# Machine Learning

## ASSIGNMENT 2

---

## Members:

Mark Magdy : 2205040                    Arsany Osama : 2205122

Rana Ashraf : 2205019

## Overview

The objective of this project is to explore different classification models and evaluate their performance using the MAGIC gamma telescope dataset. The dataset consists of two classes: gamma and hadron, with 12,332 gamma events and 6,688 hadron events.

## Goals

1. Apply various classification models, including Decision Trees, Naïve Bayes, Random Forests, and AdaBoost, to classify gamma and hadron events.
2. Perform parameter tuning for each model to optimize their performance.
3. Evaluate the models using accuracy, precision, recall, F1-score, and confusion matrix.
4. Compare the performance of different models and identify the most effective one for this classification task.

## Data Preprocessing

1) Set Column Names :

Assigning meaningful names to the features and label

```
# Set column names
data.columns = ['f1', 'f2', 'f3', 'f4', 'f5', 'f6', 'f7', 'f8', 'f9', 'f10', 'class']
#print("Data after renaming columns:\n", data)
```

## 2) Handle Missing Values :

Dropping rows containing missing values such as NaN

```
# Handle missing values (if any)
data.dropna(inplace=True)
```

## 3) Scaling :

Standard scaling ensures that all features have zero mean and unit variance, making them comparable and preventing any single feature from dominating the model.

```
data[['f1', 'f2', 'f3', 'f4', 'f5', 'f6', 'f7', 'f8', 'f9', 'f10']] = StandardScaler().fit_transform(
    data[['f1', 'f2', 'f3', 'f4', 'f5', 'f6', 'f7', 'f8', 'f9', 'f10']])
#print("Data after scaling:\n", data)
```

## 4) Data Balancing :

Since the dataset is class-imbalanced, we randomly take a subset of gamma ignoring extra class readings to make both classes equal in size (6688).

```
Balance the dataset, randomly put aside the extra readings
for the gamma "g" class to make both classes equal in size
'''
gamma_data_equal = gamma_data.sample(n=len(hadron_data), random_state=42)
balanced_data = pd.concat([gamma_data_equal, hadron_data])
```

## 5) Shuffling :

Ensures that the order of the samples does not influence the training of the model. It helps in preventing any inherent patterns in the dataset from affecting the model's learning process

```
'''
#Shuffle the combined dataset to ensure randomness
'''
balanced_data = balanced_data.sample(frac=1, random_state=42).reset_index(drop=True)

'''
```

## 6) Data Split :

We split the dataset randomly into training and testing sets, with 70% for training and 30% for testing..

```
'''
Split the dataset into training 70% and testing sets 30%
"random_state" makes the code reproducible
'''
train_data, test_data = train_test_split(balanced_data, test_size=0.3, random_state=42)
```

## 7) Prepare Training Set :

Extract the training features by dropping 'class' column and extract target variable

```
'''
Prepare  the  training data  and  Labels
'''
X_train = train_data.drop(columns=['class'])
y_train = train_data['class']
```

## 8) Prepare Testing Set :

We repeat the same process of extracting features but from the test dataset.

```
'''
Prepare the testing data and labels
'''
X_test = test_data.drop(columns=['class'])
y_test = test_data['class']

return X_train, y_train, X_test, y_test
```

# Classification Models

## I. Decision Trees

Create instance of the (DescisionTreeClassifier) class. The model is trained using the training data (x_train and y_train) by calling the (fit) method.

```
model = tree.DecisionTreeClassifier().fit(x_train , y_train)
```

The trained model is used to make predictions on the testing data (x_test).

```
predicted_y = model.predict(x_test)
```

Accuracy is calculated by comparing the predicted labels with the actual labels using the (accuracy_score) function.

```
accuracy = accuracy_score(y_test, predicted_y)
print("=" * 40)
print("Accuracy of the Decision Tree Classifier:", accuracy , " --> " , accuracy * 100 , "%")
```

```
========================================
Accuracy of the Decision Tree Classifier: 0.785447296287067  -->  78.5447296287067 %
========================================
```

The confusion matrix is computed using the (confusion_matrix) function to evaluate the performance of the model.

```
print("Confusion Matrix")
confusion_mat = confusion_matrix(y_test, predicted_y)
print(confusion_mat)
```

The confusion matrix has four important terms:

True Positives (TP): The cases where the model predicted the class correctly, and the actual class is positive (in this case, "g") = 1591

True Negatives (TN): The cases where the model predicted the class correctly, and the actual class is negative (in this case, "h") = 1561

False Positives (FP): The cases where the model predicted the class incorrectly (predicted positive, but actually negative) = 385

False Negatives (FN): The cases where the model predicted the class incorrectly (predicted negative, but actually positive) = 476

```
Confusion Matrix
[[1591  385]
 [ 476 1561]]
```

```python
print("Report")
report = classification_report(y_test, predicted_y)
print(report)
```

Classification report, including precision, recall, and F1-score, is generated using the (classification_report) function.

```
Report
              precision    recall  f1-score   support

           g       0.77      0.81      0.79      1976
           h       0.80      0.77      0.78      2037

    accuracy                           0.79      4013
   macro avg       0.79      0.79      0.79      4013
weighted avg       0.79      0.79      0.79      4013
```

## II.    AdaBoost

Initialize an AdaBoost classifier using (AdaBoostClassifier) from scikit-learn.

```
# AdaBoost classifier
adaboost = AdaBoostClassifier()
```

The AdaBoost classifier is fitted to the training data using (fit)

```
# Fit on training set using features and label
adaboost.fit(X_train, y_train)
```

Then, it predicts the labels for the testing data using (predict)

```
# Predict labels of testing features
adaboost_predictions = adaboost.predict(X_test)
```

Perform cross-validation to evaluate the performance of the AdaBoost classifier using (cross_val_score).

```
AdaBoost Cross-Validation Accuracy: 80.9034%
----------------------------------------
```

Then, it predicts the labels for the testing data using (predict)

```
# Predict labels of testing features
adaboost_predictions = adaboost.predict(X_test)
```

Calculate the accuracy of the AdaBoost classifier on the testing data using (accuracy_score) and print it

```
# Accuracy on testing features and label
accuracy = accuracy_score(y_test, adaboost_predictions)
print("AdaBoost Testing Accuracy: {:.4f}%".format(accuracy * 100))
```

```
========================================
AdaBoost Testing Accuracy: 80.3887%
========================================
```

The confusion matrix is calculated using (confusion_matrix) to evaluate the performance of the classifier in terms of true positive, true negative, false positive, and false negative predictions.

```
# Calculate confusion matrix
adaboost_conf_matrix = confusion_matrix(y_test, adaboost_predictions)
print("AdaBoost Confusion Matrix:")
print(adaboost_conf_matrix)
```

True Positive (TP): The number of instances where the model correctly predicts the positive class (gamma events in this case) = 1608

True Negative (TN): The number of instances where the model correctly predicts the negative class (hadron events in this case) = 1618

False Positive (FP): The number of instances where the model incorrectly predicts the positive class (incorrectly classified as gamma events when they are actually hadron events) = 368

False Negative (FN): The number of instances where the model incorrectly predicts the negative class (incorrectly classified as hadron events when they are actually gamma events) = 419

```
AdaBoost Confusion Matrix:
[[1608  368]
 [ 419 1618]]
```

Classification report of Default AdaBoost is generated using the (classification_report) function.

```
========================================
AdaBoost Classification Report:
              precision    recall  f1-score   support

           g       0.79      0.81      0.80      1976
           h       0.81      0.79      0.80      2037

    accuracy                           0.80      4013
   macro avg       0.80      0.80      0.80      4013
weighted avg       0.80      0.80      0.80      4013
```

Model Parameter Tuning to find the optimal combination of hyperparameters that maximizes the model's performance

Define parameter grid that specifies the hyperparameters to be tuned and the range of values to explore. In this case, we're tuning the (n_estimators) parameter for the AdaBoost classifier, which represents the number of weak learners (decision trees) in the ensemble.

By defining a parameter grid, we create a set of values that the grid search algorithm will explore to find the best combination of hyperparameters.

```
# Define the parameter grid for AdaBoost
param_grid = {'n_estimators': [50, 100, 200]}
```

Grid search is a technique that performs search over a specified parameter grid to find the optimal combination of hyperparameters using CROSS-VALIDATION, which helps prevent overfitting

We use 5-fold cross-validation (cv=5), meaning the dataset is split into 5 equal parts, and the model is trained and evaluated 5 times.

Once the parameter grid and cross-validation settings are defined, Fit Grid Search (GridSearchCV) to Training Data (X_train and y_train). This step iterates through all combinations of hyperparameters, trains each model on the training data, and evaluates its performance using cross-validation.

```
# Perform grid search with cross-validation
grid_search = GridSearchCV(AdaBoostClassifier(), param_grid, cv=5)
grid_search.fit(X_train, y_train)
```

After fitting the grid search object, we extract the best estimator using the (best_estimator_) attribute that represents the AdaBoost classifier with the optimal combination of hyperparameters

```
# Get the best classifier with tuned parameters
best_adaboost = grid_search.best_estimator_
```

we use the best estimator to predict the labels for the testing data (X_test)

```
# Predict the labels for the testing data using the best classifier
best_adaboost_predictions = best_adaboost.predict(X_test)
```

We use the (accuracy_score) function to calculate the accuracy of the TUNED AdaBoost classifier on the testing data (X_test and y_test).

```
# Calculate accuracy on testing data using the best classifier
accuracy_tuned = accuracy_score(y_test, best_adaboost_predictions)
print("Tuned AdaBoost Testing Accuracy: {:.4f}%".format(accuracy_tuned * 100))
```

The tuned AdaBoost model achieved a testing accuracy of approximately 81.86%, while the AdaBoost model without parameter tuning achieved a testing accuracy of around 80.39%. This indicates that the performance of the tuned model improved slightly compared to the default model.

```
========================================
Tuned AdaBoost Testing Accuracy: 81.8590%
========================================
```

Compute the confusion matrix for the TUNED matrix

```
# Calculate confusion matrix using the best classifier
confusion_tuned = confusion_matrix(y_test, best_adaboost_predictions)
print("Tuned AdaBoost Confusion Matrix:")
print(confusion_tuned)
```

```
==========================================
Tuned AdaBoost Confusion Matrix:
[[1644  332]
 [ 396 1641]]
==========================================
```

Notice the following differences between Tuned and Default Confusion Matrix :

1) True Positives (TP):

   The tuned AdaBoost model correctly classified 36 additional positive instances compared to the default AdaBoost model.

This indicates that the tuned model improved in identifying instances belonging to the positive class.

2) False Positives (FP):

The tuned AdaBoost model reduced the number of false positive predictions by 36 compared to the default AdaBoost model. This suggests that the tuned model produced fewer instances where negative instances were incorrectly classified as positive.

3) False Negatives (FN):

The tuned AdaBoost model decreased the number of false negative predictions by 23 compared to the default AdaBoost model. This indicates an improvement in correctly identifying positive instances that were previously misclassified as negative.

4) True Negatives (TN):

The number of true negative predictions remained relatively similar between the two models, with the tuned AdaBoost model correctly classifying three additional negative instances compared to the default AdaBoost model.

Overall, the differences in the confusion matrices demonstrate that the tuned AdaBoost model achieved improvements in sensitivity (Identifying positive instances) and specificity (Identifying negative instances) compared to the default AdaBoost model.

Compute precision, recall, and F1-score for each class based on the predictions of the TUNED AdaBoost classifier on the testing data.

```python
# Calculate precision, recall, and F1-score using the best classifier
report_tuned = classification_report(y_test, best_adaboost_predictions)
print("Tuned AdaBoost Classification Report:")
print(report_tuned)
```

Notice the following differences between Tuned and Default Classification Report:

1) Precision:

   Precision for class 'g' is 0.81, indicating a slight improvement over the default AdaBoost model of 0.79

   Precision for class 'h' is 0.83, indicating a slight improvement over the default AdaBoost model of 0.81

2) Recall:

   Recall for class 'g' is 0.83, indicating an improvement over the default AdaBoost model of 0.81

   Recall for class 'h' is 0.81, indicating a slight increase compared to the default AdaBoost model of 0.79

3) F1-score:

   The F1-score for class 'g' is 0.82, indicating an improvement over the default AdaBoost model of 0.80

   The F1-score for class 'h' is 0.82, indicating an improvement over the default AdaBoost model of 0.80

4) Accuracy:

   The overall accuracy of the tuned AdaBoost model is 0.82, indicating that 82% of all predictions were correct, which is higher than the default AdaBoost model of 80%

## III. Random Forests

The Random Forest classifier is initialized without specifying any hyperparameters, which means it will use default values.

```python
# Random Forest classifier
random_forest = RandomForestClassifier()
```

Cross-validation is performed using the `cross_val_score()` function. It evaluates the performance of the Random Forest classifier on the training data.

```python
# Cross Validation to evaluate the performance
random_forest_cv = cross_val_score(random_forest, X_train, y_train, cv=5)
print("Random Forest Cross-Validation Accuracy: {:.4f}%".format(random_forest_cv.mean()
```

Parameters of `cross_val_score()` :

1. `random_forest`: This parameter represents the Random Forest classifier object. It is the model that will be evaluated using cross-validation.

2. `X_train`: This parameter denotes the feature data of the training set.

3. `y_train`: This parameter represents the target labels of the training set.

4. `cv=5`: This parameter specifies the number of folds for cross-validation. In this case, `cv=5` indicates that a 5-fold cross-validation will be performed. It means that the training data will be divided into 5 subsets or folds, and the Random Forest classifier will be trained and evaluated 5 times, each time using a different fold as the validation set. Here's how the data is divided in each iteration of the 5-fold cross-validation:

1. Iteration 1:
   - Training set: Folds 2, 3, 4, 5
   - Testing set: Fold 1
2. Iteration 2:
   - Training set: Folds 1, 3, 4, 5

- Testing set: Fold 2
3. Iteration 3:
    - Training set: Folds 1, 2, 4, 5
    - Testing set: Fold 3
4. Iteration 4:
    - Training set: Folds 1, 2, 3, 5
    - Testing set: Fold 4
5. Iteration 5:
    - Training set: Folds 1, 2, 3, 4
    - Testing set: Fold 5

Output:
```
Random Forest Cross-Validation Accuracy: 85.7204%
```

The Random Forest classifier is fitted to the training data using the `fit()` method, trained model is then used to predict the labels of the testing data using the `predict()` method,

The accuracy of the Random Forest classifier is calculated by comparing the predicted labels with the actual labels of the testing data using the `accuracy_score()` function.

```python
# Fit on training set using features and label
random_forest.fit(X_train, y_train)

# Predict labels of testing features
random_forest_predictions = random_forest.predict(X_test)

# Accuracy on testing features and label
accuracy = accuracy_score(y_test, random_forest_predictions)
print("Random Forest Testing Accuracy: {:.4f}%".format(accuracy * 100))
```

Output :
```
Random Forest Testing Accuracy: 85.0984%
```

The confusion matrix is calculated to evaluate the performance of the Random Forest classifier in terms of true positives, true negatives, false positives, and false negatives using the `confusion_matrix()` function.

Random Forest Confusion Matrix:

[[1758  218]

 [ 380 1657]]

Overall Report:
Random Forest Classification Report:

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| g | 0.82 | 0.89 | 0.85 | 1976 |
| h | 0.88 | 0.81 | 0.85 | 2037 |
| accuracy |  |  | 0.85 | 4013 |
| macro avg | 0.85 | 0.85 | 0.85 | 4013 |
| weighted avg | 0.85 | 0.85 | 0.85 | 4013 |

Tuning:

In the code below, we used a technique called grid search that is used for hyperparameter tuning. Here's an explanation of the code:

1. `param_grid`: It is a dictionary that defines the combinations of hyperparameter values to be searched. In this case, the hyperparameters being tuned are `'n_estimators'` and `'max_depth'`.

  - `'n_estimators'` refers to the number of decision trees in the Random Forest ensemble.

  - `'max_depth'` specifies the maximum depth of each decision tree.

  The possible values for `'n_estimators'` are 50, 100, and 200, while the possible values for `'max_depth'` are `None`, 10, and 20.

2. `GridSearchCV()`: It is a class provided by scikit-learn that performs an exhaustive search over the specified hyperparameter values using cross-validation.

  - The first argument (`RandomForestClassifier()`) represents the base model to be tuned, which is the Random Forest classifier in this case.

- The second argument (`param_grid`) is the dictionary of hyperparameters and their possible values.

- The `cv=5` parameter indicates that 5-fold cross-validation will be used during the grid search.

3. `grid_search.fit(X_train, y_train)`: This line fits the grid search object to the training data (`X_train` and `y_train`) to find the best hyperparameter combination.

Then get the best estimator by `grid_search.best_estimator_` after that predict from features in X_test

```python
# Define the parameter grid for Random Forest
param_grid = {'n_estimators': [50, 100, 200], 'max_depth': [None, 10, 20]}

# Perform grid search with cross-validation
grid_search = GridSearchCV(RandomForestClassifier(), param_grid, cv=5)
grid_search.fit(X_train, y_train)

# Get the best classifier with tuned parameters
best_random_forest = grid_search.best_estimator_

# Predict the labels for the testing data using the best classifier
best_random_forest_predictions = best_random_forest.predict(X_test)
```

Accuracy:

Tuned Random Forest Testing Accuracy: 85.1981%

Confusion Matrix:

Tuned Random Forest Confusion Matrix:

[[1751  225]

 [ 369 1668]]

Report:

Tuned Random Forest Classification Report:

| | precision | recall | f1-score | support |
|---|---|---|---|---|
| 9 | 0.83 | 0.89 | 0.85 | 1976 |

| | | | | |
|---|---|---|---|---|
| h | 0.88 | 0.82 | 0.85 | 2037 |
| | | | | |
| accuracy | | | 0.85 | 4013 |
| macro avg | 0.85 | 0.85 | 0.85 | 4013 |
| weighted avg | 0.85 | 0.85 | 0.85 | 4013 |

## IV.  Naive Bayes

The Naive Bayes model is trained on the training set using the (fit) function, which takes features and labels of the training set as input.

```python
nb = GaussianNB()
nb.fit(X_train, y_train)
```

The function (predict) takes the testing features (X_test) as input and predicts the corresponding labels. We obtain the predicted labels for the testing features, which can then be compared with the actual labels (y_test) to evaluate the performance.

```python
nb_predictions = nb.predict(X_test) #predict labels of testing features
```

After training the model, it is tested on the testing set to assess its performance in real-world scenarios. The accuracy of the model on the testing data is calculated using the (accuracy_score) function.

```python
print("="*40)
#Accuracy on testing features and label
accuracy = accuracy_score(y_test, nb_predictions)
print("Naive Bayes Testing Accuracy: {:.4f}%".format(accuracy * 100))
print("="*40)
```

```
Naive Bayes Testing Accuracy: 64.0917%
```

A confusion matrix is generated to provide a detailed breakdown of the model's performance

```
#Calculate confusion matrix
nb_conf_matrix = confusion_matrix(y_test, nb_predictions)
print("Naive Bayes Confusion Matrix:")
print(nb_conf_matrix)
print("="*40)
```

True Positive (TP): The number of samples that were correctly predicted as positive (gamma events) by the classifier = 1774

False Positive (FP): The number of samples that were incorrectly predicted as positive (gamma events), when it actually belongs to the negative class (hadron events) = 202

True Negative (TN): The number of samples that were correctly predicted as negative (hadron events) by the classifier = 798

False Negative (FN): The number of samples that were incorrectly predicted as negative (hadron events) but belong to e positive class (gamma events) =1239

```
========================================
Naive Bayes Confusion Matrix:
[[1774  202]
 [1239  798]]
========================================
```

Finally, a classification report is generated, which includes precision, recall, and F1-score for each class (gamma and hadron). This report provides insights into the model's performance, including its ability to correctly classify instances of each class.

```
========================================
Naive Bayes Classification Report:
              precision    recall  f1-score   support

           g       0.59      0.90      0.71      1976
           h       0.80      0.39      0.53      2037

    accuracy                           0.64      4013
   macro avg       0.69      0.64      0.62      4013
weighted avg       0.69      0.64      0.62      4013
```

# Comparisons

1. **Naive Bayes:**
   - Algorithm Type: Probabilistic classifier based on Bayes' theorem with strong (naive) independence assumptions between features.
   - Strengths:
     - Simple and easy to implement.
     - Performs well with high-dimensional data and large datasets.
     - Fast training and prediction time.
   - Weaknesses:
     - Assumes independence between features, which might not hold true in real-world datasets.
     - Limited expressive power compared to more complex models.
     - Sensitive to irrelevant features.

2. **Decision Tree:**
   - Algorithm Type: Non-parametric supervised learning algorithm used for classification and regression.
   - Strengths:
     - Easy to understand and interpret.
     - Handles both numerical and categorical data well.
     - Automatically selects important features.
     - Robust to outliers and missing values.
   - Weaknesses:
     - achieve overfitting, especially with deep trees.
     - Can be unstable with small variations in data.
     - Tends to create biased trees if some classes dominate.

### 3. AdaBoost (Adaptive Boosting):

- Algorithm Type: Ensemble learning method that combines multiple weak learners sequentially to create a strong learner.
- Strengths:
  - Can achieve high accuracy even with simple base learners (weak learners).
  - Reduces overfitting by focusing on misclassified instances.
  - Less susceptible to the overfitting problem compared to individual weak learners.
- Weaknesses:
  - Sensitive to noisy data and outliers.
  - Can be computationally expensive, especially with a large number of weak learners.
  - Performance can degrade if weak classifiers are too complex or too weak.

### 4. Random Forest:

- Algorithm Type: Ensemble learning method that constructs a multitude of decision trees at training time and outputs the mode of the classes as the prediction.
- Strengths:
  - Excellent performance on many types of problems.
  - Handles high-dimensional data well.
  - Robust to overfitting due to ensemble averaging.
  - Automatically selects important features.
  - Can handle missing values and maintain accuracy.
- Weaknesses:
  - Less interpretable compared to individual decision trees.
  - Can be slow to evaluate and make predictions, especially with large forests.
  - Requires more memory and computational resources compared to individual decision trees.

# Questions

## I. Balancing:

How does class-imbalance affect the performance of classification models?

> It affects classification models by biasing them towards the majority class, leading to poor performance on minority classes. Accuracy becomes unreliable, and models may struggle to generalize to new data. It also causes sensitivity to changes, ignoring feature importance.

Did you explore any other techniques for balancing the dataset besides randomly downsampling the majority class?

> 1) Random Oversampling: Randomly duplicating instances from the minority class to match the size of the majority class
>
> 2) Ensemble Techniques: Using ensemble methods such as BalancedRandomForest, which train multiple classifiers on balanced subsets of the data.
>
> 3) Cluster-Based Over-sampling: Identifying clusters within the minority class and make new instances based on the characteristics of these clusters.

How did balancing the dataset affect the performance of the models?

> It improves the performance models by enhancing prediction accuracy for minority classes, reducing bias towards the majority class, and leads to more reliable accuracy.

## II.  Model Tuning:

### Grid Search vs. Random Search. What is the main difference?

Grid Search: searches through a specified subset of hyperparameters, testing all possible combinations.

Random Search: randomly selects hyperparameter values from specified distributions or ranges.

### Did you consider using grid search or random search for hyperparameter tuning instead of cross-validation? Why or why not?

Grid search and random search are techniques for hyperparameter tuning, while cross-validation is a method for estimating the performance of a model. Typically, grid search or random search is used within a cross-validation framework to find the optimal hyperparameters.

### Is it necessary that accuracy increases when cv folds number increase?

No, increasing the number of folds in cross-validation does not guarantee an increase in accuracy. The relationship between the number of folds and accuracy can be influenced by various factors, and it is not a direct correlation.

### What's the purpose of cross validation?

Cross-validation can help in detecting overfitting but it doesn't guarantee that overfitting will be completely eliminated.

### Can you explain the concept of tuning in machine learning models?

Tuning in machine learning refers to the process of finding the optimal values for hyperparameters in a model

Are hyperparameters settings or configurations that are learned from the data ?

> No, they are not learned from the data but set before training the model

How does tuning impact model accuracy compared to using default hyperparameter values?

> Tuning helps improve model accuracy by finding the hyperparameter values that work best for a specific dataset and problem. Default hyperparameter values may not be suitable for every scenario, and tuning allows us to customize the model to achieve better performance. By finding the optimal values, tuning can enhance the model's ability to capture complex patterns and make more accurate predictions.

Does cross-validation ensure that overfitting is avoided?

> it does not guarantee the elimination of overfitting. Cross-validation can indicate if a model is sensitive to specific training instances or if there is significant variation in performance across folds, which may suggest overfitting

.

## III.   Model Classification:

What makes these classification models special? (Decision Tree, AdaBoost, Random Forest, Naive Bayes)

1)  Decision Tree: Decision trees are easy to understand. They can handle both numerical and categorical data and are less affected by outliers.
2)  AdaBoost : AdaBoost is an ensemble learning method that combines multiple weak learners to create a strong learner. It focus on misclassified instances and adjusts the weights of misclassified instances, making it more robust to outliers
3)  Random Forest: Random Forest is another ensemble learning technique that builds multiple decision trees and splits over random

dimensions. It has high accuracy and is able to handle large datasets. It can detect complex relationships between features and is less sensitive to noisy data.

4) Naive Bayes: Naive Bayes is a probabilistic classifier. It is  efficient for large datasets. Naive Bayes is useful for text classification tasks, such as spam detection.

## How can you control the complexity of a decision tree to avoid overfitting?

1. Maximum Depth: Limit the depth of the tree.
2. Minimum Samples per Split: Set the minimum number of samples required to split a node.
3. Minimum Samples per Leaf: Specify the minimum number of samples in a leaf node.
4. Maximum Features: Restrict the number of features considered for each split.

## How does the number of trees in a random forest affect the model's performance and generalization ability?

More trees reduce bias, improve averaging, and reduce overfitting,stabilizes predictions, enhances robustness to noise, and allows for capturing complex patterns.

## How do Naive Bayes and Decision Tree differ in terms of their underlying principles?

Naive Bayes is a probabilistic classifier based on Bayes' theorem and assumes that the features are conditionally independent given the class label.

Decision Tree is a non-parametric supervised learning algorithm that builds a tree-like model by recursively partitioning the data based on features' values.

## What are the advantages of AdaBoost over Naive Bayes and Decision Tree?

-AdaBoost is an ensemble method that combines multiple weak classifiers to create a strong classifier.

-It can improve the overall performance by focusing on difficult-to-classify examples.

-AdaBoost is less prone to overfitting compared to Decision Trees and can handle complex datasets better.

## How does Random Forest differ from AdaBoost?

-Random Forest is also an ensemble method that combines multiple decision trees, but it uses a different approach than AdaBoost.

-Random Forest builds each tree independently, without adjusting the weights of training examples.

-It introduces randomness by using a subset of features and bootstrapping to create diverse trees.

-AdaBoost, on the other hand, assigns higher weights to misclassified examples to improve their classification.

## Which technique is more interpretable, Decision Tree, or Random Forest?

Decision Trees are generally more interpretable because they represent a series of if-then rules.

Each node in the tree corresponds to a feature and a threshold, allowing for clear interpretation of the splitting process.

Random Forest, being an ensemble of Decision Trees, provides predictions based on voting or averaging, making it less interpretable than a single Decision Tree

## IV.   Model Evaluation:

**Can you explain the meaning of accuracy, precision, recall, F-score, and the confusion matrix?**

1. Accuracy: measures the correctly classified instances out of the total instances in the dataset. Accuracy alone may not be sufficient for imbalanced datasets since it does not consider the distribution of classes.

   Accuracy=Number of Correct Predictions / Total Number of Predictions

2. Precision: measures the proportion of correctly predicted positives out of all instances predicted as positive. It represents the model's ability to correctly identify positive instances without misclassifying negative instances as positive.

   Precision=True Positives / True Positives+False Positives

3. Recall (Sensitivity): measures the proportion of correct positive predictions out of all actual positive instances in the dataset. It represents the model's ability to capture all positive instances without missing any.

   Recall=True Positives/ True Positives+False Negatives

4. F1-score: mean of precision and recall. It provides a balance between precision and recall.

   F1-score=2×Precision×Recall / Precision+Recall

5. Confusion Matrix: A confusion matrix is a table that visualizes the performance of a classification model by comparing actual and predicted classes. It provides insights into the model's performance across different classes.
   It consists of four main elements:
   - True Positives (TP): Instances correctly predicted as positive.
   - True Negatives (TN): Instances correctly predicted as negative.
   - False Positives (FP): Instances incorrectly predicted as positive
   - False Negatives (FN): Instances incorrectly predicted as negative

Are there any other classification models you could have explored for this task?

1. Logistic Regression:  a linear model used for binary classification. It models the probability of a binary outcome
2. K-Nearest Neighbors (KNN): KNN is a  method used for classification and regression tasks. It predicts the class of a point based on the majority class among its k nearest neighbors.