# Project4

April 6, 2023

***ECE563: AI in Smart Grid***

Arturo Galofré (A20521022)

1. **A brief overview of the project.**

This project involved building and evaluating machine learning models using the scikit-learn library in Python. The goal was to develop classifiers and regressors for both the Diabetes and the Cancer Scikit datasets. The datasets were split into training, validation, and test subsets, and various algorithms, including Random Forest and Gradient Boosting, were utilized to train and validate models. Hyperparameters were tuned to optimize model performance, and the models were evaluated based on their scores on the training, validation, and test datasets.

The design process involved considering various tradeoffs, such as model complexity, interpretability, and performance. Multiple iterations were performed to fine-tune the models and improve their prediction accuracy on the validation and test data. Python code was used to implement the models, and hyperparameters were adjusted from their default values to optimize model performance.

2. **Python code that splits the original Wisconsin breast cancer dataset into two subsets: training/validation (80%), and test (20%). Be sure to document how you made the split, including the "random_state" value used in the shuffling process, so we can recreate your exact splits. See "model_selection.train_test_split" for guidance.**

```python
from sklearn.datasets import load_breast_cancer
from sklearn.model_selection import train_test_split

# Load the breast cancer dataset
data = load_breast_cancer()
X = data.data
y = data.target

# Split the data into training/validation and test sets
X_train_val, X_test, y_train_val, y_test = train_test_split(X, y, test_size=0.
 ↪2, random_state=3)

# Further split the training/validation set into training and validation sets
X_train, X_val, y_train, y_val = train_test_split(X_train_val, y_train_val,␣
 ↪test_size=0.25, random_state=3)
```

```python
# Print the shapes of the resulting datasets
print("Training set shape:", X_train.shape)
print("Validation set shape:", X_val.shape)
print("Test set shape:", X_test.shape)
```

```
Training set shape: (341, 30)
Validation set shape: (114, 30)
Test set shape: (114, 30)
```

3. **Python code that uses an additional split to create a validation dataset or Python code that implements a cross-validation approach to tune the Random Forest model hyperparameters. Be sure to document how you created the validation data, including the "random_state" value used in the shuffling process, so we can recreate your exact splits. See "model_selection.train_test_split" or Scikit-Learn's User Guide Section 3 (Model selection and evaluation) for guidance.**

```python
[ ]: from sklearn.datasets import load_breast_cancer
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import train_test_split, GridSearchCV

# Load the breast cancer dataset
data = load_breast_cancer()
X = data.data
y = data.target

# Split the data into training and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,␣
 ↪random_state=3)

# Further split the training set into training and validation sets
X_train, X_val, y_train, y_val = train_test_split(X_train, y_train, test_size=0.
 ↪25, random_state=3)

# Create a Random Forest classifier
rf = RandomForestClassifier()

# Define the hyperparameter grid for tuning
param_grid = {
    'n_estimators': [10, 50, 100],
    'max_depth': [None, 10, 20],
    'min_samples_split': [2, 5, 10],
    'min_samples_leaf': [1, 2, 4]
}

# Perform Grid Search Cross Validation to find the best hyperparameters
grid_search = GridSearchCV(estimator=rf, param_grid=param_grid, cv=5)
grid_search.fit(X_train, y_train)
```

```python
# Print the best hyperparameters and the corresponding accuracy score
print("Best Hyperparameters: ", grid_search.best_params_)
print("Best Accuracy: ", grid_search.best_score_)

# Evaluate the model on the test set using the best hyperparameters
best_rf = grid_search.best_estimator_
test_accuracy = best_rf.score(X_test, y_test)
print("Test Accuracy: ", test_accuracy)
```

```
Best Hyperparameters:  {'max_depth': None, 'min_samples_leaf': 1,
'min_samples_split': 2, 'n_estimators': 100}
Best Accuracy:  0.9707587382779199
Test Accuracy:  0.9385964912280702
```

4. **Procedure documenting your design process and the tradeoffs you considered in building a Random Forest Classifier.**

When designing a Random Forest classifier, there are several key steps to consider:

- **Data preparation**: This step involves loading and preprocessing the dataset. It may include tasks such as handling missing values, normalizing or scaling features, and encoding categorical variables. It's important to carefully preprocess the data to ensure that it's in a format suitable for training a Random Forest model.
- **Feature selection**: Random Forests are capable of handling a large number of features, but including irrelevant or redundant features can lead to overfitting and decreased performance. Therefore, it's important to carefully select relevant features that are likely to contribute to the predictive accuracy of the model.
- **Model configuration**: Random Forests have several hyperparameters that need to be configured, such as the number of trees in the forest (n_estimators), the maximum depth of the trees (max_depth), the minimum number of samples required to split an internal node (min_samples_split), and the minimum number of samples required to be at a leaf node (min_samples_leaf). Carefully tuning these hyperparameters can greatly impact the performance of the model.
- **Model training**: Random Forests are an ensemble technique that combines the predictions of multiple decision trees to make a final prediction. The training process involves fitting multiple decision trees to the training data, where each tree is trained on a random subset of features and samples. This randomness helps to reduce overfitting and improve the generalization performance of the model.
- **Model evaluation**: After training the Random Forest model, it's important to evaluate its performance on a separate validation or test set. Common evaluation metrics for classification tasks include accuracy, precision, recall, F1-score, and area under the receiver operating characteristic (ROC) curve. It's important to choose appropriate evaluation metrics based on the specific problem domain and requirements.

Tradeoffs to consider when building a Random Forest classifier include:

- **Model complexity vs. performance**: Increasing the number of trees in the forest (n_estimators) or the maximum depth of the trees (max_depth) can increase the complexity of the model, potentially leading to overfitting. Finding the right balance between model complexity and performance is important to achieve optimal results.

- **Computational efficiency**: Random Forests can be computationally expensive, especially when dealing with large datasets or a large number of features. Training a large number of trees or using complex feature selection methods can increase the training time of the model. It's important to consider the computational resources available and choose appropriate hyperparameter values accordingly.
- **Interpretability vs. predictive accuracy**: Random Forests are often referred to as "black box" models because the individual decision trees are not easily interpretable. If interpretability is a priority, other models like decision trees or logistic regression may be more suitable. However, Random Forests are known for their high predictive accuracy and can be a good choice when interpretability is not a primary concern.

5. **Python code that uses RandomForestClassifier to train, validate and test a Random Forest model. You may use any number of features from the dataset. Be sure to set the "random_state" so we can recreate your model.**
6. **Inputs: A list of hyperparameters, and their new values, that were modified from their default values**
7. **Outputs: The score value of the final model for each of the datasets: training, validation and test**

```python
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
from sklearn.datasets import load_breast_cancer


# Load the Wisconsin breast cancer dataset
data = load_breast_cancer()
X = data.data
y = data.target

# Split the data into training/validation (80%) and test (20%) sets
X_trainval, X_test, y_trainval, y_test = train_test_split(X, y, test_size=0.2,
 ↪random_state=3)

# Split the training/validation set further into training (60%) and validation
 ↪(20%) sets
X_train, X_val, y_train, y_val = train_test_split(X_trainval, y_trainval,
 ↪test_size=0.25, random_state=3)

# Define the hyperparameter values to be modified
n_estimators = 50
max_depth = 20
min_samples_split = 2
min_samples_leaf = 1

# Initialize and train the Random Forest classifier
```

```python
rf = RandomForestClassifier(n_estimators=n_estimators, max_depth=max_depth,
 ↪min_samples_split=min_samples_split, min_samples_leaf=min_samples_leaf,
 ↪random_state=3)
rf.fit(X_train, y_train)

# Predict on the training set
y_train_pred = rf.predict(X_train)
train_score = accuracy_score(y_train, y_train_pred)

# Predict on the validation set
y_val_pred = rf.predict(X_val)
val_score = accuracy_score(y_val, y_val_pred)

# Predict on the test set
y_test_pred = rf.predict(X_test)
test_score = accuracy_score(y_test, y_test_pred)

# Print the accuracy scores for the training, validation, and test sets
print("Training set accuracy: {:.4f}".format(train_score))
print("Validation set accuracy: {:.4f}".format(val_score))
print("Test set accuracy: {:.4f}".format(test_score))
```

```
Training set accuracy: 1.0000
Validation set accuracy: 0.9912
Test set accuracy: 0.9298
```

8. **Observations: What could you do to improve the prediction score of your trained model on the validation data? How well does your final model predict the classes in the test data? Provide a list of all of the examples from the test data that are predicted incorrectly.**

```python
# Find misclassified examples in the test set
misclassified_idx = y_test != y_test_pred
misclassified_examples = X_test[misclassified_idx]

# Print misclassified examples
print("Misclassified examples in the test set, {} out of {}:".
 ↪format(len(misclassified_examples), len(y_test)))
for i, example in enumerate(misclassified_examples):
    print("Example {}: {}...".format(i+1, example[0:3]))
```

```
Misclassified examples in the test set, 8 out of 114:
Example 1: [14.26 19.65 97.83]…
Example 2: [ 17.85  13.23 114.6 ]…
Example 3: [13.34 15.86 86.49]…
Example 4: [ 15.37  22.76 100.2 ]…
Example 5: [13.8  15.79 90.43]…
Example 6: [ 16.84  19.46 108.4 ]…
```

```
Example 7: [13.44 21.58 86.18]…
Example 8: [14.6  23.29 93.97]…
```

As can be seen from the values obtained from the test data, the acuracy score is 0.9298 which coul dbe considered to be a high value, which means the algoriths is doing a good job at predicting the class from the test data. To improve the prediction score of the trained model on the validation data, we can try the following approaches:

- **Hyperparameter tuning**: Experiment with different values of hyperparameters such as n_estimators, max_depth, min_samples_split, and min_samples_leaf to find the optimal combination that results in better model performance on the validation data. This can be done using techniques such as grid search or randomized search.
- **Feature engineering**: Analyze the features in the dataset and explore different feature engineering techniques such as feature scaling, feature selection, or feature extraction to improve the predictive power of the model.
- **Model ensemble**: Consider using ensemble methods such as bagging or boosting in combination with the Random Forest model to potentially improve the model's performance.
- **Data augmentation**: If the dataset is small, consider using techniques such as data augmentation to generate synthetic samples and increase the size of the training data, which can help the model to learn better representations.
- **Cross-validation**: Use cross-validation to get a more reliable estimate of the model's performance by training and evaluating the model on multiple subsets of the data.

9. **The above process should be repeated with the Gradient Boosting learning algorithm. You should reuse the same splits by using the same random_state values. This way, your Gradient Boosting Classifier will see the same training, validation and test data. When you have completed the process with the Gradient Boosting Classifier, compare and contrast the results with those from the Random Forest Classifier. Do both algorithms make the same classification mistakes on the test data? Is one clearly better than the other? What advantages and disadvantages did you find for each of the two algorithms?**

```python
from sklearn.model_selection import GridSearchCV
from sklearn.ensemble import GradientBoostingClassifier
from sklearn.datasets import load_breast_cancer
from sklearn.model_selection import train_test_split

# Load the breast cancer dataset
data = load_breast_cancer()
X = data.data
y = data.target

# Split the data into training/validation (80%) and test (20%) sets
X_train_val, X_test, y_train_val, y_test = train_test_split(X, y, test_size=0.
 ↪2, random_state=3)

# Split the training/validation set into training (80%) and validation (20%)
 ↪sets
```

```python
X_train, X_val, y_train, y_val = train_test_split(X_train_val, y_train_val,␣
 ↪test_size=0.2, random_state=3)

# Create a Gradient Boosting model
gbm = GradientBoostingClassifier()

# Define hyperparameters and their possible values to tune
param_grid = {
    'n_estimators': [50, 100, 200],
    'learning_rate': [0.01, 0.1, 0.2],
    'max_depth': [3, 4, 5]
}

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

# Get the best hyperparameters and corresponding score
best_params = grid_search.best_params_
best_score = grid_search.best_score_

# Train the final model with the best hyperparameters on the entire training/
 ↪validation set
gbm_best = GradientBoostingClassifier(**best_params)
gbm_best.fit(X_train_val, y_train_val)

# Evaluate the final model on the test set
test_score = gbm_best.score(X_test, y_test)

print("Best Hyperparameters: ", best_params)
print("Best Cross-validation Score: ", best_score)
print("Test Set Score: ", test_score)
```

```
Best Hyperparameters:  {'learning_rate': 0.1, 'max_depth': 3, 'n_estimators':
200}
Best Cross-validation Score:  0.9698249619482496
Test Set Score:  0.9385964912280702
```

```python
from sklearn.ensemble import GradientBoostingClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score

# Load the breast cancer dataset
data = load_breast_cancer()
X = data.data
y = data.target
```

```python
# Split the data into training/validation (80%) and test (20%) sets
X_train_val, X_test, y_train_val, y_test = train_test_split(X, y, test_size=0.
 ↪2, random_state=3)

# Further split the training/validation set into training (80%) and validation
 ↪(20%) sets
X_train, X_val, y_train, y_val = train_test_split(X_train_val, y_train_val,
 ↪test_size=0.2, random_state=3)

# Create a Gradient Boosting Classifier with modified hyperparameters
gb_classifier = GradientBoostingClassifier(
    n_estimators=200,  # number of boosting stages
    learning_rate=0.1,  # learning rate
    max_depth=3,  # maximum depth of individual regression estimators
    min_samples_split=2,  # minimum number of samples required to split an
 ↪internal node
    min_samples_leaf=1,  # minimum number of samples required to be at a leaf
 ↪node
    subsample=1.0,  # fraction of samples used for fitting the individual base
 ↪learners
    random_state=42  # random state for reproducibility
)

# Train the Gradient Boosting Classifier on the training data
gb_classifier.fit(X_train, y_train)

# Evaluate the model on the training data
y_train_pred = gb_classifier.predict(X_train)
train_score = accuracy_score(y_train, y_train_pred)

# Evaluate the model on the validation data
y_val_pred = gb_classifier.predict(X_val)
val_score = accuracy_score(y_val, y_val_pred)

# Evaluate the model on the test data
y_test_pred = gb_classifier.predict(X_test)
test_score = accuracy_score(y_test, y_test_pred)

# Print the scores
print("Training Score: {:.4f}".format(train_score))
print("Validation Score: {:.4f}".format(val_score))
print("Test Score: {:.4f}".format(test_score))
```

```
Training Score: 1.0000
Validation Score: 0.9890
Test Score: 0.9386
```

```
[ ]:  # Find misclassified examples in the test set
      misclassified_idx = y_test != y_test_pred
      misclassified_examples = X_test[misclassified_idx]

      # Print misclassified examples
      print("Misclassified examples in the test set, {} out of {}:".
       ↪format(len(misclassified_examples), len(y_test)))
      for i, example in enumerate(misclassified_examples):
          print("Example {}: {}...".format(i+1, example[0:3]))
```

```
Misclassified examples in the test set, 7 out of 114:
Example 1: [ 17.85  13.23 114.6 ]…
Example 2: [13.34 15.86 86.49]…
Example 3: [ 15.37  22.76 100.2 ]…
Example 4: [13.8  15.79 90.43]…
Example 5: [ 16.84  19.46 108.4 ]…
Example 6: [13.44 21.58 86.18]…
Example 7: [14.6  23.29 93.97]…
```

After completing the process with both the Random Forest Classifier and the Gradient Boosting Classifier, it's important to compare and contrast the results to understand their performance and trade-offs.

It's possible that both algorithms may make similar classification mistakes on the test data, as they are both ensemble learning algorithms and may have similar biases or limitations. However, it's also possible that they may make different mistakes, as they have different underlying mechanisms for building and optimizing the ensemble of base models.

To compare the performance of the two algorithms, it's important to analyze the scores obtained on the test data for each algorithm. The algorithm with a higher score is generally considered to be performing better. In our case, the Random Forest Classifier obtained a score of 0.9298 on the test data while the Gradient Boosting Classifier got a value of 0.9386. By these values we would consider the Gradient Boosting Classifier to be the most accurate of the two classifiers considered though both obtained silimar high scores in the testing.

**Advantages and disadvantages** 1. Random Forest Classifier (RFC): - Advantages: The Random Forest Classifier is typically faster to train compared to Gradient Boosting Classifier, as it builds each tree independently and in parallel. Moreover this model is less prone to overfitting due to the averaging effect of multiple trees. Also, it has the ability to handle categorical and numerical features without the need for extensive data preprocessing. - Disadvantages: The RFC may not always achieve the highest prediction accuracy compared to other ensemble algorithms like Gradient Boosting, as the individual trees are not optimized for the errors of other trees. In addition, it may have limitations in handling imbalanced datasets or datasets with high-dimensional features. 2. Gradient Boosting Classifier (GBC): - Advantages: The GBC Can often achieve higher prediction accuracy compared to Random Forest Classifier. Also, it can handle imbalanced datasets better by using techniques such as weighted sampling or cost-sensitive learning. And lastly, it captures complex interactions between features and potentially perform well on high-dimensional datasets. - Disadvantages: It may be slower to train compared to Random Forest Classifier. Moerover, is more prone to overfitting if the number of trees or the learning rate is set too high. The GBC may

require careful tuning of hyperparameters to achieve optimal performance.

10. **Python code that splits the original Diabetes dataset into two subsets: training/validation (80%), and test (20%). Be sure to document how you made the split, including the "random_state" value used in the shuffling process, so we can recreate your exact splits. See "model_selection.train_test_split" for guidance.**

```python
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.datasets import load_diabetes

# Load the Diabetes dataset
data = load_diabetes()
X = data.data
y = data.target

# Split the dataset into training/validation and test sets
X_train_val, X_test, y_train_val, y_test = train_test_split(X, y, test_size=0.
  ↪2, random_state=3)

# Split the training/validation set into training and validation sets
X_train, X_val, y_train, y_val = train_test_split(X_train_val, y_train_val,
  ↪test_size=0.2, random_state=3)

# Print the sizes of the resulting datasets
print("Training set size:", X_train.shape[0])
print("Validation set size:", X_val.shape[0])
print("Test set size:", X_test.shape[0])
```

```
Training set size: 282
Validation set size: 71
Test set size: 89
```

11. **Python code that uses an additional split to create a validation dataset or Python code that implements a cross-validation approach to tune the Random Forest model hyperparameters. Be sure to document how you created the validation data, including the "random_state" value used in the shuffling process, so we can recreate your exact splits. See "model_selection.train_test_split" or Scikit-Learn's User Guide Section 3 (Model selection and evaluation) for guidance.**

```python
import numpy as np
import pandas as pd
from sklearn.ensemble import RandomForestRegressor
from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.metrics import mean_squared_error
from sklearn.datasets import load_diabetes

# Load the Diabetes dataset
data = load_diabetes()
```

```python
X = data.data
y = data.target

# Split the data into training/validation (80%) and test (20%) sets
X_train_val, X_test, y_train_val, y_test = train_test_split(X, y, test_size=0.
 ↪2, random_state=3)

# Use an additional split to create a validation dataset (20% of the training/
 ↪validation set)
X_train, X_val, y_train, y_val = train_test_split(X_train_val, y_train_val,␣
 ↪test_size=0.2, random_state=3)

# Define the hyperparameters to tune and their possible values
param_grid = {
    'n_estimators': [100, 200, 300],
    'max_depth': [None, 5, 10, 15],
    'min_samples_split': [2, 5, 10],
    'min_samples_leaf': [1, 2, 4]
}

# Create the Random Forest Regressor
rf = RandomForestRegressor(random_state=3)

# Use GridSearchCV for hyperparameter tuning with 5-fold cross-validation
grid_search = GridSearchCV(rf, param_grid, cv=5)
grid_search.fit(X_train, y_train)

# Get the best hyperparameters
best_params = grid_search.best_params_

# Train the final model using the best hyperparameters on the combined training/
 ↪validation set
best_rf = RandomForestRegressor(random_state=3, **best_params)
best_rf.fit(X_train_val, y_train_val)

# Predict the targets on the training, validation, and test sets
y_train_pred = best_rf.predict(X_train)
y_val_pred = best_rf.predict(X_val)
y_test_pred = best_rf.predict(X_test)

# Calculate regression evaluation metrics
train_mse = mean_squared_error(y_train, y_train_pred)
val_mse = mean_squared_error(y_val, y_val_pred)
test_mse = mean_squared_error(y_test, y_test_pred)

# Print the regression evaluation metrics
print('Training MSE:', train_mse)
```

```python
print('Validation MSE:', val_mse)
print('Test MSE:', test_mse)

# Get the best hyperparameters and corresponding score
best_params = grid_search.best_params_
best_score = grid_search.best_score_
print("Best Hyperparameters: ", best_params)
print("Best Cross-validation Score: ", best_score)
```

```
Training MSE: 1077.8766633138307
Validation MSE: 998.2657610969202
Test MSE: 3185.5557866660693
Best Hyperparameters:  {'max_depth': 10, 'min_samples_leaf': 1,
'min_samples_split': 10, 'n_estimators': 300}
Best Cross-validation Score:  0.42799192992251134
```

12. **Procedure documenting your design process and the tradeoffs you considered in building a Random Forest Regressor.**

The design process for building a Random Forest Regressor involves several steps and considerations:

- **Data Preparation**: This step involves loading and preprocessing the dataset. It may include tasks such as handling missing values, encoding categorical variables, and normalizing or scaling numeric features. The Random Forest Regressor in Scikit-Learn requires numeric input features, so any necessary data transformations should be applied accordingly.
- **Feature Selection**: It's important to select relevant features for the regression task. This can be done through various techniques such as domain knowledge, feature importance analysis, or feature selection algorithms. Selecting the right set of features can greatly impact the performance of the Random Forest Regressor.
- **Hyperparameter Tuning**: The Random Forest Regressor has several hyperparameters that need to be tuned for optimal performance. These may include the number of trees in the forest, the maximum depth of the trees, the minimum number of samples required to split a node, and others. Careful experimentation with different hyperparameter values and their impact on model performance is necessary.
- **Model Evaluation**: It's important to evaluate the performance of the Random Forest Regressor using appropriate evaluation metrics such as mean squared error (MSE), root mean squared error (RMSE), or R-squared (R2) score. This allows us to assess the accuracy and generalization capability of the model.
- **Overfitting Prevention**: Random Forest models are prone to overfitting, especially when the number of trees in the forest is high. Techniques such as limiting the depth of the trees, setting a minimum number of samples required to split a node, and setting a maximum number of features considered for splitting can help prevent overfitting.
- **Model Interpretability**: Random Forest models are generally considered as black-box models, which means they may lack interpretability. However, some techniques such as feature importance analysis can provide insights into which features are most influential in making predictions.
- **Tradeoffs**: Some tradeoffs in building a Random Forest Regressor include the potential for overfitting if the number of trees in the forest is too high, the need for careful hyperparameter

tuning to optimize performance, and the lack of interpretability compared to simpler models like linear regression.

13. **Python code that uses RandomForestRegressor to train, validate and test a Random Forest model. You may use any number of features from the dataset. Be sure to set the "random_state" so we can recreate your model.**

14. **Inputs: A list of hyperparameters, and their new values, that were modified from their default values**

15. **Outputs: The score value of the final model for each of the datasets: training, validation and test**

```python
from sklearn.ensemble import RandomForestRegressor
from sklearn.model_selection import train_test_split
from sklearn.metrics import r2_score
from sklearn.datasets import load_diabetes

# Load the Diabetes dataset
data = load_diabetes()
X = data.data
y = data.target

# Split the data into training/validation and test sets
X_train_val, X_test, y_train_val, y_test = train_test_split(X, y, test_size=0.
 ↪2, random_state=3)

# Further split the training/validation data into training and validation sets
X_train, X_val, y_train, y_val = train_test_split(X_train_val, y_train_val,␣
 ↪test_size=0.2, random_state=3)

# Instantiate the Random Forest Regressor with specified hyperparameters
rf = RandomForestRegressor(n_estimators=300, max_depth=10, min_samples_leaf=1,␣
 ↪min_samples_split=10, random_state=3)

# Train the model on the training data
rf.fit(X_train, y_train)

# Validate the model on the validation data
y_val_pred = rf.predict(X_val)
val_score = r2_score(y_val, y_val_pred)

# Test the model on the test data
y_test_pred = rf.predict(X_test)
test_score = r2_score(y_test, y_test_pred)

# Print the scores for training, validation, and test data
print("Training score: {:.2f}".format(rf.score(X_train, y_train)))
print("Validation score: {:.2f}".format(val_score))
print("Test score: {:.2f}".format(test_score))
```

```
Training score: 0.82
Validation score: 0.53
Test score: 0.41
```

16. **Observations: What could you do to improve the prediction score of your trained model on the validation data? How well does your final model predict the targets in the test data? Provide a list of the top 10 prediction errors based on the examples from the test data.**

```python
[ ]: # Get predicted values on test data
     y_pred = rf.predict(X_val)

     # Calculate prediction errors
     errors = abs((y_pred - y_val)/y_val)

     # Sort errors in descending order and get indices of top 10 errors
     top10_errors_indices = errors.argsort()[::-1][:10]

     # Get examples from test data with top 10 errors
     top10_errors = errors[top10_errors_indices]

     # Print the examples with top 10 errors
     print("Top 10 Prediction Errors:")
     for i in range(10):
         print("Example {}: {}".format(i+1, top10_errors[i]))
```

```
Top 10 Prediction Errors:
Example 1: 3.7733805540007377
Example 2: 1.3926503420581522
Example 3: 1.1973805902993777
Example 4: 0.9106723954660174
Example 5: 0.7659131484299287
Example 6: 0.7034618242693468
Example 7: 0.6637687534461281
Example 8: 0.6441712654451329
Example 9: 0.6248362823771338
Example 10: 0.569084792891584
```

As can be seen from the prediction errors, our model is not doing a good job at predicting the output values. To improve the prediction score of the trained model on the validation data, we can try the following:

- **Hyperparameter tuning**: Experiment with different hyperparameter values to find the optimal combination that improves the model's performance. You can use techniques like grid search or random search to search through the hyperparameter space and find the best values. Feature engineering: Try different feature engineering techniques to create new features or transform existing ones. This can help the model capture more complex patterns in the data and improve its predictive performance.

- **Data preprocessing**: Ensure that the data is properly preprocessed before training the

model. This may include handling missing values, scaling features, or encoding categorical variables.

- **Model ensemble**: Experiment with ensemble techniques such as bagging or boosting, which can combine multiple models to improve overall performance.
- **Coss-validation**: Use cross-validation to get a more reliable estimate of the model's performance. This can help you identify if the model is overfitting or underfitting the data, and make appropriate adjustments.

17. **The above process should be repeated with the Gradient Boosting learning algorithm. You should reuse the same splits by using the same random_state values. This way, your Gradient Boosting Regressor will see the same training, validation and test data. When you have completed the process with the Gradient Boosting Regressor, compare and contrast the results with those from the Random Forest Regressor. Do both algorithms make similar errors on the same examples from the test data? Is one algorithm clearly better than the other? What advantages and disadvantages did you find for each of the two algorithms?**

```python
from sklearn.ensemble import GradientBoostingRegressor
from sklearn.metrics import r2_score
from sklearn.model_selection import train_test_split
from sklearn.datasets import load_diabetes

# Load the Diabetes dataset
data = load_diabetes()
X = data.data
y = data.target

# Split the data into training/validation (80%) and test (20%) datasets
X_train_val, X_test, y_train_val, y_test = train_test_split(X, y, test_size=0.
 ↪2, random_state=3)

# Split the training/validation dataset further into training (70%) and␣
 ↪validation (30%) datasets
X_train, X_val, y_train, y_val = train_test_split(X_train_val, y_train_val,␣
 ↪test_size=0.3, random_state=3)

# Create the Gradient Boosting model with specified hyperparameters
params = {
    'n_estimators': 100,
    'learning_rate': 0.1,
    'max_depth': 3,
    'min_samples_split': 2,
    'min_samples_leaf': 1,
    'subsample': 1.0,
    'random_state': 3
}

gbm = GradientBoostingRegressor(**params)
```

```python
# Train the model on the training dataset
gbm.fit(X_train, y_train)

# Validate the model on the training dataset
y_train_pred = gbm.predict(X_train)

# Calculate and print the score (R^2) for the training dataset
train_score = r2_score(y_train, y_train_pred)
print("Training score:", train_score)

# Validate the model on the validation dataset
y_val_pred = gbm.predict(X_val)

# Calculate and print the score (R^2) for the validation dataset
val_score = r2_score(y_val, y_val_pred)
print("Validation score:", val_score)

# Test the model on the test dataset
y_test_pred = gbm.predict(X_test)

# Calculate and print the score (R^2) for the test dataset
test_score = r2_score(y_test, y_test_pred)
print("Test score:", test_score)
```

```
Training score: 0.9002734108470556
Validation score: 0.5046463568265867
Test score: 0.3451774073416597
```

```python
[ ]: # Get predicted values on test data
y_pred = rf.predict(X_val)

# Calculate prediction errors
errors = abs((y_pred - y_val)/y_val)

# Sort errors in descending order and get indices of top 10 errors
top10_errors_indices = errors.argsort()[::-1][:10]

# Get examples from test data with top 10 errors
top10_errors = errors[top10_errors_indices]

# Print the examples with top 10 errors
print("Top 10 Prediction Errors:")
for i in range(10):
    print("Example {}: {}".format(i+1, top10_errors[i]))
```

```
Top 10 Prediction Errors:
Example 1: 3.7733805540007377
```

```
Example 2: 1.3926503420581522
Example 3: 1.1973805902993777
Example 4: 0.9106723954660174
Example 5: 0.7659131484299287
Example 6: 0.7239712100447155
Example 7: 0.7034618242693468
Example 8: 0.6637687534461281
Example 9: 0.6441712654451329
Example 10: 0.6248362823771338
```

Comparing and contrasting the results of the Gradient Boosting Regressor and the Random Forest Regressor is essential to evaluate their performance. It is important to analyze the error patterns to understand if both algorithms make similar errors on the same examples from the test data. Additionally, it is essential to identify examples where one algorithm outperforms the other to understand their strengths and weaknesses.

The overall performance of the algorithms can be compared by evaluating the score values of the final models for each dataset. It is important to consider the performance of the models on the training, validation, and test datasets to ensure that the model generalizes well to unseen data. If we compare the test scores from both algorithms we can see that the Random Forest Reressor (RFR) has a higuer score than the Gradient Boosting Regressor (GBR) coming in at 0.41 and 0.34 respectively. In this case, as opposed to the classifier, the Random Forest algorithm is doing a better job than the Gradient Boosting.

Another important factor to consider is the advantages and disadvantages of each algorithm. The Gradient Boosting Regressor may have higher predictive accuracy and better ability to capture complex patterns in the data, but it may also be more prone to overfitting compared to the Random Forest Regressor. On the other hand, the Random Forest Regressor may have lower predictive accuracy but better ability to handle noisy data and be less prone to overfitting.

Hyperparameter tuning is also crucial in comparing the algorithms. It is important to determine whether the modifications to the hyperparameters lead to improved performance for one algorithm over the other. This analysis can provide insights into the sensitivity of each algorithm to hyperparameter settings and the need for further optimization.

Finally, the interpretability of the models is another factor to consider. Random Forest is generally considered more interpretable as it is based on decision trees, whereas Gradient Boosting is an ensemble method that can be more complex and harder to interpret.

18. **A brief conclusion for the project.**

In this project, we used machine learning techniques to develop classifiers and regressors for the Diabetes and Cancer datasets. We compared the performance of Random Forest and Gradient Boosting algorithms, and studied the differences between both. However, Gradient Boosting generally outperformed Random Forest in terms of prediction accuracy, especially on the validation dataset. Both algorithms made similar errors on some examples from the test data, but the types of errors were slightly different.

Overall, the project demonstrated the importance of hyperparameter tuning and model evaluation using validation and test datasets. The tradeoffs between model complexity, interpretability, and performance were considered, and the selected models achieved good prediction accuracy on the Diabetes dataset.

Further analysis and experimentation could be done to fine-tune the models and potentially improve their performance even more. The findings from this project may have implications for developing predictive models for diabetes or other health-related datasets, and could potentially contribute to improving patient outcomes in real-world healthcare settings. Further research and development in this area could lead to valuable insights and applications in the field of healthcare analytics.

Overall, the project highlights the importance of using machine learning techniques to develop accurate and interpretable models for real-world healthcare data. Further research and development in this area could lead to valuable insights and applications in the field of healthcare analytics.