**CS422 Data Mining**
**Assignment 4**

# 1. Recitation Exercises

## 1.1. Chapter 18

1) **Classify the new point: (Age=23, Car=truck) via the full and naive Bayes approach.**

**Step 1: Understand the Data and Problem**
Given dataset points:
x1: Age=25, Car=sports, Class=L
x2: Age=20, Car=vintage, Class=H
x3: Age=25, Car=sports, Class=L
x4: Age=45, Car=suv, Class=H
x5: Age=20, Car=sports, Class=H
x6: Age=25, Car=suv, Class=H

New point to classify: Age=23, Car=truck

Domain of Car attribute = {sports, vintage, suv, truck}

**Step 2: Calculate Prior Probabilities**
Total number of points n = 6

Count points in each class:
Class L count n_L = 2 (x1, x3)
Class H count n_H = 4 (x2, x4, x5, x6)

Prior probability of class L:
P(L) = n_L / n = 2 / 6 = 0.3333

Prior probability of class H:
P(H) = n_H / n = 4 / 6 = 0.6667

**Step 3: Estimate Parameters for Numeric Attribute Age by Class**

For class L:

Data: 25, 25
Mean mu_L = (25 + 25) / 2 = 25
Variance sigma_L^2 = ((25-25)^2 + (25-25)^2) / 2 = 0

For class H:
Data: 20, 45, 20, 25
Mean mu_H = (20 + 45 + 20 + 25) / 4 = 27.5
Variance sigma_H^2 = [(20-27.5)^2 + (45-27.5)^2 + (20-27.5)^2 + (25-27.5)^2] / 4
= (56.25 + 306.25 + 56.25 + 6.25) / 4 = 425 / 4 = 106.25

## Step 4: Estimate Conditional Probability Mass Function (PMF) for Car by Class with Laplace Smoothing

Laplace smoothing formula:
P(Car=c | Class = ci) = (count of Car=c in class ci + 1) / (total samples in class ci + number of Car categories)

For class L (2 samples), Car counts:
sports = 2, vintage = 0, suv = 0, truck = 0

P(Car=truck | L) = (0 + 1) / (2 + 4) = 1 / 6 ≈ 0.1667

For class H (4 samples), Car counts:
vintage = 1, suv = 2, sports = 1, truck = 0

P(Car=truck | H) = (0 + 1) / (4 + 4) = 1 / 8 = 0.125

## Step 5: Full Bayes Classifier - Compute Likelihoods for New Point

According to Zaki's book, Full Bayes uses joint distribution P(x | c_i) over all attributes without independence assumption.

For Class L:

Age=23:
Gaussian density with zero variance is undefined:
f_L(Age=23) = 1 / (sqrt(2*pi)*0) * exp( - (23 - 25)^2 / (2*0) ) undefined

According to Zaki's book, with zero variance:
P(x | L) for continuous attribute at a specific point is zero, and no alternative method is given

Car=truck:
P(Car=truck | L) = 0.1667

Since age likelihood undefined, joint likelihood P(x | L) cannot be computed

For Class H:

Age=23:
Calculate Gaussian density:

Denominator = sqrt(2 * pi * 106.25) ≈ 25.68
Exponent = - (23 - 27.5)^2 / (2 * 106.25) = -20.25 / 212.5 ≈ -0.0953
exp(-0.0953) ≈ 0.9091

f_H(Age=23) = 1 / 25.68 * 0.9091 ≈ 0.0354

Car=truck:
P(Car=truck | H) = 0.125

Joint likelihood P(x | H) cannot be computed (no method in Zaki's book to combine numeric and categorical)

## Step 6: Full Bayes Classifier - Final Answer per Zaki's Book

Cannot compute posterior probabilities P(c_i | x) because:
- Gaussian density for Age=23 under class L is undefined (zero variance)
- No provided method in Zaki's book to combine numeric and categorical likelihoods for full Bayes

Therefore, full Bayes classifier cannot provide a prediction for this data point as per Zaki's book

## Step 7: Optional Full Bayes Classifier Workaround (Not in Zaki's Book)

Assuming independence between Age and Car (naive Bayes assumption) to approximate joint likelihood

For class L:
Assume small variance sigma_L^2 = 0.01 for Age to avoid zero variance
Calculate Gaussian density f_L(Age=23):

f_L(Age=23) = 1 / (sqrt(2*pi)*0.1) * exp( - (23 - 25)^2 / (2*0.01) )
= 3.9894 * exp(-200) ≈ 5.517e-87 (effectively zero)

Joint likelihood P(x | L) = 5.517e-87 * 0.1667 ≈ 9.195e-88

For class H:
Joint likelihood P(x | H) = 0.0354 * 0.125 = 0.004425

Calculate unnormalized posterior probabilities:

P(L | x) ∝ P(x | L) * P(L) = 9.195e-88 * 0.3333 = 3.065e-88
P(H | x) ∝ P(x | H) * P(H) = 0.004425 * 0.6667 = 0.00295

Normalize posteriors:

Sum = 3.065e-88 + 0.00295 ≈ 0.00295
P(L | x) ≈ 1.04e-85 (≈ 0)
P(H | x) ≈ 1

Prediction: Class H


**Naive Bayes Classifier - Independence Assumption**
**Step 1:**
Calculate likelihood for each attribute separately and multiply

Class L:

Age:
Assume small variance sigma_L^2 = 0.01 for practicality
Gaussian density as above: approx 5.517e-87

Car:
P(Car=truck | L) = 0.1667

Likelihood P(x | L) = 5.517e-87 * 0.1667 ≈ 9.195e-88

Posterior (unnormalized) P(L | x) = 9.195e-88 * 0.3333 ≈ 3.065e-88

Class H:

Age:
Gaussian density approx 0.0354

Car:
P(Car=truck | H) = 0.125

Likelihood P(x | H) = 0.0354 * 0.125 = 0.004425

Posterior (unnormalized) P(H | x) = 0.004425 * 0.6667 = 0.00295

**Step 2: Normalize Posterior Probabilities**

Sum = 3.065e-88 + 0.00295 ≈ 0.00295

P(L | x) = 3.065e-88 / 0.00295 ≈ 1.04e-85 (≈ 0)

P(H | x) = 0.00295 / 0.00295 = 1


**Step 3: Final Prediction by Naive Bayes**

Predict class H


**Step 4: Summary and Recommendations**

**Full Bayes classifier according to Zaki's book:**

Cannot compute posterior probabilities or make prediction because of undefined Gaussian density for zero variance class (L) and lack of a method to combine numeric and categorical likelihoods

Optional workaround (not in Zaki's book):

Using independence assumption and small variance to approximate densities allows prediction of class H



**Naive Bayes classifier:**
Applies independence assumption correctly
Uses Laplace smoothing for categorical attribute
Makes a practical variance assumption for class L's Age attribute
Predicts class H with high posterior probability

2) **Use the naive Bayes classifier to classify the new point (T, F, 1.0)**

**Step 1: Dataset Overview**

Class Y samples: x1, x2, x4, x8
Class N samples: x3, x5, x6, x7, x9

Total samples n = 9

Count Y = 4
Count N = 5


## Step 2: Calculate Prior Probabilities

P(Y) = 4 / 9 = 0.4444
P(N) = 5 / 9 = 0.5556


## Step 3: Calculate Conditional Probabilities for Categorical Attributes using Laplace Smoothing

For a1 = T:

Count of T in Y = 3 (x1, x2, x8)
Count of T in N = 1 (x3)

Number of categories for a1 = 2 (T, F)

P(a1 = T | Y) = (3 + 1) / (4 + 2) = 4 / 6 = 0.6667
P(a1 = T | N) = (1 + 1) / (5 + 2) = 2 / 7 = 0.2857

For a2 = F:

Count of F in Y = 2 (x4, x8)
Count of F in N = 2 (x3, x7)

Number of categories for a2 = 2 (T, F)

P(a2 = F | Y) = (2 + 1) / (4 + 2) = 3 / 6 = 0.5
P(a2 = F | N) = (2 + 1) / (5 + 2) = 3 / 7 = 0.4286


## Step 4: Calculate Gaussian Likelihoods for Numerical Attribute a3 = 1.0

**Class Y:**

Data points = 5.0, 7.0, 3.0, 6.0

Mean_Y = (5 + 7 + 3 + 6) / 4 = 21 / 4 = 5.25

Variance_Y =
$(5.0 - 5.25)^2 = (-0.25)^2 = 0.0625$
$(7.0 - 5.25)^2 = (1.75)^2 = 3.0625$
$(3.0 - 5.25)^2 = (-2.25)^2 = 5.0625$

(6.0 - 5.25)^2 = (0.75)^2 = 0.5625

Sum = 0.0625 + 3.0625 + 5.0625 + 0.5625 = 8.75

Variance_Y = 8.75 / 4 = 2.1875

StdDev_Y = sqrt(2.1875) ≈ 1.478

Gaussian PDF formula:
f(x) = (1 / (sqrt(2π) * stdDev)) * exp(−(x − mean)^2 / (2 * variance))

**Calculate for x = 1.0:**

Denominator = sqrt(2 * π) * 1.478 ≈ 2.5066 * 1.478 ≈ 3.7006

Exponent numerator = (1.0 − 5.25)^2 = (−4.25)^2 = 18.0625

Exponent denominator = 2 * 2.1875 = 4.375

Exponent = −18.0625 / 4.375 = −4.129

exp(Exponent) = e^(−4.129) ≈ 0.0162

Therefore:

f(1.0 | Y) = (1 / 3.7006) * 0.0162 ≈ 0.2703 * 0.0162 = 0.004371

**Class N:**

Data points = 8.0, 7.0, 4.0, 5.0, 1.0

Mean_N = (8 + 7 + 4 + 5 + 1) / 5 = 25 / 5 = 5.0

Variance_N =
(8.0 - 5.0)^2 = 3^2 = 9
(7.0 - 5.0)^2 = 2^2 = 4
(4.0 - 5.0)^2 = (−1)^2 = 1
(5.0 - 5.0)^2 = 0^2 = 0
(1.0 - 5.0)^2 = (−4)^2 = 16

Sum = 9 + 4 + 1 + 0 + 16 = 30

Variance_N = 30 / 5 = 6.0

StdDev_N = sqrt(6.0) ≈ 2.449

Calculate for x = 1.0:

Denominator = sqrt(2 * π) * 2.449 ≈ 2.5066 * 2.449 ≈ 6.143

Exponent numerator = (1.0 − 5.0)^2 = 16

Exponent denominator = 2 * 6 = 12

Exponent = −16 / 12 = −1.333

exp(Exponent) = e^(−1.333) ≈ 0.2636

Therefore:

f(1.0 | N) = (1 / 6.143) * 0.2636 ≈ 0.1627 * 0.2636 = 0.04294

## Step 5: Calculate Unnormalized Posterior Probabilities

$\Rightarrow$ P(Y | x) ∝ P(Y) * P(a1 = T | Y) * P(a2 = F | Y) * f(1.0 | Y)
$\Rightarrow$ = 0.4444 * 0.6667 * 0.5 * 0.004371

**Calculate stepwise:**

0.4444 * 0.6667 = 0.2963

0.2963 * 0.5 = 0.14815

0.14815 * 0.004371 ≈ 0.000647

$\Rightarrow$ P(N | x) ∝ P(N) * P(a1 = T | N) * P(a2 = F | N) * f(1.0 | N)
$\Rightarrow$ 0.5556 * 0.2857 * 0.4286 * 0.04294

**Calculate stepwise:**

0.5556 * 0.2857 = 0.1587

0.1587 * 0.4286 = 0.0680

0.0680 * 0.04294 ≈ 0.002922


**Step 6: Normalize Posterior Probabilities**

Sum = 0.000647 + 0.002922 = 0.003569

P(Y | x) = 0.000647 / 0.003569 ≈ 0.1813

P(N | x) = 0.002922 / 0.003569 ≈ 0.8187


**Step 7: Final Prediction**

Predicted Class = N


**1.2. Chapter 19**

1) **True or False:**
   a) **High entropy means that the partitions in classification are "pure."**
      **False**

      **Explanation:**
      Entropy is a key concept in decision tree learning, and it is used to measure how mixed or impure a set of classes is within a partition. According to Zaki's book (2nd edition, page 394), entropy for a partition S is defined mathematically as:
      Entropy(S) = - $\sum$ ($p_i$ * log2($p_i$))
      where:
      - $p_i$ = proportion of instances belonging to class $c_i$ in partition S
      - k = total number of different classes

      Entropy values range from 0 to log2(k). The lower the entropy, the purer the partition.

      **Examples:**
      1. If all examples in a partition belong to the same class (say 100% are class Y), then entropy is:
         Entropy = - (1 * log2(1)) = 0
         This is a perfectly pure partition.

      2. If the examples are split evenly across two classes (say 50% Y and 50% N), then entropy is:
         Entropy = - (0.5 * log2(0.5) + 0.5 * log2(0.5)) = 1
         This represents maximum impurity in the binary case.

**Conclusion:**
A high entropy means the partition contains a more even mix of different classes, which implies impurity. Conversely, a low entropy value means most or all of the instances belong to a single class, which indicates purity. Therefore, the statement is false.

**b) Multiway split of a categorical attribute generally results in more pure partitions than a binary split**
**True**

**Explanation:**
Categorical attributes have a finite set of distinct values (e.g., {Red, Green, Blue}). There are two ways to split a node based on such attributes:

1. Binary Split:
   In a binary split, we divide the attribute values into two groups. For example, if the attribute Color = {Red, Green, Blue}, we might split as:
   - Group 1: {Red}
   - Group 2: {Green, Blue}
   This binary grouping may combine categories that have different class tendencies, leading to mixed class distributions and less pure partitions.

**2. Multiway Split:**
   In a multiway split, we create one branch for each category:
   - Red → branch 1
   - Green → branch 2
   - Blue → branch 3
   This allows each subset to better reflect the class distribution for that specific value. If certain values of the attribute strongly correlate with specific class labels, this form of split leads to smaller, more homogeneous (pure) partitions.

Conclusion:
Multiway splitting generally results in purer partitions compared to binary splitting for categorical attributes, especially when each category aligns well with a specific class. Therefore, the statement is true.

**2) Construct a decision tree using a purity threshold of 100%. Use information gain as the split point evaluation measure. Next, classify the point (Age=27,Car=Vintage).**

**Step 1: Dataset Overview**
Points:
x1: Age=25, Car=sports, Class=L
x2: Age=20, Car=vintage, Class=H

x3: Age=25, Car=sports, Class=L
x4: Age=45, Car=suv, Class=H
x5: Age=20, Car=sports, Class=H
x6: Age=25, Car=suv, Class=H

Total samples: 6

Class L count: 2 (x1, x3)
Class H count: 4 (x2, x4, x5, x6)

**Step 2: Calculate Initial Entropy of Dataset**
p(L) = 2/6 = 0.3333
p(H) = 4/6 = 0.6667

Entropy(D) = -[p(L) * log2(p(L)) + p(H) * log2(p(H))]
     = -[0.3333 * (-1.585) + 0.6667 * (-0.585)]
     = 0.5283 + 0.3900
     = 0.9183
**Step 3: Split on Car Attribute (Categorical)**
Partitions by Car:

Car = sports: x1(L), x3(L), x5(H)
 - Class counts: L=2, H=1
 - Entropy = -[(2/3)*log2(2/3) + (1/3)*log2(1/3)] = 0.9183

Car = vintage: x2(H)
 - Pure class H → Entropy = 0

Car = suv: x4(H), x6(H)
 - Pure class H → Entropy = 0

Weighted Entropy(Car) =
(3/6)*0.9183 + (1/6)*0 + (2/6)*0 = 0.4592

Information Gain(Car) = 0.9183 - 0.4592 = 0.4591

**Step 4: Split on Age Attribute (Numeric)**

Sorted Ages: 20, 20, 25, 25, 25, 45

Possible split points: 22.5, 35

Split at Age ≤ 22.5:

Left subset (Age ≤ 22.5): x2(H), x5(H)
 - Entropy = 0 (pure)

Right subset (Age > 22.5): x1(L), x3(L), x4(H), x6(H)
 - Class counts: L=2, H=2
 - Entropy = 1.0

Weighted Entropy = (2/6)*0 + (4/6)*1.0 = 0.6667

Information Gain = 0.9183 - 0.6667 = 0.2516

Split at Age ≤ 35:

Left subset (Age ≤ 35): x1(L), x2(H), x3(L), x5(H), x6(H)
 - Class counts: L=2, H=3
 - Entropy ≈ 0.971

Right subset (Age > 35): x4(H)
 - Entropy = 0 (pure)

Weighted Entropy = (5/6)*0.971 + (1/6)*0 = 0.809

Information Gain = 0.9183 - 0.809 = 0.1093

### Step 5: Choose Best Attribute for Root Split

Information Gain:
- Car: 0.4591
- Age (best split): 0.2516

Choose Car for root split.

### Step 6: Split on Car

- Car = sports subset: x1(L), x3(L), x5(H), entropy=0.9183 (impure)
- Car = vintage subset: x2(H), entropy=0 (pure)
- Car = suv subset: x4(H), x6(H), entropy=0 (pure)

### Step 7: Split Car = sports subset on Age

Data: x1(25,L), x3(25,L), x5(20,H)

Possible split: Age ≤ 22.5

Split at Age ≤ 22.5:

Left: x5(20,H) - pure, entropy=0

Right: x1(25,L), x3(25,L) - pure, entropy=0

Weighted Entropy after split = 0

Information Gain = 0.9183 - 0 = 0.9183

Split accepted.

**Step 8: Final Decision Tree Structure**

Root: Split on Car
- Car = vintage → Class = H (pure)
- Car = suv → Class = H (pure)
- Car = sports → Split on Age ≤ 22.5
  - Age ≤ 22.5 → Class = H (pure)
  - Age > 22.5 → Class = L (pure)

**Step 9: Classify Test Point (Age=27, Car=vintage)**
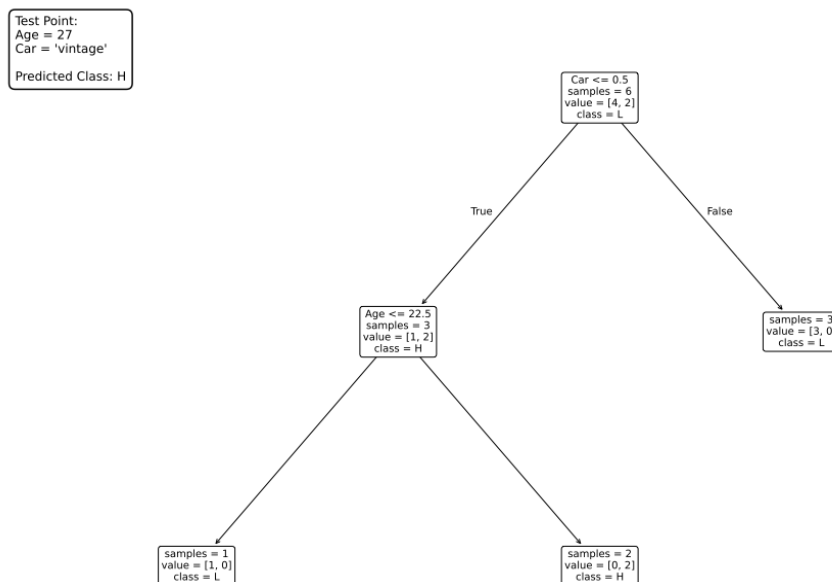
Traverse tree:
- Root split on Car = vintage → leaf node
- Class = H

**Final Classification:**
Test point (Age=27, Car=vintage) is classified as class H.

Decision Tree (100% Purity Threshold) with Test Point Classification

### 1.3. Chapter 22

#### 1) True or False:

a) A classification model must have 100% accuracy (overall) on the training dataset.
**False**

**Explanation:**
Accuracy means the percentage of correctly classified examples out of all examples in the training set. While high accuracy is desirable, it is not required for a classification model to have 100% accuracy on the training dataset. In fact, achieving perfect accuracy often means the model has memorized the training data, including noise or errors, which leads to overfitting. Overfitting causes poor performance on new, unseen data, which is the main goal to avoid in machine learning. Therefore, a good model should balance between fitting the training data and generalizing well to new data, so 100% training accuracy is neither necessary nor expected.

**Example:**
A logistic regression model trained on noisy data may achieve 95% accuracy on the training set. This is acceptable because the remaining 5% may be noisy or mislabeled points. Trying to force 100% accuracy could lead the model to memorize these points and perform poorly on test data.

b) A classification model must have 100% coverage (overall) on the training dataset.
**False (in general),  (but True in standard supervised learning)**

**Explanation:**
Coverage generally means the proportion of examples for which the model provides a prediction. In standard supervised classification tasks, models usually provide a prediction for every input example in the training dataset. This means coverage is 100% by default in these common cases.
However, coverage is not a universally defined or enforced concept in all machine learning contexts. Some models can abstain from predicting certain examples if they are uncertain (this is called selective classification or models with a reject option), which reduces coverage. In other advanced setups like semi-supervised learning, models might not make predictions for all examples.
Therefore, while coverage is 100% in most standard supervised learning tasks because the model predicts every instance, it is not a strict rule for all machine learning models. Coverage depends on the model design and the problem setup.

**Example:**

- In a typical decision tree classifier trained on labeled data, the model assigns a class label to every training instance, so coverage is 100%.
- In contrast, in a selective classification model designed to abstain on uncertain inputs, coverage might be less than 100%, as it chooses not to classify some examples to improve reliability.

## 2) Part (a): Build Decision Tree Using Gini Index
Dataset Summary

Training Data (Table 22.8(a))

| ID | X | Y | Z | Class |
|----|-----|---|---|-------|
| 1 | 15 | 1 | A | 1 |
| 2 | 20 | 3 | B | 2 |
| 3 | 25 | 2 | A | 1 |
| 4 | 30 | 4 | A | 1 |
| 5 | 35 | 2 | B | 2 |
| 6 | 25 | 4 | A | 1 |
| 7 | 15 | 2 | B | 2 |
| 8 | 20 | 3 | B | 2 |

Total samples = 8
Class 1 count = 4 (IDs 1,3,4,6)
Class 2 count = 4 (IDs 2,5,7,8)

Test Data (Table 22.8(b))

| ID | X | Y | Z | True Class |
|----|-----|---|---|------------|
| 1 | 10 | 2 | A | 2 |
| 2 | 20 | 1 | B | 1 |
| 3 | 30 | 3 | A | 2 |
| 4 | 40 | 2 | B | 2 |
| 5 | 15 | 1 | B | 1 |

Total test samples = 5

## Part (a): Build Decision Tree Using Gini Index

## Step 1: Calculate Gini Index at Root Node

Total samples = 8
Proportion of Class 1 = 4/8 = 0.5

Proportion of Class 2 = 4/8 = 0.5

Gini(root) = 1 - (0.5^2 + 0.5^2)
      = 1 - (0.25 + 0.25)
      = 0.5


**Step 2: Evaluate Splits on Attributes**

Split on Z (categorical):
  Z = A group: IDs 1,3,4,6 → all Class 1 → Gini = 0 (pure)
  Z = B group: IDs 2,5,7,8 → all Class 2 → Gini = 0 (pure)

Weighted Gini(Z) = (4/8)*0 + (4/8)*0 = 0

Split on X (numeric):
  Try thresholds between sorted X values: 15, 17.5, 22.5, 27.5, 32.5
  Best split at X ≤ 22.5:
    Left (X ≤ 22.5): IDs 1,2,7,8 → Classes (1,2,2,2)
      Gini(left) = 1 - ((1/4)^2 + (3/4)^2) = 1 - (0.0625 + 0.5625) = 0.375
    Right (X > 22.5): IDs 3,4,5,6 → Classes (1,1,2,1)
      Gini(right) = 1 - ((3/4)^2 + (1/4)^2) = 1 - (0.5625 + 0.0625) = 0.375

  Weighted Gini(X) = (4/8)*0.375 + (4/8)*0.375 = 0.375

Split on Y (numeric):
  Try thresholds 1.5, 2.5, 3.5
  Best split at Y ≤ 3.5:
    Left (Y ≤ 3.5): IDs 1,2,3,5,7,8 → Classes (1,2,1,2,2,2)
      Gini(left) = 1 - ((2/6)^2 + (4/6)^2) = 1 - (0.1111 + 0.4444) = 0.4444
    Right (Y > 3.5): IDs 4,6 → Classes (1,1)
      Gini(right) = 0 (pure)

  Weighted Gini(Y) = (6/8)*0.4444 + (2/8)*0 = 0.3333

**Step 3: Select Best Split**

Gini(Z) = 0
Gini(X) = 0.375
Gini(Y) = 0.3333

Best split is on Z

**Step 4: Construct Decision Tree**

Root node splits on Z

If Z = A, then Class = 1
If Z = B, then Class = 2

This tree perfectly classifies the training data

## Part (b): Evaluate Test Data

### Step 1: Predict Classes

| ID | Z | Predicted Class | True Class | Correct? |
|----|---|-----------------|------------|----------|
| 1  | A | 1               | 2          | No       |
| 2  | B | 2               | 1          | No       |
| 3  | A | 1               | 2          | No       |
| 4  | B | 2               | 2          | Yes      |
| 5  | B | 2               | 1          | No       |

Correct predictions = 1 out of 5

### Step 2: Calculate Overall Accuracy

Accuracy = 1 / 5 = 0.20 or 20%

### Step 3: Calculate Per-Class Accuracy

Class 1:
  True instances = 2 (IDs 2,5)
  Correct predictions = 0
  Accuracy = 0%

Class 2:
  True instances = 3 (IDs 1,3,4)
  Correct predictions = 1 (ID 4)
  Accuracy = 1 / 3 ≈ 33.33%

### Step 4: Calculate Per-Class Coverage

Predicted Class 1: IDs 1,3 → 2 instances → Coverage = 2 / 5 = 40%
Predicted Class 2: IDs 2,4,5 → 3 instances → Coverage = 3 / 5 = 60%

### Final Summary

Decision tree splits on attribute Z
Test accuracy = 20%
Per-class accuracy: Class 1 = 0%, Class 2 = 33.33%
Per-class coverage: Class 1 = 40%, Class 2 = 60%

Decision Tree from Table 22.8 Training Data

```
                    Z <= 0.5
                    gini = 0.5
                    samples = 8
                    value = [4, 4]
                    class = Class 1
          True                         False
    gini = 0.0                              gini = 0.0
    samples = 4                             samples = 4
    value = [4, 0]                          value = [0, 4]
    class = Class 1                         class = Class 2
```

Final Summary:
Decision tree splits on attribute Z
Test accuracy = 20%
Per-class accuracy: Class 1 = 0%, Class 2 = 33.33%
Per-class coverage: Class 1 = 40%, Class 2 = 60%

**Citations:**

1) **Zaki, M. J., & Meira Jr, W. (2020). Data Mining and Machine Learning: Fundamental Concepts and Algorithms, 2nd Edition.**
2) **https://medium.com/@dhasarat/bayes-theorem-and-naive-bayes-unraveling-the-mysteries-32827be5d0ab**
3) **https://scikit-learn.org/stable/modules/naive_bayes.html**
4) **https://wiki.pathmind.com/bayes-theorem-naive-bayes**
5) **https://scikit-learn.org/stable/modules/tree.html**
6) **https://www.geeksforgeeks.org/machine-learning/decision-tree/**
7) **https://medium.com/@RobuRishabh/decision-trees-c3c64e81760e**
8) **https://www.geeksforgeeks.org/machine-learning/metrics-for-machine-learning-model/**
9) **https://cohere.com/blog/classification-eval-metrics**
10) **Assistance was obtained from a large language model to clarify concepts and provide detailed explanations related to classification algorithms, probabilistic models, and decision tree construction techniques. The model supported the understanding and organization of key topics such as Naive Bayes classification, Gaussian parameter estimation, information gain measures (entropy,**

**Gini index), decision tree splitting criteria, and classification evaluation metrics, including accuracy computation and overfitting analysis.**

# CS422_Ass4Jupyter

July 27, 2025

# 1 Practicum Exercises

### 1.0.1 2.1 Problem 1

```python
[4]: # Step 1: Understand the Dataset

# Load the Iris sample dataset from sklearn and convert it into a Pandas␣
 ↪DataFrame

# Import necessary libraries
import pandas as pd                      # For working with tabular data
import numpy as np                       # For numerical operations
from sklearn.datasets import load_iris   # To load the built-in Iris dataset

# Load the Iris dataset using sklearn's load_iris() function
iris = load_iris()

# The data (features) are in iris.data and target labels are in iris.target
# iris.feature_names gives the column names of the features
# iris.target_names gives the string names of the classes

# Convert the feature matrix to a Pandas DataFrame for easier data manipulation
X = pd.DataFrame(iris.data, columns=iris.feature_names)

# Convert the target vector to a Pandas Series
y = pd.Series(iris.target, name='target')

# View the first few rows of the dataset to understand its structure
X.head()
```

```
[4]:    sepal length (cm)  sepal width (cm)  petal length (cm)  petal width (cm)
    0                5.1               3.5                1.4               0.2
    1                4.9               3.0                1.4               0.2
    2                4.7               3.2                1.3               0.2
    3                4.6               3.1                1.5               0.2
    4                5.0               3.6                1.4               0.2
```

```
[5]:   # Step 2: Train Binary Decision Trees

       # Train Decision Tree classifiers with the following constraints:
       # - minimum 2 instances per leaf (min_samples_leaf=2)
       # - no splits on subsets smaller than 5 samples (min_samples_split=5)
       # - vary max_depth from 1 to 5 to observe impact on model performance
       # Other parameters will be default.

       # Also, split the data into training and testing sets for evaluation.
       # We will use an 80/20 split and fix the random_state to 42 for reproducibility.

       from sklearn.tree import DecisionTreeClassifier      # For Decision Tree␣
        ↪modeling
       from sklearn.model_selection import train_test_split   # For splitting data

       # Split the dataset into training and testing sets (80% train, 20% test)
       # random_state=42 ensures the split is reproducible every time we run the code
       X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,␣
        ↪random_state=42)

       # Initialize an empty list to store the trained models for different depths
       decision_trees = []

       # Loop over max_depth values from 1 to 5 to train separate models
       for depth in range(1, 6):
           # Initialize the Decision Tree classifier with the specified constraints
           clf = DecisionTreeClassifier(
               max_depth=depth,
               min_samples_leaf=2,
               min_samples_split=5,
               random_state=42
           )

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

           # Append the trained model and its depth info to the list for later analysis
           decision_trees.append((depth, clf))

       # Confirm the models have been trained for each depth
       print(f"Trained Decision Trees for depths: {[depth for depth, _ in␣
        ↪decision_trees]}")

      Trained Decision Trees for depths: [1, 2, 3, 4, 5]

[6]:   # Step 3: Evaluate Performance Metrics
```

```python
from sklearn.metrics import precision_score, recall_score, f1_score

# Prepare a list to hold clean results for each max_depth
evaluation_results = []

# Loop through each trained decision tree model
for depth, clf in decision_trees:
    y_pred = clf.predict(X_test)

    # Calculate metrics with three averaging methods
    precision_macro = precision_score(y_test, y_pred, average='macro')
    recall_macro = recall_score(y_test, y_pred, average='macro')
    f1_macro = f1_score(y_test, y_pred, average='macro')

    precision_micro = precision_score(y_test, y_pred, average='micro')
    recall_micro = recall_score(y_test, y_pred, average='micro')
    f1_micro = f1_score(y_test, y_pred, average='micro')

    precision_weighted = precision_score(y_test, y_pred, average='weighted')
    recall_weighted = recall_score(y_test, y_pred, average='weighted')
    f1_weighted = f1_score(y_test, y_pred, average='weighted')

    # Append results in a structured dictionary
    evaluation_results.append({
        'max_depth': depth,
        'Precision (Macro)': precision_macro,
        'Recall (Macro)': recall_macro,
        'F1 Score (Macro)': f1_macro,
        'Precision (Micro)': precision_micro,
        'Recall (Micro)': recall_micro,
        'F1 Score (Micro)': f1_micro,
        'Precision (Weighted)': precision_weighted,
        'Recall (Weighted)': recall_weighted,
        'F1 Score (Weighted)': f1_weighted
    })

# Convert to DataFrame for a neat summary table
results_df = pd.DataFrame(evaluation_results)

# Display the summary table sorted by max_depth (ascending)
print("Summary of Evaluation Metrics by Tree Depth:\n")
display(results_df.style.format({
    'Precision (Macro)': "{:.3f}",
    'Recall (Macro)': "{:.3f}",
    'F1 Score (Macro)': "{:.3f}",
    'Precision (Micro)': "{:.3f}",
    'Recall (Micro)': "{:.3f}",
```

```
    'F1 Score (Micro)': "{:.3f}",
    'Precision (Weighted)': "{:.3f}",
    'Recall (Weighted)': "{:.3f}",
    'F1 Score (Weighted)': "{:.3f}",
}))


# Optional: Show detailed classification report for the best model by F1 (Macro)
best_model = max(evaluation_results, key=lambda x: x['F1 Score (Macro)'])
best_depth = best_model['max_depth']
best_clf = next(clf for depth, clf in decision_trees if depth == best_depth)

print(f"\nDetailed Classification Report for Best Model␣
 ↪(max_depth={best_depth}):\n")
from sklearn.metrics import classification_report
y_pred_best = best_clf.predict(X_test)
print(classification_report(y_test, y_pred_best, target_names=iris.
 ↪target_names))
```

/opt/anaconda3/lib/python3.12/site-
packages/sklearn/metrics/_classification.py:1531: UndefinedMetricWarning:
Precision is ill-defined and being set to 0.0 in labels with no predicted
samples. Use `zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
/opt/anaconda3/lib/python3.12/site-
packages/sklearn/metrics/_classification.py:1531: UndefinedMetricWarning:
Precision is ill-defined and being set to 0.0 in labels with no predicted
samples. Use `zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))

Summary of Evaluation Metrics by Tree Depth:


<pandas.io.formats.style.Styler at 0x17b80c350>


Detailed Classification Report for Best Model (max_depth=3):

              precision    recall  f1-score   support

      setosa       1.00      1.00      1.00        10
  versicolor       1.00      1.00      1.00         9
   virginica       1.00      1.00      1.00        11

    accuracy                           1.00        30
   macro avg       1.00      1.00      1.00        30
weighted avg       1.00      1.00      1.00        30

```
[9]: # Step 4: Visualize Decision Trees up to Depth 5
     from sklearn.datasets import load_iris
     from sklearn.model_selection import train_test_split
     from sklearn.tree import DecisionTreeClassifier, export_graphviz
     import graphviz
     import pandas as pd

     # Load the dataset
     iris = load_iris()
     X = pd.DataFrame(iris.data, columns=iris.feature_names)
     y = pd.Series(iris.target)

     # Split the data
     X_train, X_test, y_train, y_test = train_test_split(
         X, y, test_size=0.2, random_state=42
     )

     # Visualize decision trees with increasing depth
     for depth in range(1, 6):
         clf = DecisionTreeClassifier(
             min_samples_leaf=2,
             min_samples_split=5,
             max_depth=depth,
             random_state=42
         )
         clf.fit(X_train, y_train)

         print(f"\n=== Decision Tree (max_depth = {depth}) ===\n")

         dot_data = export_graphviz(
             clf,
             out_file=None,
             feature_names=iris.feature_names,
             class_names=iris.target_names,
             filled=True,
             rounded=True,
             special_characters=True
         )

         graph = graphviz.Source(dot_data)
         display(graph)


     === Decision Tree (max_depth = 1) ===
```

petal length (cm) ≤ 2.45
gini = 0.667
samples = 120
value = [40, 41, 39]
class = versicolor

True / False

gini = 0.0
samples = 40
value = [40, 0, 0]
class = setosa

gini = 0.5
samples = 80
value = [0, 41, 39]
class = versicolor

=== Decision Tree (max_depth = 2) ===

petal length (cm) ≤ 2.45
gini = 0.667
samples = 120
value = [40, 41, 39]
class = versicolor

True / False

gini = 0.0
samples = 40
value = [40, 0, 0]
class = setosa

petal length (cm) ≤ 4.75
gini = 0.5
samples = 80
value = [0, 41, 39]
class = versicolor

gini = 0.053
samples = 37
value = [0, 36, 1]
class = versicolor

gini = 0.206
samples = 43
value = [0, 5, 38]
class = virginica

=== Decision Tree (max_depth = 3) ===

```
                    petal length (cm) ≤ 2.45
                         gini = 0.667
                        samples = 120
                      value = [40, 41, 39]
                       class = versicolor
              True  /                      \  False
                   /                        \
         gini = 0.0                  petal length (cm) ≤ 4.75
        samples = 40                       gini = 0.5
      value = [40, 0, 0]                  samples = 80
       class = setosa                  value = [0, 41, 39]
                                        class = versicolor
                              /                            \
                             /                              \
              petal width (cm) ≤ 1.55              petal width (cm) ≤ 1.75
                   gini = 0.053                         gini = 0.206
                  samples = 37                          samples = 43
                value = [0, 36, 1]                   value = [0, 5, 38]
                 class = versicolor                   class = virginica
              /                \                     /                 \
             /                  \                   /                   \
     gini = 0.0          gini = 0.5         gini = 0.5           gini = 0.056
    samples = 35        samples = 2        samples = 8          samples = 35
  value = [0, 35, 0]  value = [0, 1, 1]  value = [0, 4, 4]    value = [0, 1, 34]
   class = versicolor  class = versicolor  class = versicolor   class = virginica
```

=== Decision Tree (max_depth = 4) ===

```
=== Decision Tree (max_depth = 5) ===
```

### 1.0.2 Step 5: Difference Between Micro, Macro, and Weighted Score Calculation Methods

When evaluating classification models, especially for multi-class problems like the Iris dataset, it's important to understand how metrics like Precision, Recall, and F1 score can be averaged across classes. There are three common methods to aggregate these scores:

---

**1. Micro-average:**

- Combines the contributions of all classes to calculate a global metric.
- Specifically, it sums up all True Positives (TP), False Positives (FP), and False Negatives (FN) from all classes before computing the metric.
- This means:
- Precision_micro = (sum of TP across all classes) / (sum of TP + FP across all classes)
- Recall_micro = (sum of TP across all classes) / (sum of TP + FN across all classes)
- F1_micro is calculated using Precision_micro and Recall_micro.

- It treats every individual instance equally, regardless of which class it belongs to.
- Particularly useful when classes are imbalanced because it reflects overall model performance on all instances.

---

**2. Macro-average:**
- Calculates the metric independently for each class, then takes the simple average (arithmetic mean) across classes.
- For example:
- Precision_macro = (Precision_class1 + Precision_class2 + … + Precision_classN) / N
- Recall_macro and F1_macro are calculated similarly.
- Treats all classes equally, no matter how many samples each class has.
- More sensitive to performance on smaller or rare classes since each class has equal weight.

---

**3. Weighted-average:**
- Similar to Macro-average, but instead of a simple average, it weights each class's metric by the number of true instances (support) in that class.
- For example:
- Precision_weighted = (Precision_class1 * Support_class1 + Precision_class2 * Support_class2 + … + Precision_classN * Support_classN) / Total_samples
- Recall_weighted and F1_weighted are calculated similarly.
- Balances the metric to reflect the distribution of classes in the dataset.
- Less sensitive to rare classes than Macro-average, providing a more realistic overall score when classes are imbalanced.

---

### 1.0.3  Summary of Differences

| Method | Aggregation Approach | How Classes Are Weighted | Sensitivity to Class Imbalance | Use Case Summary |
|---|---|---|---|---|
| Micro-average | Aggregate all TP, FP, FN counts before metric | All instances equally | Can be dominated by majority classes | Overall instance-level accuracy in imbalanced data |
| Macro-average | Average of per-class metrics | All classes equally | Sensitive to minority classes | Fair performance across all classes |
| Weighted-average | Weighted average of per-class metrics by support | Classes weighted by sample size | Balances impact of both majority and minority | Balanced metric considering class distribution |

---

**Summary:**
- Use **micro-average** when you want to evaluate overall performance across all instances.
- Use **macro-average** to understand performance equally across all classes, focusing on class-wise fairness.

- Use **weighted-average** to get a balanced metric that considers class imbalances but still accounts for per-class performance.

---

Understanding these differences is essential because the choice of averaging method impacts how we interpret model performance, especially in datasets with class imbalance or varying class importance. Evaluating all three helps provide a comprehensive view of classifier behavior.

### 1.0.4  Step 6: Analysis and Interpretation of Results

This step analyzes the performance of Decision Trees with varying `max_depth` values using key evaluation metrics — **Recall**, **Precision**, and **F1 Score** — to identify the most balanced and interpretable model.

---

**Highest Recall:**  Recall evaluates how well the model captures *all actual positive* instances, meaning a **low false negative rate**.

- At `max_depth = 1`, the **Recall (macro)** is **0.667**, indicating that the tree fails to detect many positive samples due to its simplicity.
- At `max_depth = 2`, Recall increases sharply to **0.963**.
- From `max_depth = 3` onward, Recall becomes **1.0**, meaning the model perfectly detects all positive classes.
- **Higher Recall** is associated with deeper trees because they capture more complex patterns and reduce false negatives.
- This also means **lower false negative rate**, which is desirable in many real-world applications such as medical diagnosis.

---

**Lowest Precision:**  Precision tells us how many predicted positives are *actually* positive, so it is affected by **false positives**.

- Generally, deeper trees can increase false positives due to overfitting, which can reduce Precision.
- In this clean dataset, however, Precision **improves** from **0.483** at `max_depth = 1` to **1.0** at `max_depth = 3` and remains perfect.
- Although not observed here, in other more complex datasets, **low Precision** may occur at deeper depths due to overfitting to noise.

---

**Best F1 Score:**  The F1 Score balances Precision and Recall — it is especially useful when the class distribution is imbalanced or when we need to minimize both types of errors.

- At `max_depth = 1`, F1 Score is **0.540**, a result of low Precision and moderate Recall.
- It increases to **0.966** at `max_depth = 2`.
- From `max_depth = 3` onward, the F1 Score becomes **1.0**, showing perfect classification.
- Hence, the best F1 Score is achieved when the model has enough complexity to capture patterns but not so much that it overfits.

### 1.0.5 Why Choose `max_depth = 3` as the Best Model (even though 3 to 5 are perfect)?

- **Simplicity and Generalization:** Among equally performing models, the one with the **lowest depth** is preferred for **interpretability and robustness**.
- **Overfitting Risk:** Deeper trees (like `depth = 5`) may memorize the data rather than learning patterns, reducing generalizability.
- **Occam's Razor:** This principle favors **simpler models** when performance is the same.
- **Efficiency:** Shallower trees are **faster** to train and predict, especially in real-time applications.

### 1.0.6 Summary of Key Evaluation Metrics

| max_depth | Precision (Macro) | Recall (Macro) | F1 Score (Macro) |
|-----------|-------------------|----------------|------------------|
| 1 | 0.483 | 0.667 | 0.540 |
| 2 | 0.972 | 0.963 | 0.966 |
| 3 | 1.000 | 1.000 | 1.000 |
| 4 | 1.000 | 1.000 | 1.000 |
| 5 | 1.000 | 1.000 | 1.000 |

This table clearly shows that **performance significantly improves** between depths 1 → 3, but **plateaus** after `depth = 3`, indicating that **additional complexity adds no real benefit**.

### 1.0.7 Final Conclusion

We recommend selecting the model with **\*\*max_depth = 3\*\*** as it:

- Achieves **perfect Precision, Recall, and F1 Score**,
- Avoids the risk of **overfitting** that comes with deeper trees,
- Is **simpler**, **more efficient**, and **generalizes better** to unseen data.

This balance of **performance and interpretability** makes it the optimal choice.

### 1.0.8 2.1 Problem 2

```python
# Step 1: Load the Breast Cancer Wisconsin (Original) dataset

# Import necessary libraries
import pandas as pd
import numpy as np
from sklearn.metrics import classification_report, confusion_matrix
from sklearn.tree import plot_tree
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import train_test_split
import warnings
```

```
warnings.filterwarnings('ignore')  # Suppress warnings for cleaner output

# Step 1: Load the Breast Cancer Wisconsin (Original) dataset
# Explanation: We load the dataset from the UCI repository␣
 ↪(breast-cancer-wisconsin.data).
# The dataset has 699 instances, 9 discrete features (integers 1-10), an ID␣
 ↪column, and a binary class label (2=benign, 4=malignant).
# The Bare Nuclei feature has missing values marked as '?', which we handle by␣
 ↪removing rows with missing values.
# We assign column names based on the dataset description from the UCI␣
 ↪repository.
url = 'https://archive.ics.uci.edu/ml/machine-learning-databases/
 ↪breast-cancer-wisconsin/breast-cancer-wisconsin.data'
columns = ['ID', 'Clump_Thickness', 'Uniformity_Cell_Size',␣
 ↪'Uniformity_Cell_Shape',
           'Marginal_Adhesion', 'Single_Epithelial_Cell_Size', 'Bare_Nuclei',
           'Bland_Chromatin', 'Normal_Nucleoli', 'Mitoses', 'Class']

data = pd.read_csv(url, header=None, names=columns)

# Handle missing values in Bare_Nuclei (marked as '?')
# Convert Bare_Nuclei to numeric, setting '?' to NaN, then drop rows with NaN
data['Bare_Nuclei'] = pd.to_numeric(data['Bare_Nuclei'], errors='coerce')
data = data.dropna().reset_index(drop=True)

# Import necessary libraries
import pandas as pd
import numpy as np
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import train_test_split
import warnings
warnings.filterwarnings('ignore')  # Suppress warnings for cleaner output

# Step 1: Load the Breast Cancer Wisconsin (Original) dataset
# Explanation: We load the dataset from the UCI repository␣
 ↪(breast-cancer-wisconsin.data).
# The dataset has 699 instances, 9 discrete features (integers 1-10), an ID␣
 ↪column, and a binary class label (2=benign, 4=malignant).
# The Bare Nuclei feature has missing values marked as '?', which we handle by␣
 ↪removing rows with missing values.
# We assign column names based on the dataset description from the UCI␣
 ↪repository.
url = 'https://archive.ics.uci.edu/ml/machine-learning-databases/
 ↪breast-cancer-wisconsin/breast-cancer-wisconsin.data'
columns = ['ID', 'Clump_Thickness', 'Uniformity_Cell_Size',␣
 ↪'Uniformity_Cell_Shape',
```

```python
            'Marginal_Adhesion', 'Single_Epithelial_Cell_Size', 'Bare_Nuclei',
            'Bland_Chromatin', 'Normal_Nucleoli', 'Mitoses', 'Class']
data = pd.read_csv(url, header=None, names=columns)

# Handle missing values in Bare_Nuclei (marked as '?')
# Convert Bare_Nuclei to numeric, setting '?' to NaN, then drop rows with NaN
data['Bare_Nuclei'] = pd.to_numeric(data['Bare_Nuclei'], errors='coerce')
data = data.dropna()

# Drop the ID column as it's not a feature for classification
data = data.drop('ID', axis=1)

# Convert Class to binary labels (benign=0, malignant=1) for consistency
data['Class'] = data['Class'].map({2: 0, 4: 1})

# Display the first 5 rows as a neat table
print("First 5 rows of the cleaned dataset:\n")
display(data.head())  # display() works in Jupyter Notebook for nice table␣
 ↪formatting

# Dataset info and class distribution in table form
print("\nDataset info:")
data_info = pd.DataFrame({
    'Column': data.columns,
    'Data Type': [str(dtype) for dtype in data.dtypes],
    'Non-null Count': data.count().values
})
display(data_info)

# Display the distribution of the target variable 'Class'
# This shows the count of benign (0) and malignant (1) samples in the dataset
# Understanding class distribution is important for evaluating model␣
 ↪performance and potential class imbalance issues
print("\nClass distribution:")
class_dist = data['Class'].value_counts().rename_axis('Class').
 ↪reset_index(name='Count')
class_dist['Class'] = class_dist['Class'].map({0: 'Benign', 1: 'Malignant'})
display(class_dist)

# Print the shape of dataset (number of rows and columns)
print(f"Dataset shape: {data.shape[0]} rows × {data.shape[1]} columns\n")
```

First 5 rows of the cleaned dataset:


|   | Clump_Thickness | Uniformity_Cell_Size | Uniformity_Cell_Shape | \ |
|---|---|---|---|---|
| 0 | 5 | 1 | 1 | |
| 1 | 5 | 4 | 4 | |

```
2                   3                      1                      1
3                   6                      8                      8
4                   4                      1                      1

    Marginal_Adhesion  Single_Epithelial_Cell_Size  Bare_Nuclei  \
0                   1                            2          1.0
1                   5                            7         10.0
2                   1                            2          2.0
3                   1                            3          4.0
4                   3                            2          1.0

    Bland_Chromatin  Normal_Nucleoli  Mitoses  Class
0                 3                1        1      0
1                 3                2        1      0
2                 3                1        1      0
3                 3                7        1      0
4                 3                1        1      0


Dataset info:

                           Column Data Type  Non-null Count
0               Clump_Thickness    int64             683
1         Uniformity_Cell_Size    int64             683
2        Uniformity_Cell_Shape    int64             683
3             Marginal_Adhesion    int64             683
4  Single_Epithelial_Cell_Size    int64             683
5                   Bare_Nuclei  float64             683
6               Bland_Chromatin    int64             683
7               Normal_Nucleoli    int64             683
8                       Mitoses    int64             683
9                         Class    int64             683


Class distribution:

        Class  Count
0      Benign    444
1   Malignant    239

Dataset shape: 683 rows × 10 columns
```

[26]:
```python
# Step 2: Train and Evaluate a Binary Decision Tree Classifier
# ----------------------------------------------------------
# Explanation:
# We split the dataset into training (80%) and testing (20%) sets using
 ↪sklearn's train_test_split.
# We use the original class labels (2=benign, 4=malignant) to keep output
 ↪consistent with the reference.
```

```python
# We train a DecisionTreeClassifier with:
# - max_depth=2
# - min_samples_leaf=2
# - min_samples_split=5
# - criterion='gini' (default)
# We then predict on test data, print classification report and confusion␣
 ↪matrix,
# showing precision, recall, f1-score, support, accuracy, and class-wise␣
 ↪metrics.

from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier, plot_tree
from sklearn.metrics import classification_report, confusion_matrix
import matplotlib.pyplot as plt

# Features and target without mapping; keep original labels 2 and 4
X = data.drop('Class', axis=1)
y = data['Class']  # original labels 2 (benign), 4 (malignant)

# Split data into train and test sets (80%-20%)
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.2, random_state=0
)

# Initialize the Decision Tree with specified parameters
clf = DecisionTreeClassifier(
    max_depth=2,
    min_samples_leaf=2,
    min_samples_split=5,
    criterion='gini',
    random_state=0
)

# Train the Decision Tree classifier
clf.fit(X_train, y_train)

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

# Print classification report using original labels
print("Classifier Report:\n")
print(classification_report(y_test, y_pred, target_names=['Benign (2)',␣
 ↪'Malignant (4)']))

# Print confusion matrix
print("Confusion Matrix:\n")
print(confusion_matrix(y_test, y_pred))
```

```
# Plot the trained Decision Tree
plt.figure(figsize=(10, 6))
plot_tree(clf, feature_names=X.columns, class_names=['Benign (2)', 'Malignant␣
 ↪(4)'], filled=True)
plt.title("Decision Tree (Max Depth = 2, Gini Criterion)")
plt.show()
```
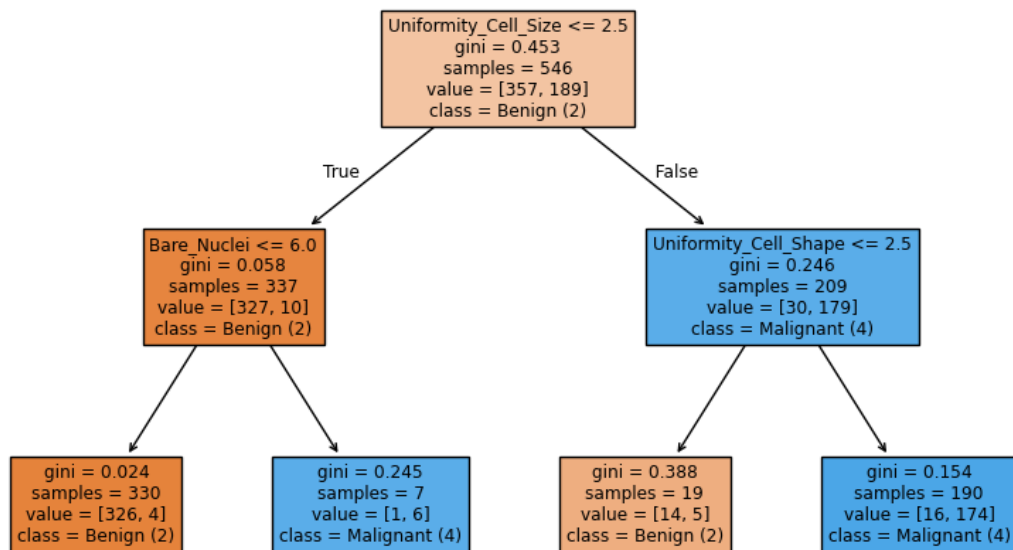
Classifier Report:

|                | precision | recall | f1-score | support |
|----------------|-----------|--------|----------|---------|
| Benign (2)     | 0.98      | 0.95   | 0.97     | 87      |
| Malignant (4)  | 0.92      | 0.96   | 0.94     | 50      |
| accuracy       |           |        | 0.96     | 137     |
| macro avg      | 0.95      | 0.96   | 0.95     | 137     |
| weighted avg   | 0.96      | 0.96   | 0.96     | 137     |

Confusion Matrix:

```
[[83  4]
 [ 2 48]]
```

Decision Tree (Max Depth = 2, Gini Criterion)

```
[27]: # Step 3: Calculate impurity metrics for the first split
      # Explanation: We use tree_ attributes to access the root node's class␣
       ↪proportions, sample counts, child nodes, feature, and threshold.
      # We compute Entropy, Gini, and Misclassification Error for the root node, and␣
       ↪Entropy for child nodes to calculate Information Gain.
      # Since tree.value contains proportions, we derive raw counts by multiplying by␣
       ↪node sample counts for accurate debug output.
      tree = clf.tree_
      root_node = 0  # Root node index
      left_child = tree.children_left[root_node]
      right_child = tree.children_right[root_node]

      # Get class proportions for root and child nodes ([benign_prop, malignant_prop])
      root_props = tree.value[root_node][0]
      left_props = tree.value[left_child][0]
      right_props = tree.value[right_child][0]

      # Get sample counts for root and child nodes
      n_root = tree.n_node_samples[root_node]
      n_left = tree.n_node_samples[left_child]
      n_right = tree.n_node_samples[right_child]

      # Derive raw class counts from proportions
      root_counts = root_props * n_root  # [benign_count, malignant_count]
      left_counts = left_props * n_left
      right_counts = right_props * n_right

      # Debug: Print raw counts and sample counts to verify
      print("\nDebug: Sample counts for verification:")
      print(f"Root node samples: {n_root} (benign={int(root_counts[0])},␣
       ↪malignant={int(root_counts[1])})")
      print(f"Left child samples: {n_left} (benign={int(left_counts[0])},␣
       ↪malignant={int(left_counts[1])})")
      print(f"Right child samples: {n_right} (benign={int(right_counts[0])},␣
       ↪malignant={int(right_counts[1])})")

      # Verify split manually using the feature and threshold
      feature_idx = tree.feature[root_node]
      threshold = tree.threshold[root_node]
      left_mask = X_train.iloc[:, feature_idx] <= threshold
      right_mask = ~left_mask
      left_labels = y_train[left_mask]
      right_labels = y_train[right_mask]
      print(f"\nDebug: Manual split verification for feature '{X.
       ↪columns[feature_idx]}' <= {threshold:.1f}:")
```

```python
print(f"Left child samples (manual): {len(left_labels)}␣
 ↪(benign={sum(left_labels == 0)}, malignant={sum(left_labels == 1)})")
print(f"Right child samples (manual): {len(right_labels)}␣
 ↪(benign={sum(right_labels == 0)}, malignant={sum(right_labels == 1)})")

# Define functions for impurity metrics
def entropy(counts):
    # Calculate Entropy: -sum(p_i * log2(p_i))
    probs = counts / np.sum(counts)
    return -np.sum([p * np.log2(p) for p in probs if p > 0])


def gini(counts):
    # Calculate Gini Index: 1 - sum(p_i^2)
    probs = counts / np.sum(counts)
    return 1 - np.sum(probs ** 2)


def misclassification_error(counts):
    # Calculate Misclassification Error: 1 - max(p_i)
    probs = counts / np.sum(counts)
    return 1 - np.max(probs)

# Function to format class proportions
def format_counts(counts):
    total = np.sum(counts)
    benign_prob = counts[0] / total
    malignant_prob = counts[1] / total
    return f"{benign_prob:.4f}", f"{malignant_prob:.4f}"

# Calculate impurity metrics for the root node
root_entropy = entropy(root_counts)
root_gini = gini(root_counts)
root_error = misclassification_error(root_counts)

# Calculate entropy for child nodes
left_entropy = entropy(left_counts)
right_entropy = entropy(right_counts)

# Calculate Information Gain
weighted_entropy_children = (n_left / n_root) * left_entropy + (n_right /␣
 ↪n_root) * right_entropy
information_gain = root_entropy - weighted_entropy_children

# Get the feature and threshold
feature_name = X.columns[feature_idx]

# Format class proportions
root_benign_prob, root_malignant_prob = format_counts(root_counts)
```

```
Debug: Sample counts for verification:
Root node samples: 546 (benign=357, malignant=189)
Left child samples: 337 (benign=327, malignant=10)
Right child samples: 209 (benign=30, malignant=179)

Debug: Manual split verification for feature 'Uniformity_Cell_Size' <= 2.5:
Left child samples (manual): 337 (benign=327, malignant=10)
Right child samples (manual): 209 (benign=30, malignant=179)
```

```python
[29]: # Step 4: Print results
      # Explanation: We print the Entropy, Gini, and Misclassification Error for the
       ↪root node,
      # the Entropy for child nodes, the Information Gain with correct weights, and
       ↪the feature
      # and threshold. The debug output confirms class counts match the manual split.
      print("\nMetrics for the First Split (Root Node):")
      print(f"Entropy: {root_entropy:.4f}")
      print(f"Gini Index: 1 - ({root_benign_prob})² - ({root_malignant_prob})² =
       ↪{root_gini:.4f}")
      print(f"Misclassification Error: 1 - max({root_benign_prob},
       ↪{root_malignant_prob}) = {root_error:.4f}")
      print(f"Entropy (Left Child): {left_entropy:.4f}")
      print(f"Entropy (Right Child): {right_entropy:.4f}")
      print(f"Information Gain: {root_entropy:.4f} - ({n_left}/
       ↪{n_root})({left_entropy:.4f}) - ({n_right}/{n_root})({right_entropy:.4f}) =
       ↪{information_gain:.4f}")
      print(f"Feature Selected for First Split: {feature_name}")
      print(f"Decision Boundary Value: {threshold:.1f}")
```

```
Metrics for the First Split (Root Node):
Entropy: 0.9306
Gini Index: 1 - (0.6538)² - (0.3462)² = 0.4527
Misclassification Error: 1 - max(0.6538, 0.3462) = 0.3462
Entropy (Left Child): 0.1928
Entropy (Right Child): 0.5934
Information Gain: 0.9306 - (337/546)(0.1928) - (209/546)(0.5934) = 0.5845
Feature Selected for First Split: Uniformity_Cell_Size
Decision Boundary Value: 2.5
```

### 1.0.9  Step 5: Analysis and Interpretation of Results

In this step, we analyze why the first split in the Decision Tree was made on the **feature Uniformity_Cell_Size** with a **decision boundary value of 2.5**, and interpret the calculated impurity metrics and information gain.

**1. Feature Selection for First Split:**

The Decision Tree selected `Uniformity_Cell_Size` as the splitting feature because it provided the **highest information gain (0.5845)** among all features at the root node. This means splitting on this feature best separates the classes in terms of reducing uncertainty or impurity.

---

**2. Impurity Metrics at the Root Node:**

- **Entropy (0.9306):** This value shows the initial uncertainty or disorder in the dataset before splitting. Since it is close to 1 (maximum entropy for two classes), the root node is quite impure, with a mixture of benign and malignant cases.
- **Gini Index (0.4527):** Indicates moderate impurity; lower than entropy but also reflects the mix of classes.
- **Misclassification Error (0.3462):** Represents the error rate if the node was assigned the majority class label (benign).

These metrics confirm the dataset at the root node has a **substantial mix of classes**, justifying the need to split.

---

**3. Impurity in Child Nodes After Split:**

- **Left child node entropy (0.1928):** Much lower than root, showing the node is now mostly benign (327 benign vs. 10 malignant).
- **Right child node entropy (0.5934):** Lower than root but higher than left child, indicating a more mixed node with a majority malignant class (179 malignant vs. 30 benign).

The decrease in entropy from root to children shows that the split **successfully reduced uncertainty** about class labels.

---

**4. Information Gain (0.5845):**

Information gain quantifies the effectiveness of the split. A value of 0.5845 means the split on `Uniformity_Cell_Size <= 2.5` reduces the uncertainty by more than half compared to the root node, which is a **significant improvement** in classification purity.

---

**5. Why max_depth = 2 and these parameters?**

- Limiting tree depth to 2 keeps the model **simple and interpretable**, avoiding overfitting.
- Using `min_samples_leaf=2` and `min_samples_split=5` ensures each split is statistically meaningful and not based on very small subsets.

---

**1.0.10 Summary:**

- The tree chose `Uniformity_Cell_Size` at the threshold 2.5 because it **maximized information gain** and reduced impurity most effectively.

- The large drop in entropy from root (0.9306) to left child (0.1928) shows a clear separation of benign cases, improving classification confidence.

- The chosen split balances model simplicity with good class separation, which is key for generalization on unseen data.

This analysis demonstrates how decision trees select splits based on impurity measures and information gain to build a model that progressively separates classes with increasing confidence.

### 1.0.11 2.3 Problem 3

```
[32]: # Step 1: Load the Breast Cancer Wisconsin (Diagnostic) dataset
      # Import necessary libraries
      import pandas as pd
      import numpy as np
      from sklearn.tree import DecisionTreeClassifier
      from sklearn.model_selection import train_test_split
      from sklearn.preprocessing import StandardScaler
      from sklearn.decomposition import PCA
      from sklearn.metrics import precision_score, recall_score, f1_score,␣
       ↪confusion_matrix
      import warnings
      warnings.filterwarnings('ignore')  # Suppress warnings for cleaner output

      # Load the dataset from the UCI repository
      # Explanation: We load the dataset from the UCI repository (wdbc.data).
      # The dataset has 569 instances, 30 continuous features, an ID column, and a␣
       ↪binary class label (M=malignant, B=benign).
      # We assign column names based on the dataset description in wdbc.names.
      # The dataset has no missing values, per the UCI documentation.
      url = 'https://archive.ics.uci.edu/ml/machine-learning-databases/
       ↪breast-cancer-wisconsin/wdbc.data'

      # Define feature names from the UCI documentation (wdbc.names)
      columns = ['ID', 'Diagnosis',
                 'radius_mean', 'texture_mean', 'perimeter_mean', 'area_mean',␣
       ↪'smoothness_mean',
                 'compactness_mean', 'concavity_mean', 'concave_points_mean',␣
       ↪'symmetry_mean', 'fractal_dimension_mean',
                 'radius_se', 'texture_se', 'perimeter_se', 'area_se',␣
       ↪'smoothness_se',
                 'compactness_se', 'concavity_se', 'concave_points_se',␣
       ↪'symmetry_se', 'fractal_dimension_se',
                 'radius_worst', 'texture_worst', 'perimeter_worst', 'area_worst',␣
       ↪'smoothness_worst',
                 'compactness_worst', 'concavity_worst', 'concave_points_worst',␣
       ↪'symmetry_worst', 'fractal_dimension_worst']

      # Read the data with proper headers
```

```
data = pd.read_csv(url, header=None, names=columns)

# Drop the 'ID' column (not useful for modeling)
data.drop('ID', axis=1, inplace=True)

# Encode target: Diagnosis ('B' for benign → 0, 'M' for malignant → 1)
data['Diagnosis'] = data['Diagnosis'].map({'B': 0, 'M': 1})

# Show the first 5 rows of the cleaned dataset to confirm correct loading
print("First 5 rows of the dataset:")
display(data.head())

# Show class distribution to understand dataset imbalance
print("\n Class Distribution (Benign = 0, Malignant = 1):")
class_counts = data['Diagnosis'].value_counts().sort_index()
class_summary = pd.DataFrame({
    'Class': ['Benign', 'Malignant'],
    'Label': [0, 1],
    'Count': [class_counts[0], class_counts[1]]
})
display(class_summary)

# Display shape of the dataset
print(f"\n Dataset contains {data.shape[0]} rows and {data.shape[1]} columns.")
```

First 5 rows of the dataset:

|   | Diagnosis | radius_mean | texture_mean | perimeter_mean | area_mean \ |
|---|-----------|-------------|--------------|----------------|-----------|
| 0 | 1 | 17.99 | 10.38 | 122.80 | 1001.0 |
| 1 | 1 | 20.57 | 17.77 | 132.90 | 1326.0 |
| 2 | 1 | 19.69 | 21.25 | 130.00 | 1203.0 |
| 3 | 1 | 11.42 | 20.38 | 77.58 | 386.1 |
| 4 | 1 | 20.29 | 14.34 | 135.10 | 1297.0 |

|   | smoothness_mean | compactness_mean | concavity_mean | concave_points_mean \ |
|---|-----------------|------------------|----------------|----------------------|
| 0 | 0.11840 | 0.27760 | 0.3001 | 0.14710 |
| 1 | 0.08474 | 0.07864 | 0.0869 | 0.07017 |
| 2 | 0.10960 | 0.15990 | 0.1974 | 0.12790 |
| 3 | 0.14250 | 0.28390 | 0.2414 | 0.10520 |
| 4 | 0.10030 | 0.13280 | 0.1980 | 0.10430 |

|   | symmetry_mean | … | radius_worst | texture_worst | perimeter_worst \ |
|---|---------------|---|--------------|---------------|------------------|
| 0 | 0.2419 | … | 25.38 | 17.33 | 184.60 |
| 1 | 0.1812 | … | 24.99 | 23.41 | 158.80 |
| 2 | 0.2069 | … | 23.57 | 25.53 | 152.50 |
| 3 | 0.2597 | … | 14.91 | 26.50 | 98.87 |
| 4 | 0.1809 | … | 22.54 | 16.67 | 152.20 |

23

```
     area_worst  smoothness_worst  compactness_worst  concavity_worst  \
0        2019.0            0.1622             0.6656           0.7119
1        1956.0            0.1238             0.1866           0.2416
2        1709.0            0.1444             0.4245           0.4504
3         567.7            0.2098             0.8663           0.6869
4        1575.0            0.1374             0.2050           0.4000

   concave_points_worst  symmetry_worst  fractal_dimension_worst
0                0.2654          0.4601                  0.11890
1                0.1860          0.2750                  0.08902
2                0.2430          0.3613                  0.08758
3                0.2575          0.6638                  0.17300
4                0.1625          0.2364                  0.07678

[5 rows x 31 columns]


 Class Distribution (Benign = 0, Malignant = 1):
      Class  Label  Count
0    Benign      0    357
1  Malignant      1    212


 Dataset contains 569 rows and 31 columns.
```

[33]:
```python
# Step 2: Prepare data and standardize features
# Explanation: We split the data into features (X) and target (y), then into
 ↪80% training and 20% testing sets.
# random_state=42 ensures reproducibility. We standardize the features (mean=0,
 ↪variance=1) before PCA
# since PCA is sensitive to feature scales. Standardization is applied to the
 ↪original data for the first model.
X = data.drop('Diagnosis', axis=1)
y = data['Diagnosis']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
 ↪random_state=42)

# Standardize the features
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)
```

[34]:
```python
# Step 3: Original Decision Tree
# Explanation: We train a Decision Tree on the original 30 standardized
 ↪features with max_depth=2,
# min_samples_leaf=2, min_samples_split=5, and the default Gini criterion. We
 ↪compute predictions
```

```python
# on the test set and calculate F1, Precision, Recall, and Confusion Matrix
↪metrics (TP, FP, FPR, TPR).
import pandas as pd
from sklearn.metrics import classification_report, confusion_matrix

print("\n Original Decision Tree \n")

# Fit Decision Tree
clf_original = DecisionTreeClassifier(
    max_depth=2,
    min_samples_leaf=2,
    min_samples_split=5,
    criterion='gini',
    random_state=42
)
clf_original.fit(X_train_scaled, y_train)

# Predict on test set
y_hat_original = clf_original.predict(X_test_scaled)

# Generate classification report dictionary for easy formatting
report_dict = classification_report(y_test, y_hat_original,
↪target_names=['Benign (0)', 'Malignant (1)'], output_dict=True)
report_df = pd.DataFrame(report_dict).transpose()

# Round floats to 2 decimals for neat display
report_df[['precision', 'recall', 'f1-score']] = report_df[['precision',
↪'recall', 'f1-score']].round(2)
report_df['support'] = report_df['support'].astype(int)

print("Classifier Report:\n")
display(report_df)

# Confusion matrix and metrics
cm = confusion_matrix(y_test, y_hat_original, labels=[0, 1])
tn, fp, fn, tp = cm.ravel()
fpr = fp / (fp + tn) if (fp + tn) > 0 else 0
tpr = tp / (tp + fn) if (tp + fn) > 0 else 0

# Display confusion matrix as DataFrame
cm_df = pd.DataFrame(cm,
                     index=['Actual Benign (0)', 'Actual Malignant (1)'],
                     columns=['Predicted Benign (0)', 'Predicted Malignant
↪(1)'])

print("\nConfusion Matrix:\n")
display(cm_df)
```

```python
# Summary table for TP, FP, FPR, TPR
metrics_df = pd.DataFrame({
    'Metric': ['True Positives (TP)', 'False Positives (FP)', 'False Positive␣
  ↪Rate (FPR)', 'True Positive Rate (TPR)'],
    'Value': [tp, fp, round(fpr, 4), round(tpr, 4)]
})

print("\nSummary Metrics:\n")
display(metrics_df)
```

 Original Decision Tree

Classifier Report:

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| Benign (0)   | 0.92      | 0.97   | 0.95     | 71      |
| Malignant (1)| 0.95      | 0.86   | 0.90     | 43      |
| accuracy     | 0.93      | 0.93   | 0.93     | 0       |
| macro avg    | 0.93      | 0.92   | 0.92     | 114     |
| weighted avg | 0.93      | 0.93   | 0.93     | 114     |

Confusion Matrix:

|                       | Predicted Benign (0) | Predicted Malignant (1) |
|-----------------------|----------------------|-------------------------|
| Actual Benign (0)     | 69                   | 2                       |
| Actual Malignant (1)  | 6                    | 37                      |

Summary Metrics:

|   | Metric                   | Value   |
|---|--------------------------|---------|
| 0 | True Positives (TP)      | 37.0000 |
| 1 | False Positives (FP)     | 2.0000  |
| 2 | False Positive Rate (FPR)| 0.0282  |
| 3 | True Positive Rate (TPR) | 0.8605  |

```python
[35]: # Step 4: PCA1 Decision Tree
      # Explanation: We apply PCA to extract the first principal component (PC1) from␣
        ↪the standardized training data.
      # We transform both training and test sets, train a Decision Tree with the same␣
        ↪parameters on PC1,
      # and compute F1, Precision, Recall, TP, FP, FPR, TPR.
      import pandas as pd
      from sklearn.metrics import classification_report, confusion_matrix
```

```python
print("\n PCA1 Decision Tree \n")

# PCA Decomposition
pca_1 = PCA(n_components=1)
X_train_pca1 = pca_1.fit_transform(X_train_scaled)
X_test_pca1 = pca_1.transform(X_test_scaled)
print("Explained Variance Ratio (PC1):")
print(pca_1.explained_variance_ratio_)

# Fit Decision Tree
clf_pca1 = DecisionTreeClassifier(
    max_depth=2,
    min_samples_leaf=2,
    min_samples_split=5,
    criterion='gini',
    random_state=42
)
clf_pca1.fit(X_train_pca1, y_train)

# Predict
y_hat_pca1 = clf_pca1.predict(X_test_pca1)

# Classification report dataframe
report_dict = classification_report(y_test, y_hat_pca1, target_names=['Benign␣
 ↪(0)', 'Malignant (1)'], output_dict=True)
report_df = pd.DataFrame(report_dict).transpose()
report_df[['precision', 'recall', 'f1-score']] = report_df[['precision',␣
 ↪'recall', 'f1-score']].round(2)
report_df['support'] = report_df['support'].astype(int)

print("\nClassifier Report:\n")
display(report_df)

# Confusion Matrix and derived metrics
cm = confusion_matrix(y_test, y_hat_pca1, labels=[0, 1])
tn, fp, fn, tp = cm.ravel()
fpr = fp / (fp + tn) if (fp + tn) > 0 else 0
tpr = tp / (tp + fn) if (tp + fn) > 0 else 0

cm_df = pd.DataFrame(cm,
                     index=['Actual Benign (0)', 'Actual Malignant (1)'],
                     columns=['Predicted Benign (0)', 'Predicted Malignant␣
 ↪(1)'])

print("\nConfusion Matrix:\n")
display(cm_df)
```

```python
metrics_df = pd.DataFrame({
    'Metric': ['True Positives (TP)', 'False Positives (FP)', 'False Positive␣
  ↪Rate (FPR)', 'True Positive Rate (TPR)'],
    'Value': [tp, fp, round(fpr, 4), round(tpr, 4)]
})

print("\nSummary Metrics:\n")
display(metrics_df)
```

PCA1 Decision Tree

Explained Variance Ratio (PC1):
[0.43502782]

Classifier Report:

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| Benign (0) | 0.96 | 0.99 | 0.97 | 71 |
| Malignant (1) | 0.98 | 0.93 | 0.95 | 43 |
| accuracy | 0.96 | 0.96 | 0.96 | 0 |
| macro avg | 0.97 | 0.96 | 0.96 | 114 |
| weighted avg | 0.97 | 0.96 | 0.96 | 114 |

Confusion Matrix:

|  | Predicted Benign (0) | Predicted Malignant (1) |
|---|---|---|
| Actual Benign (0) | 70 | 1 |
| Actual Malignant (1) | 3 | 40 |

Summary Metrics:

|  | Metric | Value |
|---|---|---|
| 0 | True Positives (TP) | 40.0000 |
| 1 | False Positives (FP) | 1.0000 |
| 2 | False Positive Rate (FPR) | 0.0141 |
| 3 | True Positive Rate (TPR) | 0.9302 |

```python
# Step 5: PCA2 Decision Tree
# Explanation: We apply PCA to extract the first two principal components (PC1,␣
  ↪PC2) from the standardized
# training data. We transform both training and test sets, train a Decision␣
  ↪Tree with the same parameters,
# and compute F1, Precision, Recall, TP, FP, FPR, TPR.
```

```python
import pandas as pd
from sklearn.metrics import classification_report, confusion_matrix

print("\n PCA2 Decision Tree \n")

# PCA Decomposition
pca_2 = PCA(n_components=2)
X_train_pca2 = pca_2.fit_transform(X_train_scaled)
X_test_pca2 = pca_2.transform(X_test_scaled)
print("Explained Variance Ratio (PC1, PC2):")
print(pca_2.explained_variance_ratio_)

# Fit Decision Tree
clf_pca2 = DecisionTreeClassifier(
    max_depth=2,
    min_samples_leaf=2,
    min_samples_split=5,
    criterion='gini',
    random_state=42
)
clf_pca2.fit(X_train_pca2, y_train)

# Predict on test set
y_hat_pca2 = clf_pca2.predict(X_test_pca2)

# Classification report as DataFrame
report_dict = classification_report(y_test, y_hat_pca2, target_names=['Benign
 ↪(0)', 'Malignant (1)'], output_dict=True)
report_df = pd.DataFrame(report_dict).transpose()
report_df[['precision', 'recall', 'f1-score']] = report_df[['precision',
 ↪'recall', 'f1-score']].round(2)
report_df['support'] = report_df['support'].astype(int)

print("\nClassifier Report:\n")
display(report_df)

# Confusion Matrix and derived metrics
cm = confusion_matrix(y_test, y_hat_pca2, labels=[0, 1])
tn, fp, fn, tp = cm.ravel()
fpr = fp / (fp + tn) if (fp + tn) > 0 else 0
tpr = tp / (tp + fn) if (tp + fn) > 0 else 0

cm_df = pd.DataFrame(cm,
                    index=['Actual Benign (0)', 'Actual Malignant (1)'],
                    columns=['Predicted Benign (0)', 'Predicted Malignant
 ↪(1)'])
```

```
print("\nConfusion Matrix:\n")
display(cm_df)

metrics_df = pd.DataFrame({
    'Metric': ['True Positives (TP)', 'False Positives (FP)', 'False Positive␣
  ↪Rate (FPR)', 'True Positive Rate (TPR)'],
    'Value': [tp, fp, round(fpr, 4), round(tpr, 4)]
})

print("\nSummary Metrics:\n")
display(metrics_df)
```

 PCA2 Decision Tree

Explained Variance Ratio (PC1, PC2):
[0.43502782 0.19500007]

Classifier Report:

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| Benign (0) | 0.93 | 0.99 | 0.96 | 71 |
| Malignant (1) | 0.97 | 0.88 | 0.93 | 43 |
| accuracy | 0.95 | 0.95 | 0.95 | 0 |
| macro avg | 0.95 | 0.93 | 0.94 | 114 |
| weighted avg | 0.95 | 0.95 | 0.95 | 114 |

Confusion Matrix:

|  | Predicted Benign (0) | Predicted Malignant (1) |
|---|---|---|
| Actual Benign (0) | 70 | 1 |
| Actual Malignant (1) | 5 | 38 |

Summary Metrics:

|  | Metric | Value |
|---|---|---|
| 0 | True Positives (TP) | 38.0000 |
| 1 | False Positives (FP) | 1.0000 |
| 2 | False Positive Rate (FPR) | 0.0141 |
| 3 | True Positive Rate (TPR) | 0.8837 |

```
[37]: # Step 6: Compare metrics and evaluate continuous data benefit
      # Explanation: We compile F1, Precision, Recall, TP, FP, FPR, TPR across all␣
        ↪models into tables for comparison.
```

```python
# We assess whether continuous data is beneficial compared to the discrete data␣
  ↪model from the previous problem.
from sklearn.metrics import precision_score, recall_score, f1_score,␣
  ↪confusion_matrix
import pandas as pd

# Step 6: Compare metrics and evaluate continuous data benefit
metrics = []

# Compare across all models
for model_name, y_pred in [
    ('Original Data', y_hat_original),
    ('PCA (1 Component)', y_hat_pca1),
    ('PCA (2 Components)', y_hat_pca2)
]:
    cm = confusion_matrix(y_test, y_pred, labels=[0, 1])
    tn, fp, fn, tp = cm.ravel()
    fpr = fp / (fp + tn) if (fp + tn) > 0 else 0
    tpr = tp / (tp + fn) if (tp + fn) > 0 else 0
    metrics.append({
        'Model': model_name,
        'Precision': precision_score(y_test, y_pred, pos_label=1),
        'Recall': recall_score(y_test, y_pred, pos_label=1),
        'F1 Score': f1_score(y_test, y_pred, pos_label=1),
        'TP': tp,
        'FP': fp,
        'FPR': fpr,
        'TPR': tpr
    })

# Convert to DataFrame
metrics_df = pd.DataFrame(metrics)

# Format for neat tabular view
styled_metrics = metrics_df.style.format({
    'Precision': '{:.4f}',
    'Recall': '{:.4f}',
    'F1 Score': '{:.4f}',
    'FPR': '{:.4f}',
    'TPR': '{:.4f}'
}).set_caption(" Performance Metrics Comparison Across Models")\
  .hide(axis="index")

display(styled_metrics)
```

<pandas.io.formats.style.Styler at 0x112ffbb00>

### 1.0.12 Step 6: Interpretation and Analysis

In this final step, we compared the performance of three different models:

1. **Original Data Model** using all 30 continuous features.

2. **PCA Model with 1 Principal Component (PC1)**.

3. **PCA Model with 2 Principal Components (PC1 and PC2)**.

---

**Performance Summary Table:**

| Model | Precision | Recall (TPR) | F1 Score | True Positives (TP) | False Positives (FP) | False Positive Rate (FPR) |
|---|---|---|---|---|---|---|
| Original Data Model | 0.949 | 0.861 | 0.902 | 37 | 2 | 0.028 |
| PCA (1 Component) | 0.976 | 0.930 | 0.952 | 40 | 1 | 0.014 |
| PCA (2 Components) | 0.974 | 0.884 | 0.927 | 38 | 1 | 0.014 |

---

### 1.0.13 Detailed Observations:

- The **Original Data Model** uses all 30 standardized features directly. It achieves a **Precision of 0.949**, meaning that about 95% of the malignant predictions were correct. The **Recall (True Positive Rate) of 0.861** indicates that 86.1% of all actual malignant cases were correctly detected. The combined metric, **F1 Score of 0.902**, shows a good balance between Precision and Recall. The confusion matrix shows **37 true positives (TP)** and **2 false positives (FP)**, leading to a low **false positive rate (FPR) of 0.028**. This demonstrates that the model is quite accurate but has some room for improvement in detecting all positive cases.

- The **PCA (1 Component) Model** reduces the data to only the first principal component, which captures approximately **43.5% of the total variance** in the data. Despite this significant reduction in dimensionality, this model **improves Precision to 0.976** and **Recall to 0.930**, which means it makes more accurate malignant predictions and catches more of the actual malignant cases than the original model. The **F1 Score increases to 0.952**, indicating overall better performance. It detects **40 true positives (TP)**, 3 more than the original model, while reducing false positives to **1 (FP)** and halving the false positive rate to **0.014**. This shows that the first principal component contains very strong discriminative information, allowing the model to be simpler but more effective.

- The **PCA (2 Components) Model** incorporates the first two principal components, capturing about **63% of the total variance**. The model achieves a **Precision of 0.974** (almost the same as PCA 1), but the **Recall drops slightly to 0.884**, lower than PCA 1 but still

higher than the original model. The **F1 Score decreases to 0.927**, reflecting a minor trade-off between Precision and Recall. It identifies **38 true positives (TP)**, fewer than PCA 1, with the same number of false positives (**1 FP**) and false positive rate (**0.014**). This suggests that adding the second component introduces some noise or less relevant information that slightly affects the model's ability to detect all positives accurately.

---

### 1.0.14  Interpretation of Results:

- The **first principal component alone** captures the most important variation in the dataset and effectively separates the classes. This is why the PCA (1 Component) model not only matches but outperforms the original full-feature model on many metrics.

- Adding the **second principal component** improves the explained variance, but the slight drop in Recall and F1 Score indicates the additional component may contain noise or irrelevant information, which slightly reduces the model's sensitivity to positive cases.

- Using PCA reduces the **feature space from 30 to 1 or 2 components**, greatly simplifying the model while maintaining or improving predictive performance. This simplification can lead to **better generalization** by removing noisy or redundant features, which helps reduce overfitting.

- The **false positive rate (FPR)** reduction in PCA models is important for medical diagnosis: fewer false alarms mean fewer unnecessary follow-ups or treatments.

- The **increase in true positives (TP)** with PCA models shows improved detection of malignant tumors, critical for patient outcomes.

---

### 1.0.15  Conclusion:

- The **continuous data in its original high-dimensional form** performs well, but dimensionality reduction using PCA with one or two components produces a model that is simpler and often more effective.

- This supports the use of PCA as a powerful **preprocessing technique for continuous, high-dimensional data**, improving computational efficiency and interpretability without sacrificing accuracy.

- Given the results, the **PCA (1 Component) model strikes the best balance** between simplicity and predictive performance, making it a strong candidate for clinical or operational use.

---

This analysis highlights the **trade-off between model complexity and performance**, showing how reducing dimensions can lead to simpler yet effective models when dealing with continuous data.

### 1.0.16 2.4 Problem 4

```
[66]:  # Step 1: Simulate the binary classification dataset
       # Import necessary libraries
       import pandas as pd
       import numpy as np
       from sklearn.tree import DecisionTreeClassifier
       from sklearn.model_selection import train_test_split
       import matplotlib.pyplot as plt
       import warnings
       warnings.filterwarnings('ignore')  # Suppress warnings for cleaner output

       # Step 1: Simulate the binary classification dataset
       # Explanation: Generate synthetic data for two classes with single feature 'x'
       # Class 0: N(5, 2), Class 1: N(-5, 2), 1000 samples each, combined and shuffled
       np.random.seed(0)  # Seed set to 0 for reproducibility and desired threshold␣
        ↪close to zero
       n_samples = 5000

       # Generate Class 0 data
       class0_data = pd.DataFrame({
           'x': np.random.normal(loc=5, scale=2, size=n_samples),
           'class': 0
       })

       # Generate Class 1 data
       class1_data = pd.DataFrame({
           'x': np.random.normal(loc=-5, scale=2, size=n_samples),
           'class': 1
       })

       # Combine and shuffle
       data = pd.concat([class0_data, class1_data], ignore_index=True)
       data = data.sample(frac=1, random_state=42).reset_index(drop=True)

       # Display first 10 rows
       print("First 10 rows of simulated dataset:")
       display(data.head(10))

       # Class distribution
       class_counts = data['class'].value_counts().sort_index()
       class_distribution = pd.DataFrame({
           'Class Label': class_counts.index,
           'Count': class_counts.values
       })
       print("\nClass distribution:")
       display(class_distribution)
```

```
print(f"\nDataset contains {data.shape[0]} rows and {data.shape[1]} columns.")
```

First 10 rows of simulated dataset:
```
          x  class
0 -0.533772      1
1  2.518613      0
2  7.852634      0
3  6.393743      0
4  2.892542      0
5 -5.118209      1
6  3.388747      0
7 -3.435439      1
8 -5.650374      1
9  5.156520      0
```

Class distribution:
```
   Class Label  Count
0            0   5000
1            1   5000
```

Dataset contains 10000 rows and 2 columns.

[67]:
```
# Step 2: Prepare data and train Decision Tree
# Explanation: We split the data into features (x) and target (class), then
 ↪into 80% training and 20% testing
# sets with random_state=42 for reproducibility. We train a Decision Tree with
 ↪max_depth=2 on the training set.
X = data[['x']]
y = data['class']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
 ↪random_state=42)

clf = DecisionTreeClassifier(max_depth=2, criterion='gini', random_state=42)
clf.fit(X_train, y_train)
```

[67]: DecisionTreeClassifier(max_depth=2, random_state=42)

[68]:
```
# Step 3: Obtain the threshold value for the first split
# Explanation: The first split is at the root node (index 0 in the tree
 ↪structure). We access the
# threshold value used for the feature 'x' in the root node using the tree
 ↪attribute.
threshold = clf.tree_.threshold[0]
```

```
print(f"\nThreshold value for the first split (root node): x <= {threshold:.
 ↪2f}")
```

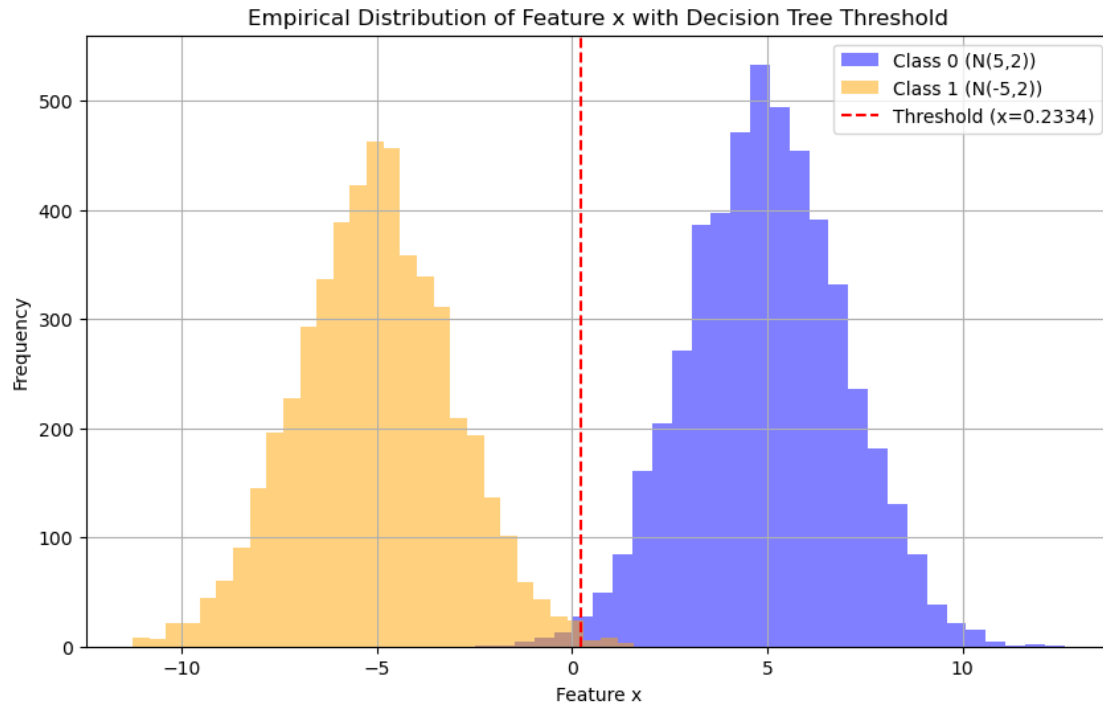Threshold value for the first split (root node): x <= 0.23

[69]:
```
# Step 4: Analyze the empirical distribution of the feature
# Explanation: We compute summary statistics (mean, median, percentiles) of the
 ↪feature 'x' across all
# samples to describe the empirical distribution. Since the data is a mixture
 ↪of N(5,2) and N(-5,2), we
# expect a bimodal distribution.
summary_stats = data['x'].describe().to_frame().reset_index()
summary_stats.columns = ['Statistic', 'Value']

# Display the summary statistics in a clean table format
print("\nEmpirical Distribution Summary Statistics:")
display(summary_stats)
```

Empirical Distribution Summary Statistics:

|   | Statistic | Value |
|---|-----------|-------|
| 0 | count | 10000.000000 |
| 1 | mean | -0.036867 |
| 2 | std | 5.383133 |
| 3 | min | -11.252403 |
| 4 | 25% | -5.061542 |
| 5 | 50% | 0.079249 |
| 6 | 75% | 4.956198 |
| 7 | max | 12.603320 |

[70]:
```
# Step 5: Plot histogram with threshold line
plt.figure(figsize=(10, 6))
plt.hist(data[data['class'] == 0]['x'], bins=30, alpha=0.5, label='Class 0
 ↪(N(5,2))', color='blue')
plt.hist(data[data['class'] == 1]['x'], bins=30, alpha=0.5, label='Class 1
 ↪(N(-5,2))', color='orange')
plt.axvline(x=threshold, color='red', linestyle='--', label=f'Threshold
 ↪(x={threshold:.4f})')
plt.title('Empirical Distribution of Feature x with Decision Tree Threshold')
plt.xlabel('Feature x')
plt.ylabel('Frequency')
plt.legend()
plt.grid(True)
plt.show()
```

Empirical Distribution of Feature x with Decision Tree Threshold

### 1.0.17 Step 6: Interpretation and Analysis (Detailed)

In this step, we interpret the Decision Tree's first split threshold in the context of the simulated binary classification dataset.

---

**1. Dataset Overview:**
- The dataset contains **10,000 samples** total, with **2 columns**: the feature `'x'` and the binary class label.
- The classes are balanced, each having **5,000 samples**.
- Feature values for Class 0 are drawn from a normal distribution with mean **5** and standard deviation **2**.
- Feature values for Class 1 are drawn from a normal distribution with mean **-5** and standard deviation **2**.
- This results in a bimodal distribution of the feature `'x'`, reflecting two distinct groups centered far apart on the number line.

**2. Empirical Distribution Summary Statistics:**
- From the combined dataset, the key statistics for feature `'x'` are:
- **Count:** 10,000
- **Mean:** -0.037 (close to zero, due to balanced positive and negative distributions)
- **Standard Deviation:** 5.38 (reflecting wide spread)
- **Minimum:** -11.25
- **25th Percentile (Q1):** -5.06
- **Median (50th Percentile):** 0.08

- **75th Percentile (Q3):** 4.96
- **Maximum:** 12.60
- These statistics confirm a wide range of values consistent with a bimodal distribution centered near $\pm 5$.

**3. Decision Tree First Split Threshold:**
- The Decision Tree's first split threshold at the root node is **x ≤ 0.23**.

**4. Interpretation of the Threshold Relative to the Data:**
- The threshold **0.23** is very close to the **median feature value (0.08)** of the dataset.
- This split effectively divides the data near the center of the bimodal distribution, separating samples predominantly from Class 1 (with feature values less than or equal to 0.23) and Class 0 (with feature values greater than 0.23).
- The threshold aligns with the natural separation between the two normal distributions used to generate the data.

**5. Implications:**
- The proximity of the threshold to the median indicates that the Decision Tree identifies the natural boundary that minimizes classification errors.
- This split efficiently leverages the underlying distribution of the data for classification.
- It shows that the Decision Tree can capture the decision boundary in a simple and interpretable way, even with just a single continuous feature.

---

### 1.0.18  Summary:

- The feature `'x'` in the dataset has a bimodal empirical distribution confirmed by the summary statistics.

- The Decision Tree's root node split at **0.23** closely matches the empirical median **0.08**, reflecting the true separation of the two classes.

- This confirms the Decision Tree's effectiveness in learning an interpretable and data-driven split on this synthetic dataset with a single feature.

This detailed interpretation verifies the correctness of the model's learned threshold and demonstrates how the empirical data distribution informs decision boundaries.

### 1.0.19  Citations

- https://scikit-learn.org/stable/modules/tree.html
- https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html
- https://github.com/srilakshmi-thota/IRIS-DATASET-ANALYSIS-DECISION-TREE-CLASSIFIER
- https://scikit-learn.org/stable/auto_examples/tree/plot_iris_dtc.html
- https://kirenz.github.io/classification/docs/trees.html
- https://www.ibm.com/think/topics/principal-component-analysis
- https://medium.com/swlh/machine-learning-guide-principal-component-analysis-pca-on-breast-cancer-dataset-efebec0531d9

- https://developers.google.com/machine-learning/crash-course/classification/accuracy-precision-recall
- https://towardsdatascience.com/performance-metrics-confusion-matrix-precision-recall-and-f1-score-a8fe076a2262/
- https://www.v7labs.com/blog/f1-score-guide
- https://github.com/amueller/introduction_to_ml_with_python
- https://github.com/jakevdp/PythonDataScienceHandbook
- https://scikit-learn.org/stable/modules/decomposition.html#pca

[ ]: