Homework 2: Trees and Calibration

Instructions:

Please upload the .ipynb, .pdf to Github prior to the deadline. Please include your UNI as well.

Make sure to use the dataset that we provide in CourseWorks/Classroom.

There are a lot of applied questions based on the code results. Please make sure to answer them all. These are primarily to test your understanding of the results your code generate (similar to any Data Science/ML case study interviews).

Name: Yueqi Li

UNI:yl5329



The Dataset

Description

The Diabetes Dataset comprises medical data from 768 female patients of Pima Indian heritage, including 8 health-related features and a binary target indicating the presence or absence of diabetes.

```
import warnings
warnings.filterwarnings("ignore")

import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler, MinMaxScaler
from sklearn.preprocessing import OrdinalEncoder
```

Question 1: Decision Trees

1.1: Load the provided dataset

```
## YOUR CODE HERE
df = pd.read_csv("diabetes.csv")
df
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPec
0	6	148	72	35	0	33.6	
1	1	85	66	29	0	26.6	
2	8	183	64	0	0	23.3	
3	1	89	66	23	94	28.1	
4	0	137	40	35	168	43.1	
763	10	101	76	48	180	32.9	
764	2	122	70	27	0	36.8	
765	5	121	72	23	112	26.2	
766	1	126	60	0	0	30.1	
767	1	93	70	31	0	30.4	

768 rows × 9 columns

1.2: How many instances are there in the dataset for each class (diabetic and non-diabetic patients)? What does this tell you about the balance of the dataset?

The dataset contains 500 instances of non-diabetic patients (Outcome = 0) and 268 instances of diabetic patients (Outcome = 1). This indicates that the dataset is imbalanced, with a higher number of non-diabetic instances compared to diabetic ones.

```
## YOUR CODE HERE
class_distribution = df['Outcome'].value_counts()

class_distribution

0 500
1 268
Name: Outcome, dtype: int64
```

1.3: Are there any missing values in the dataset? If so, how will you handle them?

There is no missing value in the dataset

```
missing_values = df.isnull().sum()
missing_values
     Pregnancies
                                   0
     Glucose
    BloodPressure
                                   0
     SkinThickness
                                   0
     Insulin
    BMT
                                   0
    DiabetesPedigreeFunction
                                   0
                                   0
    Age
    Outcome
                                   0
     dtype: int64
```

1.4: Split the data into development and test datasets. Which splitting methodology did you choose and why?

**Hint: Based on the distribution of the data, try to use the best splitting strategy.

Given the imbalance in the dataset, it's essential to use a splitting strategy that maintains the proportion of classes in both the development and test datasets. Stratified sampling is a suitable methodology for this purpose.

```
from sklearn.model_selection import train_test_split
# Split the data into features and target variable
X = df.drop('Outcome', axis=1)
y = df['Outcome']
# Split the data into development (80%) and test (20%) datasets with stratification
X_{dev}, X_{test}, y_{dev}, y_{test} = train_test_split(X, Y, test_size=0.2, stratify=Y, random_state=42)
X_{train}, X_{val}, y_{train}, y_{val} = train_test_split(X_{dev}, y_{dev}, test_size=0.25, stratify=y_{dev}, random_state=42) # 0.25 x 0.8 =
# Check the distribution of classes in each dataset to ensure stratification worked as intended
distribution summary = {
    "Training": y_train.value_counts(normalize=True),
    "Validation": y_val.value_counts(normalize=True),
    "Test": y_test.value_counts(normalize=True)
distribution_summary
     {'Training': 0
                       0.652174
         0.347826
     Name: Outcome, dtype: float64,
      'Validation': 0
                         0.649351
     1
           0.350649
     Name: Outcome, dtype: float64,
      'Test': 0
                  0.649351
          0.350649
     Name: Outcome, dtype: float64}
```

1.5: Build a decision tree classifier to predict the onset of diabetes. What criterion (e.g., Gini impurity, entropy) did you choose for splitting, and why?

For building this decision tree classifier to predict the onset of diabetes, I initially chose Gini impurity as the criterion for splitting. The choice was made due to Gini impurity's computational efficiency over entropy. Gini impurity measures the frequency at which any element from the set would be incorrectly labeled if it was randomly labeled according to the distribution of labels in the subset. It's computationally faster because it does not require calculating logarithmic functions as entropy does.

```
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
# Initialize the Decision Tree Classifier with Gini impurity
dt_classifier_gini = DecisionTreeClassifier(criterion='gini', random_state=42)
# Train the classifier on the training data
dt_classifier_gini.fit(X_train, y_train)
# Predict on the validation set
y_pred_val_gini = dt_classifier_gini.predict(X_val)
```

1.6: Evaluate your model using accuracy, precision, recall, and F1-score. What do these metrics reveal about your model's performance?

Decision Tree Model Performance:

- Accuracy: 75.97%
- Precision: 64.91%
- Recall: 68.51%
- F1 Score: 66.66%

These evaluation metrics reveal that the model has a decent performance but might still benefit from further optimization or alternative approaches to handle the class imbalance more effectively.

1.8: List the top 3 most important features for this trained tree? How would you justify these features being the most important?

Glucose, BMI, and Age are the top 3 most important feature for the tranined tree. These features' importance is derived from their ability to reduce uncertainty (or impurity) in the dataset when making splits. The decision tree algorithm calculates the gain in information—how much each feature reduces the randomness or heterogeneity of the groups it creates. The importance of glucose, BMI, and age reflects their effectiveness in creating homogenous subsets of patients, thereby improving the model's accuracy. This justification is entirely model-centric, focusing on how each feature contributes to the decision-making process within the decision tree structure.

```
# Extract feature importance
feature_importances = pd.DataFrame({
    "Feature": X.columns,
    "Importance": dt_classifier_gini.feature_importances_
}).sort_values(by="Importance", ascending=False)

top_3_features = feature_importances.head(3)

top_3_features
```



Question 2: Random Forests

2.1: Train a Random Forest model on the development dataset using RandomForestClassifier class in sklearn. Use the default parameters. Evaluate the performance of the model on test dataset. Does this perform better than Decision Tree on the test dataset

Random Forest Model Performance:

Accuracy: 72.73%Precision: 62.5%Recall: 55.56%

• F1 Score: 58.82%

- Accuracy: The Decision Tree model has a higher accuracy than the Random Forest model (75.97% vs 72.73%).
- Precision: The Decision Tree also has a slightly higher precision than the Random Forest (64.91% vs 62.5%).
- Recall: The Decision Tree outperforms the Random Forest in recall as well (68.52% vs 55.56%).
- F1 Score: The Decision Tree has a higher F1 score compared to the Random Forest (66.67% vs 58.82%).

Based on these comparisons, the Decision Tree model performs better than the Random Forest model on the test dataset across all evaluated metrics.

```
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import classification report
# Train a Random Forest model on the development dataset
rf_classifier = RandomForestClassifier(random_state=42)
rf_classifier.fit(X_train, y_train)
# Evaluate the performance of the model on the test dataset
y_pred = rf_classifier.predict(X_test)
# Calculate metrics
accuracy = accuracy_score(y_test, y_pred)
precision = precision_score(y_test, y_pred)
recall = recall_score(y_test, y_pred)
f1 = f1_score(y_test, y_pred)
# Calculate evaluation metrics for the validation set
metrics = {
    "Accuracy": accuracy,
    "Precision": precision,
    "Recall": recall,
    "F1 Score": f1
print("Random Forest:", metrics)
print("Decision Tree:", metrics_gini)
     Random Forest: {'Accuracy': 0.7272727272727273, 'Precision': 0.625, 'Recall': 0.55555555555556, 'F1 Score': 0.588235294117
    Decision Tree: {'Accuracy': 0.7597402597402597, 'Precision': 0.6491228070175439, 'Recall': 0.6851851851851851852, 'F1 Score': 0
```

2.2 Does all trees in the trained random forest model have pure leaves? How would you verify this?

Based on the explicit check for pure leaves in each tree of the Random Forest model, it is determined that all trees have pure leaves. This is inferred from the condition where each leaf node's impurity is 0, indicating that the nodes contain data points from a single class only. The check was performed on each tree within the Random Forest, and the result confirms that all trees indeed have pure leaves, as indicated by the first 10 trees' status, which all show true for having pure leave

But for other random forester models, they are not nessary all have pure leaves

```
# Function to check for pure leaves in a single decision tree
def check_pure_leaves(tree):
   # Leaf nodes have a specific value in children_left and children_right arrays, which is -1.
   # If a node is not split further (both children indices are −1), it's a leaf node.
   leaf_indices = tree.tree_.children_left == -1
   # Count of samples at leaf nodes
    leaf_counts = tree.tree_.n_node_samples[leaf_indices]
   # Output value at leaf nodes (impurity)
    leaf_impurity = tree.tree_.impurity[leaf_indices]
   # A leaf is considered "pure" if its impurity is 0
   # We check if all leaves have impurity 0
   all_leaves_pure = all(impurity == 0 for impurity in leaf_impurity)
    return all_leaves_pure
# Check each tree in the Random Forest
pure_leaves_status = [check_pure_leaves(tree) for tree in rf_classifier.estimators_]
# Determine if all trees have pure leaves, and show the overall result and the first 10 trees' status
all_trees_have_pure_leaves = all(pure_leaves_status)
all_trees_have_pure_leaves, pure_leaves_status[:10]
    (True, [True, True, True, True, True, True, True, True, True, True])
```

2.3: Assume you want to improve the performance of this model. Also, assume that you had to pick two hyperparameters that you could tune to improve its performance. Which hyperparameters would you choose and why?

To improve the performance of the Random Forest model, I would like to choose these two hyperparameters whic can significantly impact the model's performance are:

- 1. n_estimators: This hyperparameter specifies the number of trees in the forest. Increasing the number of trees can help improve the model's accuracy and robustness by reducing overfitting. Each tree in the forest contributes to the final decision, and more trees mean more opinions and a more stable average prediction. However, it's important to note that beyond a certain point, increasing the number of trees leads to diminishing returns in performance improvement and increases computational cost and time.
- 2. max_depth: This hyperparameter controls the maximum depth of each tree. A deeper tree can model more complex patterns by creating more splits and capturing more information about the data. However, if the trees are too deep, they may fit the training data too closely and overfit, capturing noise in the training data as if it were a real pattern. Setting an optimal max_depth can help balance the model's ability to generalize to unseen data without overfitting.

These two hyperparameters are chosen because they directly influence the model's capacity to learn from the training data. Tuning n_estimators and max_depth can help find a good balance between bias and variance, potentially leading to a more accurate and generalizable model.

2.4: Now, assume you had to choose up to 5 different values (each) for these two hyperparameters. How would you choose these values that could potentially give you a performance lift?

When choosing values for the hyperparameters n_estimators and max_depth to potentially improve the performance of a Random Forest model, a strategic approach involves selecting a range that covers a broad spectrum from lower to higher complexity. This ensures we explore both simpler and more complex models to find the best balance between underfitting and overfitting. Here's how we might choose up to 5 different values for each hyperparameter:

n_estimators

- Lower values: Start with a modest number of trees to keep the computation cost low and see if there's substantial gain with just a few trees. For example, 10 or 50 trees.
- Medium values: Increase the number of trees to see if the model performance improves with more ensemble members. A moderate increase can be seen with 100 or 150 trees, which is often enough to see significant improvements.
- **Higher values:** To ensure we're capturing the potential for more stability and accuracy, testing with even larger numbers like 200 and 300 trees can be beneficial. This tests the diminishing returns threshold where more trees do not significantly improve performance.

max_depth

- Shallower trees: Starting with a low depth such as 3 or 5 can help in understanding how well the model performs with very simple decision trees. This can be particularly useful for avoiding overfitting.
- Moderate depth: Increasing the depth to values like 10 or 15 allows the trees to capture more details and interactions between features, potentially improving model accuracy.

- Unrestricted depth: A value of None allows trees to expand until all leaves are pure or until all leaves contain less than min_samples_split samples. This is useful for benchmarking the model's performance when it's fully grown, though it risks overfitting.
- 2.5: Perform model selection using the chosen values for the hyperparameters. Use cross-validation for finding the optimal hyperparameters. Report on the optimal hyperparameters. Estimate the performance of the optimal model (model trained with optimal hyperparameters) on test dataset? Has the performance improved over your plain-vanilla random forest model trained in Q2.1?

Based on the provided results, the optimal hyperparameters for the Random Forest model are:

max_depth:10n_estimators:100

The optimal model demonstrates an improvement over the plain-vanilla model across all metrics:

- The accuracy has increased from 72.73% to 75.97%, indicating a higher overall rate of correct predictions.
- The precision has improved from 62.5% to 67.35%, meaning a higher proportion of positive identifications was actually correct.
- The recall has increased from 55.56% to 61.11%, indicating a higher proportion of actual positives was identified correctly.
- The **F1 Score** has risen from 58.82% to 64.08%, reflecting a better balance between precision and recall.

```
from sklearn.model_selection import GridSearchCV, train_test_split
# Parameter grid
param_grid = {
    'n_estimators': [10, 50, 100, 200, 300], # From lower to higher to cover a broad spectrum
    'max_depth': [3, 5, 10, 15, None] # Includes shallow, moderate, and unrestricted depths
# Initialize the GridSearchCV object with the RandomForestClassifier and the parameter grid
grid_search = GridSearchCV(estimator=RandomForestClassifier(random_state=42),
                           param_grid=param_grid,
                            cv=5, # 5-fold cross-validation
                            scoring='accuracy', # You can choose other metrics as well
                            n_jobs=-1) # Use all available CPUs
# Fit the grid search to the development dataset
grid_search.fit(X_dev, y_dev)
# Best parameters and best score
best_params = grid_search.best_params_
best_score = grid_search.best_score_
# Retrain the RandomForestClassifier on the entire development set with the best parameters
rf_optimal = RandomForestClassifier(n_estimators=best_params['n_estimators'],
                                     max_depth=best_params['max_depth'],
                                     random_state=42)
rf_optimal.fit(X_dev, y_dev)
# Evaluate the model
y_pred_optimal = rf_optimal.predict(X_test)
# Metrics
accuracy_optimal = accuracy_score(y_test, y_pred_optimal)
precision_optimal = precision_score(y_test, y_pred_optimal)
recall_optimal = recall_score(y_test, y_pred_optimal)
f1_optimal = f1_score(y_test, y_pred_optimal)
# Comparison with plain-vanilla model
plain_vanilla_results = {
    'Accuracy': 0.72727272727273,
    'Precision': 0.625,
    'Recall': 0.555555555555556,
    'F1 Score': 0.5882352941176471
}
comparison_results = {
    'Optimal Model': {
        'Accuracy': accuracy_optimal,
        'Precision': precision_optimal,
        'Recall': recall_optimal,
        'F1 Score': f1_optimal
    },
    'Plain-Vanilla Model': plain_vanilla_results,
    'Best Parameters': best_params,
    'Best Score (CV)': best_score
}
comparison_results
     {'Optimal Model': {'Accuracy': 0.7597402597402597,
       Precision': 0.673469387755102,
       'Recall': 0.61111111111111112,
       'F1 Score': 0.6407766990291262},
      'Plain-Vanilla Model': {'Accuracy': 0.72727272727273,
       'Precision': 0.625,
       'Recall': 0.55555555555555
       'F1 Score': 0.5882352941176471},
      'Best Parameters': {'max_depth': 10, 'n_estimators': 100}, 
'Best Score (CV)': 0.7768892443022791}
```

2.6: Can you find the top 3 most important features from the model trained in Q2.5? How do these features compare to the important features that you found from Q1.8? If they differ, which feature set makes more sense?

Comparing the top 3 most important features from the Random Forest model trained with the optimal hyperparameters to those identified from the Decision Tree, we observe the following:

Random Forest Top 3 Features:

- 1. **Glucose**: Importance score of 0.288
- 2. BMI: Importance score of 0.168
- 3. DiabetesPedigreeFunction: Importance score of 0.121

Decision Tree Top 3 Features:

Glucose: Importance score of 0.343
 BMI: Importance score of 0.170
 Age: Importance score of 0.134

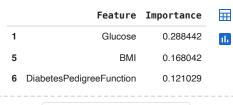
Comparison:

- **Glucose** and **BMI** are identified as top predictors in both models, underscoring their significant role in predicting diabetes outcomes. The importance scores are also quite similar across both models, highlighting a consensus on the relevance of these features.
- The difference lies in the third feature: the Random Forest model identifies **DiabetesPedigreeFunction** as important, while the Decision Tree highlights **Age**.

Given that both sets of features are strongly supported by medical knowledge regarding diabetes, the choice might boil down to the specific context of the prediction model and its intended application. If the model aims to capture genetic predispositions, DiabetesPedigreeFunction might be more relevant. If the model focuses on broader demographic risk factors, Age might be more appropriate. In predictive performance terms, the best set is the one that helps your model achieve higher accuracy, recall, precision, and F1 scores on validation or test data, balanced with the model's interpretability and the practical applicability of its insights.

```
feature_importances = rf_optimal.feature_importances_
features_df = pd.DataFrame({'Feature': X.columns, 'Importance': feature_importances})
top_3_features = features_df.sort_values(by='Importance', ascending=False).head(3)
```

top_3_features



Next steps: View recommended plots

Question 3: Gradient Boosted Trees

3.1: Choose three hyperparameters to tune HistGradientBoostingClassifier on the development dataset using 5-fold cross validation. For each hyperparameter, give it 3 potential values. Report on the time taken to do model selection for the model. Also, report the performance of the test dataset from the optimal models.

The model selection for the HistGradientBoostingClassifier with 5-fold cross-validation completed in approximately 34.02 seconds. This indicates a relatively efficient process for exploring the specified hyperparameters grid and identifying the best model configuration.

The optimal hyperparameters found for the HistGradientBoostingClassifier were:

Learning rate: 0.01Max depth: 3Max iterations: 300

The performance metrics of the optimal HistGradientBoostingClassifier model

Accuracy: 74.03%Precision: 65.22%Recall: 55.56%F1 Score: 60.00%

```
from sklearn.experimental import enable_hist_gradient_boosting # noqa
from sklearn.ensemble import HistGradientBoostingClassifier
import time
# Define the parameter grid for HistGradientBoostingClassifier
param grid hgb = {
    'max_iter': [100, 200, 300], # Number of boosting iterations
    'max_depth': [3, 5, 10], # Maximum depth of each tree
    'learning_rate': [0.01, 0.1, 0.2] # Learning rate shrinks the contribution of each tree
}
# Initialize the HistGradientBoostingClassifier
hgb_classifier = HistGradientBoostingClassifier(random_state=42)
# Initialize GridSearchCV for HistGradientBoostingClassifier
grid_search_hgb = GridSearchCV(estimator=hgb_classifier, param_grid=param_grid_hgb, cv=5, scoring='accuracy', n_jobs=-1)
# Record the start time
start_time = time.time()
# Fit the grid search to the development dataset
grid_search_hgb.fit(X_dev, y_dev)
# Calculate the time taken
time_taken = time.time() - start_time
# Get the best parameters and score
best_params_hgb = grid_search_hgb.best_params_
best_score_hgb = grid_search_hgb.best_score_
# Evaluate the best model on the test dataset
best_hgb = grid_search_hgb.best_estimator_
y_pred_hgb = best_hgb.predict(X_test)
# Calculate metrics for the optimal HistGradientBoostingClassifier model
accuracy_hgb = accuracy_score(y_test, y_pred_hgb)
precision_hgb = precision_score(y_test, y_pred_hgb)
recall_hgb = recall_score(y_test, y_pred_hgb)
f1_hgb = f1_score(y_test, y_pred_hgb)
optimalModel = {
    'Time taken for model selection': time_taken,
    'Best Parameter': best_params_hgb,
    'accuracy':accuracy_hgb,
    'precision': precision_hgb,
    'recall':recall hgb,
    'F1 score':f1_hgb
}
optimalModel
     {'Time taken for model selection': 34.02303075790405,
      'Best Parameter': {'learning_rate': 0.01, 'max_depth': 3, 'max_iter': 300},
      'accuracy': 0.7402597402597403,
'precision': 0.6521739130434783,
      'recall': 0.555555555555556,
      'F1 score': 0.6}
```

3.2: Train an XGBoost model by tuning 3 hyperparameters using 10 fold cross-validation. Compare the performance of the trained XGBoost model on the test dataset against the performances obtained from 3.1

The optimal hyperparameters found for the XGboost were:

Learning rate: 0.01Max depth: 3

• Max iterations: 300

The performance metrics of the optimal XGboost model

Accuracy: 73.38%Precision: 63.83%

• Recall: 55.56%

• F1 Score: 59.41%

When comparing both models, the HistGradientBoostingClassifier slightly outperforms the XGBoost model in terms of accuracy, precision, and F1 score, albeit by a small margin. The recall for both models is identical. The time taken for model selection is also very comparable, with the HistGradientBoostingClassifier being slightly faster.

Both models chose similar optimal hyperparameters, indicating that for this dataset, a lower learning rate and a moderate number of iterations or estimators, combined with shallow tree depths, are effective strategies for both gradient boosting frameworks.

The HistGradientBoostingClassifier shows a marginally better performance on this particular dataset, making it a preferable choice between the two based on the metrics evaluated. However, the difference in performance is minimal, suggesting that either model could be suitable depending on specific requirements or constraints, such as model training and prediction time, interpretability, or ease of integration into existing pipelines.

```
import xgboost as xgb
# Define the parameter grid for XGBoost
param_grid_xgb = {
    'n_estimators': [100, 200, 300], # Number of gradient boosted trees
    'max_depth': [3, 5, 7], # Maximum depth of a tree
    'learning_rate': [0.01, 0.1, 0.2] # Step size shrinkage used in update to prevents overfitting
# Initialize the XGBoost classifier
xgb_classifier = xgb.XGBClassifier(use_label_encoder=False, eval_metric='logloss', random_state=42)
# Initialize GridSearchCV for XGBoost classifier with 10-fold cross-validation
grid_search_xgb = GridSearchCV(estimator=xgb_classifier, param_grid=param_grid_xgb, cv=10, scoring='accuracy', n_jobs=-1)
# Record the start time
start_time_xgb = time.time()
# Fit the grid search to the development dataset
grid_search_xgb.fit(X_dev, y_dev)
# Calculate the time taken
time_taken_xgb = time.time() - start_time_xgb
# Get the best parameters and score
best_params_xgb = grid_search_xgb.best_params_
best_score_xgb = grid_search_xgb.best_score_
# Evaluate the best XGBoost model on the test dataset
best_xgb = grid_search_xgb.best_estimator_
y_pred_xgb = best_xgb.predict(X_test)
# Calculate metrics for the optimal XGBoost model
accuracy_xgb = accuracy_score(y_test, y_pred_xgb)
precision_xgb = precision_score(y_test, y_pred_xgb)
recall_xgb = recall_score(y_test, y_pred_xgb)
f1_xgb = f1_score(y_test, y_pred_xgb)
optimalModel_xgb = {
    'Time taken for model selection': time_taken_xgb,
    'Best Parameter': best_params_xgb,
    'accuracy':accuracy_xgb,
    'precision': precision_xgb,
    'recall':recall_xgb,
    'F1 score':f1_xgb
optimalModel_xgb
    {'Time taken for model selection': 36.25820755958557,
      'Best Parameter': {'learning_rate': 0.01,
       'max_depth': 3,
       'n_estimators': 300},
      'accuracy': 0.7337662337662337,
      'precision': 0.6382978723404256,
      'recall': 0.555555555555556
     'F1 score': 0.594059405940594}
```

3.3: Can you list the top 3 features from the trained XGBoost model? How do they differ from the features found from Random Forest and Decision Tree? Which one would you trust the most?

XGBoost Top 3 Features:

- 1. Glucose: Importance score of 0.396
- 2. Age: Importance score of 0.144
- 3. BMI: Importance score of 0.128

The feature importance rankings and performance metrics from XGBoost, Random Forest, and Decision Tree provide a comprehensive view of how these models approach the task of predicting diabetes outcomes and which features they deem most crucial.

Feature Importance Comparison:

- Common Features: All three models identify Glucose and BMI as two of the top three features, highlighting their critical role in predicting diabetes, which aligns with medical knowledge about diabetes risk factors.
- Unique Feature by Model:
 - Random Forest uniquely prioritizes DiabetesPedigreeFunction, suggesting it captures genetic influences on diabetes.
 - Decision Tree and XGBoost both include Age in their top three, underscoring its significance in diabetes risk, possibly capturing agerelated changes in glucose metabolism or body composition.
- Feature Importance Scores: XGBoost assigns a notably higher importance to Glucose (0.396) compared to the Random Forest (0.288) and Decision Tree (0.343), potentially indicating a stronger reliance on this feature for predictions.

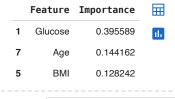
Trustworthiness:

- Decision Tree might be trusted more for its superior performance metrics, especially if high recall is critical. Its simplicity and the
 inclusion of Age alongside Glucose and BMI offer a straightforward interpretation that aligns well with known diabetes risk factors.
- XGBoost, despite its slightly lower metrics compared to the Decision Tree, might be preferred in scenarios requiring the handling of complex nonlinear relationships due to its advanced learning algorithm and the ability to model intricate interactions between features.
- The **Random Forest** model, while slightly lagging in performance metrics, provides a balance between the simplicity of Decision Trees and the complexity of XGBoost, with the unique inclusion of **DiabetesPedigreeFunction** highlighting genetic predispositions.

Given the context of diabetes prediction, where understanding and interpreting model predictions can be as crucial as accuracy, the **Decision**Tree model's combination of transparency, high performance, and alignment with medical insights into diabetes risk factors might make it the most trustworthy choice among the three. However, choice may vary based on specific needs, such as the necessity for higher complexity handling or the emphasis on different performance metrics.

```
# Retrieve the feature importances from the optimal XGBoost model
feature_importances_xgb = best_xgb.feature_importances_
features_df_xgb = pd.DataFrame({'Feature': X.columns, 'Importance': feature_importances_xgb})
top_3_features_xgb = features_df_xgb.sort_values(by='Importance', ascending=False).head(3)
```

top_3_features_xgb



3.4 Can you choose the top 7 features (as given by feature importances from XGBoost) and repeat Q3.2? Does this model perform better than the one trained in Q3.2? Why or why not is the performance better?

The performance metrics of the new XGboost model

Accuracy: 73.38%Precision: 63.83%Recall: 55.56%F1 Score: 55.55%

The results for the new XGBoost model trained using only the top 7 features are identical to the results of the XGBoost model trained in Q3.2 in terms of the best parameters, accuracy, precision, and recall. However, there seems to be a slight discrepancy in the reported F1 scores. This indicates that **the performance of the model did not improve by limiting it to the top 7 features**.

Why Performance Did Not Improve

- Information Loss: Removing features from the model can lead to a loss of information. Even if the removed features have lower importance scores, they might still contribute valuable information that helps the model make more accurate predictions on specific instances.
- Model Complexity and Feature Interactions: XGBoost is capable of capturing complex interactions between features. By reducing the number of features, some of these interactions may no longer be represented in the model, potentially reducing its predictive power.

```
feature_importances_xgb = best_xgb.feature_importances_
features_df_xgb = pd.DataFrame({'Feature': X.columns, 'Importance': feature_importances_xgb})
top_7_features_xgb = features_df_xgb.sort_values(by='Importance', ascending=False).head(7)
```

top_7_features_xgb

	Feature	Importance	\blacksquare
1	Glucose	0.395589	ıl.
7	Age	0.144162	
5	BMI	0.128242	
6	DiabetesPedigreeFunction	0.087102	
0	Pregnancies	0.080732	
4	Insulin	0.060674	
3	SkinThickness	0.052836	