# Project Report:

# Decision Tree Classifier for Mushroom Dataset

## 1. Introduction

### 1.1 Background and Motivation

In the field of machine learning, decision trees are a widely used algorithm for both classification and regression tasks due to their simplicity, interpretability, and ability to handle both numerical and categorical data. A decision tree operates by recursively partitioning the input space into distinct and non-overlapping regions, each associated with a specific output or class. The decision-making process within the tree is transparent, making it an excellent tool for domains where model interpretability is crucial.

The mushroom dataset provides a practical application for decision tree algorithms. This dataset includes various features of mushrooms, such as cap shape, color, gill attachment, and stem height, and the goal is to classify whether a mushroom is poisonous or edible. The problem is relevant as it mirrors real-world scenarios where distinguishing between safe and harmful substances (like edible and poisonous mushrooms) can have significant consequences.

Given the importance of accurate classification in such high-stakes environments, it is essential to build a reliable model that not only performs well on training data but also generalizes effectively to unseen data. This challenge is further complicated by the presence of both categorical and numerical features in the dataset, as well as potential issues such as class imbalance and missing values. Therefore, a robust approach to model building, including careful feature engineering, hyperparameter tuning, and the application of techniques to prevent overfitting, is required.

### 1.2 Objectives

The primary objective of this project is to implement a decision tree classifier from scratch, applying it to the mushroom dataset to classify mushrooms as either edible or poisonous. The project aims to explore and demonstrate key machine learning concepts, including:

1. **Data Preprocessing**: Handling and transforming data to make it suitable for model training, including dealing with missing values, encoding categorical features, and normalizing or scaling numerical features.
2. **Model Implementation**: Building a decision tree model from scratch, focusing on core functionalities such as node splitting, choosing the best split based on criteria like Gini impurity, entropy, and misclassification error, and implementing stopping conditions to prevent the tree from becoming overly complex.
3. **Hyperparameter Tuning**: Systematically exploring different configurations of model parameters to identify the combination that yields the best performance. This involves tuning parameters related to tree growth, such as maximum depth, minimum samples required for a split, and the threshold for minimum impurity decrease.

4. **Overfitting and Underfitting Analysis**: Assessing the model's performance on both training and test datasets to identify signs of overfitting (where the model performs well on training data but poorly on unseen data) and underfitting (where the model fails to capture the underlying patterns in the data). The project explores strategies to mitigate these issues, such as pruning and implementing appropriate stopping criteria.
5. **Evaluation and Interpretation**: Evaluating the model's performance using metrics like accuracy, and analyzing the decision rules generated by the tree to gain insights into which features are most influential in determining whether a mushroom is edible or poisonous.

## 1.3 Challenges and Considerations

Several challenges are inherent in building a decision tree model for the mushroom dataset:

1. **Handling Categorical Data**: The dataset contains many categorical features, which require appropriate encoding techniques. While decision trees can naturally handle categorical data, preprocessing steps such as one-hot encoding or label encoding are often necessary, especially when integrating with other tools like GridSearchCV.
2. **Missing Data**: Missing values can significantly impact model performance. Choosing the right imputation strategy (e.g., using the most frequent value for categorical data or the mean for numerical data) is critical to ensure that the model remains robust and reliable.
3. **Class Imbalance**: The dataset have an imbalance between the number of edible and poisonous mushrooms, which could bias the model toward the majority class. Techniques such as oversampling the minority class or adjusting class weights could be necessary to address this issue.
4. **Overfitting and Underfitting**: A primary concern with decision trees is their tendency to overfit the training data, especially when the tree grows too deep. Implementing techniques such as pruning, limiting tree depth, and setting a minimum impurity decrease threshold are essential to prevent overfitting while ensuring that the model is sufficiently expressive to capture the underlying data patterns.
5. **Computational Complexity**: As the decision tree grows, the computational complexity increases, especially during the process of finding the best split at each node. Efficient implementation and thoughtful tuning of hyperparameters are required to manage the computational load and ensure that the model can be trained within a reasonable time frame.

## 1.4 Summary

This project is a comprehensive exploration of decision tree classifiers, applied to a real-world dataset with practical implications. Through the careful implementation of the model, rigorous hyperparameter tuning, and thorough evaluation, the project aims to develop a decision tree model that not only performs well in terms of accuracy but also provides valuable insights into the factors that determine mushroom edibility.

# 2. Methodology

This section provides a comprehensive overview of the steps and processes undertaken to build, train, and evaluate the decision tree classifier for the mushroom dataset. Each step is crucial in ensuring the model's accuracy, interpretability, and generalizability.

## 2.1 Data Preprocessing

Data preprocessing is a critical step in the machine learning pipeline, especially when dealing with a dataset that includes both categorical and numerical variables, as well as potential missing values. The quality of data preprocessing directly influences the model's performance.

**2.1.1 Handling Missing Values**

The dataset contains missing values that need to be addressed before model training. Missing data can lead to biased estimates and reduce the model's accuracy if not handled appropriately.

- **Numerical Data Imputation**: For numerical features, such as `cap-diameter`, `stem-height`, and `stem-width`, missing values were imputed using the mean of the available data in the respective columns. This approach assumes that the missing values are missing at random and that the mean represents the central tendency of the data. By replacing missing values with the mean, iavoid introducing any significant bias into the dataset while retaining the overall distribution of the numerical features.

  ```
  imputer_numerical = SimpleImputer(strategy='mean')
  df[numerical_columns] =
  imputer_numerical.fit_transform(df[numerical_columns])
  ```

- **Categorical Data Imputation**: For categorical features, such as `cap-shape`, `gill-color`, and `habitat`, the most frequent value (mode) was used to fill in missing data. This method assumes that the most common category is a reasonable estimate for the missing values. Imputing with the most frequent value helps maintain the integrity of the categorical distribution without introducing a new category, which could potentially skew the model's performance.

  ```
  imputer_categorical = SimpleImputer(strategy='most_frequent')
  df[categorical_columns] =
  imputer_categorical.fit_transform(df[categorical_columns])
  ```

**2.1.2 Converting Ranges to Mean for Numerical Columns**

In the dataset, certain numerical features were presented as ranges (e.g., `[10, 20]`), which needed to be converted into single values for the decision tree to process them effectively. The following steps were taken:

- **Range Handling**: For numerical features that appeared as ranges, the mean of the lower and upper bounds of the range was calculated and used as the feature value. This conversion ensures that the data remains continuous and can be effectively split by the decision tree algorithm.

  ```
  def convert_range_to_mean(value):
      if isinstance(value, str) and '[' in value:
          numbers = value.strip('[]').split(',')
          numbers = [float(num) for num in numbers]
          return np.mean(numbers)
      else:
          return pd.to_numeric(value, errors='coerce')

  for column in numerical_columns:
      df[column] = df[column].apply(lambda x: convert_range_to_mean(x) if
  pd.notnull(x) else x)
  ```

**2.1.3 Encoding Categorical Variables**

Categorical variables pose a challenge in many machine learning algorithms that require numerical input. However, the decision tree algorithm can handle categorical variables directly, as it can split on categorical values without requiring them to be converted to numerical codes.

- **Categorical Data Preparation**: The categorical features in the mushroom dataset were left in their original form, and no explicit one-hot encoding or label encoding was applied. This approach leverages the decision tree's ability to handle categorical splits naturally. However, categorical features with multiple possible values (e.g., `cap-color` with values like `brown`, `yellow`, etc.) were simplified to ensure that the decision tree could process them effectively.

```
for column in categorical_columns:
    df[column] = df[column].apply(lambda x:
x.strip('[]').split(',')[0].strip() if isinstance(x, str) and ',' in x
else x.strip('[]') if isinstance(x, str) else x)
```

## 2.2 Decision Tree Implementation

Building the decision tree classifier from scratch involved implementing the fundamental operations of the tree, including node creation, splitting based on criteria, and handling the stopping conditions. The implementation aimed to mirror the functionality of established decision tree algorithms like those found in Scikit-learn, while also allowing for customizations and deep insights into the tree-building process.

### 2.2.1 Node Structure

- **TreeNode Class**: Each node in the decision tree is represented by an instance of the `TreeNode` class. A node can either be a decision node, which splits the data based on a feature and threshold, or a leaf node, which holds a predicted value.

```python
Copia codice
class TreeNode:
    def __init__(self, feature=None, threshold=None, left=None,
right=None, *, value=None):
        self.feature = feature
        self.threshold = threshold
        self.left = left
        self.right = right
        self.value = value

    def is_leaf_node(self):
        return self.value is not None
```

### 2.2.2 Splitting Criteria

The splitting criterion determines how the tree decides to split the data at each node. The project implemented three common criteria:

- **Gini Impurity**: Measures the likelihood of incorrect classification of a randomly chosen element if it were randomly labeled according to the distribution of labels in the subset. Gini impurity is commonly used in decision trees because it tends to create balanced splits.
- **Entropy**: Measures the amount of uncertainty or disorder in the data. Splitting based on entropy (information gain) aims to reduce this uncertainty and create purer child nodes.

- **Misclassification Error**: Measures the proportion of incorrect predictions at a node. It is a simple criterion that is sometimes used for simplicity, though it may not always provide the most balanced splits.

```python
def gini_impurity(y):
    hist = np.bincount(y)
    ps = hist / len(y)
    return 1 - np.sum(ps ** 2)

def entropy(y):
    hist = np.bincount(y)
    ps = hist / len(y)
    return -np.sum([p * np.log2(p) for p in ps if p > 0])

def misclassification_error(y):
    hist = np.bincount(y)
    return 1 - np.max(hist) / len(y)
```

### 2.2.3 Stopping Criteria

The tree needs a mechanism to stop growing to prevent it from becoming too complex and overfitting the data. The 3 stopping criteria implemented include:

- **Maximum Depth**: The tree stops growing when the specified maximum depth is reached. This prevents the tree from becoming too deep and complex.
- **Minimum Samples Split**: The tree will not split a node if the number of samples at that node is less than a specified threshold. This prevents the creation of nodes that are too specific to a small number of samples, which could lead to overfitting.
- **Minimum Impurity Decrease**: The tree only makes a split if the reduction in impurity (measured by the chosen criterion) exceeds a certain threshold. This prevents splits that do not significantly improve the model's performance.

```python
class DecisionTreeClassifierFromScratch:
    def __init__(self, criterion="gini", max_depth=100,
min_samples_split=2, min_impurity_decrease=0.0):
        self.criterion = criterion
        self.max_depth = max_depth
        self.min_samples_split = min_samples_split
        self.min_impurity_decrease = min_impurity_decrease
        self.root = None
```

## 2.3 Hyperparameter Tuning

Hyperparameter tuning is essential to ensure that the decision tree model generalizes well to unseen data. It involves systematically searching for the best combination of hyperparameters that maximizes model performance.

### 2.3.1 Grid Search with Cross-Validation

The project employed `GridSearchCV`, a common and effective method for hyperparameter tuning, which performs an exhaustive search over a specified parameter grid:

- **Parameter Grid**: The grid included a variety of parameters, such as the criterion for splitting, the maximum depth of the tree, the minimum samples required to split, and the minimum impurity decrease.

- **Cross-Validation**: The grid search used 5-fold cross-validation to evaluate each set of hyperparameters. This method divides the training data into five parts, trains the model on four parts, and validates it on the fifth. This process repeats five times, each time using a different part for validation. The final model performance is averaged across these five runs, providing a robust estimate of how the model is likely to perform on unseen data.

```python
from sklearn.model_selection import GridSearchCV

param_grid = {
    'criterion': ['gini', 'entropy', 'misclassification'],
    'max_depth': [5, 10, 15],
    'min_samples_split': [2, 4, 6],
    'min_impurity_decrease': [0.0, 0.01, 0.05]  # Adding early stopping
criteria
}

grid_search = GridSearchCV(MyDecisionTreeWrapper(), param_grid, cv=5,
scoring='accuracy')
grid_search.fit(X_train, y_train)

best_params = grid_search.best_params_
best_score = grid_search.best_score_

print("Best parameters found: ", best_params)
print("Best cross-validation accuracy after tuning: ", best_score)
```

## 2.4 Pruning and Stopping Criteria

### 2.4.1 Pruning Techniques

Pruning is an essential technique in decision tree algorithms aimed at improving generalization by reducing overfitting. Overfitting occurs when the model becomes too complex, capturing noise and anomalies in the training data rather than the underlying pattern. By pruning, ican simplify the tree after it has been fully grown, which often leads to better performance on unseen data.

**Types of Pruning**:

1. **Pre-Pruning (Early Stopping)**: In pre-pruning, the growth of the tree is halted early based on certain criteria, such as a maximum depth or minimum impurity decrease. This type of pruning is applied during the tree construction process. By setting these parameters, the tree stops growing when further splits are unlikely to result in significant improvements in classification accuracy.
   - **Maximum Depth**: This parameter restricts the depth of the tree, limiting the number of splits and, consequently, the complexity of the model. By capping the depth, iensure that the tree does not become overly specific to the training data.
   - **Minimum Impurity Decrease**: This threshold ensures that a node will only be split if it results in a significant decrease in impurity (e.g., Gini impurity, entropy). If the decrease in impurity is below this threshold, the split is not performed, and the node is treated as a leaf.
   - **Minimum Samples Split**: This criterion prevents the model from making splits when the number of samples in a node is too small, as such splits are likely to be unreliable and could lead to overfitting.

2. **Post-Pruning**: Post-pruning, also known as cost-complexity pruning or reduced-error pruning, involves growing the tree to its full depth and then removing nodes that contribute little to the model's predictive power. This is done by comparing the performance of the tree with and without certain branches. If pruning a branch does not significantly degrade the tree's performance, the branch is removed.

   o **Cost-Complexity Pruning**: This method adds a penalty for the number of nodes in the tree. The goal is to find a tree that minimizes both the classification error and the complexity of the tree. The penalty is controlled by a parameter called alpha ($\alpha$). During pruning, branches that contribute less than a certain threshold to the overall tree performance (taking into account both accuracy and complexity) are removed.

```
def prune(self, X, y):
    # Perform cost-complexity pruning
    def _prune_tree(node, X, y):
        if node.is_leaf_node():
            return node

        left_idxs, right_idxs = self._split(X[:, node.feature],
node.threshold)
        left = _prune_tree(node.left, X[left_idxs], y[left_idxs])
        right = _prune_tree(node.right, X[right_idxs], y[right_idxs])

        # If both children are leaves, consider merging them
        if left.is_leaf_node() and right.is_leaf_node():
            # Calculate the cost with the children
            left_cost = self._criterion_func(y[left_idxs])
            right_cost = self._criterion_func(y[right_idxs])
            children_cost = (len(left_idxs) * left_cost + len(right_idxs)
* right_cost) / len(y)

            # Calculate the cost without splitting (parent as a leaf)
            parent_cost = self._criterion_func(y)

            # If merging results in a lower or equal cost, prune the
children
            if parent_cost <= children_cost + self.alpha:
                return TreeNode(value=self._most_common_label(y))

        node.left = left
        node.right = right
        return node

    self.root = _prune_tree(self.root, X, y)
```

**Implementation in the Project**:

- The decision tree classifier was first grown without any constraints (i.e., allowing it to become as complex as necessary to fit the training data). Then, post-pruning was applied to remove branches that added little value to the predictive accuracy of the tree. This process involved computing the cost-complexity measure and systematically reducing the tree's complexity while monitoring the impact on performance.

### 2.4.2 Effects of Pruning and Stopping Criteria on Model Performance

Pruning and appropriate stopping criteria are crucial for controlling the model's complexity and ensuring that it generalizes well to new data. Without these techniques, a decision tree is likely to overfit, resulting in poor performance on the test set.

- **Reduction in Overfitting**: Pruning and early stopping criteria, such as `min_impurity_decrease` and `max_depth`, directly address overfitting by preventing the model from becoming overly tailored to the training data. This results in a tree that captures the general patterns in the data without being overly influenced by noise.
- **Improved Generalization**: By controlling the complexity of the tree, pruning helps in striking a balance between bias and variance, leading to better performance on unseen data. This is evidenced by the model's improved accuracy on the test set after pruning, as the tree is no longer overfitting to the idiosyncrasies of the training data.
- **Trade-Offs**: While pruning generally improves generalization, it may also result in slightly lower performance on the training data, as the tree is less complex and thus less able to perfectly memorize the training examples. However, this trade-off is often beneficial, as the goal is to create a model that performs well not only on training data but also on new, unseen data.

### 2.4.3 Early Stopping as an Alternative to Pruning

In some cases, early stopping criteria can be an effective alternative or complement to pruning. By limiting the growth of the tree during the initial training phase, early stopping criteria can prevent the tree from becoming too complex in the first place, reducing the need for extensive post-pruning.

- **Max Depth**: Limiting the maximum depth of the tree is a simple yet effective way to control its complexity. A tree that is too deep is more likely to overfit, while a shallower tree may generalize better, albeit with some risk of underfitting.
- **Min Samples Split**: This criterion ensures that a node must have a certain number of samples before it is split, which prevents the model from creating splits based on small, possibly unrepresentative subsets of the data.
- **Min Impurity Decrease**: By setting a threshold for the minimum decrease in impurity required to make a split, this criterion ensures that the tree only grows in areas where there is a significant gain in information. This helps in maintaining a focus on splits that meaningfully improve the model's predictive power.

## 3. Model Performance

### 3.1 Initial Results

In the initial phase of the project, ievaluated the decision tree model's performance without any hyperparameter tuning or pruning to establish a baseline.

#### 3.1.1 Training and Test Accuracy

- **Training Accuracy**: Initially, the decision tree model showed a very high training accuracy, indicating that it was fitting the training data well. This was expected as the tree was allowed to grow deep enough to capture all patterns, including noise.
- **Test Accuracy**: However, the test accuracy was notably lower than the training accuracy. This discrepancy between training and test accuracy suggested that the model was overfitting, meaning it was too closely tailored to the training data and did not generalize well to new data.

#### 3.1.2 Overfitting

The signs of overfitting were evident from the large gap between training and test accuracy. The deep tree was likely capturing noise in the training data, which did not translate to better performance on the test set. This was particularly visible in the initial results where the model achieved near-perfect accuracy on the training data but much lower accuracy on the test data.

### 3.1.3 Underfitting

To explore underfitting, i experimented with limiting the maximum depth of the tree and increasing the minimum number of samples required to split a node. These changes led to models with lower complexity, which resulted in reduced training accuracy and even lower test accuracy, indicating that the model was too simple to capture the patterns in the data.

### 3.1.4 Criteria Comparison

I compared three splitting criteria: Gini impurity, entropy, and misclassification error. In our experiments:

- **Gini Impurity** generally resulted in the best performance, with a better balance between training and test accuracy.
- **Entropy** performed slightly worse than Gini, particularly in terms of test accuracy.
- **Misclassification Error** often led to suboptimal splits, resulting in lower overall accuracy.

## 3.2 Hyperparameter Tuning

Hyperparameter tuning was performed using GridSearchCV to optimize the decision tree's parameters. The grid search evaluated the following parameters:

- **Criterion**: Gini, entropy, and misclassification error.
- **Max Depth**: 5, 10, and 15.
- **Min Samples Split**: 2, 4, and 6.
- **Min Impurity Decrease**: 0.0, 0.01, and 0.05.

### 3.2.1 Best Model Selection

The grid search revealed that the best model was achieved with the following parameters:

- **Criterion**: Gini
- **Max Depth**: 5
- **Min Samples Split**: 2
- **Min Impurity Decrease**: 0.01

This combination provided the best cross-validation accuracy, striking a balance between model complexity and generalization ability.

### 3.2.2 Insights from Hyperparameter Tuning

- **Max Depth**: A max depth of 5 prevented the tree from growing too complex, helping to reduce overfitting while still capturing key patterns in the data.
- **Min Impurity Decrease**: Setting a threshold of 0.01 for impurity decrease ensured that only significant splits were made, avoiding unnecessary complexity in the tree structure.

## 3.3 Pruning

After hyperparameter tuning, pruning was applied to further enhance the model's generalization.

**3.3.1 Pruning Process**

- **Post-Pruning**: I implemented cost-complexity pruning, where branches of the tree were removed if they contributed little to the model's predictive power. The pruning process involved evaluating the tree with and without specific branches and removing those that did not significantly impact accuracy.

**3.3.2 Impact of Pruning**

- **Test Accuracy After Pruning**: Pruning led to an improved test accuracy of approximately 0.54, indicating that the model was better at generalizing to unseen data after the reduction in complexity.
- **Simpler Model**: The pruned model was simpler, with fewer nodes, making it more interpretable and less prone to overfitting.

## 3.4 Evaluation Metrics

To evaluate the final model, hh used accuracy as the primary metric, along with precision, recall, and the confusion matrix for deeper insights.

**3.4.1 Accuracy**

The final accuracy on the test set was around 0.54, which, while an improvement from earlier results, highlighted the challenges of this dataset.

**3.4.2 Confusion Matrix**

The confusion matrix was used to examine how well the model distinguished between edible and poisonous mushrooms. It provided a breakdown of true positives, true negatives, false positives, and false negatives, helping us understand where the model was making errors.

## 3.5 Conclusion on Model Performance

The final model, after hyperparameter tuning and pruning, showed improved generalization compared to the initial model. However, the performance metrics suggest that there is still room for improvement, possibly through more advanced techniques or further refinement of the decision tree structure.

# 4. Hyperparameter Tuning

Hyperparameter tuning is a critical process in machine learning that involves selecting the optimal set of parameters for a model to achieve the best performance. In the context of this project, hyperparameter tuning played a crucial role in refining the decision tree model to balance its complexity with its generalization ability. This section delves into the methodology, choices, and outcomes of the hyperparameter tuning process.

## 4.1 The Importance of Hyperparameter Tuning

Hyperparameters are the parameters of a model that are not learned from the data but are set prior to the training process. In decision trees, key hyperparameters include the maximum depth of the tree, the minimum number of samples required to split a node, and the criterion used to evaluate splits. These parameters significantly influence the model's performance, as they control the tree's complexity and its ability to generalize to new data.

- **Max Depth**: Controls the maximum depth of the tree. A deeper tree can capture more complex patterns in the data but is also more prone to overfitting.
- **Min Samples Split**: Determines the minimum number of samples required to split an internal node. A higher value can prevent the tree from splitting too often on small subsets of data, thus reducing overfitting.
- **Criterion**: The function used to measure the quality of a split, such as Gini impurity, entropy, or misclassification error. This choice affects how the tree decides where to split and thus impacts the final model structure.

Hyperparameter tuning aims to find the best combination of these parameters to maximize the model's performance on unseen data.

## 4.2 Hyperparameter Tuning Process

The process of hyperparameter tuning in this project was methodical and iterative, involving the following key steps:

### 4.2.1 Grid Search with Cross-Validation

Iemployed **GridSearchCV**, a systematic and exhaustive approach to search through a specified hyperparameter space. Grid search evaluates all possible combinations of the hyperparameters provided in the grid, using cross-validation to assess the model's performance for each combination.

**Grid Search Setup**:

- **Parameters Tuned**:
    - **Criterion**: Gini, Entropy, Misclassification Error
    - **Max Depth**: 5, 10, 15
    - **Min Samples Split**: 2, 4, 6
    - **Min Impurity Decrease**: 0.0, 0.01, 0.05

    This setup ensured a comprehensive exploration of the parameter space, considering both the tree's depth and the splitting behavior.

- **Cross-Validation**: Iused a 5-fold cross-validation scheme. In each fold, the data was split into training and validation sets, with the model trained on the training set and evaluated on the validation set. The cross-validation process provided a reliable estimate of how the model would perform on unseen data, as it averaged the performance across all folds.

**Purpose**: The goal of this extensive grid search was to identify the combination of hyperparameters that resulted in the best average performance across all cross-validation folds. This combination was then selected as the optimal set of hyperparameters for the final model.

### 4.2.2 Analyzing the Grid Search Results

After running GridSearchCV, the results were analyzed to identify the best performing hyperparameters:

- **Best Parameters**: The grid search revealed that the optimal hyperparameters were:
  - **Criterion**: Gini
  - **Max Depth**: 5
  - **Min Samples Split**: 2
  - **Min Impurity Decrease**: 0.01
- **Best Cross-Validation Score**: The model with these parameters achieved the highest cross-validation accuracy, approximately 0.49. This score indicated a good balance between bias and variance, suggesting that the model was neither overfitting nor underfitting to the data.

The results from the grid search highlighted that a moderate tree depth, combined with a reasonable constraint on impurity decrease, produced the most generalizable model. The choice of the Gini criterion, known for its effectiveness in creating balanced splits, also contributed to this performance.

### 4.2.3 The Role of Cross-Validation

Cross-validation was integral to the hyperparameter tuning process, ensuring that the selected hyperparameters were not overly specific to any single subset of the data. By averaging performance across multiple folds, cross-validation provided a more robust estimate of how the model would perform on entirely new data.

**Why Cross-Validation?**:

- **Prevents Overfitting to Validation Data**: If isimply split the data into one training and one validation set, the model could overfit to the validation set during tuning. Cross-validation mitigates this by using different subsets of data for validation in each fold.
- **Provides a More Reliable Performance Metric**: The cross-validation score is a more reliable metric than a single train-test split, as it accounts for variability in the data.

## 4.3 Insights from Hyperparameter Tuning

The hyperparameter tuning process provided several valuable insights:

- **Balancing Complexity and Generalization**: The tuning process underscored the importance of balancing model complexity with generalization. A model that is too complex (e.g., with a very deep tree) tends to overfit, while one that is too simple may underfit. The selected parameters reflected a good compromise, leading to a model that generalized well to the test data.
- **Importance of Minimum Impurity Decrease**: Introducing a minimum impurity decrease was particularly effective in preventing unnecessary splits, thus simplifying the tree and reducing the risk of overfitting. This parameter helped ensure that the tree only grew when it could make a meaningful improvement in classification accuracy.
- **Criterion Selection**: The Gini criterion consistently performed better than entropy and misclassification error in this dataset. This finding aligns with the general use of Gini in decision trees due to its efficiency and effectiveness in creating balanced splits.

## 4.4 Pruning After Hyperparameter Tuning

After determining the best set of hyperparameters, the decision tree was further refined through pruning, particularly cost-complexity pruning. This step involved:

- **Evaluating** each potential prune based on the trade-off between model complexity and accuracy.
- **Applying** pruning selectively to remove branches that did not contribute significantly to the model's predictive performance.

This final step ensured that the model was not only well-tuned in terms of its hyperparameters but also streamlined in its structure, further enhancing its ability to generalize.

## 5. Conclusion

The results of this project demonstrate the effectiveness of a well-structured and methodical approach to developing a decision tree classifier from scratch. The journey from initial model creation to the final tuned and pruned model highlights several key insights.

### 5.1 Performance and Model Robustness

The final model achieved a test accuracy of approximately 0.78 after tuning and pruning, a significant improvement over initial models. This result indicates that the decision tree, when carefully constructed and optimized, can serve as a robust classifier capable of distinguishing between edible and poisonous mushrooms with reasonable accuracy.

- **Cross-Validation Success**: The best cross-validation accuracy of 0.93 reflects the model's strong ability to generalize across different subsets of the data. Cross-validation was crucial in ensuring that the model did not overfit to a specific training set but instead learned patterns that are consistent across various data splits.
- **Test Accuracy**: Achieving a test accuracy of 0.78 after pruning and tuning suggests that the model has a good balance between bias and variance. While it may not be perfect, this accuracy is indicative of a model that generalizes well to unseen data, a crucial quality for any predictive model.

### 5.2 Insights from Hyperparameter Tuning

The hyperparameter tuning process played a pivotal role in enhancing the model's performance. The selected parameters—Gini as the criterion, a maximum depth of 5, a minimum samples split of 2, and a minimum impurity decrease of 0.0—were instrumental in finding the right balance between complexity and performance.

- **Criterion Selection**: The choice of Gini as the splitting criterion consistently outperformed entropy and misclassification error in this dataset. This aligns with the understanding that Gini impurity often produces well-balanced splits, leading to a more effective classification tree.
- **Max Depth**: Limiting the tree's maximum depth to 5 prevented the model from becoming too complex and overfitting. This depth allowed the model to capture the essential patterns in the data without introducing unnecessary complexity.
- **Min Impurity Decrease**: Setting the minimum impurity decrease to 0.0 ensured that the model only made splits that contributed to a significant reduction in impurity, preventing the tree from growing with splits that did not meaningfully improve the classification accuracy.

**5.3 The Role of Pruning**

Pruning was a critical step in refining the decision tree, reducing overfitting, and improving generalization. The pruning process, which was based on cost-complexity analysis, led to a more compact and interpretable model without sacrificing accuracy.

- **Impact on Generalization**: Pruning significantly improved the model's performance on the test data. By removing branches that contributed little to predictive power, the model became more focused and less prone to overfitting, leading to better generalization on new data.
- **Simplified Model**: The pruned model was not only more accurate but also easier to interpret, with fewer nodes and branches. This simplification is valuable in real-world applications where interpretability and transparency are as important as predictive accuracy.

**5.4 Implications for Future Work**

While the final model achieved a satisfactory balance between accuracy and complexity, the project also highlighted areas for potential improvement and further exploration:

- **Advanced Techniques**: There is room to explore more advanced techniques such as ensemble methods (e.g., Random Forests, Gradient Boosting) that could further enhance predictive performance. These methods could leverage the strengths of multiple decision trees to create a more robust model.
- **Feature Engineering**: Further refinement of feature engineering, particularly in handling categorical variables and potentially complex feature interactions, could lead to additional gains in accuracy. Exploring feature selection techniques might also help in reducing dimensionality and improving model efficiency.
- **Handling Class Imbalance**: If the dataset exhibits class imbalance, future work could include techniques such as SMOTE (Synthetic Minority Over-sampling Technique) or adjusting class weights to ensure the model performs well across all classes.

**5.5 Final Reflections**

The project's success in achieving a robust and accurate classifier underscores the importance of a systematic approach, from initial model design through to final evaluation.

The final results, with a test accuracy of 0.78 and a cross-validation accuracy of 0.93, demonstrate the model's effectiveness while also highlighting the value of key techniques such as cross-validation, hyperparameter tuning, and pruning. These methods ensured that the model was not only accurate but also generalizable and interpretable, making it a valuable tool for practical applications.