# ML LAB 3

Name: Dhanya Prabhu

Date: 19/08/2025

SRN: PES2UG23CS169

Section: C

Mushrooms.csv:

```
Internal Nodes:       5
 PS C:\Dhanya\PESENGINEERING\sem5\ML\Lab_2\code\pytorch_implementation> python test.py --ID EC_C_PES2UG23CS169_Lab3 --data mushrooms.csv --print-tree
Running tests with PYTORCH framework
===============================================================
 target column: 'class' (last column)
Original dataset info:
Shape: (8124, 23)
Columns: ['cap-shape', 'cap-surface', 'cap-color', 'bruises', 'odor', 'gill-attachment', 'gill-spacing', 'gill-size', 'gill-color', 'stalk-shape', 'stalk-root', 'stalk-surface-above-ring', 'stalk-surface-below-ring', 'stalk-color
-above-ring', 'stalk-color-below-ring', 'veil-type', 'veil-color', 'ring-number', 'ring-type', 'spore-print-color', 'population', 'habitat', 'class']

First few rows:

cap-shape: ['x' 'b' 's' 'f' 'k'] -> [5 0 4 2 3]

cap-surface: ['s' 'y' 'f' 'g'] -> [2 3 0 1]

cap-color: ['n' 'y' 'w' 'g' 'e'] -> [4 9 8 3 2]

class: ['p' 'e'] -> [1 0]

Processed dataset shape: torch.Size([8124, 23])
Number of features: 22
Features: ['cap-shape', 'cap-surface', 'cap-color', 'bruises', 'odor', 'gill-attachment', 'gill-spacing', 'gill-size', 'gill-color', 'stalk-shape', 'stalk-root', 'stalk-surface-above-ring', 'stalk-surface-below-ring', 'stalk-colo
r-above-ring', 'stalk-color-below-ring', 'veil-type', 'veil-color', 'ring-number', 'ring-type', 'spore-print-color', 'population', 'habitat']
Target: class
Framework: PYTORCH
Data type: <class 'torch.Tensor'>

===============================================================
DECISION TREE CONSTRUCTION DEMO
===============================================================
Total samples: 8124
Training samples: 6499
Testing samples: 1625

Constructing decision tree using training data...

🌐 Decision tree construction completed using PYTORCH!

🌳 DECISION TREE STRUCTURE
===============================================================
Root [odor] (gain: 0.9083)
├── = 0:
│   ├── Class 0
├── = 1:
│   ├── Class 1
├── = 2:
│   ├── Class 1
├── = 3:
│   ├── Class 0
├── = 4:
│   ├── Class 1
├── = 5:
│   ├── [spore-print-color] (gain: 0.1469)
│   ├── = 0:
│   │   ├── Class 0
│   ├── = 1:
│   │   ├── Class 0
│   ├── = 2:
│   │   ├── Class 0
│   ├── = 3:
│   │   ├── Class 0
│   ├── = 4:
│   │   ├── Class 0
│   ├── = 5:
│   │   ├── Class 1
│   ├── = 7:
│   │   ├── [habitat] (gain: 0.2217)
│   │   ├── = 0:
│   │   │   ├── [gill-size] (gain: 0.7642)
│   │   │   ├── = 0:
│   │   │   │   ├── Class 0
│   │   │   ├── = 1:
│   │   │   │   ├── Class 1
│   │   ├── = 1:
│   │   │   ├── Class 0
│   │   ├── = 2:
│   │   │   ├── [cap-color] (gain: 0.7300)
│   │   │   ├── = 1:
│   │   │   │   ├── Class 0
│   │   │   ├── = 4:
│   │   │   │   ├── Class 0
│   │   │   ├── = 8:
│   │   │   │   ├── Class 1
│   │   │   ├── = 9:
│   │   │   │   ├── Class 1
│   │   ├── = 4:
│   │   │   ├── Class 0
│   │   ├── = 6:
│   │   │   ├── Class 0
│   ├── = 8:
│   │   ├── Class 0
├── = 6:
│   ├── Class 1
├── = 7:
│   ├── Class 1
├── = 8:
│   ├── Class 1

📊 OVERALL PERFORMANCE METRICS
===============================================================
Accuracy:              1.0000 (100.00%)
Precision (weighted): 1.0000
Recall (weighted):     1.0000
F1-Score (weighted):  1.0000
Precision (macro):     1.0000
Recall (macro):        1.0000
F1-Score (macro):      1.0000

🌱 TREE COMPLEXITY METRICS
===============================================================
Maximum Depth:         4
Total Nodes:           29
Leaf Nodes:            24
Internal Nodes:        5
```

Tictactoe.csv:

```
● PS C:\Dhanya\PESENGINEERING\sem5\ML\Lab_2\code\pytorch_implementation> python test.py --ID EC_C_PES2UG23CS169_Lab3 --data tictactoe.csv --print-tree
   Running tests with PYTORCH framework
   ================================================================
    target column: 'Class' (last column)
   Original dataset info:
   Shape: (958, 10)
   Columns: ['top-left-square', 'top-middle-square', 'top-right-square', 'middle-left-square', 'middle-middle-square', 'middle-right-square', 'bottom-left-square', 'bottom-middle-square', 'bottom-right-square', 'Class']

   First few rows:

   top-left-square: ['x' 'o' 'b'] -> [2 1 0]

   top-middle-square: ['x' 'o' 'b'] -> [2 1 0]

   top-right-square: ['x' 'o' 'b'] -> [2 1 0]

   Class: ['positive' 'negative'] -> [1 0]

   Processed dataset shape: torch.Size([958, 10])
   Number of features: 9
   Features: ['top-left-square', 'top-middle-square', 'top-right-square', 'middle-left-square', 'middle-middle-square', 'middle-right-square', 'bottom-left-square', 'bottom-middle-square', 'bottom-right-square']
   Target: Class
   Framework: PYTORCH
   Data type: <class 'torch.Tensor'>

   ================================================================
   DECISION TREE CONSTRUCTION DEMO
   ================================================================
   Total samples: 958
   Training samples: 766
   Testing samples: 192

   Constructing decision tree using training data...

   🌳 Decision tree construction completed using PYTORCH!

   🌲 DECISION TREE STRUCTURE
   ================================================================
   Root [middle-middle-square] (gain: 0.0834)
   ├── = 0:
   │   ├── [bottom-left-square] (gain: 0.1056)
   │   │   ├── = 0:
   │   │   │   ├── [top-right-square] (gain: 0.9024)
   │   │   │   │   ├── = 1:
   │   │   │   │   │   ├── Class 0
   │   │   │   │   ├── = 2:
   │   │   │   │   │   ├── Class 1
   │   │   ├── = 1:
   │   │   │   ├── [top-right-square] (gain: 0.2782)
   │   │   │   │   ├── = 0:
   │   │   │   │   │   ├── Class 0
   │   │   │   │   ├── = 1:
   │   │   │   │   │   ├── Class 0
   │   │   │   │   ├── = 2:
   │   │   │   │   │   ├── [top-left-square] (gain: 0.1767)
   │   │   │   │   │   │   ├── = 0:
   │   │   │   │   │   │   │   ├── [bottom-right-square] (gain: 0.9183)
   │   │   │   │   │   │   │   │   ├── = 1:
   │   │   │   │   │   │   │   │   │   ├── Class 0
   │   │   │   │   │   │   │   │   ├── = 2:
   │   │   │   │   │   │   │   │   │   ├── Class 1
   │   │   │   │   │   │   ├── = 1:
   │   │   │   │   │   │   │   ├── [top-middle-square] (gain: 0.6058)
   │   │   │   │   │   │   │   │   ├── = 0:
   │   │   │   │   │   │   │   │   │   ├── [middle-left-square] (gain: 0.9183)
   │   │   │   │   │   │   │   │   │   │   ├── = 1:
   │   │   │   │   │   │   │   │   │   │   │   ├── Class 0
   │   │   │   │   │   │   │   │   │   │   ├── = 2:
   │   │   │   │   │   │   │   │   │   │   │   ├── Class 1
   │   │   │   │   │   │   │   │   ├── = 1:
   │   │   │   │   │   │   │   │   │   ├── Class 1
   │   │   │   │   │   │   │   │   ├── = 2:
   │   │   │   │   │   │   │   │   │   ├── Class 0
   │   │   │   │   │   │   ├── = 2:
   │   │   │   │   │   │   │   ├── [top-middle-square] (gain: 0.3393)
   │   │   │   │   │   │   │   │   ├── = 0:
   │   │   │   │   │   │   │   │   │   ├── [middle-left-square] (gain: 0.9183)
   │   │   │   │   │   │   │   │   │   │   ├── = 0:
   │   │   │   │   │   │   │   │   │   │   │   ├── Class 0
```

```
                                   ─ = 1:
                                   ├── Class 1
                                   ─ = 2:
                                   ├── Class 0
                              ─ = 1:
                              [middle-left-square] (gain: 0.9183)
                              ─ = 0:
                              ├── Class 1
                              ─ = 1:
                              ├── Class 1
                              ─ = 2:
                              ├── Class 0
                         ─ = 2:
                         ├── Class 1
              ─ = 2:
              [top-right-square] (gain: 0.1225)
              ─ = 0:
              ├── Class 1
              ─ = 1:
              [middle-right-square] (gain: 0.1682)
              ─ = 0:
              ├── Class 1
              ─ = 1:
              [bottom-right-square] (gain: 0.9403)
              ─ = 0:
              ├── Class 1
              ─ = 1:
              ├── Class 0
              ─ = 2:
              ├── Class 1
              ─ = 2:
              [top-left-square] (gain: 0.9183)
              ─ = 0:
              ├── Class 1
              ─ = 1:
              ├── Class 0
              ─ = 2:
              ├── Class 1
         ─ = 2:
         ├── Class 1
   ─ = 1:
   [top-right-square] (gain: 0.0223)
   ─ = 0:
   [bottom-left-square] (gain: 0.2247)
   ─ = 0:
   ├── Class 0
   ─ = 1:
   ├── Class 0
   ─ = 2:
   [middle-right-square] (gain: 0.1159)
   ─ = 0:
   [top-left-square] (gain: 0.1771)
   ─ = 0:
   [middle-left-square] (gain: 0.9183)
   ─ = 0:
   ├── Class 1
   ─ = 1:
   ├── Class 1
   ─ = 2:
   ├── Class 0
   ─ = 1:
   [bottom-right-square] (gain: 0.9710)
   ─ = 1:
   ├── Class 0
   ─ = 2:
   ├── Class 1
   ─ = 2:
   ├── Class 1
   ─ = 1:
   [middle-left-square] (gain: 0.9887)
   ─ = 0:
   ├── Class 1
   ─ = 1:
   ├── Class 0
   ─ = 2:
   ├── Class 1
   ─ = 2:
   [bottom-middle-square] (gain: 0.2400)
   ├── Class 1
```

```
                              ─ = 0:
                              [top-left-square] (gain: 1.0000)
                              ─ = 1:
                              ├── Class 0
                              ─ = 2:
                              ├── Class 1
                         ─ = 1:
                         ├── Class 0
                    ─ = 2:
                    [bottom-right-square] (gain: 0.9710)
                    ─ = 1:
                    ├── Class 0
                    ─ = 2:
                    ├── Class 1
         ─ = 1:
         [bottom-left-square] (gain: 0.4759)
         ─ = 0:
         ├── Class 0
         ─ = 1:
         ├── Class 0
         ─ = 2:
         [top-middle-square] (gain: 0.1974)
         ─ = 0:
         ├── Class 1
         ─ = 1:
         [top-left-square] (gain: 0.3436)
         ─ = 0:
         [bottom-middle-square] (gain: 0.9183)
         ─ = 1:
         ├── Class 0
         ─ = 2:
         ├── Class 1
         ─ = 1:
         ├── Class 0
         ─ = 2:
         [bottom-middle-square] (gain: 0.5917)
         ─ = 0:
         ├── Class 1
         ─ = 1:
         ├── Class 0
         ─ = 2:
         ├── Class 1
         ─ = 2:
         [bottom-right-square] (gain: 0.1245)
         ─ = 0:
         [middle-left-square] (gain: 0.9183)
         ─ = 1:
         ├── Class 0
         ─ = 2:
         ├── Class 1
         ─ = 1:
         [bottom-middle-square] (gain: 0.5613)
         ─ = 0:
         [top-left-square] (gain: 1.0000)
         ─ = 1:
         ├── Class 0
         ─ = 2:
         ├── Class 1
         ─ = 1:
         ├── Class 1
         ─ = 2:
         ├── Class 0
         ─ = 2:
         [bottom-middle-square] (gain: 0.6122)
         ─ = 0:
         ├── Class 0
         ─ = 1:
         [middle-right-square] (gain: 0.9183)
         ─ = 1:
         ├── Class 1
         ─ = 2:
         ├── Class 0
         ─ = 2:
         ├── Class 1
   ─ = 2:
   [bottom-right-square] (gain: 0.0777)
   ─ = 0:
   [top-left-square] (gain: 0.3462)
   ─ = 0:
```

```
                         ─ = 0:
                         ├── Class 0
                    ─ = 1:
                    ├── Class 0
               ─ = 2:
               [top-middle-square] (gain: 0.7008)
               ─ = 0:
               ├── Class 0
               ─ = 1:
               [middle-right-square] (gain: 0.7219)
               ─ = 0:
               ├── Class 0
               ─ = 1:
               ├── Class 1
               ─ = 2:
               ├── Class 0
               ─ = 2:
               ├── Class 0
          ─ = 1:
          [top-left-square] (gain: 0.5439)
          ─ = 0:
          ├── Class 0
          ─ = 1:
          ├── Class 0
          ─ = 2:
          [top-middle-square] (gain: 0.4687)
          ─ = 0:
          [bottom-middle-square] (gain: 0.9183)
          ─ = 0:
          ├── Class 1
          ─ = 1:
          ├── Class 0
          ─ = 2:
          ├── Class 0
          ─ = 1:
          [middle-right-square] (gain: 0.9183)
          ─ = 0:
          ├── Class 1
          ─ = 1:
          ├── Class 1
          ─ = 2:
          ├── Class 0
          ─ = 2:
          ├── Class 1
     ─ = 2:
     [middle-right-square] (gain: 0.4731)
     ─ = 0:
     [top-middle-square] (gain: 0.6464)
     ─ = 0:
     ├── Class 1
     ─ = 1:
     ├── Class 0
     ─ = 2:
     [top-left-square] (gain: 0.8113)
     ─ = 1:
     ├── Class 0
     ─ = 2:
     ├── Class 1
     ─ = 1:
     [middle-left-square] (gain: 0.3995)
     ─ = 0:
     [bottom-middle-square] (gain: 0.8113)
     ─ = 0:
     ├── Class 1
     ─ = 1:
     ├── Class 0
     ─ = 2:
     ├── Class 1
     ─ = 1:
     ├── Class 0
     ─ = 2:
     [top-middle-square] (gain: 0.8113)
     ─ = 1:
     ├── Class 0
     ─ = 2:
     ├── Class 1
   ├── Class 1
```

```
   ─ = 2:
   [bottom-right-square] (gain: 0.0269)
   ─ = 0:
   [top-left-square] (gain: 0.1239)
   ─ = 0:
   ├── Class 1
   ─ = 1:
   [bottom-middle-square] (gain: 0.1033)
   ─ = 0:
   [middle-left-square] (gain: 0.1605)
   ─ = 0:
   ├── Class 1
   ─ = 1:
   [bottom-left-square] (gain: 1.0000)
   ─ = 1:
   ├── Class 0
   ─ = 2:
   ├── Class 1
   ─ = 2:
   [middle-right-square] (gain: 0.5917)
   ─ = 0:
   ├── Class 0
   ─ = 1:
   ├── Class 1
   ─ = 2:
   ├── Class 1
   ─ = 1:
   ├── Class 1
   ─ = 2:
   [top-middle-square] (gain: 0.4591)
   ─ = 0:
   [middle-right-square] (gain: 0.9183)
   ─ = 0:
   ├── Class 0
   ─ = 1:
   ├── Class 1
   ─ = 2:
   ├── Class 0
   ─ = 1:
   [top-right-square] (gain: 0.6122)
   ─ = 0:
   ├── Class 1
   ─ = 1:
   ├── Class 0
   ─ = 2:
   [middle-right-square] (gain: 0.9183)
   ─ = 0:
   ├── Class 1
   ─ = 1:
   ├── Class 1
   ─ = 2:
   ├── Class 0
   ─ = 2:
   ├── Class 1
   ─ = 2:
   ├── Class 1
   ─ = 1:
   [top-left-square] (gain: 0.0713)
   ─ = 0:
   [middle-right-square] (gain: 0.0760)
   ─ = 0:
   [bottom-left-square] (gain: 0.2455)
   ─ = 0:
   ├── Class 1
   ─ = 1:
   [bottom-middle-square] (gain: 0.9710)
   ─ = 0:
   ├── Class 0
   ─ = 2:
   ├── Class 1
   ─ = 2:
   ├── Class 1
   ─ = 1:
   [top-right-square] (gain: 0.7207)
   ─ = 0:
   ├── Class 1
```

```
                         ├── Class 0
                    ─ = 2:
                    [middle-left-square] (gain: 0.3060)
                    ─ = 0:
                    ├── Class 1
                    ─ = 1:
                    ├── Class 1
                    ─ = 2:
                    [top-middle-square] (gain: 1.0000)
                    ─ = 0:
                    ├── Class 1
                    ─ = 2:
                    ├── Class 0
               ─ = 2:
               [middle-left-square] (gain: 0.2294)
               ─ = 0:
               [top-middle-square] (gain: 0.5000)
               ─ = 0:
               ├── Class 0
               ─ = 1:
               ├── Class 1
               ─ = 2:
               [top-right-square] (gain: 1.0000)
               ─ = 0:
               ├── Class 0
               ─ = 1:
               ├── Class 1
               ─ = 1:
               [bottom-left-square] (gain: 0.3219)
               ─ = 0:
               ├── Class 1
               ─ = 1:
               [top-right-square] (gain: 1.0000)
               ─ = 0:
               ├── Class 1
               ─ = 2:
               ├── Class 0
               ─ = 2:
               ├── Class 1
               ─ = 2:
               ├── Class 1
          ─ = 1:
          [bottom-left-square] (gain: 0.1145)
          ─ = 0:
          ├── Class 1
          ─ = 1:
          [middle-left-square] (gain: 0.3407)
          ─ = 0:
          [bottom-middle-square] (gain: 0.9183)
          ─ = 1:
          ├── Class 0
          ─ = 2:
          ├── Class 1
          ─ = 1:
          ├── Class 0
          ─ = 2:
          [bottom-middle-square] (gain: 0.6500)
          ─ = 0:
          ├── Class 1
          ─ = 1:
          ├── Class 0
          ─ = 2:
          ├── Class 1
          ─ = 2:
          [top-right-square] (gain: 0.1913)
          ─ = 0:
          ├── Class 1
          ─ = 1:
          [top-middle-square] (gain: 0.3774)
          ─ = 0:
          ├── Class 1
          ─ = 1:
          ├── Class 0
          ─ = 2:
          [middle-right-square] (gain: 0.8113)
          ─ = 0:
          ├── Class 1
          ─ = 1:
          ├── Class 0
```

```
                         ├── Class 1
                    ─ = 1:
                    ├── Class 0
                    ─ = 2:
                    ├── Class 0
               ─ = 2:
               [middle-left-square] (gain: 0.9710)
               ─ = 0:
               ├── Class 0
               ─ = 1:
               ├── Class 1
               ─ = 2:
               ├── Class 0
          ─ = 1:
          [middle-left-square] (gain: 0.4846)
          ─ = 0:
          ├── Class 1
          ─ = 1:
          [bottom-left-square] (gain: 0.9183)
          ─ = 1:
          ├── Class 0
          ─ = 2:
          ├── Class 1
          ─ = 2:
          [bottom-left-square] (gain: 0.0760)
          ─ = 0:
          ├── Class 0
          ─ = 1:
          [top-right-square] (gain: 0.9183)
          ─ = 1:
          ├── Class 1
          ─ = 2:
          ├── Class 0
          ─ = 2:
          [top-right-square] (gain: 0.9183)
          ─ = 1:
          ├── Class 0
          ─ = 2:
          ├── Class 1
     ─ = 2:
     [middle-left-square] (gain: 0.3425)
     ─ = 0:
     [top-right-square] (gain: 0.9710)
     ─ = 1:
     ├── Class 1
     ─ = 2:
     ├── Class 0
     ─ = 1:
     [top-right-square] (gain: 0.9183)
     ─ = 0:
     ├── Class 0
     ─ = 1:
     ├── Class 0
     ─ = 2:
     ├── Class 1
     ─ = 2:
     ├── Class 1
   ─ = 2:
   ├── Class 1
```

📊 OVERALL PERFORMANCE METRICS
============================================
Accuracy:              0.8730 (87.30%)
Precision (weighted):  0.8741
Recall (weighted):     0.8730
F1-Score (weighted):   0.8734
Precision (macro):     0.8590
Recall (macro):        0.8638
F1-Score (macro):      0.8613

🌳 TREE COMPLEXITY METRICS
============================================
Maximum Depth:          7
Total Nodes:            281
Leaf Nodes:             180
Internal Nodes:         101

Nursery.csv:

```
PS C:\Dhanya\PESENGINEERING\sem5\ML\Lab_2\code\pytorch_implementation> python test.py --ID EC_C_PES2UG23CS169_Lab3 --data Nursery.csv --print-tree
Running tests with PYTORCH framework
=========================================================
 target column: 'class' (last column)
Original dataset info:
Shape: (12960, 9)
Columns: ['parents', 'has_nurs', 'form', 'children', 'housing', 'finance', 'social', 'health', 'class']

First few rows:

parents: ['usual' 'pretentious' 'great_pret'] -> [2 1 0]

has_nurs: ['proper' 'less_proper' 'improper' 'critical' 'very_crit'] -> [3 2 1 0 4]

form: ['complete' 'completed' 'incomplete' 'foster'] -> [0 1 3 2]

class: ['recommend' 'priority' 'not_recom' 'very_recom' 'spec_prior'] -> [2 1 0 4 3]

Processed dataset shape: torch.Size([12960, 9])
Number of features: 8
Features: ['parents', 'has_nurs', 'form', 'children', 'housing', 'finance', 'social', 'health']
Target: class
Framework: PYTORCH
Data type: <class 'torch.Tensor'>


=========================================================
DECISION TREE CONSTRUCTION DEMO
=========================================================
Total samples: 12960
Training samples: 10368
Testing samples: 2592

Constructing decision tree using training data...
```

```
                    = 2:
                    ├── Class 1
                = 2:
                ├── [social] (gain: 0.1579)
                    = 0:
                    ├── [housing] (gain: 0.1963)
                        = 0:
                        ├── [finance] (gain: 0.4934)
                            = 0:
                            ├── Class 4
                            = 1:
                            ├── [form] (gain: 0.6058)
                                = 0:
                                ├── Class 4
                                = 1:
                                ├── Class 4
                                = 2:
                                ├── Class 1
                                = 3:
                                ├── Class 1
                        = 1:
                        ├── [form] (gain: 0.1555)
                            = 0:
                            ├── [children] (gain: 0.8631)
                                = 0:
                                ├── Class 4
                                = 1:
                                ├── Class 1
                                = 2:
                                ├── Class 1
                                = 3:
                                ├── Class 1
                            = 1:
                            ├── Class 1
                            = 2:
                            ├── Class 1
                            = 3:
                            ├── Class 1
                        = 2:
                        ├── [children] (gain: 0.5185)
                            = 0:
                            ├── [form] (gain: 0.7219)
                                = 0:
                                ├── Class 4
                                = 1:
                                ├── Class 4
                                = 2:
                                ├── Class 1
                                = 3:
                                ├── Class 4
                            = 1:
                            ├── [form] (gain: 0.9710)
                                = 0:
                                ├── Class 4
                                = 1:
                                ├── Class 4
                                = 3:
                                ├── Class 1
                            = 2:
                            ├── Class 1
                            = 3:
                            ├── Class 1
                    = 1:
                    ├── Class 1
                    = 2:
                    ├── [housing] (gain: 0.1933)
                        = 0:
                        ├── [finance] (gain: 0.4243)
                            = 0:
                            ├── Class 4
                            = 1:
                            ├── [children] (gain: 0.4228)
                                = 0:
                                ├── Class 4
                                = 1:
                                ├── Class 1
                                = 2:
                                ├── Class 1
                                = 3:
```

```
                            ├── Class 1
                        = 1:
                        ├── [children] (gain: 0.1793)
                            = 0:
                            ├── [form] (gain: 0.9183)
                                = 0:
                                ├── Class 4
                                = 1:
                                ├── Class 1
                                = 2:
                                ├── Class 1
                                = 3:
                                ├── Class 1
                            = 1:
                            ├── Class 1
                            = 2:
                            ├── Class 1
                            = 3:
                            ├── Class 1
                        = 2:
                        ├── [children] (gain: 0.4667)
                            = 0:
                            ├── [form] (gain: 0.6500)
                                = 0:
                                ├── Class 4
                                = 1:
                                ├── Class 4
                                = 2:
                                ├── Class 1
                                = 3:
                                ├── Class 1
                            = 1:
                            ├── [form] (gain: 1.0000)
                                = 0:
                                ├── Class 4
                                = 1:
                                ├── Class 4
                                = 2:
                                ├── Class 1
                                = 3:
                                ├── Class 1
                            = 2:
                            ├── Class 1
                            = 3:
                            ├── Class 1
                = 2:
                ├── [social] (gain: 0.1983)
                    = 0:
                    ├── [parents] (gain: 0.1465)
                        = 0:
                        ├── Class 1
                        = 1:
                        ├── [housing] (gain: 0.2147)
                            = 0:
                            ├── [finance] (gain: 0.4408)
                                = 0:
                                ├── Class 4
                                = 1:
                                ├── [children] (gain: 0.4353)
                                    = 0:
                                    ├── Class 4
                                    = 1:
                                    ├── Class 1
                                    = 2:
                                    ├── Class 1
                                    = 3:
                                    ├── Class 1
                            = 1:
                            ├── [form] (gain: 0.0948)
                                = 0:
                                ├── [children] (gain: 0.7219)
                                    = 0:
                                    ├── Class 4
                                    = 1:
                                    ├── Class 1
                                    = 2:
                                    ├── Class 1
```

```
                        = 2:
                        ├── Class 1
                        = 3:
                        ├── Class 1
                    = 1:
                    ├── Class 1
                    = 2:
                    ├── Class 1
                    = 3:
                    ├── Class 1
                = 2:
                ├── [children] (gain: 0.4054)
                    = 0:
                    ├── [form] (gain: 0.8631)
                        = 0:
                        ├── Class 4
                        = 1:
                        ├── Class 4
                        = 2:
                        ├── Class 1
                        = 3:
                        ├── Class 4
                    = 1:
                    ├── [form] (gain: 0.9852)
                        = 0:
                        ├── Class 4
                        = 1:
                        ├── Class 4
                        = 2:
                        ├── Class 1
                        = 3:
                        ├── Class 1
                    = 2:
                    ├── Class 1
                    = 3:
                    ├── Class 1
            = 2:
            ├── [housing] (gain: 0.2021)
                = 0:
                ├── [finance] (gain: 0.5127)
                    = 0:
                    ├── Class 4
                    = 1:
                    ├── [children] (gain: 0.4345)
                        = 0:
                        ├── Class 4
                        = 1:
                        ├── Class 1
                        = 2:
                        ├── Class 1
                        = 3:
                        ├── Class 1
                = 1:
                ├── [form] (gain: 0.1589)
                    = 0:
                    ├── [children] (gain: 0.8631)
                        = 0:
                        ├── Class 4
                        = 1:
                        ├── Class 1
                        = 2:
                        ├── Class 1
                        = 3:
                        ├── Class 1
                    = 1:
                    ├── Class 1
                    = 2:
                    ├── Class 1
                    = 3:
                    ├── Class 1
                = 2:
                ├── [children] (gain: 0.3632)
                    = 0:
                    ├── [form] (gain: 0.9183)
                        = 0:
                        ├── Class 4
                        = 1:
                        ├── Class 4
                        = 2:
```

```
                    ├── Class 4
                    = 2:
                    ├── Class 1
                    = 3:
                    ├── Class 4
                = 1:
                ├── [form] (gain: 0.9852)
                    = 0:
                    ├── Class 4
                    = 1:
                    ├── Class 4
                    = 2:
                    ├── Class 1
                    = 3:
                    ├── Class 1
                = 2:
                ├── Class 1
                = 3:
                ├── Class 1
    ├── = 1:
    ├── [parents] (gain: 0.4439)
        = 0:
        ├── [housing] (gain: 0.1910)
            = 0:
            ├── [finance] (gain: 0.4530)
                = 0:
                ├── Class 1
                = 1:
                ├── [children] (gain: 0.4591)
                    = 0:
                    ├── Class 1
                    = 1:
                    ├── Class 1
                    = 2:
                    ├── Class 3
                    = 3:
                    ├── Class 3
            = 1:
            ├── [form] (gain: 0.2011)
                = 0:
                ├── [children] (gain: 0.9710)
                    = 0:
                    ├── Class 1
                    = 1:
                    ├── Class 3
                    = 2:
                    ├── Class 3
                = 1:
                ├── Class 3
                = 2:
                ├── Class 3
                = 3:
                ├── Class 3
            = 2:
            ├── [children] (gain: 0.4729)
                = 0:
                ├── [form] (gain: 0.6500)
                    = 0:
                    ├── Class 1
                    = 1:
                    ├── Class 1
                    = 2:
                    ├── Class 3
                    = 3:
                    ├── Class 1
                = 1:
                ├── [form] (gain: 1.0000)
                    = 0:
                    ├── Class 1
                    = 1:
                    ├── Class 1
                    = 2:
                    ├── Class 3
                    = 3:
                    ├── Class 3
                = 2:
                ├── Class 3
                = 3:
                ├── Class 3
```

```
            = 2:
            ├── Class 1
    = 2:
    ├── [parents] (gain: 0.1553)
        = 0:
        ├── Class 1
        = 1:
        ├── [housing] (gain: 0.2299)
            = 0:
            ├── [finance] (gain: 0.4139)
                = 0:
                ├── Class 4
                = 1:
                ├── [children] (gain: 0.4997)
                    = 0:
                    ├── Class 4
                    = 1:
                    ├── Class 4
                    = 2:
                    ├── Class 1
                    = 3:
                    ├── Class 1
            = 1:
            ├── [form] (gain: 0.2422)
                = 0:
                ├── [children] (gain: 1.0000)
                    = 0:
                    ├── Class 4
                    = 2:
                    ├── Class 1
                    = 3:
                    ├── Class 1
                = 1:
                ├── Class 1
                = 2:
                ├── Class 1
                = 3:
                ├── Class 1
            = 2:
            ├── [children] (gain: 0.3854)
                = 0:
                ├── [form] (gain: 0.9710)
                    = 1:
                    ├── Class 4
                    = 2:
                    ├── Class 1
                    = 3:
                    ├── Class 4
                = 1:
                ├── [form] (gain: 0.9852)
                    = 0:
                    ├── Class 4
                    = 1:
                    ├── Class 4
                    = 2:
                    ├── Class 1
                    = 3:
                    ├── Class 1
                = 2:
                ├── Class 1
                = 3:
                ├── Class 1
        = 2:
        ├── [housing] (gain: 0.1933)
            = 0:
            ├── [finance] (gain: 0.3888)
                = 0:
                ├── Class 4
                = 1:
                ├── [children] (gain: 0.4039)
                    = 0:
                    ├── Class 4
                    = 1:
                    ├── Class 1
                    = 2:
                    ├── Class 1
                    = 3:
                    ├── Class 1
            = 1:
```

```
                ├── [children] (gain: 0.8631)
                    = 0:
                    ├── Class 4
                    = 1:
                    ├── Class 1
                    = 2:
                    ├── Class 1
                    = 3:
                    ├── Class 1
                = 1:
                ├── Class 1
                = 2:
                ├── Class 1
                = 3:
                ├── Class 1
            = 2:
            ├── [children] (gain: 0.4323)
                = 0:
                ├── [form] (gain: 0.8113)
                    = 0:
                    ├── Class 4
                    = 1:
                    ├── Class 4
                    = 2:
                    ├── Class 1
                    = 3:
                    ├── Class 4
                = 1:
                ├── [form] (gain: 0.9183)
                    = 0:
                    ├── Class 4
                    = 2:
                    ├── Class 1
                    = 3:
                    ├── Class 1
                = 2:
                ├── Class 1
                = 3:
                ├── Class 1
    = 3:
    ├── [parents] (gain: 0.2121)
        = 0:
        ├── [social] (gain: 0.4863)
            = 0:
            ├── Class 1
            = 1:
            ├── [housing] (gain: 0.2658)
                = 0:
                ├── [finance] (gain: 0.3976)
                    = 0:
                    ├── Class 1
                    = 1:
                    ├── [children] (gain: 0.4997)
                        = 0:
                        ├── Class 1
                        = 1:
                        ├── Class 1
                        = 2:
                        ├── Class 3
                        = 3:
                        ├── Class 3
                = 1:
                ├── [form] (gain: 0.0874)
                    = 0:
                    ├── [children] (gain: 0.6500)
                        = 0:
                        ├── Class 1
                        = 1:
                        ├── Class 3
                        = 2:
                        ├── Class 3
                        = 3:
                        ├── Class 3
                    = 1:
                    ├── Class 3
                    = 2:
                    ├── Class 3
                    = 3:
                    ├── Class 3
```

```
            └ = 3:                                              └ = 3:
              ├ Class 3                                           ├ Class 1
        └ = 2:                                              ├ = 1:
          ├ [children] (gain: 0.3084)                         ├ [form] (gain: 0.8631)
          ├ = 0:                                              ├ = 0:
            ├ [form] (gain: 0.9183)                             ├ Class 1
            ├ = 1:                                            ├ = 1:
              ├ Class 1                                         ├ Class 1
            ├ = 2:                                            ├ = 2:
              ├ Class 3                                         ├ Class 3
            ├ = 3:                                            ├ = 3:
              ├ Class 1                                         ├ Class 3
          ├ = 1:                                             ├ = 2:
            ├ [form] (gain: 0.9183)                            ├ Class 3
            ├ = 0:                                            ├ = 3:
              ├ Class 1                                         ├ Class 3
            ├ = 2:                                         ├ = 1:
              ├ Class 3                                       ├ [children] (gain: 0.1481)
            ├ = 3:                                           ├ = 0:
              ├ Class 3                                         ├ [form] (gain: 0.8524)
          ├ = 2:                                               ├ = 0:
            ├ Class 3                                            ├ Class 1
          ├ = 3:                                               ├ = 1:
            ├ Class 3                                            ├ Class 3
      ├ = 2:                                                   ├ = 2:
        ├ Class 1                                               ├ Class 3
  ├ = 1:                                                       ├ = 3:
    ├ [social] (gain: 0.1591)                                   ├ Class 3
    ├ = 0:                                                   ├ = 1:
      ├ [housing] (gain: 0.1614)                               ├ Class 3
      ├ = 0:                                                 ├ = 2:
        ├ [finance] (gain: 0.4315)                             ├ Class 3
        ├ = 0:                                               ├ = 3:
          ├ Class 4                                            ├ Class 3
        ├ = 1:                                            ├ = 2:
          ├ [children] (gain: 0.4353)                        ├ [children] (gain: 0.4086)
          ├ = 0:                                             ├ = 0:
            ├ Class 4                                          ├ [form] (gain: 0.8524)
          ├ = 1:                                              ├ = 0:
            ├ Class 1                                           ├ Class 1
          ├ = 2:                                              ├ = 1:
            ├ Class 1                                           ├ Class 1
          ├ = 3:                                              ├ = 2:
            ├ Class 1                                           ├ Class 3
      ├ = 1:                                                   ├ = 3:
        ├ [children] (gain: 0.1793)                            ├ Class 1
        ├ = 0:                                               ├ = 1:
          ├ [form] (gain: 0.9183)                              ├ [form] (gain: 0.9928)
          ├ = 0:                                               ├ = 0:
            ├ Class 4                                            ├ Class 1
          ├ = 1:                                               ├ = 1:
            ├ Class 1                                            ├ Class 1
          ├ = 2:                                               ├ = 2:
            ├ Class 1                                            ├ Class 3
          ├ = 3:                                               ├ = 3:
            ├ Class 1                                            ├ Class 3
        ├ = 1:                                               ├ = 2:
          ├ Class 1                                            ├ Class 3
        ├ = 2:                                               ├ = 3:
          ├ Class 1                                            ├ Class 3
        ├ = 3:
          ├ Class 1
      ├ = 2:
        ├ [children] (gain: 0.4353)
        ├ = 0:                                    ▓ OVERALL PERFORMANCE METRICS
          ├ [form] (gain: 0.8113)                 ========================================
          ├ = 0:                                  Accuracy:              0.9867 (98.67%)
            ├ Class 4                              Precision (weighted): 0.9876
          ├ = 1:                                  Recall (weighted):     0.9867
            ├ Class 4                              F1-Score (weighted):   0.9872
          ├ = 2:                                  Precision (macro):     0.7604
            ├ Class 1                              Recall (macro):        0.7654
          ├ = 3:                                  F1-Score (macro):      0.7628
            ├ Class 4
        ├ = 1:                                    🌳 TREE COMPLEXITY METRICS
          ├ [form] (gain: 1.0000)                 ========================================
          ├ = 0:                                  Maximum Depth:         7
            ├ Class 4                              Total Nodes:           952
          ├ = 1:                                  Leaf Nodes:            680
                                                  Internal Nodes:        272
```

1. Performance Comparision:
   - Mushrooms achieves perfect classification.
   - Nursery has a very high performance and is almost perfect.
   - TicTacToe is strong but lower to others
2. Tree:
   - Mushroom: Depth – 4

     Nodes – 29

     Root – odor

   - TicTacToe : Depth – 7
     Nodes – 281

Root - middle middle square

- Nursery: Depth – 7

  Nodes – 952

  Root – odor

- Tree complexity: Nursery dataset creates the largest tree as it has many categorical attributes and 5 class targets. Mushroom tree is very small and separable. TicTacToe is in between.

3. Dataset-Specific Insights:
   - Mushrooms
     - Feature Importance: odor is dominant (94% of splits).
     - Class Distribution: Balanced (52% edible, 48% poisonous).
     - Decision Patterns: If odor is foul/fishy, mushroom is poisonous; otherwise edible.
     - Overfitting: Minimal (small tree, perfect separation).
   - Nursery
     - Feature Importance: health (52%), has_nurs-.
     - Class Distribution: Almost balanced across 5 categories.
     - Decision Patterns: First split on health; followed by has_nurs.
     - Overfitting: Deep tree (391 nodes) hints at some overfitting, but accuracy remains high.
   - TicTacToe
     - Feature Importance: bottom-left, top-left, and middle-middle.
     - Class Distribution: Imbalanced (65% one class, 35% other).
     - Decision Patterns: Root splits often capture winning/losing square positions.
     - Overfitting: Larger tree relative to dataset size; risk of memorizing rare board states.

**4. C**omparative Analysis Report

a) Algorithm Performance

- Highest Accuracy: Mushrooms (100%) because of a single dominant feature (odor).

- Dataset Size Effect: Larger datasets (Nursery, TicTacToe) → deeper, more complex trees.

- Number of Features: Too many features (Nursery) means more complexity, risk of overfitting. Few decisive features (Mushrooms) means clean separations.

b) Data Characteristics Impact

- Class Imbalance:

  TicTacToe (65–35 imbalance) lowers performance slightly

  Mushroom balanced classes → perfect accuracy.

  Nursery has mild imbalance → tree depth grows to capture rare classes.

- Binary vs Multi-valued Features: Binary features (TicTacToe) lead to more branching combinations; Multi-valued features (Mushroom's odor, Nursery's health) allow quick separation.

c) Practical Applications

- Mushrooms: Useful in bioinformatics / food safety. High interpretability and reliability.

- Nursery: Mimics admission/recommendation systems, where social & financial attributes matter.

- TicTacToe: Demonstrates game state classification; relevant to board game AI or reinforcement learning benchmarks.

d) Improvements

- Mushrooms: Already perfect; nothing needed.

- Nursery: Pruning the tree to reduce overfitting.

- TicTacToe: Balance dataset, prune tree, or use ensemble methods to generalize better.