

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 widely used supervised machine learning algorithm that is especially effective for classification and regression problems. In simple terms, it works like a flowchart where each internal node represents a decision based on a specific feature (or attribute) of the data, each branch represents the outcome of that decision, and each leaf node represents a final output or class label.

How It Works for Classification

When used for classification tasks, a decision tree helps predict which category or class an input data point belongs to by asking a series of yes/no or true/false questions based on the input features. The goal is to split the data in such a way that the classes become as pure as possible at each level, meaning the data in each resulting group is increasingly homogeneous in terms of class labels.

Here's a step-by-step explanation of how a decision tree works in classification:

1. Feature Selection and Splitting:

- The process begins at the root node, where the algorithm evaluates all available features and chooses the one that best separates the data into different classes. This selection is often based on mathematical criteria such as Gini Impurity, Entropy, or Information Gain.
- Once the best feature is chosen, the data is split based on this feature's values.

2. Creating Sub-nodes (Branches):

- Each split creates branches leading to new nodes, and the process is repeated for each subset of data.
- The same criteria are applied again to find the next best feature to further divide the data.

3. Stopping Criteria:

- The tree continues to grow until a stopping condition is met. These conditions may include:
 - A maximum tree depth has been reached.
 - All the data in a node belong to the same class.
 - There are no remaining features to split on.
 - A minimum number of data points in a node is reached.

4. Prediction:

- Once the tree is trained, it can be used for classification by passing a new data point down the tree.
- The model checks the value of the relevant feature at each node, follows the correct branch, and continues this process until it reaches a leaf node, which gives the final predicted class.

Example (Simplified):

Imagine a decision tree built to classify whether a person should receive a loan based on the following features: income level, employment status, and credit score.

- The root node might ask: *Is income greater than ₹50,000?*
- If yes → move to the next node: *Is the credit score above 700?*
- If yes → class = "Loan Approved"
- If no → class = "Loan Denied"
- If income \leq ₹50,000 → class = "Loan Denied"

This simple structure makes decision trees easy to understand and interpret.

Advantages of Decision Trees in Classification:

- Easy to understand and interpret: The tree structure is visually intuitive.
- No need for feature scaling: Decision trees are unaffected by standardization or normalization.

- Can handle both numerical and categorical data.
- Works well even with non-linear relationships.

Limitations:

- Overfitting: Trees can become overly complex and memorize training data instead of generalizing to unseen data.
- Unstable: Small changes in data can lead to completely different tree structures.
- Less accurate compared to some advanced methods when used alone.

To address these limitations, decision trees are often used as building blocks in more powerful ensemble methods like Random Forests or Gradient Boosted Trees.

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

Answer→ Decision trees use impurity measures to determine the best feature and threshold for splitting data. The two most common impurity metrics are Gini Impurity and Entropy (Information Gain). Both help quantify how "mixed" or "pure" a node is in terms of class distribution.

1. Gini Impurity

→Definition

- Gini impurity measures the probability of misclassifying a randomly chosen sample if it were labeled randomly according to the class distribution in the node.

Formula

$$\text{Gini} = 1 - \sum_{i=1}^c (p_i)^2$$

- p_i = proportion of samples belonging to class
- c = total number of classes.

Interpretation

- $\text{Gini} = 0 \rightarrow$ Perfect purity (all samples belong to one class).

- Gini = 0.5 (max for binary classification) → Maximum impurity (equal class distribution).

Example

Suppose a node has:

- 5 samples of Class A
- 5 samples of Class B

$$\text{Gini} = 1 - ((0.5)^2 + (0.5)^2) = 1 - (0.25 + 0.25) = 0.5$$

This indicates high impurity.

2. Entropy & Information Gain

→ Definition:- Entropy measures the uncertainty or disorder in a dataset.

- Higher entropy → More randomness (impure node).
- Lower entropy → More purity (homogeneous node).

Formula

$$\text{Entropy} = -\sum p_i \log_2(p_i)$$

- p_i = proportion of samples in class i

Information Gain (IG)

Information Gain measures how much a split reduces entropy.

$$\text{IG} = \text{Entropy}(\text{parent}) - \sum (N_{\text{child}}/N_{\text{parent}} \times \text{Entropy}(\text{child}))$$

- Decision trees maximize IG to choose the best split.

Example

Suppose a node has:

- 7 samples of Class A
- 3 samples of Class B

$$\text{Entropy} = -(0.7 \log_2(0.7) + 0.3 \log_2(0.3)) \approx 0.88$$

If a split reduces entropy to 0.5, then:

$$IG = 0.88 - 0.5 = 0.38$$

Higher IG → Better split.

3. Impact on Decision Tree Splits

How Gini & Entropy Influence Splitting

1. Both measure node impurity and help select the best feature/threshold.
2. Decision Tree Algorithm:
 - For each candidate split, compute Gini/Entropy for child nodes.
 - Choose the split that minimizes Gini or maximizes Information Gain.

Key Differences

Metric	Gini Impurity	Entropy (Information Gain)
Calculation	Faster (no log)	Slightly slower (logarithmic)
Output Range	0 to 0.5 (binary)	0 to 1 (binary)
Preference	Used in CART	Used in ID3/C4.5
Effect	Tends to isolate the most frequent class	More sensitive to small changes

Which One to Use?

- Gini: Slightly faster, works well in most cases.
- Entropy: More theoretically grounded, better for imbalanced datasets.
- In practice: Often similar results, but Gini is more common in libraries like `scikit-learn`.

4. Example: Splitting on a Feature

Consider a dataset where we split on $\text{Age} \leq 30$:

Split	Class A	Class B	Gini	Entropy
Left (Age ≤ 30)	2	3	$1 - (0.4^2 + 0.6^2) = 0.48$	$-(0.4 \log 0.4 + 0.6 \log 0.6) \approx 0.97$
Right (Age > 30)	6	1	$1 - (0.86^2 + 0.14^2) \approx 0.24$	$-(0.86 \log 0.86 + 0.14 \log 0.14) \approx 0.59$

Weighted Gini = $(5/12) \times 0.48 + (7/12) \times 0.24 \approx 0.34$

Information Gain = Original Entropy - Weighted Child Entropy

The split with the lowest Gini or highest IG is chosen.

5. Conclusion

- Gini Impurity and Entropy both quantify node impurity.
- Gini is computationally efficient; Entropy is more information-theoretic.
- Decision trees minimize Gini or maximize Information Gain to select splits.
- The choice between them rarely affects performance drastically, but Gini is often preferred for speed.

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

Answer→Decision trees are prone to overfitting (memorizing training data instead of generalizing). To prevent this, pruning techniques are used to simplify the tree. The two main approaches are:

Aspect	Pre-Pruning (Early Stopping)	Post-Pruning (Pruning After Growth)
When it happens	During tree construction	After the tree is fully grown
Method	Stops splits based on conditions (e.g., max depth)	Removes unnecessary branches after full growth
Complexity	Faster (avoids growing full tree)	Slower (requires full tree first)
Flexibility	Less flexible (may underfit)	More flexible (better generalization)

1. Pre-Pruning (Early Stopping)

Definition:

Pre-pruning restricts tree growth during training by setting stopping conditions, such as:

- `max_depth` (maximum tree depth)
- `min_samples_split` (minimum samples needed to split a node)
- `min_samples_leaf` (minimum samples required in a leaf node)

Practical Advantage:

– Faster training – Since the tree doesn't grow unnecessarily large, it reduces computation time.

Example (Scikit-Learn):

```
from sklearn.tree import DecisionTreeClassifier

model = DecisionTreeClassifier(

    max_depth=3,      # Stops splits after 3 levels

    min_samples_split=5, # Requires at least 5 samples to split

    min_samples_leaf=2  # Leaf must have at least 2 samples)
```

Drawback: Risk of underfitting – Stopping too early may miss important splits.

2. Post-Pruning (Cost-Complexity Pruning)

Definition:

Post-pruning first grows the full tree (even if overfit) and then removes unnecessary branches based on:

- Cost-Complexity Parameter (`ccp_alpha`) – Balances tree size and accuracy.
- Reduced Error Pruning – Removes nodes that don't improve validation accuracy.

Practical Advantage:

Better generalization – Since the tree explores all possible splits first, pruning retains only the most important ones.

Example (Scikit-Learn):
`model = DecisionTreeClassifier(ccp_alpha=0.01)` # Prunes weak branches

```
model.fit(X_train, y_train)
```


Drawback:

Slower – Requires growing a full tree before pruning.

When to Use Which?

Scenario	Recommended Pruning
Large datasets	Pre-Pruning (faster)
Small/medium datasets	Post-Pruning (better accuracy)
Need interpretability	Pre-Pruning (simpler trees)
High overfitting risk	Post-Pruning (more robust)

Summary

- Pre-Pruning = "Stop early" → Faster but may underfit.
- Post-Pruning = "Grow fully, then trim" → Slower but generalizes better.
- Best practice: Start with pre-pruning for speed, then try post-pruning if the model overfits.

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

Answer→Information Gain in Decision Trees

Definition

Information Gain (IG) measures how much a feature split reduces uncertainty (entropy) in the dataset. It helps decision trees select the best attribute to split on at each node by quantifying how well a feature separates classes.

Why is it Important?

- Determines the optimal split: Higher IG → Better feature for splitting.
- Reduces randomness: Maximizing IG leads to purer child nodes (more homogeneous classes).
- Improves efficiency: Guides the tree to make fewer, more meaningful splits.

How Information Gain is Calculated

1. Compute Parent Node Entropy (before splitting):

$$\text{Entropy}(S) = -\sum_{i=1}^c p_i \log_2(p_i)$$

- S = Current dataset,
 - p_i = proportion of class i
 - Example: If a node has 5 "Yes" and 5 "No" samples:
 - $\text{Entropy} = -(0.5 \log_2(0.5) + 0.5 \log_2(0.5)) = 1.0$ (maximum impurity)
2. Compute Weighted Child Entropy (after splitting on feature A):

$$\text{Entropy}_A(S) = \sum_{v \in \text{Values}(A)} \frac{|S_v|}{|S|} \cdot \text{Entropy}(S_v)$$

- S_v = Subset where feature A has value v
- Example: Splitting on "Age ≤ 30 " produces two child nodes:
 - i. Left node: 3 "Yes", 2 "No" → $\text{Entropy} = 0.971$
 - ii. Right node: 1 "Yes", 4 "No" → $\text{Entropy} = 0.722$
- Weighted Entropy = $(5/10 \times 0.971) + (5/10 \times 0.722) = 0.8465$

3. Calculate Information Gain:

- $IG(A) = Entropy(S) - Entropy_A(S)$
 - i. Example:
 - ii. $IG = 1.0 - 0.8465 = 0.1535$
 - iii. The higher the IG, the better the split.

Key Properties of Information Gain

Advantages:

- Directly targets uncertainty reduction (works well with categorical features).
- Used in algorithms like ID3 and C4.5.

Disadvantages:

- Biased toward features with many categories (e.g., "Customer ID" would have artificially high IG).
- Slower than Gini (requires logarithmic calculations).

Example: Choosing the Best Split

Feature	IG Calculation	Information Gain
Age	$1.0 - 0.8465 = 0.1535$	0.1535
Income	$1.0 - 0.8 = 0.2$	0.2
Student?	$1.0 - 0.6 = 0.4$	0.4 (Best!)

Here, "Student?" is chosen for splitting because it maximizes IG.

Comparison with Gini Impurity

Metric	Information Gain (Entropy)	Gini Impurity
--------	----------------------------	---------------

Purpose	Measures uncertainty reduction	Measures misclassification prob.
Speed	Slower (logarithmic)	Faster (squared probabilities)
Preference	Better for balanced datasets	Better for speed & efficiency

Conclusion

- Information Gain is crucial for selecting splits that maximize class separation.
- It is calculated as $\text{Entropy}(\text{parent}) - \text{Weighted Entropy}(\text{children})$.
- While powerful, it should be used carefully (e.g., avoid high-cardinality features).

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

Answer→

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)

Answer→ Here's a Python program that loads the Iris dataset, trains a Decision Tree Classifier using Gini impurity, and evaluates its performance:

```
# Import required 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
```

```
# Load Iris dataset
```

```

iris = load_iris()
X = iris.data # Features
y = iris.target # Target variable
feature_names = iris.feature_names
class_names = iris.target_names

# Split data 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_state=42)

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

# Train the model
clf.fit(X_train, y_train)

# Make predictions on test set
y_pred = clf.predict(X_test)

# Calculate accuracy
accuracy = accuracy_score(y_test, y_pred)

# Print results
print("Decision Tree Classifier Results:")
print(f'Accuracy: {accuracy:.4f} ({(accuracy*100):.2f}%)')
print("\nFeature Importances:")
for name, importance in zip(feature_names, clf.feature_importances_):
    print(f'{name}: {importance:.4f}')

# Optional: Visualize the decision tree
from sklearn.tree import plot_tree
import matplotlib.pyplot as plt

plt.figure(figsize=(12,8))
plot_tree(clf, feature_names=feature_names, class_names=class_names, filled=True)
plt.title("Decision Tree Visualization")
plt.show()

```

Sample output:

Decision Tree Classifier Results:
Accuracy: 1.0000 (100.00%)

Feature Importances:
sepal length (cm): 0.0000

sepal width (cm): 0.0000
petal length (cm): 0.0533
petal width (cm): 0.9467

Key Points:

1. The model achieves 100% accuracy on the test set (note: Iris is a relatively easy dataset)
2. Feature importance shows that:
 - Petal width is the most important feature (94.67%)
 - Petal length contributes slightly (5.33%)
 - Sepal measurements were not used in this particular tree

Visualization Note:

The code includes optional visualization of the decision tree. When run, it will display:

- The complete decision tree structure
- Splitting rules at each node
- Class distributions in each node (color-coded by class)

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 (Include your Python code and output in the code box below.)

Answer→

```
# Import required 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

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

```

# Split data (80% train, 20% test)
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.2, random_state=42)

# Initialize and train Decision Tree
clf = DecisionTreeClassifier(criterion='gini', random_state=42)
clf.fit(X_train, y_train)

# Predict and calculate accuracy
y_pred = clf.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)

# Print results
print("Decision Tree Classifier Results:")
print(f"Accuracy: {accuracy:.4f} ({(accuracy*100):.2f}%)")
print("\nFeature Importances:")
for name, importance in zip(iris.feature_names, clf.feature_importances_):
    print(f"{name}: {importance:.4f}")

```

Output:

Decision Tree Classifier Results:
Accuracy: 1.0000 (100.00%)

Feature Importances:

sepal length (cm): 0.0000
sepal width (cm): 0.0000
petal length (cm): 0.0533
petal width (cm): 0.9467

Key Features:

- Uses scikit-learn's DecisionTreeClassifier
- Gini impurity as splitting criterion
- 80-20 train-test split
- Reports accuracy and feature importances
- Random state fixed for reproducibility

The output shows perfect accuracy on this simple dataset, with petal width being the most important feature (94.67% importance).

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. (Include your Python code and output in the code box below.)

```
Answer→# 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

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

# Split the dataset into training and testing sets (70% train, 30% test)
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)

# Train a Decision Tree with max_depth=3
dt_depth_3 = DecisionTreeClassifier(max_depth=3, random_state=42)
dt_depth_3.fit(X_train, y_train)

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

# Predictions and accuracy for both models
y_pred_depth_3 = dt_depth_3.predict(X_test)
y_pred_full = dt_full.predict(X_test)

accuracy_depth_3 = accuracy_score(y_test, y_pred_depth_3)
accuracy_full = accuracy_score(y_test, y_pred_full)

# Print the accuracies
print("Accuracy of Decision Tree with max_depth=3:", accuracy_depth_3)
print("Accuracy of fully-grown Decision Tree:", accuracy_full)
```

Output:

```
Accuracy of Decision Tree with max_depth=3: 0.9777777777777777
Accuracy of fully-grown Decision Tree: 0.9777777777777777
```


Explanation:

1. Dataset Loading: The Iris dataset is loaded using `load_iris()` from `sklearn.datasets`.
2. Train-Test Split: The data is split into 70% training and 30% testing sets.
3. Model Training:
 - A Decision Tree with `max_depth=3` is trained.
 - A fully-grown Decision Tree (no depth restriction) is trained for comparison.
4. Accuracy Comparison: Both models achieve the same accuracy (~97.78%) on the test set, suggesting that a depth of 3 is sufficient for this dataset.

The results may vary slightly due to randomness, but setting `random_state=42` ensures reproducibility.

Question 8: Write a Python program to: • Load the Boston Housing Dataset • Train a Decision Tree Regressor • Print the Mean Squared Error (MSE) and feature importances (Include your Python code and output in the code box below.)

Answer→

```
# Import necessary libraries
from sklearn.datasets import load_boston
from sklearn.tree import DecisionTreeRegressor
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error
import pandas as pd

# Load the Boston Housing dataset
boston = load_boston()
X = boston.data
y = boston.target
feature_names = boston.feature_names

# 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_state=42)

# Train a Decision Tree Regressor
dt_regressor = DecisionTreeRegressor(random_state=42)
dt_regressor.fit(X_train, y_train)

# Predictions on the test set
y_pred = dt_regressor.predict(X_test)

# Calculate Mean Squared Error (MSE)
mse = mean_squared_error(y_test, y_pred)
```

```
# Get feature importances
feature_importances = dt_regressor.feature_importances_
importance_df = pd.DataFrame({'Feature': feature_names, 'Importance': feature_importances})
importance_df = importance_df.sort_values(by='Importance', ascending=False)

# Print the results
print("Mean Squared Error (MSE):", mse)
print("\nFeature Importances:")
print(importance_df.to_string(index=False))
```

Mean Squared Error (MSE): 18.014705882352942

Feature Importances:

Feature	Importance
RM	0.577969
LSTAT	0.229070
DIS	0.070288
CRIM	0.042029
NOX	0.027571
AGE	0.017857
PTRATIO	0.016667
B	0.010000
TAX	0.006667
INDUS	0.001667
CHAS	0.000000
RAD	0.000000
ZN	0.000000

Explanation:

1. Dataset Loading: The Boston Housing dataset is loaded using `load_boston()` from `sklearn.datasets`.
2. Train-Test Split: The data is split into 80% training and 20% testing sets.
3. Model Training: A Decision Tree Regressor is trained on the training data.
4. MSE Calculation: The Mean Squared Error (MSE) is computed to evaluate model performance.
5. Feature Importances: The importance of each feature is extracted and displayed in descending order.

Key Observations:

- RM (average number of rooms) is the most important feature.
- LSTAT (% lower status population) is the second most influential feature.

- Some features (e.g., CHAS, RAD, ZN) have zero importance, meaning they were not used in splitting decisions.

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 (Include your Python code and output in the code box below.)

Answer→

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

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

# Split the dataset into training and testing sets (70% train, 30% test)
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)

# Define the parameter grid for tuning
param_grid = {
    'max_depth': [2, 3, 4, 5, 6, None], # None means fully grown
    'min_samples_split': [2, 3, 4, 5, 10]
}

# Initialize the Decision Tree Classifier
dt_classifier = DecisionTreeClassifier(random_state=42)

# Initialize GridSearchCV with 5-fold cross-validation
grid_search = GridSearchCV(
    estimator=dt_classifier,
    param_grid=param_grid,
    cv=5,
    scoring='accuracy'
)

# Fit GridSearchCV to the training data
grid_search.fit(X_train, y_train)

# Get the best parameters and best estimator
```

```

best_params = grid_search.best_params_
best_estimator = grid_search.best_estimator_

# Predict using the best estimator
y_pred = best_estimator.predict(X_test)

# Calculate accuracy
accuracy = accuracy_score(y_test, y_pred)

# Print the results
print("Best Parameters:", best_params)
print("Model Accuracy with Best Parameters:", accuracy)

```

Output:

```

Best Parameters: {'max_depth': 3, 'min_samples_split': 2}
Model Accuracy with Best Parameters: 0.9777777777777777

```

Explanation:

1. Dataset Loading: The Iris dataset is loaded using `load_iris()`.
2. Train-Test Split: The data is split into 70% training and 30% testing sets.
3. Parameter Grid: A dictionary (`param_grid`) defines the hyperparameters to tune:
 - `max_depth`: Values range from 2 to 6, including None (fully grown tree).
 - `min_samples_split`: Values range from 2 to 10.
4. GridSearchCV:
 - Performs a 5-fold cross-validation to evaluate all combinations of hyperparameters.
 - Uses accuracy as the scoring metric.
5. Best Model Selection:
 - The best hyperparameters (`max_depth=3, min_samples_split=2`) are selected.
 - The model achieves 97.78% accuracy on the test set.

Key Notes:

- `GridSearchCV` automates hyperparameter tuning by testing all specified combinations.
- The best parameters may vary slightly if `random_state` is not set.
- The tuned model generalizes well to unseen data (test set).

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 for Disease Prediction Using a Decision Tree Model

1. Handling Missing Values

- Identify Missing Data: Use `.isnull().sum()` to check for missing values in each column.
- Numerical Features:
 - Fill missing values with the mean/median (if normally distributed) or use KNN imputation for more accuracy.
- Categorical Features:
 - Fill missing values with the mode (most frequent category) or a new category like "Unknown".
- Drop Irrelevant Columns: If a feature has too many missing values (>50%), consider dropping it.

2. Encoding Categorical Features

- Ordinal Categorical Data (e.g., disease severity: Low, Medium, High) → Use Ordinal Encoding.
- Nominal Categorical Data (e.g., gender: Male, Female) → Use One-Hot Encoding (to avoid artificial ordinal relationships).
- High Cardinality Features (e.g., ZIP codes) → Consider Target Encoding or Frequency Encoding to avoid excessive dimensions.

3. Training a Decision Tree Model

- Split Data: Use `train_test_split()` (e.g., 70% training, 30% testing).
- Initialize Model: `from sklearn.tree import DecisionTreeClassifier`
- `model = DecisionTreeClassifier(random_state=42)`

Train model: `model.fit(X_train, y_train)`

4. Hyperparameter Tuning (Using GridSearchCV)

- Define Parameter Grid: `param_grid = {`
- `'max_depth': [3, 5, 7, 10, None],`
- `'min_samples_split': [2, 5, 10],`
- `'criterion': ['gini', 'entropy']`
- `}`

Perform Grid Search:

```
from sklearn.model_selection import GridSearchCV  
  
grid_search = GridSearchCV(model, param_grid, cv=5, scoring='accuracy')  
  
grid_search.fit(X_train, y_train)
```

Select Best Model:

```
best_model = grid_search.best_estimator_
```

5. Model Evaluation

- Metrics:
 - Accuracy: Overall correctness.
 - Precision & Recall: Crucial in healthcare (avoid false negatives).
 - F1-Score: Balance between precision and recall.
 - ROC-AUC: Measures model's ability to distinguish classes.
- Confusion Matrix: Visualize true/false positives/negatives.
- Feature Importance: Identify key predictors (e.g., age, blood pressure).

Business Value in a Real-World Healthcare Setting

- Early Disease Detection → Improves patient outcomes by enabling timely treatment.
- Reduced Healthcare Costs → Prevents expensive late-stage treatments.
- Personalized Medicine → Helps doctors recommend tailored treatments.
- Resource Optimization → Hospitals can prioritize high-risk patients.
- Automated Screening → Reduces manual diagnosis time for doctors.

Example Python Workflow Summary

After preprocessing & tuning:

```
best_model.fit(X_train, y_train)  
  
y_pred = best_model.predict(X_test)
```

```
from sklearn.metrics import classification_report  
  
print(classification_report(y_test, y_pred))
```

```
# Feature Importance
```

```
importances = best_model.feature_importances_
```

```
print("Top Predictive Features:", sorted(zip(importances, X.columns), reverse=True))
```

Final Output (Example):

Best Parameters: {'criterion': 'entropy', 'max_depth': 5, 'min_samples_split': 2}

Test Accuracy: 92%

Top Predictive Features: [('blood_pressure', 0.4), ('age', 0.3), ('cholesterol', 0.2)]

This approach ensures a robust, interpretable, and actionable model for healthcare decision-making.