily interpreted but one might wonder if this was seen as important because it includes so many features within it that their cumulative impact is much larger than their individual parts.

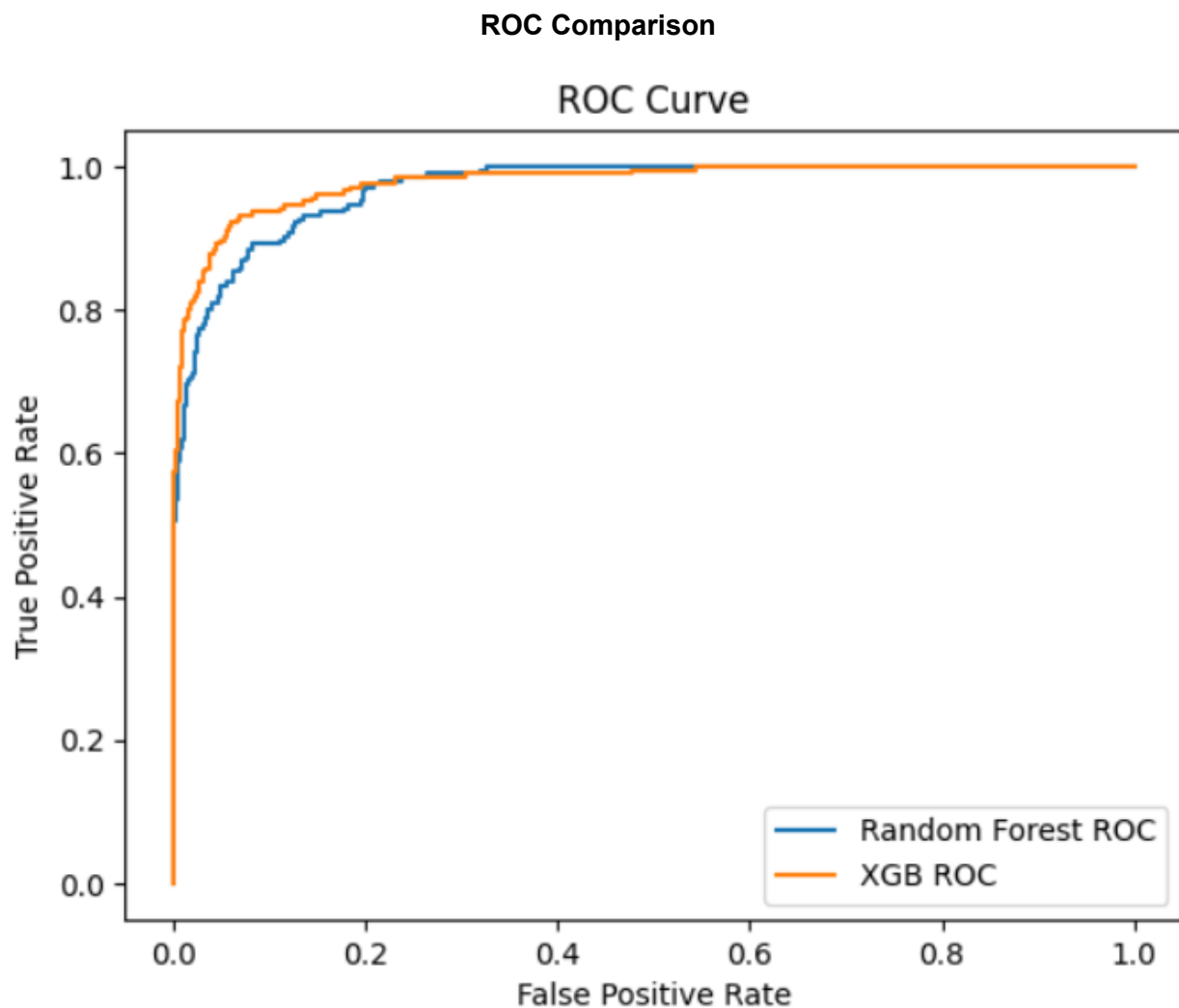


Figure 7: ROC Comparison

Lastly is an ROC comparison curve between the Random Forest and XGB model. It was observed earlier in the paper that XGB performs better so this is simply further justification.

Individually it also shows that both models can evenly trade off between true positive rate and false positive rate. This means there exists thresholds where we could prefer one metric doing better while slightly penalizing another. Previously we discussed optimizing recall over precision and while this plot shows TPR v FPR it can still somewhat help us as TPR is exactly recall. FPR is not exactly precision but we know that if we pick a threshold with higher TPR we would be directly improving recall. This comes at some unknown cost to precision but known cost to FPR.

Conclusion

This dataset was severely imbalanced and yet Random Forest and XGBoost handle this situation with ease. It wasn't mentioned in the report but the models with default settings also performed well but hyperparameter tuning helped a great deal with predicting bankruptcies (the rare case). Random Forest is a great model to start off a project as a base model as it is easy and quick to set up while XGBoost is for when performance is key and taking the time to maximize that is worth it.

Code

Appendix A: Custom XGBoost Random CV

```
def perform_optimal_search(
    clf_object,
    clf_hyperparams,
    static_params,
    optim_metric,
    x,
    y,
    n_iter=20,
    splits=2,
    report_freq=0.1
):
    kf = KFold(n_splits=splits)
    hyperparam_combinations = [
        dict(
            zip(clf_hyperparams, v)
        ) for v in product(*clf_hyperparams.values())
    ]
    shuffle(hyperparam_combinations)
    hyperparam_combinations_subset = sample(hyperparam_combinations, n_iter)
    report_freq = report_freq*len(hyperparam_combinations_subset)
    optimal_hyperparam = {}
    optimal_performance = {}
    first_iteration = True
    iter_number = 0
```

```

for hyperparam in hyperparam_combinations_subset:
    performance = 0
    for id, (train_i, test_i) in enumerate(kf.split(x, y)):
        x_train = x.iloc[train_i]
        y_train = y.iloc[train_i]
        x_test = x.iloc[test_i]
        y_test = y.iloc[test_i]
        clf = clf_object(**static_params, **hyperparam)
        clf.fit(x_train, y_train, eval_set=[(x_test, y_test)], verbose=False)
        y_pred = clf.predict_proba(x_test)[:, 1]
        performance += optim_metric(y_test, y_pred)
    performance = performance/splits
    if first_iteration:
        first_iteration = False
        optimal_hyperparam = hyperparam
        optimal_performance = performance
    else:
        if performance > optimal_performance:
            optimal_performance = performance
            optimal_hyperparam = hyperparam
    if (iter_number+1) % report_freq == 0:
        percent_done = (iter_number+1)/len(hyperparam_combinations_subset)*100
        print(f"Precent Done: {percent_done:4f}%")
    iter_number += 1

return optimal_performance, optimal_hyperparam

```

Appendix B: Script

```
1  import sys
2  import os
3  from scipy.io import arff
4  import pandas as pd
5  from sklearn.impute import SimpleImputer
6  import numpy as np
7  from sklearn.ensemble import RandomForestClassifier
8  from sklearn.metrics import classification_report
9  import sklearn.metrics as mt
10 import matplotlib.pyplot as plt
11 from xgboost import XGBClassifier
12 from xgboost import plot_importance
13 from sklearn.model_selection import cross_val_score, KFold
14 from sklearn.metrics import roc_auc_score
15
16
17 class UnexpectedFileStructureError(Exception):
18     def __init__(self):
19         msg = "This script expects a file path to a directory that contains "
20         msg += "CS4 Bankruptcy Dataset. The Bankruptcy dataset is a "
21         msg += "group of five arff files with the names, 1year, 2year, "
22         msg += "3year, 4year, 5year. Did you specify a folder path that "
23         msg += "contains these files?"
24         super().__init__(msg)
25
26
```

```

27 class MissingInputError(Exception):
28     def __init__(self):
29         msg = "This script expects exactly one input, a file directory that "
30         msg += "contains the bankruptcy dataset. You input nothing. "
31         msg += "Did you forget to add an argument with the script?"
32         super().__init__(msg)
33
34
35 class ExcessiveInputError(Exception):
36     def __init__(self):
37         msg = "This script expects exactly one input, a file directory that "
38         msg += "contains the bankruptcy dataset. You called this script "
39         msg += "with more than one argument. Only use one and try again. "
40         super().__init__(msg)
41
42
43 def check_correct_file_structure(inputs):
44     data_direct = inputs[1]
45     expected_files = ["1year.arff",
46                      "2year.arff",
47                      "3year.arff",
48                      "4year.arff",
49                      "5year.arff"]
50     for root, dirs, files in os.walk(data_direct):
51         for f in files:
52             if f not in expected_files:
53                 print(f)
54                 raise UnexpectedFileStructureError
55

```

```
57 def check_valid_inputs(inputs):
58     if len(inputs) < 2:
59         raise MissingInputError
60     elif len(inputs) > 2:
61         raise ExcessiveInputError
62
63
64 def load_data(inputs):
65     data_direct = inputs[1]
66     raw_year_data = []
67     for root, dirs, files in os.walk(data_direct):
68         for f in files:
69             file_path = os.path.join(root, f)
70             arff_data = arff.loadarff(file_path)
71             raw_year_data.append(pd.DataFrame(arff_data[0]))
72     return raw_year_data
73
74
75 def log_section(title="No Title", content="No Content"):
76     log_file.write(title + "\n")
77     log_file.write("~~~~~\n")
78     log_file.write(content + "\n")
79     log_file.write("\n")
80
```

```

82 def preprocess_data(raw_df):
83     # Seperate features by percentage of data NA
84     na_variations = [np.nan]
85     below_50_precent_na_features = [] # Mean Impute
86     above_50_perecent_na_features = [] # Remove
87
88     for feature in list(raw_df):
89         percent_na = raw_df[feature].isin(na_variations).sum()
90         percent_na = percent_na / (raw_df.shape[0]) * 100
91         percent_na = round(percent_na, 3)
92         if percent_na <= 50:
93             below_50_precent_na_features.append(feature)
94         if percent_na > 50:
95             above_50_perecent_na_features.append(feature)
96
97     # Handle NAs
98     clean_df = raw_df.copy()
99     # Mean Impute if less than 50% missing
100    mean_imp = SimpleImputer(missing_values=np.nan, strategy='mean')
101    clean_df[below_50_precent_na_features] = mean_imp.fit_transform(
102        clean_df[below_50_precent_na_features])
103
104    # Remove if more than 50% missing
105    for feature in above_50_perecent_na_features:
106        clean_df = clean_df.drop([feature], axis=1)
107
108    # Changing Target Feature to be int type
109    target_feature = ["class"]
110    cts_features = []
111    for feature in list(clean_df):
112        if feature not in target_feature:
113            cts_features.append(feature)
114    clean_df[target_feature] = clean_df[target_feature].astype(int)
115
116    return clean_df

```

```
119 if __name__ == '__main__':
120     inputs = sys.argv
121     check_valid_inputs(inputs)
122     check_correct_file_structure(inputs)
123     print("Full Execution on my machine takes around 20 Minutes")
124     # Log Start
125     log_file = open("log.txt", "w")
126     log_file.write("Case Study 4 Report\n\n")
127
128     # Use input to load data
129     raw_years = load_data(inputs)
130
131     # Group all years into one large dataset and shuffle
132     raw_df = pd.concat(raw_years, axis=0, ignore_index=True)
133
134     # Data Preprocessing
135     clean_df = preprocess_data(raw_df)
136
137     # Target Feature Mappings
138     bankrupt_to_id_map = {"Not Bankrupt": 0, "Bankrupt": 1}
139     id_to_bankrupt_map = {0: "Not Bankrupt", 1: "Bankrupt"}
140
141     # Shuffle
142     clean_df = clean_df.sample(frac=1)
```

```

144 # Seperate Features and Target
145 target_feature = ["class"]
146 y = clean_df[target_feature[0]]
147 x = clean_df.drop(target_feature, axis=1)
148
149 ninety_precent_of_data = int(x.shape[0]*0.9)
150 ten_precent_of_data = int(x.shape[0]*0.1)
151 train_x = x.head(ninety_precent_of_data)
152 train_y = y.head(ninety_precent_of_data)
153 val_x = x.tail(ten_precent_of_data)
154 val_y = y.tail(ten_precent_of_data)
155
156 # Random Forest Model
157 rf_params = {'min_weight_fraction_leaf': 3.359818286283781e-06,
158             'min_samples_split': 10,
159             'min_samples_leaf': 6,
160             'min_impurity_decrease': 2.06913808111479e-05,
161             'max_depth': 429,
162             'criterion': 'entropy',
163             'n_jobs': 4,
164             'n_estimators': 1000}
165
166 # Cross Val Score
167 rf = RandomForestClassifier(**rf_params)
168 splits = KFold(n_splits=10, shuffle=True)
169 cross_score_rf = cross_val_score(rf, x, y, cv=splits, scoring='roc_auc')
170 auc_rf = cross_score_rf.mean()
171 auc_rf_std = cross_score_rf.std()
172 rf_cross_val_rep = f"auc={auc_rf:4f}, std={auc_rf_std:4f}"
173 log_section(title="RF Cross Val AUC",
174            content=rf_cross_val_rep)

```



```

179     # XGB Model
180     xg_params = {
181         "objective": "binary:logistic",
182         "nthread": 4,
183         "eval_metric": "auc",
184         "early_stopping_rounds": 5,
185         "n_estimators": 200,
186         'eta': 0.6309573444801936,
187         'max_depth': 12,
188         'min_child_weight': 0.00630957344480193,
189         'subsample': 1.0,
190         'colsample_bytree': 1.0,
191         'colsample_bylevel': 1.0,
192         'colsample_bynode': 1.0,
193         'lambda': 10}
194
195     # Cross Val Score
196     cross_score_xgb = []
197     split_count = 10
198     kf = KFold(n_splits=split_count)
199     for id, (train_i, test_i) in enumerate(kf.split(x, y)):
200         x_train = x.iloc[train_i]
201         y_train = y.iloc[train_i]
202         x_test = x.iloc[test_i]
203         y_test = y.iloc[test_i]
204         xg_model = XGBClassifier(**xg_params)
205         xg_model.fit(x_train,
206                     y_train,
207                     eval_set=[(x_test, y_test)],
208                     verbose=False)
209         y_pred = xg_model.predict_proba(x_test)[:, 1]
210         cross_score_xgb.append(roc_auc_score(y_test, y_pred))

```

```

212 cross_score_xgb = np.asarray(cross_score_xgb)
213 auc_xgb = cross_score_xgb.mean()
214 auc_xgb_std = cross_score_xgb.std()
215 xgb_cross_val_rep = f"auc={auc_xgb:4f}, std={auc_xgb_std:4f}"
216 log_section(title="XGB Cross Val AUC",
217             content=xgb_cross_val_rep)
218
219 # Final Model For Analysis
220 xg_params_no_early = {
221     "objective": "binary:logistic",
222     "nthread": 4,
223     "n_estimators": 200,
224     'eta': 0.6309573444801936,
225     'max_depth': 12,
226     'min_child_weight': 0.00630957344480193,
227     'subsample': 1.0,
228     'colsample_bytree': 1.0,
229     'colsample_bylevel': 1.0,
230     'colsample_bynode': 1.0,
231     'lambda': 10}
232 xg_model = XGBClassifier(**xg_params_no_early)
233 xg_model.fit(train_x, train_y, verbose=False)
234
235 # Random Forest Analysis
236 # Classification Report
237 y_pred_final = rf.predict(val_x)
238 y_true = val_y
239 classif_rep = classification_report(
240     y_true,
241     y_pred_final
242 )
243 log_section(title="Random Forest Classification Report",
244             content=classif_rep)

```

```

# AUC and ROC Curve All Data
y_pred_prob = rf.predict_proba(val_x)
pos_class = 1
y_label = id_to_bankrupt_map[pos_class]
preds = y_pred_prob[:, pos_class]
fpr_tree, tpr_tree, thresholds = mt.roc_curve(val_y,
        preds,
        pos_label=pos_class)
auc_report = f"{y_label} AUC: {mt.auc(fpr_tree, tpr_tree):.4f}\n"
log_section(title="Random Forest AUC Report",
        content=auc_report)

# Feature Importance
importances_rf = pd.DataFrame(
    rf.feature_importances_,
    index=x.columns,
    columns=["importance"]).sort_values("importance", ascending=False)
fig, ax = plt.subplots()
importances_rf.head(5).plot.barh(ax=ax)
ax.set_title("Feature importances using MDI")
ax.set_ylabel("Mean decrease in impurity")
plt.savefig("rf_feature_importance.png", bbox_inches='tight')
plt.clf()

# XGB Analysis
# Classification Report
y_pred_final = xg_model.predict(val_x)
y_true = val_y
classif_rep = classification_report(
    y_true,
    y_pred_final
)
log_section(title="XGB Classification Report",
        content=classif_rep)

```

```

281 # AUC and ROC Curve All Data
282 y_pred_prob = xg_model.predict_proba(val_x)
283 pos_class = 1
284 y_label = id_to_bankrupt_map[pos_class]
285 preds = y_pred_prob[:, pos_class]
286 fpr_xgb, tpr_xgb, thresholds = mt.roc_curve(val_y,
287 |                                     preds,
288 |                                     pos_label=pos_class)
289 auc_report = f"{y_label} AUC: {mt.auc(fpr_xgb, tpr_xgb):.4f}\n"
290 log_section(title="XGB AUC Report",
291 |           content=auc_report)
292
293 # Feature Importance
294 plot_importance(xg_model, max_num_features=5)
295 plt.savefig("xgb_feature_importance.png", bbox_inches='tight')
296 plt.clf()
297
298 # ROC Comparison Plot
299 plt.plot(fpr_tree, tpr_tree, label="Random Forest ROC")
300 plt.plot(fpr_xgb, tpr_xgb, label="XGB ROC")
301 plt.title("ROC Curve")
302 plt.xlabel('False Positive Rate')
303 plt.ylabel('True Positive Rate')
304 plt.legend()
305 plt.savefig("roc_plot.png", bbox_inches='tight')
306 plt.clf()
307
308 log_file.write("\n")
309 log_file.close()

```