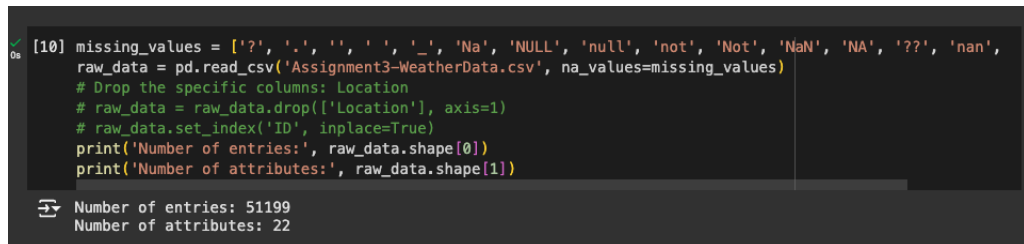# I.    Problem Statement

The goal of this report is to accurately predict the likelihood of rainfall on the following day, which is represented by the **"RainTomorrow"** attribute.

# II.    Data Preprocessing

## 1. Dataset

The dataset contains nearly a decade of daily weather observations collected from various locations throughout Australia. There are **51199** entries and **22** attributes in this dataset (see Figure 1). It is easy to see that the dataset contains a mixture of categorical and numerical variables (see Figure 2), where categorical variables have data type **"object"** and numerical variables have data type **"float64".**

**Figure 1**

```
[10] missing_values = ['?', '.', '', ' ', '_', 'Na', 'NULL', 'null', 'not', 'Not', 'NaN', 'NA', '??', 'nan',
     raw_data = pd.read_csv('Assignment3-WeatherData.csv', na_values=missing_values)
     # Drop the specific columns: Location
     # raw_data = raw_data.drop(['Location'], axis=1)
     # raw_data.set_index('ID', inplace=True)
     print('Number of entries:', raw_data.shape[0])
     print('Number of attributes:', raw_data.shape[1])

    Number of entries: 51199
    Number of attributes: 22
```

**Figure 2**

```
RangeIndex: 51199 entries, 0 to 51198
Data columns (total 22 columns):
 #   Column         Non-Null Count   Dtype
---  ------         --------------   -----
 0   Location       51199 non-null   object
 1   MinTemp        50959 non-null   float64
 2   MaxTemp        51085 non-null   float64
 3   Rainfall       50667 non-null   float64
 4   Evaporation    29227 non-null   float64
 5   Sunshine       26720 non-null   float64
 6   WindGustDir    47878 non-null   object
 7   WindGustSpeed  47901 non-null   float64
 8   WindDir9am     47564 non-null   object
 9   WindDir3pm     49795 non-null   object
 10  WindSpeed9am   50695 non-null   float64
 11  WindSpeed3pm   50199 non-null   float64
 12  Humidity9am    50531 non-null   float64
 13  Humidity3pm    49840 non-null   float64
 14  Pressure9am    46138 non-null   float64
 15  Pressure3pm    46136 non-null   float64
 16  Cloud9am       31795 non-null   float64
 17  Cloud3pm       30551 non-null   float64
 18  Temp9am        50851 non-null   float64
 19  Temp3pm        50161 non-null   float64
 20  RainToday      50667 non-null   object
 21  RainTomorrow   51199 non-null   int64
dtypes: float64(16), int64(1), object(5)
memory usage: 8.6+ MB
```

## 1.1.   Numerical Attribute

The dataset consists of 16 numeric attributes that each offer unique insights. By examining key metrics like the average (mean), how spread out the values are (standard deviation), and specific points in the distribution (like the 25th, 50th, and 75th percentiles), we can get a clearer picture of how the data is distributed. Analyzing the variance and the range also helps us understand the extent of variation in the data (see Figure 3).

**Figure 3**

| | count | mean | std | min | 25% | 50% | 75% | max | variance | iqr_size | skewness | kurtosis | nulls_count | outliers_count | nulls_percent | outliers_percent |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| MinTemp | 50959.0 | 12.19 | 6.42 | -8.5 | 7.60 | 12.0 | 16.9 | 33.9 | 41.18 | 9.30 | 0.02 | -0.49 | 240 | 22 | 0.47 | 0.04 |
| MaxTemp | 51085.0 | 23.23 | 7.15 | -4.1 | 17.90 | 22.6 | 28.3 | 47.3 | 51.10 | 10.40 | 0.23 | -0.26 | 114 | 140 | 0.22 | 0.27 |
| Rainfall | 50667.0 | 2.32 | 8.14 | 0.0 | 0.00 | 0.0 | 0.8 | 268.6 | 66.19 | 0.80 | 8.90 | 135.20 | 532 | 9142 | 1.04 | 18.04 |
| Evaporation | 29227.0 | 5.45 | 4.08 | 0.0 | 2.60 | 4.8 | 7.4 | 72.2 | 16.68 | 4.80 | 3.04 | 24.36 | 21972 | 698 | 42.91 | 2.39 |
| Sunshine | 26720.0 | 7.63 | 3.78 | 0.0 | 4.90 | 8.5 | 10.6 | 14.5 | 14.29 | 5.70 | -0.50 | -0.83 | 24479 | 0 | 47.81 | 0.00 |
| WindGustSpeed | 47901.0 | 39.99 | 13.53 | 6.0 | 31.00 | 39.0 | 48.0 | 130.0 | 183.14 | 17.00 | 0.84 | 1.27 | 3298 | 1034 | 6.44 | 2.16 |
| WindSpeed9am | 50695.0 | 14.00 | 8.88 | 0.0 | 7.00 | 13.0 | 19.0 | 130.0 | 78.88 | 12.00 | 0.77 | 1.38 | 504 | 613 | 0.98 | 1.21 |
| WindSpeed3pm | 50199.0 | 18.63 | 8.77 | 0.0 | 13.00 | 19.0 | 24.0 | 83.0 | 76.99 | 11.00 | 0.61 | 0.69 | 1000 | 857 | 1.95 | 1.71 |
| Humidity9am | 50531.0 | 68.84 | 19.02 | 1.0 | 57.00 | 70.0 | 83.0 | 100.0 | 361.92 | 26.00 | -0.48 | -0.03 | 668 | 515 | 1.30 | 1.02 |
| Humidity3pm | 49840.0 | 51.49 | 20.83 | 0.0 | 37.00 | 52.0 | 66.0 | 100.0 | 433.97 | 29.00 | 0.03 | -0.52 | 1359 | 0 | 2.65 | 0.00 |
| Pressure9am | 46138.0 | 1017.65 | 7.13 | 980.5 | 1012.90 | 1017.6 | 1022.4 | 1040.6 | 50.82 | 9.50 | -0.10 | 0.27 | 5061 | 424 | 9.88 | 0.92 |
| Pressure3pm | 46136.0 | 1015.25 | 7.07 | 979.0 | 1010.40 | 1015.2 | 1020.0 | 1037.9 | 49.92 | 9.60 | -0.04 | 0.14 | 5063 | 330 | 9.89 | 0.72 |
| Cloud9am | 31795.0 | 4.44 | 2.88 | 0.0 | 1.00 | 5.0 | 7.0 | 9.0 | 8.31 | 6.00 | -0.23 | -1.54 | 19404 | 0 | 37.90 | 0.00 |
| Cloud3pm | 30551.0 | 4.51 | 2.72 | 0.0 | 2.00 | 5.0 | 7.0 | 8.0 | 7.38 | 5.00 | -0.23 | -1.46 | 20648 | 0 | 40.33 | 0.00 |
| Temp9am | 50851.0 | 16.99 | 6.51 | -6.2 | 12.25 | 16.7 | 21.6 | 39.0 | 42.36 | 9.35 | 0.09 | -0.37 | 348 | 84 | 0.68 | 0.17 |
| Temp3pm | 50161.0 | 21.69 | 6.97 | -4.4 | 16.60 | 21.1 | 26.5 | 46.7 | 48.54 | 9.90 | 0.24 | -0.16 | 1038 | 226 | 2.03 | 0.45 |
| RainTomorrow | 51199.0 | 0.22 | 0.42 | 0.0 | 0.00 | 0.0 | 0.0 | 1.0 | 0.17 | 0.00 | 1.32 | -0.25 | 0 | 11471 | 0.00 | 22.40 |

### 1.2. Categorical Attribute

The dataset consists of five categorical attributes, including wind direction, which do not possess an inherent ordinal relationship. In this situation, Label Encoding was chosen for its simplicity and strong performance, despite its drawbacks (see Figure 4). While Label Encoding assigns numerical values to each category, it can create an artificial ordinal relationship where none exists, potentially leading the model to misinterpret the categorical data as having a ranked structure, which is particularly problematic for features like wind direction.

Nonetheless, Label Encoding outperformed both One-Hot Encoding and Binary Encoding in this case. One-Hot Encoding notably increased the feature space, adding unnecessary computational complexity without enhancing the results. On the other hand, although Binary Encoding is more compact, it failed to effectively capture the distinctions between categories. Ultimately,

Label Encoding offered a streamlined and efficient approach, enabling the model to process the data effectively and achieve the best performance in this scenario.
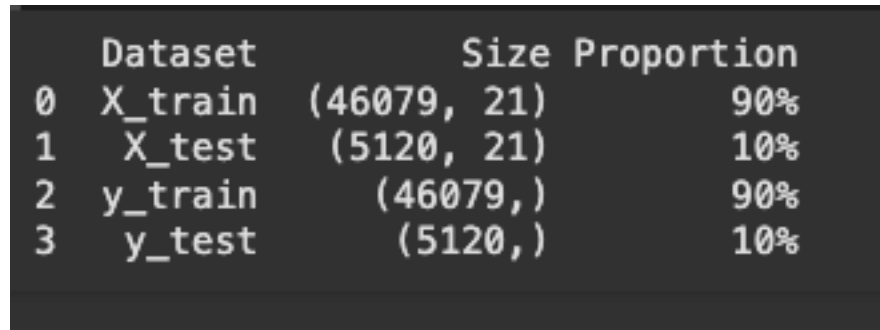
## 2. Data Splitting

Most machine learning models require a split between training, validation, and test sets. In cases of limited data, we simply divide it into training and validation sets. The training set, typically comprising a large portion of the data, helps the model learn a wide range of patterns. On the other hand, the validation set is used to re-evaluate the model and detect problems like overfitting and underfitting.

The question is "Why include a separate test set if we already have a validation test?". The test set is carefully chosen to closely resemble the actual data distribution and characteristics; therefore, it will evaluate how the model might perform in real-world deployment. When selecting the best model, we will rely on the evaluation metrics from this test set.

In this project, I will split the data into training and test sets before preprocessing to avoid data leakage (see Figure 4). This is an error where information from the test set unintentionally influences the training set. Data leakage can lead to overly optimistic results, as the model may gain prior insights about the test data; thus, its performance looks better than it would on unseen data. Later, in the modeling phase, cross-validation will further split the training set to create a validation set for fine-tuning.

Due to the complexity and size of the dataset (**51,199 entries)**, a 90/10

split is chosen.

**Figure 4**

```
    Dataset              Size Proportion
0  X_train  (46079, 21)           90%
1   X_test   (5120, 21)           10%
2  y_train     (46079,)           90%
3   y_test      (5120,)           10%
```

3. **Handling Missing Values**

Effective handling of missing values is necessary to ensure the robustness

and accuracy of predictive models. After experimenting with many imputation

methods, I identified the most effective approach for both numerical and

categorical attributes to train my model.

For numerical attributes, I initially experimented with mean, median, and

K-Nearest Neighbors (KNN) imputation. However, the best results were

achieved using **Multiple Imputation by Chained Equations (MICE),**

implemented via **IterativeImputer** (see Figure 5)**.** This method covers

complex relationships by iteratively predicting each feature based on others.

As a result, it helps preserve the data's underlying patterns and significantly

improves model performance.

**Figure 5**

```
from sklearn.experimental import enable_iterative_imputer  # Enables IterativeImputer in sklearn
from sklearn.impute import IterativeImputer

# Set up the IterativeImputer for multiple imputation by chained equations
mice_imputer = IterativeImputer(random_state=42)

# Impute only numerical columns in X_train using IterativeImputer
X_train[numerical] = mice_imputer.fit_transform(X_train[numerical])
X_test[numerical] = mice_imputer.transform(X_test[numerical])

# Check for any remaining missing values in X_train
print("Missing values in X_train after imputation:", X_train.isnull().sum().sum())
# Verify imputation by checking for any remaining missing values
print("Missing values in X_test after imputation:", X_train.isnull().sum().sum())
```

For categorical attributes, I initially tried dropping rows with missing values, given the relatively low missing rates in some columns. However, this method led to a slight decline in performance. Mode imputation or filling missing values with the most frequent category proved more effective, maintaining the distribution and structure of the categorical features, and resulting in a stronger predictive model (see Figure 6).

**Figure 6**

```
[29] # Impute missing values in each specified column with its mode
     for col in categorical_columns:
         X_train[col].fillna(X_train[col].mode()[0], inplace=True)
         X_test[col].fillna(X_train[col].mode()[0], inplace=True)
```

Despite high missing rates in columns like Sunshine (47.8%) and Evaporation (42,9%) (see Figure 7), I chose to retain these features. These features likely contain valuable information for predicting the target variable. Moreover, when I tried to remove these columns, the performance of the models decreased as well.
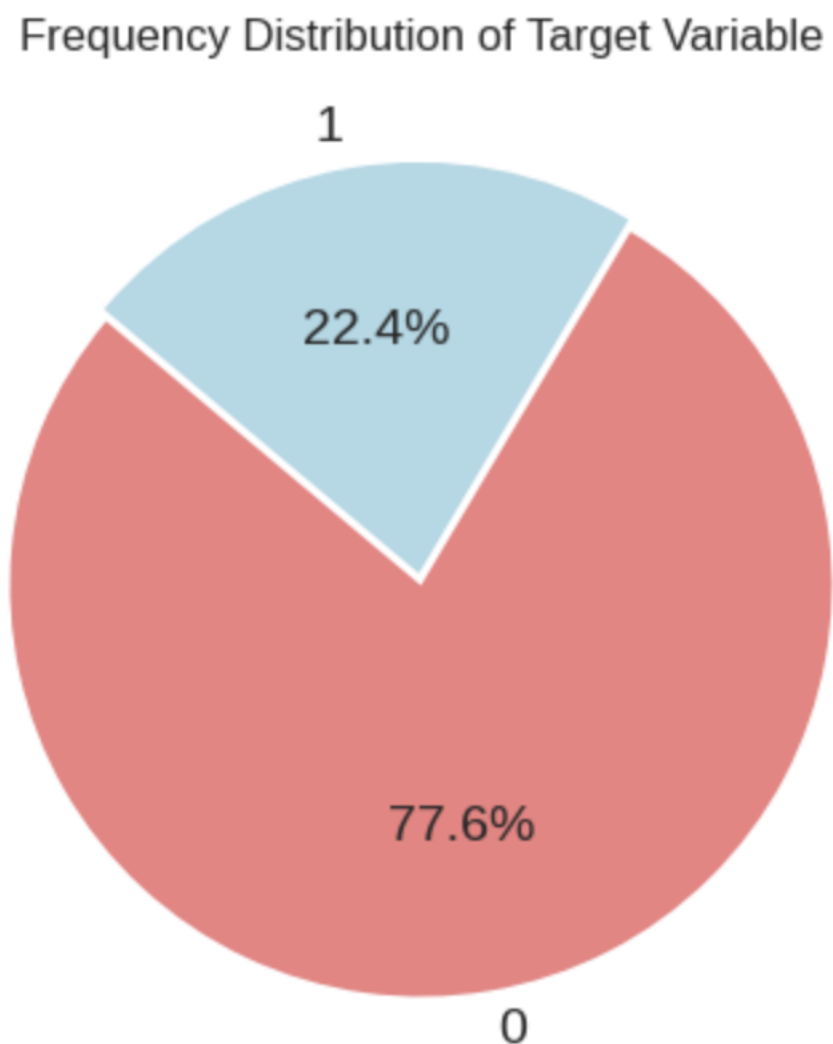
**Figure 7**

|  | nulls_count | nulls_percent |
| --- | --- | --- |
| Sunshine | 24479 | 47.811481 |
| Evaporation | 21972 | 42.914901 |
| Cloud3pm | 20648 | 40.328913 |
| Cloud9am | 19404 | 37.899178 |
| Pressure3pm | 5063 | 9.888865 |
| Pressure9am | 5061 | 9.884959 |
| WindGustSpeed | 3298 | 6.441532 |
| Humidity3pm | 1359 | 2.654349 |
| Temp3pm | 1038 | 2.027383 |
| WindSpeed3pm | 1000 | 1.953163 |
| Humidity9am | 668 | 1.304713 |
| Rainfall | 532 | 1.039083 |
| WindSpeed9am | 504 | 0.984394 |
| Temp9am | 348 | 0.679701 |
| MinTemp | 240 | 0.468759 |
| MaxTemp | 114 | 0.222661 |

4. **Dealing with Imbalanced Dataset**

Another problem that I must cover in this project is handling class imbalance. Class imbalance can lead to biased models that favor the majority

class, which will reduce predictive accuracy for the minority class. The dataset in this project exhibits class imbalance in the target variable **RainTomorrow**, where the instances of 1 (Rain) are significantly fewer than instances of 0 (no rain) (see Figure 8).

**Figure 8**

Frequency Distribution of Target Variable



To solve this problem, I experimented with different techniques, including **Synthetic Minority Over-sampling Technique (SMOTE)** and SMOTE for categorical and numerical attributes (**SMOTE-NC**). These techniques have

shown quite good results in balancing the dataset. However, the best results were obtained by upsampling the minority class directly using **bootstrapping**.

In this approach, I separated the training data into majority and minority classes. Performing bootstrapping to upsample the minority class by using the **resample** function from Scikit-Learn (see Figure 9). It created synthetic samples until it matched the size of the majority class. This balanced the dataset and allowed the model to learn from an even distribution of outcomes.

**Firgure 9**

```python
from sklearn.utils import resample

# Separate majority and minority classes
X_minority = X_train[y_train == 1]
y_minority = y_train[y_train == 1]
X_majority = X_train[y_train == 0]
y_majority = y_train[y_train == 0]

# Upsample minority class using bootstrapping
X_minority_upsampled, y_minority_upsampled = resample(X_minority, y_minority,
                                                      replace=True,
                                                      n_samples=len(X_majority),
                                                      random_state=42)

# Combine the upsampled data with the majority class
X_train_smote = pd.concat([X_majority, X_minority_upsampled])
y_train_smote = pd.concat([y_majority, y_minority_upsampled])
```

"Why we just apply this method for the training set?". The upsampling was applied exclusively to the training set to prevent data leakage. In real-world scenarios, the data is not always balanced; therefore, it's a good idea when not modifying the test set. Balancing only the training set allows the model to learn from an equal class distribution while still being evaluated on the original, imbalanced test set.

5. **Handling Outliers**

Outliers – values that are much higher or lower than the rest of the data, can affect the performance of many machine learning models. This is especially true for the models that calculate distances between data points, like SVM and KNN, where outliers can cause the model to prioritize unusual values over typical patterns. In this dataset, weather conditions fluctuate naturally, which leads to some extreme values, while real, could mislead the model's understanding of normal weather patterns.

Dropping is the common method for handling outliers; however, removing them entirely could lead to a loss of valuable information since a data point might be an outlier in one feature but still provide important context for other features. Therefore, I chose a method that "softens" the impact of these extreme values while keeping them in the dataset.

Firstly, I identified outliers using the 1.5x IQR rule (Interquartile Range) which flags values that fall below the lower bound (Q1 – 1.5 * IQR) or above the upper bound (Q3 + 1.5 * IQR) as potential outliers. Secondly, instead of removing these values, I applied Winsorization to mitigate the influence of these extreme values. In detail, values below the 5th percentile were set to the 5th percentile, and values above the 95th percentile were capped at the 95th percentile.

By applying Winsorization, the outliers don't overpower the model while their presence and subtle influence remain. This approach keeps the dataset's diversity intact and makes sure the model focuses on more typical patterns.

## 6. Feature Scaling

Feature scaling is important to enhance model performance, especially for algorithms sensitive to feature magnitude, such as SVM and KNN. In this dataset, features vary widely in their range, thus, scaling ensures that no single feature dominates due to its scale. For features containing outliers, **RobustScaler** was applied, as it scales data based on the median and interquartile range, making it less affected by extreme values. For the remaining features, **StandardScaler** was applied to center the data around a mean of zero with a standard deviation of one (see Figure 10). This process brings all features to a similar scale, benefiting the performance and stability of machine learning models.

**Figure 10**

```
from sklearn.preprocessing import RobustScaler, StandardScaler

# List of columns to exclude from scaling
robust_cols = outliers_summary[
    (outliers_summary.outliers_count_after > 0)  # Only columns with outliers
].index.difference(categorical_columns).tolist()

# Filter standard_cols to include only numeric columns that exist in X_train_smote
standard_cols = [col for col in numerical if col in X_train_smote.columns and col not in categorical_columns]

print('Standard columns for scaling:', standard_cols)

# Apply RobustScaler to columns with outliers
robust_scaler = RobustScaler()
X_train_smote[robust_cols] = robust_scaler.fit_transform(X_train_smote[robust_cols])
X_test[robust_cols] = robust_scaler.transform(X_test[robust_cols])

# Apply StandardScaler to the remaining numerical columns
standard_scaler = StandardScaler()
X_train_smote[standard_cols] = standard_scaler.fit_transform(X_train_smote[standard_cols])
X_test[standard_cols] = standard_scaler.transform(X_test[standard_cols])

# Display a summary of the processed data
stats_summary(X_train_smote).round(2)
```

## 7. Feature Selection

One of the necessary steps to improve model accuracy and efficiency is feature selection, which reduces noise and focuses on the most relevant

predictors. In this project, I used **Recursive Feature Elimination (RFE)** with a **RandomForestClassifier** as the estimator. RFE works by removing the least important features recursively, as determined by the classifier until only the most significant predictors remain. This method helps the model concentrate on the features with the strongest impact on predicting the target variable, **RainTomorrow**.

To implement RFE, I configured it to select the top 15 features from the resampled training data. After fitting RFE on the training set, the top features were identified and retained for both the training and testing set (see Figure 11).

**Figure 11**



```
[141] # Import necessary libraries for feature selection
      from sklearn.feature_selection import SelectKBest, chi2, RFE
      from sklearn.ensemble import RandomForestClassifier
      from sklearn.preprocessing import MinMaxScaler

      # Wrapper Method: Recursive Feature Elimination (RFE) with RandomForestClassifier
      selector_rfe = RFE(estimator=RandomForestClassifier(n_estimators=100, random_state=0), n_features_to_select=15, step=1)
      # Fit RFE only on the resampled training data
      X_train_selected = selector_rfe.fit_transform(X_train_smote, y_train_smote)
      selected_features_rfe = X_train.columns[selector_rfe.get_support()]
      print("Top features by RFE with RandomForest:", selected_features_rfe)


      # Transform the test data to retain only the selected features
      X_test_selected = selector_rfe.transform(X_test)

Top features by RFE with RandomForest: Index(['Location', 'MinTemp', 'MaxTemp', 'Rainfall', 'Evaporation', 'Sunshine',
       'WindGustSpeed', 'Humidity9am', 'Humidity3pm', 'Pressure9am',
       'Pressure3pm', 'Cloud9am', 'Cloud3pm', 'Temp9am', 'Temp3pm'],
      dtype='object')
```

## III. Problem Approach

To solve this prediction problem effectively, I implemented a structured approach that balances model evaluation, tuning, and optimization. Firstly, cross-validation on various classifiers to establish baselines and identify potential models. Secondly, refining these models through hyperparameter

tuning and evaluating their performance on multiple metrics to ensure balanced predictions. Finally, I combined the strengths of top models in a soft-voting ensemble, therefore, creating a robust and accurate predictor that will be optimized for real-world data.

1. **Initial Cross-Validation on Base Classifiers**

    To establish a baseline, I first applied a range of classifiers from simple to complex using default parameters in Scikit-Learn. This first step provides a general and broad understanding of which models work well with the dataset. Furthermore, it revealed the strengths and weaknesses of each model. In addition, using **StratifiedKFold** with 5 folds for cross-validation to make sure each fold retained the original class distribution, which is crucial for this imbalanced dataset. Moreover, this approach allowed me to examine each model's generalization capability and detect overfitting or underfitting.

    Moving to the evaluation part, I've chosen **F1-score** as a primary metric due to its balance of Precision and Recall. Additionally, **F1-score** is more suitable for imbalanced data than Accuracy, which is a common measure for classification problems. Models will be evaluated based on the fitting time and the mean F1-scores for both training and validation tests. Finally, we made the shortlist of models that best generalize to new and unseen data due to that evaluation.

2. **Hyperparameter Tuning of Promising Models**

After initial testing, promising classifiers were further refined by tuning their hyperparameters to maximize **F1-score**. While I experimented with traditional methods like **GridSearchCV**, these approaches proved less efficient. It can be computationally expensive as they exhaustively test all parameter combinations. On the other hand, **HalvingRandomSearchCV** has given a better efficiency because it narrows down hyperparameter options by discarding weaker configurations gradually; thus, allowing for a more concentrated exploration of the best candidates.

The combination of this method and **StratifiedKFold** cross-validation provided a reliable and robust search for optimal model parameters, finding a compromise between thoroughness and computational efficiency.

3. **Final Evaluation of Best Models**

Each optimized model was evaluated on the test set using multiple metrics to assess its performance across different criteria. Metrics considered included:

- **Accuracy:** for overall correctness, though it can be misleading with imbalanced data

- **Precision and Recall:** measure the model's ability to correctly identify positive cases.

- **F1-score:** for a balanced assessment of precision and recall

- **Confusion Matrix:** a breakdown of correct and incorrect predictions by class

- **ROC – AUC:** evaluate the model's ability to differentiate across classes at different thresholds.

4. **Ensemble with Soft Voting for Enhanced Performance**

To boost model performance, I applied a **Soft Voting** ensemble approach by combining multiple top-performing models. This approach **averages** the probability predictions of individual models, weighting them according to each model's confidence in its predictions. **HalvingRandomSearchCV** was used to optimize these weights, fine-tuning the ensemble to maximize F1-score. The ensemble was then thoroughly examined using the same measures to ensure it outperformed the individual models. The final model, an ensemble of the top classifiers, was used to predict the target variable in the unseen data. It used the same preprocessing steps as were used for the training set to guarantee consistency and reliability.

## IV. Techniques Used, Summary of Results, and Parameter Settings

1. **Cross-Validation on Base Classifiers for Initial Testing**

To begin with, I evaluated six classifiers using their default settings in Scikit-Learn and analyzing their learning curves to observe performance trends as the training size increased (see Figure 12). This allowed me to identify overfitting or underfitting behaviors. Using **StratifiedKFold** cross-

validation with 5 folds is crucial for evaluating models on an imbalanced
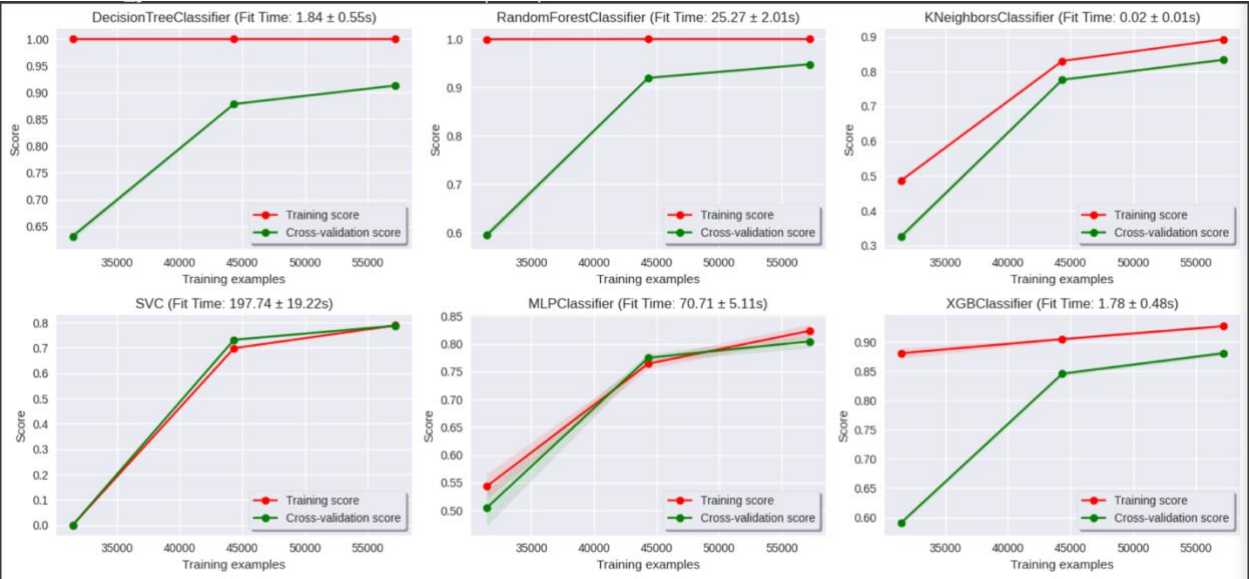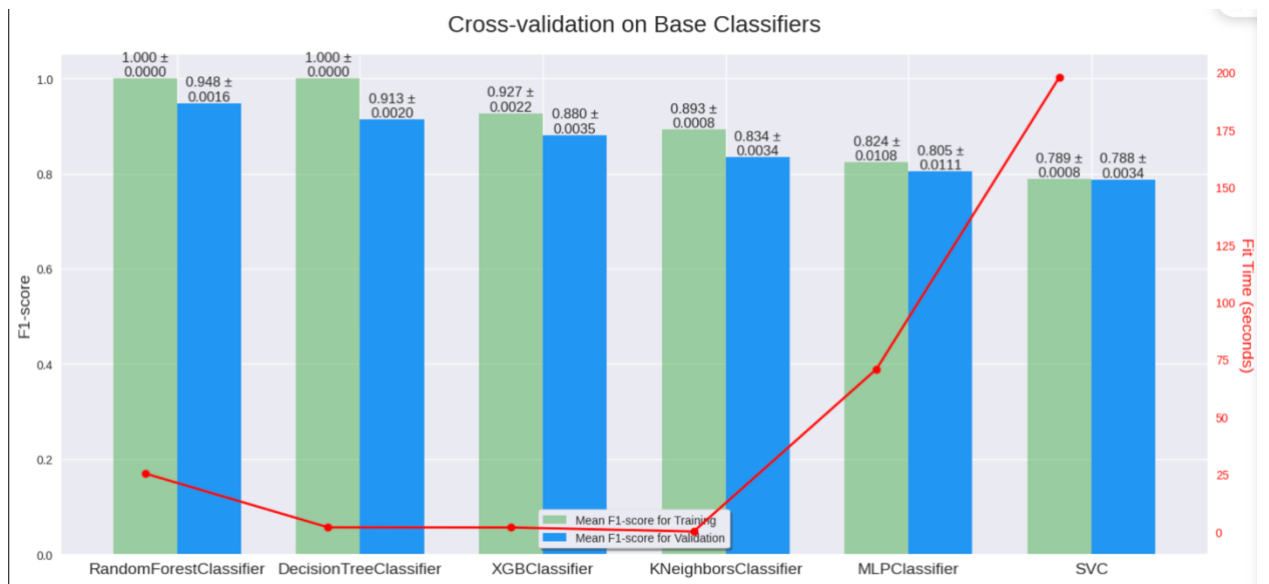
dataset (see Figures 13 & 14).

**Figure 12**



**Figure 13**

| | model_name | mean_train_score | mean_test_score | mean_fit_time | std_train_score | std_test_score | std_fit_time |
|---|---|---|---|---|---|---|---|
| 0 | RandomForestClassifier | 0.999944 | 0.947844 | 25.272283 | 0.000017 | 0.001561 | 2.013824 |
| 1 | DecisionTreeClassifier | 0.999944 | 0.913082 | 1.839152 | 0.000017 | 0.001979 | 0.552657 |
| 2 | XGBClassifier | 0.926703 | 0.880494 | 1.783009 | 0.002156 | 0.003482 | 0.477926 |
| 3 | KNeighborsClassifier | 0.892578 | 0.833969 | 0.019677 | 0.000822 | 0.003436 | 0.005849 |
| 4 | MLPClassifier | 0.823531 | 0.804535 | 70.706239 | 0.010784 | 0.011075 | 5.108980 |
| 5 | SVC | 0.788716 | 0.787512 | 197.740573 | 0.000819 | 0.003416 | 19.223913 |

**Figure 14**

Cross-validation on Base Classifiers

- **Random Forest Classifiers**: This is an ensemble learning method that builds multiple decision trees on random subsets of data and features. Next, it will aggregate their predictions to enhance accuracy and reduce variation. This "wisdom of the crowd" approach improves model robustness by mitigating biases inherent in individual trees. **Random Forest** achieved the best overall performance in testing, with a mean training F1-score of **1.000 ± 0.000** and a mean test F1-score of 0.**948 ± 0.0016**. In addition, the fitting time of **25.27 ± 2.01** seconds was moderate, which balances computational efficiency with high accuracy. As a result, **RandomForest** is a strong candidate for further tuning.

- **Decision Tree Classifier**: **Decision Tree** models are interpretable algorithms that separate data based on feature values and create a tree-like structure of decisions that lead to classification outcomes. Each split aims to maximize homogeneity within branches using metrics like Gini impurity or entropy. The Decision Tree classifier

achieved a perfect mean training F1-score of **1.000 ± 0.0000**, but its mean test F1-score of **0.913 ± 0.0020** indicates slight overfitting, as it performs worse on unseen data. However, with a low fitting time of **1.84 ± 0.55** seconds, Decision Trees are computationally efficient.

- **XGBoost**: a popular gradient-boosting algorithm, builds decision trees sequentially, each tree correcting errors from the previous ones. This iterative process improves accuracy, especially in datasets with complex patterns, by focusing on difficult cases. It achieved a mean training F1-score of **0.927 ± 0.0022** and a mean test F1-score of **0.880 ± 0.0035**, which shows good generalization. XGBoost was also highly efficient, with a fitting time of **1.78 ± 0.48** seconds, making it an attractive option for datasets with complex patterns.

- **K-Nearest Neighbors (KNN):** This is a straightforward, instance-based learning algorithm that classifies new data points based on the majority class of their closest neighbors. It relies on similarity metrics like Euclidean distance, which makes it sensitive to the scale of features. It showed moderate accuracy with a mean training F1-score of **0.893 ± 0.0008** and a mean test F1-score of **0.834 ± 0.0034**. Its fitting time was the shortest among all models at **0.02 ± 0.01** seconds, which makes it highly efficient computationally. However, KNN's simplicity limits its ability to capture complex patterns.

- **Multi-layer Perceptron (MLP):** The MLP classifier, a type of artificial neural network (ANN), consists of layers of interconnected neurons

that learn non-linear relationships in data. Each neuron takes weighted inputs, applies an activation function, and forwards the output to the next layer, thus allowing the network to learn complicated patterns. It achieved a mean training F1-score of **0.824 ± 0.0108** and a mean test F1-score of **0.805 ± 0.0111**. It had one of the highest fitting times at **70.71 ± 5.11** seconds, which indicates a significant computational cost. Although MLP showed potential in learning complex patterns, it may require additional tuning.

- **Support Vector Machine (SVC):** This is a powerful algorithm that aims to determine the best hyperplane for dividing classes by maximizing the margin between data points. This method performs well in high-dimensional areas and can handle nonlinear data using kernel functions. It achieved the lowest performance with a mean training F1-score of **0.789 ± 0.0008** and a mean test F1-score of **0.788 ± 0.0034**. The fitting time was the longest at **197.74 ± 19.22** seconds. SVC's performance suggests it may not be the best fit for this dataset.

The baseline testing revealed that Random Forest provided the strongest overall performance, with high F1-scores on both training and validation sets, as well as manageable fitting times. Additionally, XGBoost and Decision Tree also showed good validation performance, although Decision Tree showed signs of overfitting. Finally, Random Forest, XGBoost, and Decision Tree were considered for further tuning to maximize their performance.

2.  **Hyperparameter Tuning**

 After initial testing, three classifiers were considered as the most promising: RandomForestClassifier, XGBoostClassifier, and MLPClassifier. While the Decision Tree initially showed potential, further trials revealed that MLPClassifier outperformed the Decision Tree with its ability to capture complex and non-linear patterns. Therefore, I focused on tuning Random Forest, XGBoost, and MLP for optimal performance, leveraging a powerful and resource-efficient search strategy: **HalvingRandomSearchCV.**

2.1.  **Why HalvingRandomSearchCV ?**

 HalvingRandomSearchCV is an efficient and iterative approach that begins with a wide range of random configurations and gradually hones in on the most promising options. In each round, less promising configurations are phased out, allowing for the allocation of more resources—such as additional training samples or longer training epochs—to the best-performing ones. This strategy enables a comprehensive yet focused exploration of the hyperparameter space, making the search faster and more effective compared to exhaustive or purely random methods. In contrast to GridSearchCV, which meticulously evaluates every single configuration, and RandomSearchCV, which samples randomly across the entire space without any follow-up, HalvingRandomSearchCV cleverly combines random sampling with

successive halving. This results in a more targeted search that is both computationally efficient and outcome-driven.

## 2.2. Tuning each model

For each classifier, I established customized parameter distributions that reflect the unique characteristics and functionalities of its architecture and behavior:

- **XGBoostClassifier**

  XGBoost's hyperparameter tuning focused on striking a balance between depth and learning rate to effectively capture complex patterns while also preventing overfitting.

  - **n_estimators:** Set to [100, 200] to maintain an optimal number of boosting rounds.

  - **max_depth:** Tested values of [5, 6] aimed at controlling tree complexity, with deeper trees able to capture finer details, albeit with an increased risk of overfitting.

  - **learning_rate:** Values of [0.01, 0.1] were evaluated, allowing XGBoost to take smaller steps during updates, which enhances stability and lowers the risk of overfitting.

- **MLPClassifier**

  For the MLP, tuning concentrated on the network's architecture and solver to optimize learning efficiency and stability:

- hidden_layer_sizes: Configurations like (50,), (100,), and (50, 50) were tested to balance model complexity with computational efficiency.
- activation: Activation functions like relu and tanh were explored to evaluate how well the network learned non-linear relationships.
- solver: Optimizers adam and sgd were considered, with adam generally providing faster convergence for deep networks.

- **RandomForestClassifier**

  For Random Forest, tuning emphasized the number of trees and feature selection strategies to improve robustness and manage class imbalance:

  - n_estimators: Explored values of [100, 200], where more trees can lead to better performance but increase computation.
  - class_weight: Set to 'balanced' to handle the class imbalance effectively.
  - max_features: Tested sqrt and log2, controlling the number of features used in each split to prevent overfitting and enhance generalization.

3. **Evaluation of Best Estimators**

Following the hyperparameter tuning process, a comprehensive evaluation was conducted on the top three models: XGBoostClassifier, MLPClassifier, and RandomForestClassifier. This section presents a comparative analysis of their performance metrics, specifically focusing on Precision, Recall, and F1-scores across various classes, supplemented by relevant visual representations. The evaluation aims to elucidate the strengths and weaknesses of each model, as well as their effectiveness in addressing class imbalance within the dataset.

## 3.1. Classification Reports

The classification reports offer a comprehensive breakdown of Precision, Recall, and F1 scores for each class (0 and 1) (see Figure 15). These metrics emphasize how effectively each model differentiates between the majority class (Class 0) and the minority class (Class 1).

**Figure 15**

```
Classification Report for XGBoost

             precision    recall  f1-score     support
0             0.921954  0.826579  0.871666  3973.000000
1             0.557766  0.757629  0.642514  1147.000000
accuracy      0.811133  0.811133  0.811133     0.811133
macro avg     0.739860  0.792104  0.757090  5120.000000
weighted avg  0.840367  0.811133  0.820330  5120.000000

Best settings: {'n_estimators': 200, 'max_depth': 5, 'learning_rate': 0.1}

========================================================

Classification Report for MLPClassifier

             precision    recall  f1-score     support
0             0.925904  0.773723  0.843000  3973.000000
1             0.500556  0.785527  0.611469  1147.000000
accuracy      0.776367  0.776367  0.776367     0.776367
macro avg     0.713230  0.779625  0.727235  5120.000000
weighted avg  0.830616  0.776367  0.791132  5120.000000

Best settings: {'solver': 'adam', 'hidden_layer_sizes': (100,), 'activation': 'tanh'}

========================================================

Classification Report for RandomForestClassifier

             precision    recall  f1-score  support
0             0.879829  0.934307  0.906250  3973.00
1             0.710322  0.557977  0.625000  1147.00
accuracy      0.850000  0.850000  0.850000     0.85
macro avg     0.795076  0.746142  0.765625  5120.00
weighted avg  0.841856  0.850000  0.843243  5120.00

Best settings: {'n_estimators': 100, 'max_features': 'log2', 'class_weight': 'balanced'}

========================================================
```

- **XGBoostClassifier**

  - **Best parameters**: n_estimators=200, max_depth=5, learning_rate=0.1

  - **Performance summary**: The model achieved a high precision of **0.922** for Class 0, demonstrating a low false positive rate. However, its precision for Class 1 was notably lower at **0.557**, which points to some difficulties in predicting the minority class.

On the bright side, the recall values of **0.826** for Class 0 and **0.757** for Class 1 reflect strong sensitivity across both classes.

- **F1-Scores**: XGBoost's F1-scores are **0.871** (Class 0) and **0.642** (Class 1), showing balanced performance, although there is some drop in Class 1 due to precision limitations.

- **MLPClassifier**

  - **Best parameters**: solver='adam', hidden_layer_sizes=(50,), activation='tanh'

  - **Performance summary**: The MLPClassifier demonstrated impressive precision for Class 0 at **0.925**, while the precision for Class 1 was somewhat lower, sitting at **0.500**. However, the recall for Class 1 outperformed that of XGBoost, reaching **0.785**. This indicates that the MLP effectively identified positive instances, although it did result in a slightly elevated false positive rate.

  - **F1-Scores**: MLPClassifier's F1-scores were **0.843** (Class 0) and **0.611** (Class 1), showing good balance but slightly lower performance for the minority class.

- **RandomForestClassifier**

  - **Best parameters**: n_estimators=100, max_features='sqrt', class_weight='balanced'

  - **Performance summary**: The Random Forest model showed impressive precision for both classes, achieving a score of
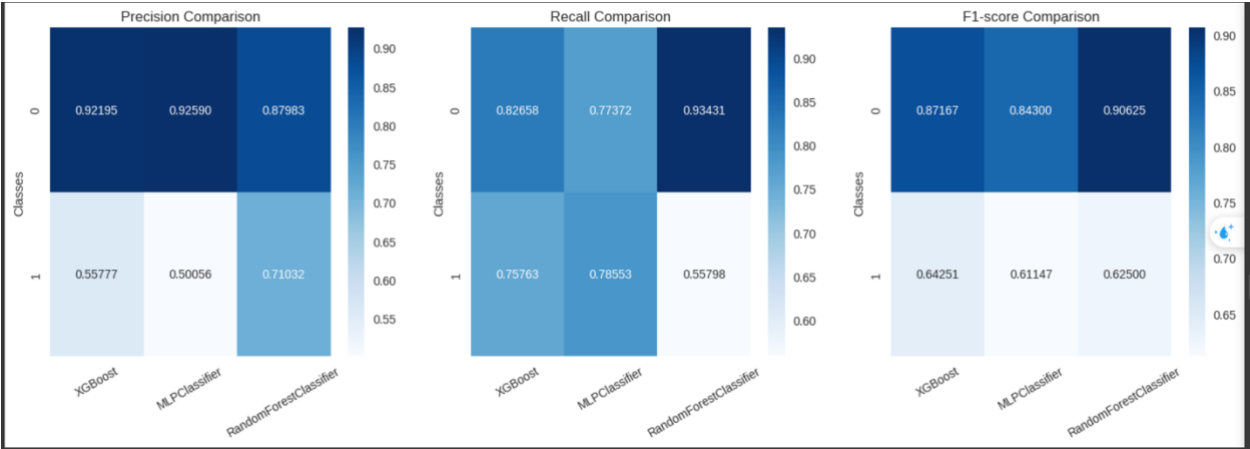
**0.879** for Class 0 and **0.710** for Class 1. Additionally, the recall was notably high for Class 0 at **0.934**, while it was somewhat lower for Class 1 at **0.557**. This indicates a well-balanced overall performance.

- **F1-Scores**: Random Forest achieved F1-scores of **0.906** (Class 0) and **0.625** (Class 1), indicating robust handling of the majority class and reliable generalization, though with a slight trade-off in sensitivity for Class 1.

### 3.2. Precision, Recall, and F1-Score Comparisons

The comparisons of Precision, Recall, and F1-score offer a visual perspective on the performance of each model for specific classes (see Figure 16).
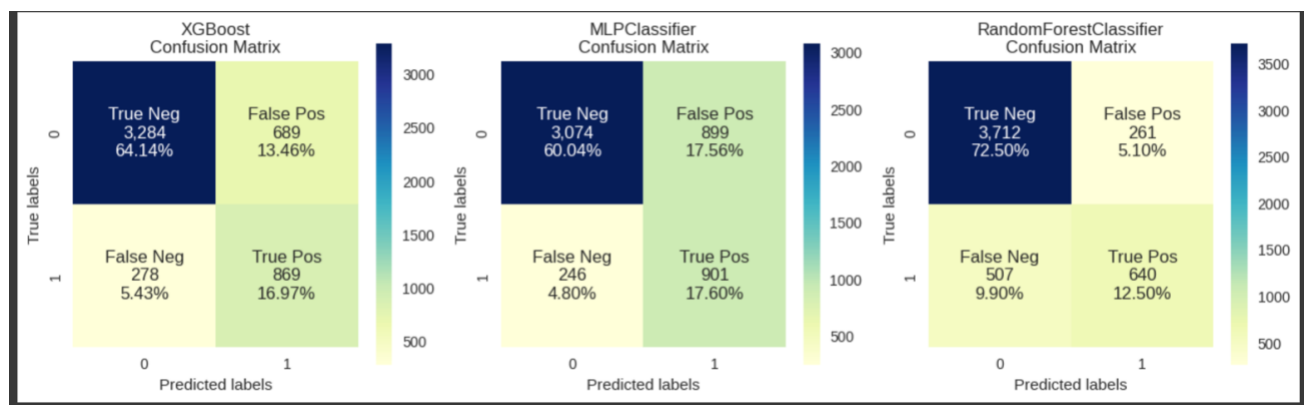
**Figure 16**



- All models demonstrated strong precision for Class 0, while Random Forest stood out with the highest precision for Class 1 at **0.71**. This indicates its effective management of class imbalance.

- **Recall**: The MLPClassifier achieved the highest recall for Class 1 at **0.785**, showcasing its superior sensitivity to positive cases. In contrast, both XGBoost and Random Forest prioritized maintaining strong recall for Class 0.

- **F1-Score**: Both XGBoost and Random Forest delivered impressive F1-scores for Class 0, with XGBoost standing out for its exceptional performance in Class 1 as well. This demonstrates a well-rounded effectiveness across both classes.

### 3.3. Confusion Matrices

The confusion matrices illustrate each model's distribution of true positives, false positives, true negatives, and false negatives (see Figure 17):
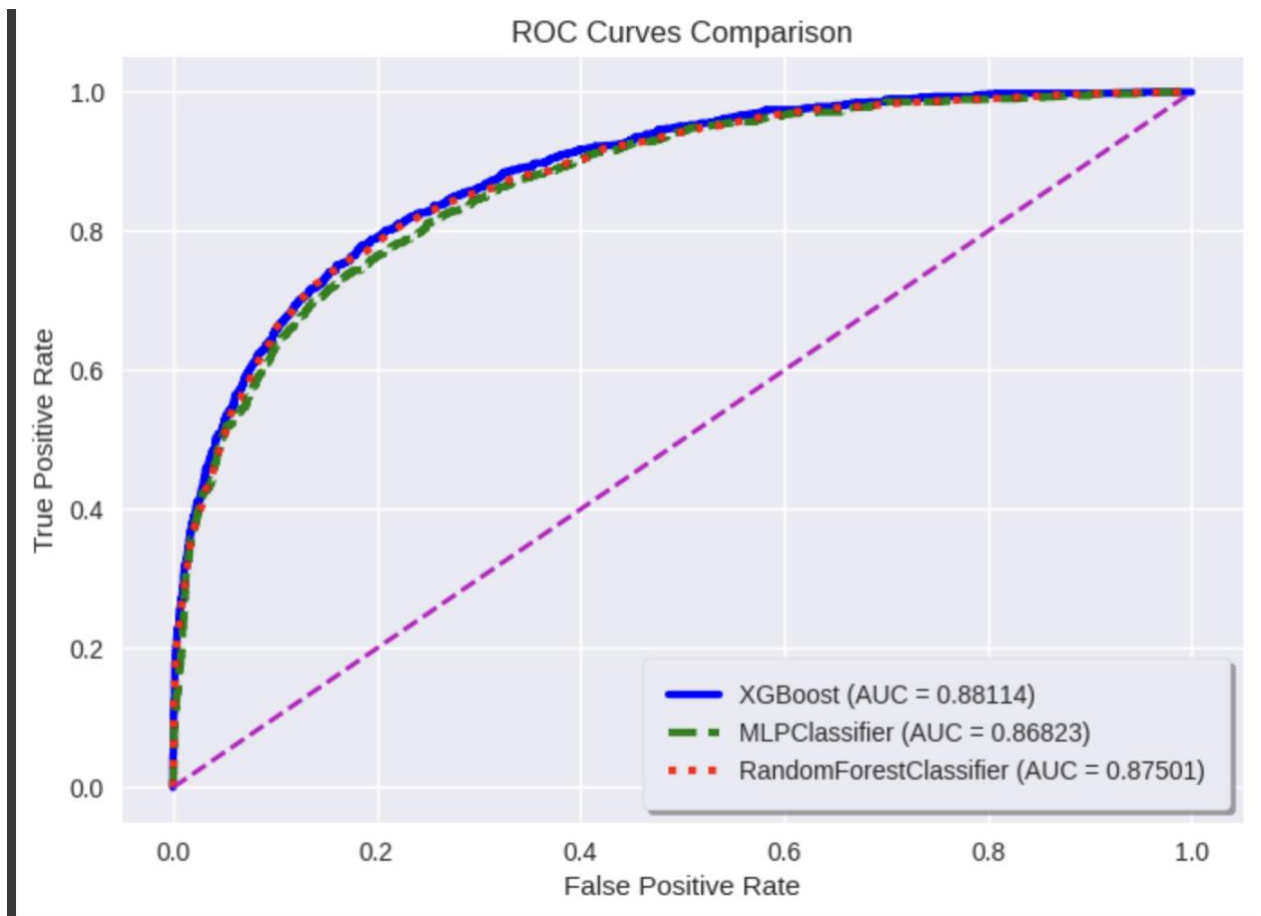
**Figure 17**



- **XGBoostClassifier**: XGBoost performed well in minimizing false positives and negatives for Class 0. However, it struggled a bit with Class 1, showing a higher rate of false negatives. This suggests that it has some challenges when it comes to accurately identifying instances from the minority class.

- o **MLPClassifier**: The MLP model exhibited a moderate number of false positives for Class 1. This led to a higher recall rate, which means it was good at identifying instances of the minority class. However, this came at the cost of slightly lower precision, suggesting that MLP tends to focus more on sensitivity—the ability to correctly identify positive cases—rather than being precise in its predictions.

- o **RandomForestClassifier**: The Random Forest model performed quite well, especially when it came to Class 0, where it had the least false positives. For Class 1, it faced some moderate false negatives. Overall, this distribution of errors shows that the Random Forest's use of the class_weight='balanced' parameter did a good job of handling the class imbalance. As a result, it maintained a more consistent performance across both classes.

### 3.4. ROC Curves

The ROC curve comparison provides a summary of each model's discriminatory power, as indicated by the Area Under the Curve (AUC) (see Figure 18):

**Figure 18**

ROC Curves Comparison

- XGBoost (AUC = 0.88114)
- MLPClassifier (AUC = 0.86823)
- RandomForestClassifier (AUC = 0.87501)

- o **XGBoostClassifier:** The AUC score of **0.881** shows that our model does a great job of telling the difference between the classes. It's close to the perfect score, which is a promising sign of its effectiveness.

- o **MLPClassifier**: The AUC was **0.868**, indicating that the model showed strong performance in class separation. While the MLP performed well, it was a bit behind XGBoost in terms of effectiveness.

- o **RandomForestClassifier:** The AUC score of **0.875** for the Random Forest model highlights its strong ability to distinguish

between classes. It performs well across all categories. This indicates a reliable balance in its predictive power.

4. **Soft Voting**

To improve model performance, I utilized a **VotingClassifier** that combines the strengths of three finely tuned models: XGBoost, MLPClassifier, and RandomForestClassifier. Adopting a soft voting strategy, this ensemble method averages the probability predictions from each model, taking into account their relative contributions. The aim was to identify a configuration that maximizes the F1-score while leveraging the distinctive strengths of each model.

**4.1. Configuration and Weight Tuning**

- **Voting Type:** The **VotingClassifier** was configured to 'soft', enabling the model to utilize weighted probabilities for its class predictions. This approach improves its adaptability, allowing it to effectively capture the nuances of both classes.

- **Weight Selection**: To figure out the best approach for our final predictions, I tried out different weight combinations. I utilized **HalvingRandomSearchCV** to experiment with sets like **[4, 1, 2],** aiming to find the perfect balance among the models (see Figure 19).

**Figure 19**

```
[155] from sklearn.ensemble import VotingClassifier

      # Setup Voting Classifier with optimized estimators
      optimized_estimators = [(name, searcher.best_estimator_) for name, searcher in searchers.items()]
      # Get the number of estimators dynamically
      soft_voting_clf = VotingClassifier(estimators=optimized_estimators, voting='soft', n_jobs=-1)

      param_grid = {
          'weights': [
              [2, 1, 2],  # Slightly prioritize XGBoost and RandomForest
              [3, 1, 2],  # Give XGBoost a higher weight than the others
              [3, 2, 2],  # Balance between XGBoost and RandomForest with some weight for MLP
              [4, 1, 2],  # Higher priority for XGBoost
              [3, 1, 3]   # Balance with equal weight to XGBoost and RandomForest, with MLP slightly less
          ]
      }
      voting_searcher = HalvingRandomSearchCV(
          soft_voting_clf, param_grid, resource='n_samples', factor=3, cv=5,
          scoring='f1', return_train_score=True, random_state=42, verbose=3
      ).fit(X_train_selected, y_train_smote)
      best_voting_clf = voting_searcher.best_estimator_
```

- **XGBoost** was assigned a weight of 4, given its strong F1 score and high precision for the majority class.

- **RandomForestClassifier** received a weight of 2, reflecting its balanced performance across both classes.

- **MLPClassifier** was assigned a weight of 1, as it demonstrated high recall for the minority class, making it a valuable addition for class balance.

## 4.2. Insights

From our cross-validation on the best estimators, the F1-scores show that the VotingClassifier has led to some significant improvements (see Figure 20). In fact, the final model reached a validation F1-score of 0.721. This approach really takes advantage of XGBoost's sharp precision and the balanced performance of RandomForest, which makes a noticeable difference in our results.

The class probabilities for this sample reveal how important the assigned weights are in shaping the final prediction. You can clearly see that XGBoost's higher weight gives it a significant edge, which helps its precision align closely with what the overall ensemble predicts. Meanwhile, both MLP and RandomForest play a crucial role by enhancing sensitivity to the minority class, helping to create a more balanced decision overall (see Figure 21).
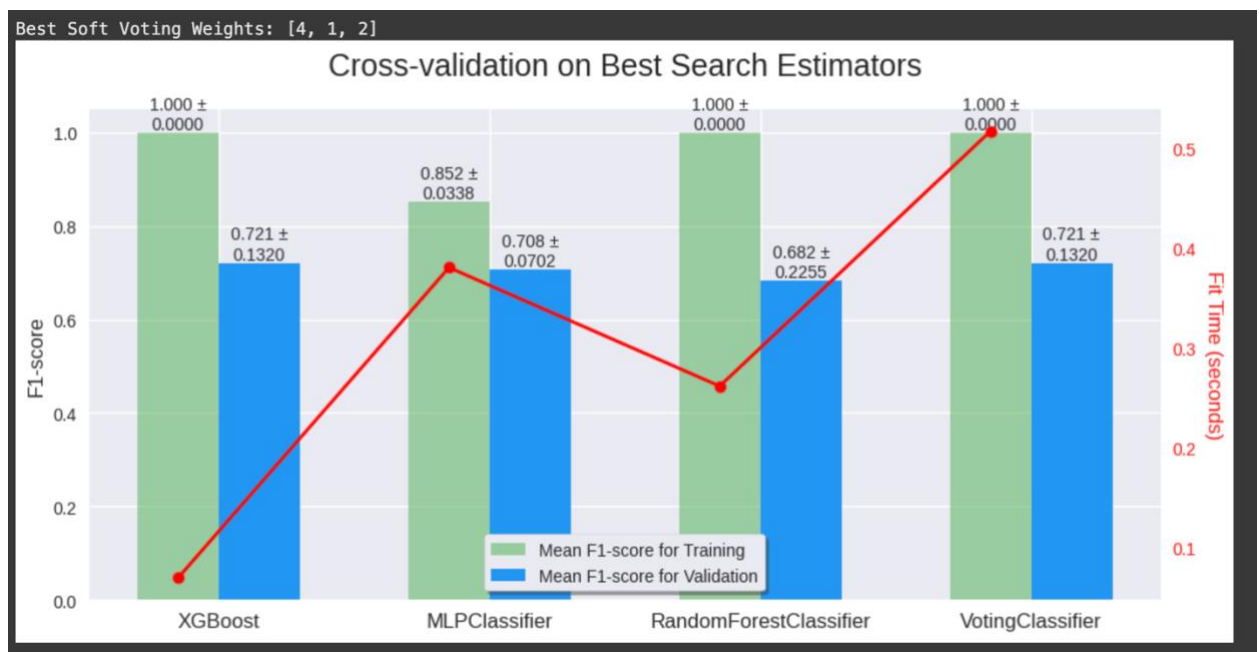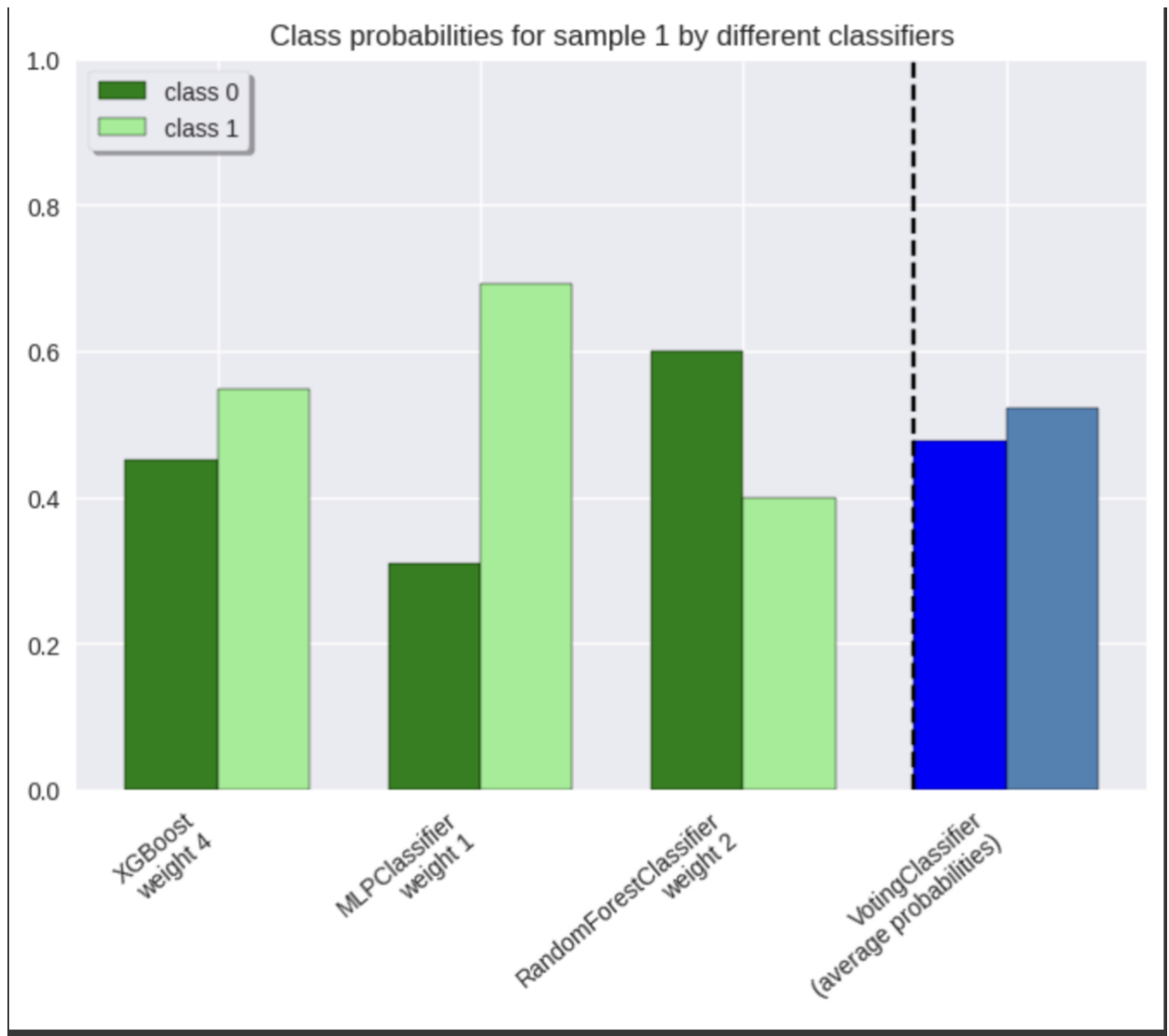
**Figure 20**



**Figure 21**

Class probabilities for sample 1 by different classifiers

### 4.3. Performance Evaluation

The classification report, confusion matrix, and ROC curve provide a comprehensive summary of the **VotingClassifier's** performance (see Figure 22):
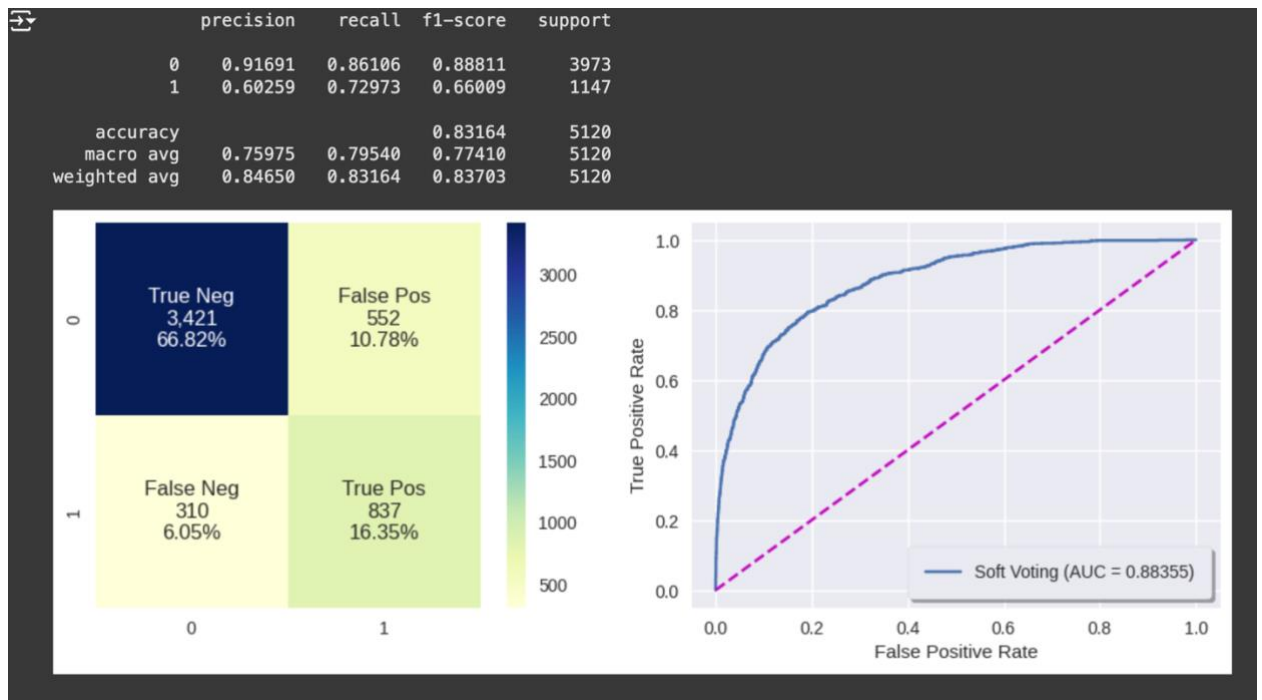
o **Precision:** The **VotingClassifier** really performed well, particularly with Class 0, as it managed to lower the number of false positives and boost the overall accuracy. It also improved Class 1 precision compared to the individual models, which

means it did a better job of correctly identifying the minority class while also reducing false positives.

- o **Recall:** The ensemble's recall for Class 1 has seen a positive change, leading to fewer missed positive instances. This means we're getting better at identifying the cases that matter most. The **VotingClassifier** does a great job of balancing precision and recall, which is essential when dealing with class imbalance.

- o **F1-Score:** The **VotingClassifier** achieved a weighted F1-score of **0.88811**, which shows that it performed better than the individual models. This result emphasizes how well it balanced the handling of both classes, leading to more effective overall results.

- o **Confusion Matrix:** The confusion matrix shows that the **VotingClassifier** has enhanced its classification accuracy. It has successfully reduced the number of false negatives for Class 1, which means it's more reliable in real-world situations where both classes are important. This improvement is crucial for ensuring that we can trust the classifier's results when it matters most.

- o **ROC Curve and AUC:** The ROC curve, which shows an AUC of 0.88355, demonstrates that the **VotingClassifier** does a great

job of distinguishing between the positive and negative classes, effectively highlighting its strong discrimination ability.

**Figure 22**



```
              precision   recall  f1-score   support

         0     0.91691   0.86106   0.88811      3973
         1     0.60259   0.72973   0.66009      1147

  accuracy                         0.83164      5120
 macro avg     0.75975   0.79540   0.77410      5120
weighted avg   0.84650   0.83164   0.83703      5120
```

5. **Prediction on Unknown Dataset and Submission on Kaggle**

In the final phase, we employed the **SoftVotingClassifier** to make predictions on the unknown dataset. To maintain consistency, we ensured that this dataset underwent the same preprocessing steps as our training data. We carefully imputed any missing values, encoded categorical variables, managed outliers through winsorization, and scaled and selected features based on our previous analysis. This meticulous approach allowed the model to interpret the unknown dataset just as it did with the training data, paving the way for reliable and accurate predictions. Finally, we formatted the predictions and saved them into a CSV file,

ensuring it was ready for submission while also aligning perfectly with the required output format.

The predictions submitted reached a public score of 0.81126 on Kaggle, showcasing the model's strong performance and ability to generalize effectively to unseen data (see Figure 23).



| ✓ | **predictions (78).csv**<br>Complete · 2h ago | **0.81126** | ☐ |

## V.    Justification for the Selected Classifier

The **SoftVotingClassifier** stands out by combining the strengths of each model. Its performance metrics were impressive across the board, showcasing its knack for generalizing effectively to new, unseen data. It also does a great job of tackling challenges like class imbalance and overfitting; therefore, it's a robust choice for various predictive tasks.

- **F1-Score:** Achieved the highest F1-score among all tested models, indicating balanced performance in both precision and recall.

- **Precision and Recall:** Excelled in precision and recall, especially in handling the minority class, which is essential for imbalanced datasets.

- **ROC-AUC**: Achieved an AUC of 0.88964, demonstrating significant ability to distinguish between the classes effectively.

- **Mitigation of Overfitting:** The ensemble method helped to mitigate overfitting by averaging the predictions, effectively smoothing out the noise and bias present in individual models.

- **Error Reduction**: Showcased exceptional management of false positives and false negatives, leading to a decrease in classification errors and resulting in more dependable predictions.

By combining the unique strengths of XGBoost's accuracy, the intricate pattern recognition of MLP, and the reliability of RandomForest, the **VotingClassifier** significantly enhances overall performance. This strategy embodies the concept of the "wisdom of the crowd," where the shared decision-making of multiple models surpasses the capabilities of any single model. As a result, the **SoftVotingClassifier** stands out as a dependable and well-rounded option for practical applications.