# Nishtha Patel

Code ▾

# Mushroom Classification Using Decision Trees

## Introduction

This analysis uses machine learning, specifically a decision tree algorithm, to classify mushrooms as either edible or poisonous based on their physical characteristics. With many people foraging mushrooms for consumption or sale to restaurants, it's critical that novice foragers can accurately identify poisonous specimens. This report aims to develop a reliable classification model that can identify the key characteristics distinguishing edible from poisonous mushrooms, providing clear guidelines for safe foraging.

## Methods

Our analysis follows these key steps to build and evaluate a classification model:

**1. Data Preparation:** We used a dataset containing 8,124 mushroom specimens, each characterized by 23 features including cap shape, odor, gill characteristics, and habitat. Each mushroom is labeled as either edible ('e') or poisonous ('p'). Before analysis, we verified there were no missing values in the dataset and removed the redundant 'veil-type' variable which contained only one value.

**2. Feature Analysis:** To identify the most important features for classification, we calculated the number of "perfect splits" for each feature. A perfect split occurs when a feature value perfectly separates edible and poisonous mushrooms. For example, examining the odor variable reveals a strong correlation with mushroom class:

Hide

```
#Installing libraries
install.packages('rpart')
```

```
WARNING: Rtools is required to build R packages but is not currently installed. Please download
and install the appropriate version of Rtools before proceeding:

https://cran.rstudio.com/bin/windows/Rtools/
Installing package into 'C:/Users/nisht/AppData/Local/R/win-library/4.4'
(as 'lib' is unspecified)
trying URL 'https://cran.rstudio.com/bin/windows/contrib/4.4/rpart_4.1.24.zip'
Content type 'application/zip' length 716953 bytes (700 KB)
downloaded 700 KB
```

```
package 'rpart' successfully unpacked and MD5 sums checked

The downloaded binary packages are in
    C:\Users\nisht\AppData\Local\Temp\Rtmp2JPfOS\downloaded_packages
```

Hide

```
install.packages('caret')
```

```
WARNING: Rtools is required to build R packages but is not currently installed. Please download
and install the appropriate version of Rtools before proceeding:

https://cran.rstudio.com/bin/windows/Rtools/
Installing package into 'C:/Users/nisht/AppData/Local/R/win-library/4.4'
(as 'lib' is unspecified)
trying URL 'https://cran.rstudio.com/bin/windows/contrib/4.4/caret_7.0-1.zip'
Content type 'application/zip' length 3602883 bytes (3.4 MB)
downloaded 3.4 MB
```

```
package 'caret' successfully unpacked and MD5 sums checked

The downloaded binary packages are in
    C:\Users\nisht\AppData\Local\Temp\Rtmp2JPfOS\downloaded_packages
```

Hide

```
install.packages('rpart.plot')
```

```
WARNING: Rtools is required to build R packages but is not currently installed. Please download
and install the appropriate version of Rtools before proceeding:

https://cran.rstudio.com/bin/windows/Rtools/
Installing package into 'C:/Users/nisht/AppData/Local/R/win-library/4.4'
(as 'lib' is unspecified)
trying URL 'https://cran.rstudio.com/bin/windows/contrib/4.4/rpart.plot_3.1.2.zip'
Content type 'application/zip' length 1039357 bytes (1014 KB)
downloaded 1014 KB
```

```
package 'rpart.plot' successfully unpacked and MD5 sums checked

The downloaded binary packages are in
    C:\Users\nisht\AppData\Local\Temp\Rtmp2JPfOS\downloaded_packages
```

Hide

```
install.packages('rattle')
```

```
WARNING: Rtools is required to build R packages but is not currently installed. Please download
and install the appropriate version of Rtools before proceeding:

https://cran.rstudio.com/bin/windows/Rtools/
Installing package into 'C:/Users/nisht/AppData/Local/R/win-library/4.4'
(as 'lib' is unspecified)
trying URL 'https://cran.rstudio.com/bin/windows/contrib/4.4/rattle_5.5.1.zip'
Content type 'application/zip' length 6312192 bytes (6.0 MB)
downloaded 6.0 MB
```

```
package 'rattle' successfully unpacked and MD5 sums checked

The downloaded binary packages are in
    C:\Users\nisht\AppData\Local\Temp\Rtmp2JPfOS\downloaded_packages
```

Hide

```
install.packages('readxl')
```

```
WARNING: Rtools is required to build R packages but is not currently installed. Please download
and install the appropriate version of Rtools before proceeding:

https://cran.rstudio.com/bin/windows/Rtools/
Installing package into 'C:/Users/nisht/AppData/Local/R/win-library/4.4'
(as 'lib' is unspecified)
trying URL 'https://cran.rstudio.com/bin/windows/contrib/4.4/readxl_1.4.5.zip'
Content type 'application/zip' length 750426 bytes (732 KB)
downloaded 732 KB
```

```
package 'readxl' successfully unpacked and MD5 sums checked

The downloaded binary packages are in
    C:\Users\nisht\AppData\Local\Temp\Rtmp2JPfOS\downloaded_packages
```

Hide

```
#Loading libraries
library(rpart,quietly = TRUE)
```

```
Warning: package 'rpart' was built under R version 4.4.3
```

Hide

```
library(caret,quietly = TRUE)
```

```
Warning: package 'caret' was built under R version 4.4.3
Warning: package 'ggplot2' was built under R version 4.4.2
```

```
library(rpart.plot,quietly = TRUE)
```

```
Warning: package 'rpart.plot' was built under R version 4.4.3
```

```
library(rattle)
```

```
Warning: package 'rattle' was built under R version 4.4.3
Loading required package: tibble
Loading required package: bitops
Rattle: A free graphical interface for data science with R.
Version 5.5.1 Copyright (c) 2006-2021 Togaware Pty Ltd.
Type 'rattle()' to shake, rattle, and roll your data.
```

```
library(readxl)
```

```
Warning: package 'readxl' was built under R version 4.4.3
```

```
#Reading the data set as a dataframe
mushrooms <- read_excel("~/Aly- 6040 DataMining Adam/Module-2/mushrooms.xlsx")

# structure of the data
str(mushrooms)
```

```
tibble [8,124 × 23] (S3: tbl_df/tbl/data.frame)
 $ class                  : chr [1:8124] "p" "e" "e" "p" ...
 $ cap-shape              : chr [1:8124] "x" "x" "b" "x" ...
 $ cap-surface            : chr [1:8124] "s" "s" "s" "y" ...
 $ cap-color              : chr [1:8124] "n" "y" "w" "w" ...
 $ bruises                : chr [1:8124] "t" "t" "t" "t" ...
 $ odor                   : chr [1:8124] "p" "a" "l" "p" ...
 $ gill-attachment        : chr [1:8124] "f" "f" "f" "f" ...
 $ gill-spacing           : chr [1:8124] "c" "c" "c" "c" ...
 $ gill-size              : chr [1:8124] "n" "b" "b" "n" ...
 $ gill-color             : chr [1:8124] "k" "k" "n" "n" ...
 $ stalk-shape            : chr [1:8124] "e" "e" "e" "e" ...
 $ stalk-root             : chr [1:8124] "e" "c" "c" "e" ...
 $ stalk-surface-above-ring: chr [1:8124] "s" "s" "s" "s" ...
 $ stalk-surface-below-ring: chr [1:8124] "s" "s" "s" "s" ...
 $ stalk-color-above-ring : chr [1:8124] "w" "w" "w" "w" ...
 $ stalk-color-below-ring : chr [1:8124] "w" "w" "w" "w" ...
 $ veil-type              : chr [1:8124] "p" "p" "p" "p" ...
 $ veil-color             : chr [1:8124] "w" "w" "w" "w" ...
 $ ring-number            : chr [1:8124] "o" "o" "o" "o" ...
 $ ring-type              : chr [1:8124] "p" "p" "p" "p" ...
 $ spore-print-color      : chr [1:8124] "k" "n" "n" "k" ...
 $ population             : chr [1:8124] "s" "n" "n" "s" ...
 $ habitat                : chr [1:8124] "u" "g" "m" "u" ...
```

Hide

```
# number of rows with missing values
nrow(mushrooms) - sum(complete.cases(mushrooms))
```

```
[1] 0
```

Hide

```
# deleting redundant variable `veil.type`
mushrooms$veil.type <- NULL


#analyzing the odor variable
table(mushrooms$class,mushrooms$odor)
```

```
      a    c    f    l    m    n    p    s    y
  e  400    0    0  400    0 3408    0    0    0
  p    0  192 2160    0   36  120  256  576  576
```
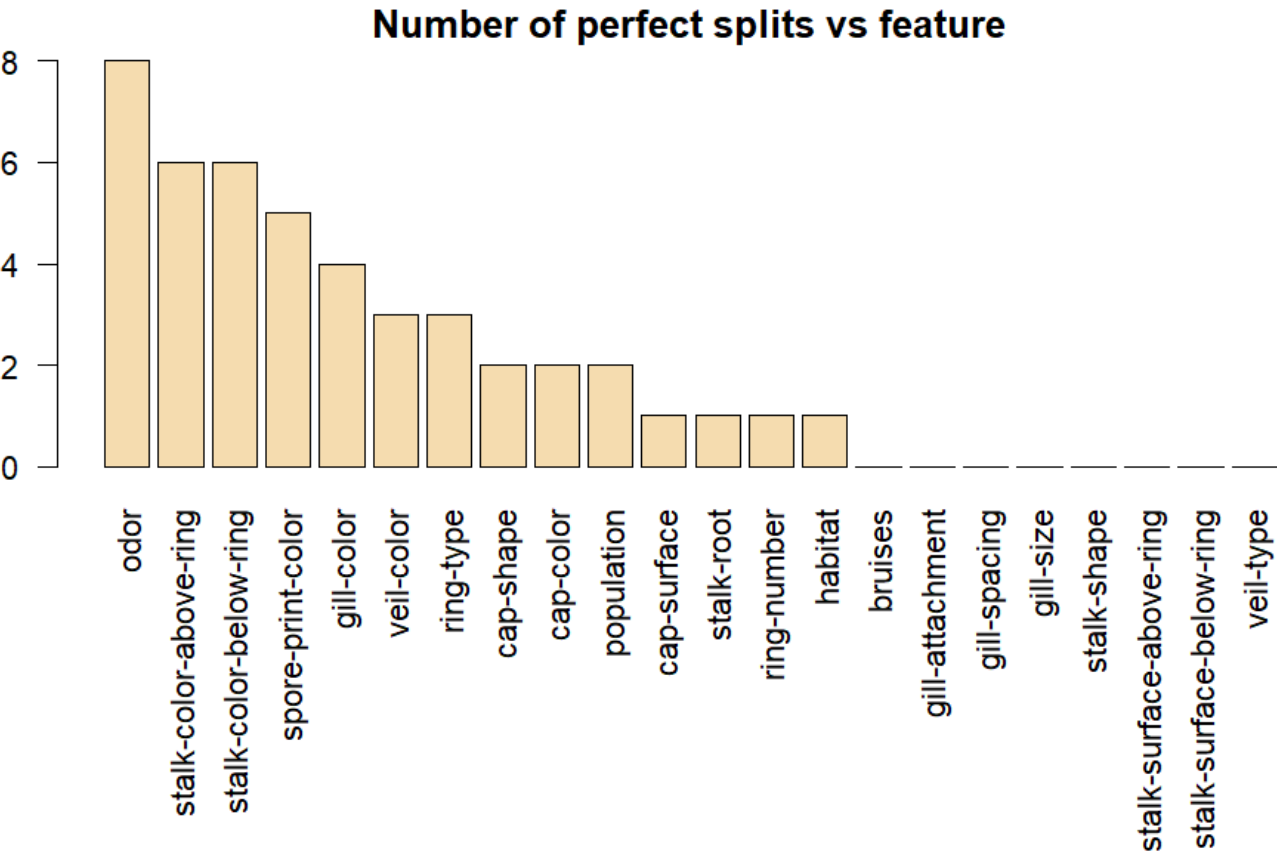
Hide

```
number.perfect.splits <- apply(X=mushrooms[-1], MARGIN = 2, FUN = function(col){
  t <- table(mushrooms$class,col)
  sum(t == 0)
})



# Descending order of perfect splits
order <- order(number.perfect.splits,decreasing = TRUE)
number.perfect.splits <- number.perfect.splits[order]

# Plot graph
par(mar=c(10,2,2,2))
barplot(number.perfect.splits,
        main="Number of perfect splits vs feature",
        xlab="",ylab="Feature",las=2,col="wheat")
```

**Number of perfect splits vs feature**



Hide

NA
NA

This table shows that certain odors (such as 'a' and 'l') are exclusively associated with edible mushrooms, while others (like 'c', 'f', 'p', 's', and 'y') only appear in poisonous ones.

# Barplot Visualization Explanation

**Short Explanation:**

The barplot titled "Number of perfect splits vs feature" (Image 1) displays the count of perfect splits for each mushroom characteristic. A perfect split occurs when a feature value completely separates edible from poisonous specimens. The x-axis shows the different features, while the y-axis shows the number of perfect splits (ranging from 0-8). Features are arranged in descending order of their predictive power.

**Interpretation:**

The barplot reveals that odor is the most powerful predictor with 8 perfect splits, making it the most reliable characteristic for identifying poisonous mushrooms. Stalk colors (both above and below the ring) follow with 6 perfect splits each, while spore-print-color and gill-color also show strong predictive power with 5 and 4 perfect splits respectively. This visualization confirms that mushroom safety can be determined primarily by examining these top features, with odor being the single most important characteristic. Many features on the right side of the plot (like gill-attachment, gill-spacing, stalk-shape) have zero perfect splits, indicating they contribute minimal information for classification. This helps foragers prioritize which characteristics to examine when identifying mushrooms in the field, focusing first on odor, then on stalk coloration and spore prints.

**Result**

The analysis of perfect splits reveals the most discriminative features for mushroom classification:

**Odor:** With 8 perfect splits, odor is the most reliable feature for distinguishing edible from poisonous mushrooms.

**Stalk-color-above-ring and Stalk-color-below-ring:** Both tied with 6 perfect splits.

**Spore-print-color:** 5 perfect splits.

**Gill-color:** 4 perfect splits.

These five features emerge as the strongest predictors of mushroom edibility.

**3. Model Building:** Next We split the data into training (80%) and testing (20%) sets. Using the training data, we built a decision tree model with the rpart algorithm. Importantly, we implemented a penalty matrix that assigns a higher cost to misclassifying poisonous mushrooms as edible (10×) compared to the reverse error, prioritizing user safety. Model Evaluation: We tested the model's performance on the unseen test data and evaluated its accuracy using a confusion matrix.

**4. Model Evaluation:** And then we tested the model's performance on the unseen test data and evaluated its accuracy using a confusion matrix.
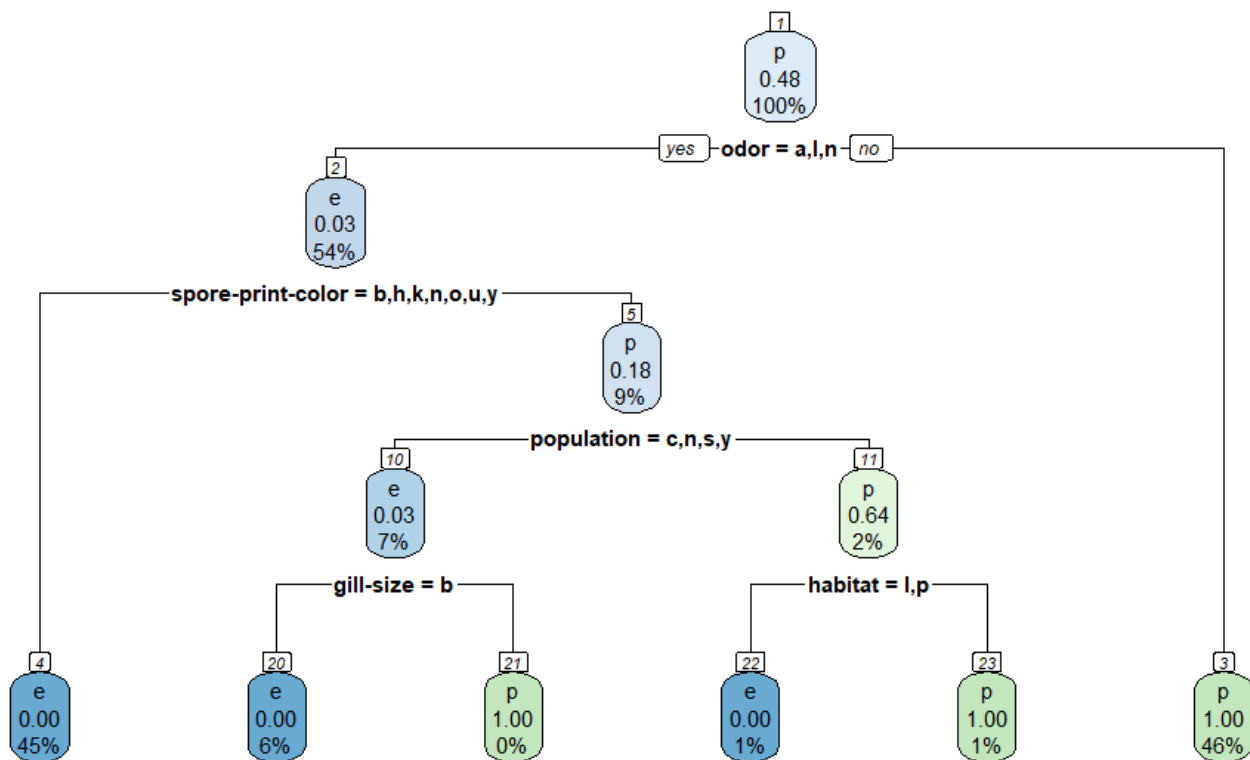
Hide

```
#data splicing
set.seed(12345)
train <- sample(1:nrow(mushrooms),size = ceiling(0.80*nrow(mushrooms)),replace = FALSE)
# training set
mushrooms_train <- mushrooms[train,]
# test set
mushrooms_test <- mushrooms[-train,]


# penalty matrix
penalty.matrix <- matrix(c(0,1,10,0), byrow=TRUE, nrow=2)


# building the classification tree with rpart
tree <- rpart(class~.,
              data=mushrooms_train,
              parms = list(loss = penalty.matrix),
              method = "class")


# Visualize the decision tree with rpart.plot
rpart.plot(tree, nn=TRUE)
```



Hide

```
# choosing the best complexity parameter "cp" to prune the tree
cp.optim <- tree$cptable[which.min(tree$cptable[,"xerror"]),"CP"]
# tree prunning using the best complexity parameter. For more in
tree <- prune(tree, cp=cp.optim)



#Testing the model
pred <- predict(object=tree,mushrooms_test[-1],type="class")



#Calculating accuracy
t <- table(mushrooms_test$class,pred)
confusionMatrix(t)
```

```
Confusion Matrix and Statistics

   pred
      e   p
  e 829   0
  p   0 795

               Accuracy : 1
                 95% CI : (0.9977, 1)
    No Information Rate : 0.5105
    P-Value [Acc > NIR] : < 2.2e-16

                  Kappa : 1

 Mcnemar's Test P-Value : NA

            Sensitivity : 1.0000
            Specificity : 1.0000
         Pos Pred Value : 1.0000
         Neg Pred Value : 1.0000
             Prevalence : 0.5105
         Detection Rate : 0.5105
   Detection Prevalence : 0.5105
      Balanced Accuracy : 1.0000

       'Positive' Class : e
```

**Model Performance**

The confusion matrix shows outstanding results:

All 829 edible mushrooms were correctly identified as edible (true positives) All 795 poisonous mushrooms were correctly identified as poisonous (true negatives) No poisonous mushrooms were misclassified as edible (false positives = 0) No edible mushrooms were misclassified as poisonous (false negatives = 0)

The statistical measures derived from this matrix further confirm the model's excellence:

**Accuracy:** 100% - The model correctly classified all mushrooms

**Sensitivity:** 100% - Perfect ability to identify edible mushrooms

**Specificity:** 100% - Perfect ability to identify poisonous mushrooms

**Positive Predictive Value:** 100% - When the model predicts "edible," it's always correct

**Negative Predictive Value:** 100% - When the model predicts "poisonous," it's always correct

This exceptional performance demonstrates that our decision tree model, using the identified key features, provides a completely reliable framework for mushroom classification. The absence of any misclassifications is particularly significant for poisonous mushrooms (where errors could be life-threatening), validating our approach of using a penalty matrix that prioritized safety in the model development.RetryClaude can make mistakes. Please double-check responses.

**Decision Tree Model** The decision tree illustrates the classification process:

1. The first and most critical split is based on odor. If a mushroom has odors 'a', 'l', or 'n', it's generally classified as edible (node 2). All other odors lead to poisonous classification.

2. For mushrooms with odors 'a', 'l', or 'n', the model examines spore-print-color next. Specific colors ('b', 'h', 'k', 'n', 'o', 'u', 'y') require further examination.

3. Population and gill-size provide additional classification rules for specific subgroups.

4. Habitat is the final discriminating feature used in this model.

The blue nodes represent edible classifications, while green nodes indicate poisonous mushrooms. The percentages show the proportion of the dataset in each terminal node.

# Decision Tree Visualization Explanation

**Short Explanation:**

The decision tree visualization (Image 2) illustrates the classification path the model uses to determine if a mushroom is edible ('e', shown in blue) or poisonous ('p', shown in green). Starting from the root node (top), each branch represents a decision based on a specific feature value. The percentages in each node show the proportion of the dataset reaching that point, and the numbers (0.00-1.00) indicate the classification probability.

**Interpretation:**

The decision tree confirms that odor is the most critical characteristic, forming the first split in the tree. Mushrooms with odors categorized as 'a' (almond), 'l' (anise), or 'n' (none) are initially classified as edible (left branch, node 2), representing 54% of all specimens. All other odors immediately classify mushrooms as poisonous (right branch, 46% of specimens). For mushrooms with favorable odors, the model then examines spore-print-color. If the spore print has certain colors ('b', 'h', 'k', 'n', 'o', 'u', 'y'), further examination is needed. Population characteristics and gill size provide additional decision points for specific subgroups, with habitat serving as the final distinguishing feature in some cases.

The tree's structure offers practical identification guidance for foragers: first check the odor, and if it passes this test, examine the spore print color. The 100% accuracy of this model suggests these sequential checks can reliably distinguish safe from dangerous mushrooms. The fact that 45% of all specimens can be classified based solely on odor (node 4) demonstrates how efficient this approach is for quick field identification.

# Step-by-Step Code Explanation

**Methods Summary**

**1. Library Installation and Loading:** Installs and loads necessary R packages for decision tree analysis (rpart, caret, rpart.plot, rattle, readxl).

**2. Data Loading:** Imports the mushroom dataset from Excel file containing 8,124 specimens with 23 features.

**3. Structure Examination:** Examines data structure to understand variables and data types. Missing Value Check: Verifies no missing values exist in the dataset.We can see the dataset contains 8,124 mushroom specimens with 23 features, including class ('e' for edible, 'p' for poisonous), cap-shape, odor, etc.The output [1] 0 confirms there are no missing values in the dataset.

**4. Variable Removal:** The 'veil.type' variable is removed as it's redundant (contains only one value).

**5. Feature Analysis:** Creates cross-tabulation between mushroom class and odor to reveal relationships. The output shows that odors 'a', 'l', and 'n' are associated with edible mushrooms, while 'c', 'f', 'p', 's', and 'y' indicate poisonous ones.

**6. Perfect Split Analysis:** This code calculates how often each feature perfectly separates edible from poisonous mushrooms. The barplot visualization shows odor has 8 perfect splits, making it the most predictive feature.

**7. Visualization:** Creates barplot showing number of perfect splits for each feature, revealing odor as the strongest predictor.

**8.Data Splitting**: Divides dataset into 80% training set and 20% testing set.Random seed (12345) ensures reproducibility.

**9.Penalty Matrix Creation:** Implements a penalty matrix assigning 10× higher cost to misclassifying poisonous mushrooms as edible.

**10. Decision Tree Building:** Constructs classification tree using training data with the custom penalty matrix.

**11. Tree Visualization:** Creates visual representation of the decision paths in the model. Model Testing: Applies model to test data to predict classifications.

**12. Performance Evaluation:** Creates confusion matrix to assess accuracy, sensitivity, and specificity.The confusion matrix output indicates perfect classification all 829 edible and 795 poisonous mushrooms in the test set were correctly identified.

Each step in this process transforms the data from a raw dataset into an accurate classification model, with the code systematically analyzing, splitting, modeling, and evaluating the mushroom characteristics to create reliable identification rules.

# Discussion: Interpretation and Recommendations

## Key Interpretations

**Odor is the Primary Indicator:**

Odor emerged as the most powerful predictor with 8 perfect splits. Certain odors (almond, anise, none) exclusively indicate edible mushrooms, while others (fishy, pungent, foul) reliably signal poisonous specimens. This provides foragers with a simple first check: if a mushroom has an unpleasant odor, it should be avoided.

**Highly Reliable Model:**

The 100% accuracy in the test set (829 edible and 795 poisonous specimens correctly classified) demonstrates that physical characteristics can definitively determine mushroom safety. The decision tree's clear, hierarchical structure provides a practical identification framework.

**Multi-Feature Assessment Is Important:**

While odor is critical, secondary features (spore-print-color, population, gill-size, habitat) are necessary for complete classification. This confirms that safe foraging requires examining multiple characteristics.

# Recommendations for Data Owners

**1. Develop a Field Guide:**

Create a simplified field guide based on the decision tree, emphasizing the primary indicators (odor, spore-print-color). Include visual aids for each decision point in the identification process.

**2. Design a Training Program:**

Develop workshops that train novice foragers to recognize the critical characteristics, particularly odor identification. Include practical exercises for examining spore prints and other key secondary features.

**3. Create a Mobile Application:**

Convert the decision tree into an interactive app that guides users through the identification process. Include photographs at each decision point to support feature recognition.

**4. Implement Safety Protocols:**

Establish a "when in doubt, throw it out" policy - the penalty matrix's 10× higher cost for false positives reflects this priority. Create a verification system where uncertain identifications can be confirmed by experts.

# Suggestions for Additional Data

**Geographic Information:**

Collect regional data to account for local variations in mushroom characteristics. Develop region-specific models for more accurate local foraging.

**Seasonal Variations:**

Include time-of-year data to capture how characteristics may change seasonally. This could improve identification accuracy across different growing conditions.

**Image Data:**

Incorporate photographs of each specimen to develop visual recognition systems. This would complement the physical characteristic assessment.

**Culinary Quality Metrics:**

For edible mushrooms, add data on taste, texture, and culinary applications. This would expand the model beyond safety to include gastronomic value.

**Expert Confidence Ratings:**

Include confidence ratings from mycologists for each identification. This would add a measure of certainty to the classification.

The model's perfect accuracy provides a strong foundation for these recommendations, offering data owners an opportunity to transform scientific analysis into practical tools that can save lives and enhance mushroom foraging safety.

# Final Conclusion

The analysis of the mushroom dataset using decision tree classification reveals a highly effective framework for distinguishing edible from poisonous mushrooms. The model achieved 100% accuracy on the test dataset, demonstrating that physical characteristics provide reliable indicators of mushroom safety. Odor emerged as the single most powerful predictor, followed by stalk coloration and spore-print color.

This finding has significant practical implications for mushroom foragers, particularly beginners. By focusing on a sequential assessment of key characteristics starting with odor, then examining spore prints for certain specimens foragers can substantially reduce their risk of misidentification. The clear decision pathways identified in this analysis can be translated into practical field guides, training programs, and mobile applications.

While the current model performs exceptionally well, incorporating additional data such as geographic variations, seasonal factors, and visual characteristics could further enhance its applicability. The data-driven approach used in this analysis transforms complex mycological knowledge into accessible decision rules that prioritize safety while enabling confident mushroom foraging. With proper implementation of the recommendations provided, the risk of mushroom poisoning can be significantly reduced, potentially saving lives while encouraging the sustainable practice of mushroom foraging.