DECISION TREE ASSIGNMENT

## Question 1: What is a Decision Tree, and how does it work in the context of classification?

Answer:

A Decision Tree is a flowchart-like structure used in machine learning for classification tasks, where it makes decisions based on a series of questions derived from the input data.

Overview of Decision Trees

A Decision Tree is a supervised learning algorithm that mimics human decision-making by breaking down a dataset into smaller subsets based on feature values. It consists of nodes and branches that represent decisions and their possible consequences. The main components of a Decision Tree include: Root Node: Represents the entire dataset and the first decision point. Decision Nodes: Internal nodes that represent tests on attributes (features) of the data.

Leaf Nodes: Terminal nodes that provide the final output or class label (e.g., "spam" or "not spam").

Branches: Lines connecting nodes, representing the outcome of a test and leading to the next node or leaf.

## How Decision Trees Work in Classification

Starting Point: The process begins at the root node, where the algorithm evaluates the dataset and selects the best feature to split the data based on a specific criterion (e.g., Gini impurity or information gain).

Splitting the Data: The tree asks a series of yes/no questions to divide the dataset into subsets. Each question corresponds to a feature, and the answers determine which branch to follow.

Recursive Process: This splitting continues recursively, creating new decision nodes and branches until a stopping criterion is met (e.g., maximum depth of the tree or minimum number of samples in a node).

Reaching Leaf Nodes: Once the tree can no longer split the data meaningfully, it reaches the leaf nodes, which provide the final classification for the input data based on the majority class of the samples in that node.

## Applications and Advantages

Classification Tasks: Decision Trees are widely used for tasks such as spam detection, customer segmentation, and medical diagnosis, where the output variable is categorical.

Interpretability: They are easy to understand and visualize, making them accessible for both beginners and professionals in machine learning.

No Feature Scaling Required: Decision Trees do not require normalization or scaling of features, simplifying the preprocessing steps

**Question 2: Explain the concepts of Gini Impurity and Entropy as impurity measures. How do they impact the splits in a Decision Tree?**

Answer:

**Gini Impurity**

Definition: Gini Impurity measures the likelihood of misclassifying a randomly chosen element from the dataset if it were labeled according to the distribution of labels in the subset. It is calculated using the formula:

[ Gini = 1 - \sum (p_i^2) ]

where $p_i$ $p$ iis the probability of a data point belonging to class i

Interpretation: The Gini Impurity value ranges from 0 to 1:

0 indicates a pure node (all instances belong to a single class).

1 indicates maximum impurity (instances are evenly distributed across classes).

Impact on Splits: In decision trees, the algorithm selects the feature that results in the lowest Gini Impurity after the split. This means that the chosen feature will create child nodes that are as pure as possible, leading to better classification accuracy.

**Entropy**

Definition: Entropy is a measure of the disorder or uncertainty in a dataset. It quantifies the impurity of a node by measuring the unpredictability of the class labels. The formula for calculating entropy is:

## ⌄   E n t r o p y

$-\sum (p_i \cdot \log_2 (p_i))$  Entropy$=-\sum(p_i \cdot \log_2(p_i))$

where $p_i$ $p$ iis the probability of a data point belonging to class i.

Interpretation: Similar to Gini Impurity, entropy values range from 0 to 1:

0 indicates a pure node (all instances belong to the same class).

1 indicates maximum disorder (instances are evenly distributed among multiple classes).

Impact on Splits: The decision tree algorithm uses entropy to determine the best feature for splitting. The feature that results in the highest information gain (the reduction in

entropy) is selected for the split. This helps in creating nodes that are more homogeneous in terms of class distribution.

**Comparison and Practical Implications**

Similarities: Both Gini Impurity and Entropy serve as effective metrics for measuring impurity in decision trees. They often yield similar results when determining the best splits.

Differences: Gini Impurity is generally faster to compute than entropy because it does not involve logarithmic calculations. However, both measures can lead to different tree structures in some cases, depending on the dataset and the distribution of classes.

Choosing Between Them: The choice between Gini Impurity and Entropy may depend on the specific application and the desired characteristics of the decision tree. Gini Impurity tends to create more balanced trees, while entropy may lead to deeper trees with more splits.

Understanding these concepts is crucial for data scientists and machine learning practitioners aiming to build robust decision tree models for various applications. By effectively utilizing Gini Impurity and Entropy, one can enhance the accuracy and interpretability of classification models.

**Question 3: What is the difference between Pre-Pruning and Post-Pruning in Decision Trees? Give one practical advantage of using each.**

Answer

**Pre-Pruning (Early Stopping)**

Definition: Pre-pruning halts the growth of the decision tree during its construction when certain conditions are met (e.g., maximum depth, minimum samples per split).

Goal: Prevent the tree from becoming too complex and overfitting the training data.

Practical Advantage: It leads to faster training and simpler models, which are easier to interpret and deploy in real-time systems

**Post-Pruning (Reduced Error Pruning)**

Definition: Post-pruning allows the tree to grow fully and then removes branches that do not contribute significantly to accuracy, often using validation data.

Goal: Improve generalization by simplifying the tree after evaluating its performance.

Practical Advantage: It often results in better accuracy on unseen data because pruning decisions are based on actual performance metrics

**Key Differences**

Timing: Pre-pruning occurs during tree construction, while post-pruning is applied after the tree is fully grown.

Efficiency: Pre-pruning is faster and more efficient, especially for large datasets, as it avoids unnecessary splits. Post-pruning is computationally intensive as it evaluates multiple pruned versions of the tree.

Risk of Underfitting: Pre-pruning may lead to underfitting if the tree's growth is stopped too early. Post-pruning minimizes this risk by starting with a fully grown tree.

Suitability: Pre-pruning is better suited for large datasets, while post-pruning is often preferred for smaller datasets where computational cost is less of a concern.

**Question 4: What is Information Gain in Decision Trees, and why is it important for choosing the best split?** Answer:

At its core, Information Gain (IG) is about reducing uncertainty. You might be wondering, "How does it work in a decision tree?" It's simple — each time we split the data, we want to make things clearer, not more confusing. Information Gain helps measure exactly how much clarity (or certainty) a split provides.

Here's the deal: Information Gain is the difference between the uncertainty (entropy) before and after a split. In other words, it's like saying, "Did this split help us make better predictions?" The higher the Information Gain, the better the split.

**IG(S, A)**: This is the Information Gain of splitting the dataset SSS based on attribute A. It's what we're trying to maximize at each node.

**Entropy(S)**: This is the entropy of the entire dataset SSS before the split. Remember, higher entropy means more uncertainty, and our goal is to reduce this uncertainty.

**\*\*|Sv||S|\*\***: This term represents the proportion of data points that take on a specific value vvv for attribute A. In simpler terms, it's the fraction of the dataset that goes down one path of the decision tree.

**Entropy(S_v)**: This is the entropy of the subset Sv, which includes only the data points that match value v of the attribute A. We calculate entropy for each subset after the split.

Essentially, Information Gain tells us how much uncertainty is removed after splitting the data based on attribute A.

### Why Is It Important for Choosing the Best Split?

Prioritizes informative features: Features that result in purer child nodes (lower entropy) yield higher information gain.

Guides tree construction: The feature with the highest IG is selected for splitting, ensuring the tree grows in a way that maximizes classification accuracy.

Reduces uncertainty: By choosing splits that reduce entropy the most, the tree becomes more efficient and interpretable

**Question 5: What are some common real-world applications of Decision Trees, and what are their main advantages and limitations?**

**Real-World Applications**

Medical Diagnosis: Decision trees are employed to predict diseases by analyzing symptoms, test results, and patient history. For instance, they can help doctors determine the likelihood of a patient having a specific condition based on various health indicators.

Customer Churn Prediction: Businesses use decision trees to identify customers who are likely to stop using their services. By analyzing customer behavior and demographics, companies can take proactive measures to retain these customers.

Loan Approval in Banking: Financial institutions utilize decision trees to assess loan applications. They evaluate factors such as income, credit history, and employment status to determine the eligibility of applicants.

Fraud Detection: Decision trees help in detecting fraudulent activities by analyzing transaction patterns. They can identify unusual behaviors that may indicate fraud, allowing companies to take preventive actions.

Student Performance Prediction: Educational institutions use decision trees to predict students' academic performance based on factors like study habits, attendance, and past results. This helps in identifying students who may need additional support.

**Advantages of Decision Trees**

Interpretability: Decision trees provide a clear and visual representation of the decision-making process, making it easy for stakeholders to understand how decisions are made.

Flexibility: They can handle both categorical and numerical data, making them versatile for various applications across different industries.

Minimal Data Preparation: Decision trees require less data preprocessing compared to other machine learning models, allowing for quicker implementation.

Non-Parametric: They do not assume any underlying distribution of the data, making them suitable for a wide range of datasets.

**Limitations of Decision Trees**

Overfitting: Decision trees can easily become too complex and overfit the training data, leading to poor generalization on unseen data. This is particularly a concern with deep trees.

Instability: Small changes in the data can lead to different tree structures, making them sensitive to variations in the dataset.

Bias towards Dominant Classes: Decision trees can be biased towards classes that have a larger number of instances, which may lead to suboptimal performance on minority classes.

Limited Predictive Power: While they are useful for interpretation, decision trees may not always provide the best predictive accuracy compared to ensemble methods like Random Forests.

In summary, Decision Trees are powerful tools for decision-making across various domains, offering clear advantages in interpretability and flexibility, while also facing challenges such as overfitting and instability. Their applications in real-world scenarios demonstrate their value in enhancing decision-making processes.

Dataset Info:

● Iris Dataset for classification tasks (sklearn.datasets.load_iris() or provided CSV).

● Boston Housing Dataset for regression tasks (sklearn.datasets.load_boston() or provided CSV).

Question 6: Write a Python program to:

● Load the Iris Dataset

● Train a Decision Tree Classifier using the Gini criterion

● Print the model's accuracy and feature importances

```python
# Import necessary libraries
from sklearn.datasets import load_iris
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
import pandas as pd

# Load the Iris dataset
iris = load_iris()
X = iris.data  # Features
y = iris.target  # Target labels

# Split the dataset into training and testing sets (80% train, 20% test)
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random

# Initialize the Decision Tree Classifier with the Gini criterion
clf = DecisionTreeClassifier(criterion='gini', random_state=42)

# Train the classifier on the training data
clf.fit(X_train, y_train)

# Predict the labels for the test set
y_pred = clf.predict(X_test)
```

```python
# Calculate and print the accuracy of the model
accuracy = accuracy_score(y_test, y_pred)
print(f"Model Accuracy: {accuracy:.2f}")

# Print the feature importances
feature_importances = clf.feature_importances_
feature_names = iris.feature_names
importance_df = pd.DataFrame({
    'Feature': feature_names,
    'Importance': feature_importances
}).sort_values(by='Importance', ascending=False)

print("\nFeature Importances:")
print(importance_df)
```

```
Model Accuracy: 1.00

Feature Importances:
            Feature  Importance
2  petal length (cm)    0.906143
3   petal width (cm)    0.077186
1   sepal width (cm)    0.016670
0  sepal length (cm)    0.000000
```

Question 7: Write a Python program to:

● Load the Iris Dataset

● Train a Decision Tree Classifier with max_depth=3 and compare its accuracy to a fully-grown tree.

```python
# Import necessary libraries
from sklearn.datasets import load_iris
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
from sklearn.tree import export_text

# Load the Iris dataset
iris = load_iris()
X = iris.data  # Features
y = iris.target  # Target labels

# Split the dataset into training and testing sets (80% train, 20% test)
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random

# Train a Decision Tree Classifier with max_depth=3
dtree_limited = DecisionTreeClassifier(max_depth=3, random_state=42)
dtree_limited.fit(X_train, y_train)

# Predict and calculate accuracy for the limited-depth tree
y_pred_limited = dtree_limited.predict(X_test)
```

```
accuracy_limited = accuracy_score(y_test, y_pred_limited)

# Train a fully-grown Decision Tree Classifier (no depth limit)
dtree_full = DecisionTreeClassifier(random_state=42)
dtree_full.fit(X_train, y_train)

# Predict and calculate accuracy for the fully-grown tree
y_pred_full = dtree_full.predict(X_test)
accuracy_full = accuracy_score(y_test, y_pred_full)

# Display results
print("Accuracy of Decision Tree with max_depth=3:", accuracy_limited)
print("Accuracy of Fully-Grown Decision Tree:", accuracy_full)

# Optional: Display the structure of the trees
print("\nDecision Tree with max_depth=3:")
print(export_text(dtree_limited, feature_names=iris.feature_names))

print("\nFully-Grown Decision Tree:")
print(export_text(dtree_full, feature_names=iris.feature_names))
```

```
Accuracy of Decision Tree with max_depth=3: 1.0
Accuracy of Fully-Grown Decision Tree: 1.0

Decision Tree with max_depth=3:
|--- petal length (cm) <= 2.45
|   |--- class: 0
|--- petal length (cm) >  2.45
|   |--- petal length (cm) <= 4.75
|   |   |--- petal width (cm) <= 1.65
|   |   |   |--- class: 1
|   |   |--- petal width (cm) >  1.65
|   |   |   |--- class: 2
|   |--- petal length (cm) >  4.75
|   |   |--- petal width (cm) <= 1.75
|   |   |   |--- class: 1
|   |   |--- petal width (cm) >  1.75
|   |   |   |--- class: 2


Fully-Grown Decision Tree:
|--- petal length (cm) <= 2.45
|   |--- class: 0
|--- petal length (cm) >  2.45
|   |--- petal length (cm) <= 4.75
|   |   |--- petal width (cm) <= 1.65
|   |   |   |--- class: 1
|   |   |--- petal width (cm) >  1.65
|   |   |   |--- class: 2
|   |--- petal length (cm) >  4.75
|   |   |--- petal width (cm) <= 1.75
|   |   |   |--- petal length (cm) <= 4.95
|   |   |   |   |--- class: 1
|   |   |   |--- petal length (cm) >  4.95
|   |   |   |   |--- petal width (cm) <= 1.55
|   |   |   |   |   |--- class: 2
|   |   |   |--- petal width (cm) >  1.55
```

```
|   |   |   |   |   |   |--- petal length (cm) <= 5.45
|   |   |   |   |   |   |    |--- class: 1
|   |   |   |   |   |   |--- petal length (cm) >  5.45
|   |   |   |   |   |   |    |--- class: 2
|   |   |--- petal width (cm) >  1.75
|   |   |   |--- petal length (cm) <= 4.85
|   |   |   |   |--- sepal width (cm) <= 3.10
|   |   |   |   |   |--- class: 2
|   |   |   |   |--- sepal width (cm) >  3.10
|   |   |   |   |   |--- class: 1
|   |   |   |--- petal length (cm) >  4.85
|   |   |   |   |--- class: 2
```

Question 8: Write a Python program to:

● Load the California Housing dataset from sklearn

● Train a Decision Tree Regressor

● Print the Mean Squared Error (MSE) and feature importances

```python
# Import necessary libraries
from sklearn.datasets import fetch_california_housing
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeRegressor
from sklearn.metrics import mean_squared_error
import numpy as np

# Step 1: Load the California Housing dataset
data = fetch_california_housing(as_frame=True)  # Load dataset as a pandas Data
X = data.data  # Features
y = data.target  # Target variable (median house value)

# Step 2: Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random

# Step 3: Train a Decision Tree Regressor
regressor = DecisionTreeRegressor(random_state=42)  # Initialize the regressor
regressor.fit(X_train, y_train)  # Train the model

# Step 4: Make predictions on the test set
y_pred = regressor.predict(X_test)

# Step 5: Calculate the Mean Squared Error (MSE)
mse = mean_squared_error(y_test, y_pred)

# Step 6: Print the feature importances
feature_importances = regressor.feature_importances_

# Display results
print("Mean Squared Error (MSE):", mse)
print("\nFeature Importances:")
for feature, importance in zip(data.feature_names, feature_importances):
```

```
        print(f"{feature}: {importance:.4f}")
```

```
Mean Squared Error (MSE): 0.495235205629094

Feature Importances:
MedInc: 0.5285
HouseAge: 0.0519
AveRooms: 0.0530
AveBedrms: 0.0287
Population: 0.0305
AveOccup: 0.1308
Latitude: 0.0937
Longitude: 0.0829
```

Question 9: Write a Python program to:

● Load the Iris Dataset

● Tune the Decision Tree's max_depth and min_samples_split using GridSearchCV

● Print the best parameters and the resulting model accuracy

```python
# Import necessary libraries
from sklearn.datasets import load_iris
from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import accuracy_score

# Load the Iris dataset
iris = load_iris()
X = iris.data  # Features
y = iris.target  # Target labels

# Split the dataset into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random

# Define the Decision Tree classifier
dt = DecisionTreeClassifier(random_state=42)

# Define the hyperparameter grid for tuning
param_grid = {
    'max_depth': [2, 3, 4, 5, 6, None],  # Possible depths of the tree
    'min_samples_split': [2, 5, 10, 20]  # Minimum samples required to split a
}

# Use GridSearchCV to find the best hyperparameters
grid_search = GridSearchCV(estimator=dt, param_grid=param_grid, cv=5, scoring='
grid_search.fit(X_train, y_train)

# Get the best parameters and the best score
best_params = grid_search.best_params_
best_model = grid_search.best_estimator_
best_score = grid_search.best_score_
```

```
# Evaluate the best model on the test set
y_pred = best_model.predict(X_test)
test_accuracy = accuracy_score(y_test, y_pred)

# Print the results
print("Best Parameters:", best_params)
print("Best Cross-Validation Accuracy:", best_score)
print("Test Set Accuracy:", test_accuracy)
```

```
Best Parameters: {'max_depth': 3, 'min_samples_split': 2}
Best Cross-Validation Accuracy: 0.9523809523809523
Test Set Accuracy: 0.9777777777777777
```

Question 10: Imagine you're working as a data scientist for a healthcare company that wants to predict whether a patient has a certain disease. You have a large dataset with mixed data types and some missing values. Explain the step-by-step process you would follow to:

● Handle the missing values

● Encode the categorical features

● Train a Decision Tree model

● Tune its hyperparameters

● Evaluate its performance And describe what business value this model could provide in the real-world setting.

Answer Step-by-Step Process

**1. Handle Missing Values**

Identify missing data: Use .isnull().sum() to check missing values per column.

Imputation strategies:

Numerical features: Use mean, median, or model-based imputation (e.g., KNN imputer).

Categorical features: Use mode imputation or create a new category like "Unknown".

Drop columns/rows: If a column has >50% missing data or rows are sparse, consider dropping them.

**2. Encode Categorical Features**

Label Encoding: For ordinal features (e.g., severity levels: mild < moderate < severe).

One-Hot Encoding: For nominal features (e.g., gender, ethnicity, symptoms).

Use pd.get_dummies() or OneHotEncoder from scikit-learn.

Avoid dummy variable trap: Drop one column if using linear models, but not necessary for Decision Trees.

### 3.Train a Decision Tree Model

Split data: Use train_test_split() to divide into training and testing sets (e.g., 80/20).

Model training:

python from sklearn.tree import DecisionTreeClassifier model = DecisionTreeClassifier(random_state=42) model.fit(X_train, y_train)

### 4.Tune Hyperparameters

Key hyperparameters:

max_depth: Prevents overfitting by limiting tree depth.

min_samples_split: Minimum samples required to split a node.

min_samples_leaf: Minimum samples required at a leaf node.

criterion: Choose between "gini" or "entropy".

### 5.Evaluate Performance

Metrics:

Accuracy: Overall correctness.

Precision & Recall: Especially important in healthcare (false positives vs false negatives).

F1 Score: Harmonic mean of precision and recall.

ROC-AUC: Measures model's ability to distinguish between classes.

### Business Value in Real-World Healthcare

Early Detection: Helps identify high-risk patients before symptoms escalate.

Resource Optimization: Prioritizes patients for testing, treatment, or specialist referral.

Personalized Care: Enables tailored treatment plans based on predicted risk.

Cost Reduction: Reduces unnecessary procedures and hospitalizations.

Compliance & Reporting: Supports regulatory reporting and quality assurance.