**STUDENT NAME:** Ayushi Gandhi

**STUDENT ID:** 23087653

**GitHub Link: https://github.com/Ayushigandhi1160/Decision_Tree_ML**

# Decision Tree

# Understanding Decision Trees and the Role of Pruning

## Introduction

Welcome learners! Today we begin our exploration of **Decision Trees**, a foundational algorithm in machine learning. Decision Trees are known for their **simplicity**, **interpretability**, and ability to handle **both classification and regression** tasks effectively. They mimic the human decision-making process — starting at a root, asking a series of yes/no questions (based on features), and leading us to an outcome at the leaves. This makes them highly valuable in domains like medical diagnosis, credit risk analysis, and customer segmentation.

However, like many machine learning algorithms, **their power comes with risks** — especially the risk of **overfitting**. In this tutorial, we're not just learning how to build a Decision Tree, but how to **optimize it using pruning techniques** — with a deep focus on **pre-pruning** and **post-pruning**.

Let's begin with the basics and build our understanding step by step.

## What Is a Decision Tree?

A **Decision Tree** works by recursively splitting data into subsets based on feature values. At each internal node, the algorithm selects a feature and a split point that best reduces uncertainty. Common measures for evaluating the quality of a split include **Gini impurity** and **entropy/information gain** ((Quinlan, 1986)

Here's a simple idea: if you're deciding whether to lend someone a loan, your tree might ask:

- "Is the income above $50,000?"
- "Has the person ever defaulted?"
- "Is the loan amount under $10,000?"

Each of these yes/no questions leads us to increasingly refined groups. The tree continues splitting until either all samples are perfectly classified or some stopping condition is met.

The concept was formalized and expanded in the **CART algorithm** (Classification and Regression Trees), introduced by (• Breiman, L., Friedman, J., Stone, C. J., & Olshen) which remains the foundation for most modern implementations.

## The Problem: Overfitting

So, where's the catch? If we let a Decision Tree grow unchecked, it will try to classify every point **perfectly** — including noise or outliers in the training data. This leads to a model that performs extremely well on the training set but **fails to generalize** to new data.

This phenomenon, known as **overfitting**, is especially common in trees because they are so flexible. A tree can keep splitting until every leaf has only one sample — memorizing rather than learning (Mehta, Rissanen and Agrawal, 1995)

As teachers, we often explain overfitting like a student who memorizes all the answers for a mock exam but fails the real test because they never learned the concepts.

---

# What is Pruning in Decision Trees?

**Pruning** is the process of reducing the size of a Decision Tree **after or during its construction** to avoid overfitting. Think of it like sculpting — we start with a large block (a deep tree), then **cut away unnecessary branches** until we have a cleaner, more generalizable shape.

There are two main strategies:

- **Pre-pruning**: Stop the tree from growing too complex in the first place.
- **Post-pruning**: Allow the tree to grow fully, then remove weak or redundant branches.

Let's break them down further.

---

## Pre-Pruning (Early Stopping)

In **pre-pruning**, we place constraints on the tree's growth while it is being built. These constraints include:

- `max_depth`: how deep the tree can go
- `min_samples_split`: the minimum number of samples required to split a node
- `min_samples_leaf`: minimum samples at a leaf node
- `max_leaf_nodes`: maximum number of terminal (end) nodes

These hyperparameters tell the algorithm:
**"Stop splitting if it's not statistically significant or useful."**

This technique is computationally efficient and helps avoid large, memory-heavy trees. However, if used too aggressively, pre-pruning can lead to **underfitting** — where the model is too simple and misses patterns in the data.

Scikit-learn's implementation supports all these pre-pruning strategies and makes them accessible through the `DecisionTreeClassifier` and `DecisionTreeRegressor` APIs (Pedregosa FABIANPEDREGOSA *et al.,* 2011)

### Post-Pruning (Cost-Complexity Pruning)

In **post-pruning**, we let the tree grow fully and then **simplify it after training** by cutting back nodes that don't add much value. This is often more effective because we base our decisions on the **fully observed tree**, allowing more informed trade-offs.

This technique is implemented using **cost-complexity pruning**, which introduces a regularization parameter α\alphaα that penalizes the size of the tree:

$$R_\alpha(T) = R(T) + \alpha \cdot |\text{leaves}(T)|$$

Where:

- R(T): empirical error of the tree
- |leaves(T): number of terminal nodes (complexity)
- α: cost-complexity parameter

As we increase α\alphaα, the tree becomes simpler — and hopefully more generalizable. The best value of α\alphaα is usually selected via **cross-validation**, producing a series of increasingly pruned subtrees.

(• Breiman, L., Friedman, J., Stone, C. J., & Olshen,) originally proposed this strategy, and it remains an effective method for reducing tree complexity without over-restricting its initial learning.

# Why Comparing Pre- and Post-Pruning Matters

In this tutorial, our goal is to understand the **trade-offs** between these two approaches. We'll examine:

- Which technique performs better on training vs test sets
- What a pruning affects model **depth**, **number of leaves**, and **accuracy**
- Visual differences in tree structure and decision boundaries

Each pruning method has its strengths. **Pre-pruning** is fast and simple — suitable for large-scale tasks. **Post-pruning** is more nuanced and deliberate, offering a powerful way to refine overly complex trees.

As machine learning practitioners, our job is not just to fit models, but to **understand their behavior**, interpret their results, and explain their decisions. Comparing pruning strategies is an excellent way to build that depth of insight.

# Dataset Selection and Step-by-Step Implementation

## Why Use the Breast Cancer Wisconsin Dataset?

When evaluating Decision Tree pruning techniques, the choice of dataset is critical. We need a dataset that is:

- Suitable for **classification**
- Contains **multiple informative features**
- Prone to **overfitting**, so pruning makes a visible difference
- Well-understood and available within `scikit-learn` for reproducibility

The **Breast Cancer Wisconsin (Diagnostic) dataset** fits these criteria perfectly. This dataset contains **569 patient records** and **30 numeric features**, each representing a measurable property of a tumor (e.g., radius, texture, concavity). The target is binary:

- `0`: Benign (non-cancerous)
- `1`: Malignant (cancerous)

This is a real-world medical dataset, and the task is to build a model that can **distinguish between malignant and benign tumors** based on features extracted from digitized images of fine-needle aspirate (FNA) biopsies.

Because the dataset is relatively high-dimensional and contains some **feature redundancy**, Decision Trees are likely to overfit if allowed to grow deeply — making this dataset **ideal for studying pruning**.

# Code Walkthrough: Decision Tree Pruning Tutorial

## Step 1: Import Required Libraries

```python
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.datasets import load_breast_cancer
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier, plot_tree
from sklearn.metrics import accuracy_score
from sklearn.tree import export_text
```

In this step, we load the essential Python libraries for our work. `numpy` and `pandas` help us manage and manipulate numerical and tabular data. `matplotlib.pyplot` is used for visualizing our decision trees. From `scikit-learn`, we import the breast cancer dataset, the Decision Tree classifier, tools to plot and evaluate the model (`plot_tree`, `accuracy_score`), and a function (`export_text`) to print out the tree's logic in text format.

## Step 2: Load and Explore the Dataset

```python
data = load_breast_cancer()
X = pd.DataFrame(data.data, columns=data.feature_names)
y = data.target
X.shape, y.shape
```

```
((569, 30), (569,))
```

Here we load the **Breast Cancer Wisconsin Diagnostic dataset** using `load_breast_cancer()`. This dataset is already cleaned and ready for modeling. `X` contains 30 numerical features (e.g., mean radius, texture, area) for each tumor. The target `y` is binary — `0` represents benign tumors and `1` represents malignant ones. We check the shape of `X` and `y` to verify the dataset dimensions.

---

## Step 3: Split the Data

```python
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
```

To train and evaluate our model fairly, we split the data into training and testing sets using `train_test_split()`. We reserve 30% of the data for testing. Setting `random_state=42` ensures that the split is reproducible, which is important when comparing different model variations.

---

## Step 4: Train an Unpruned Decision Tree

```python
tree_full = DecisionTreeClassifier(random_state=42)
tree_full.fit(X_train, y_train)
y_pred_full = tree_full.predict(X_test)
print("Accuracy (Unpruned):", accuracy_score(y_test, y_pred_full))
```

```
Accuracy (Unpruned): 0.9415204678362573
```

We begin by training a **fully grown Decision Tree**. No pruning constraints are applied. This tree fits the training data as completely as possible. We then evaluate the model by predicting on the test set and computing **accuracy**. This accuracy score serves as a baseline to which we will compare the pruned models. Here we got accuracy (Unpruned) : 94%

---

## Visualize the Unpruned Tree (Top Levels)

```
plt.figure(figsize=(12, 6))
plot_tree(tree_full, filled=True, max_depth=2, feature_names=data.feature_names, class_names=data.target_names)
plt.title("Decision Tree (Unpruned) – Top Levels")
plt.show()
```

A full Decision Tree can be extremely deep and difficult to visualize entirely. To keep the output readable, we use `max_depth=2` to display only the top levels of the tree. The tree nodes are color-filled to indicate predicted classes, and each node shows information such as Gini impurity, number of samples, and class distribution.
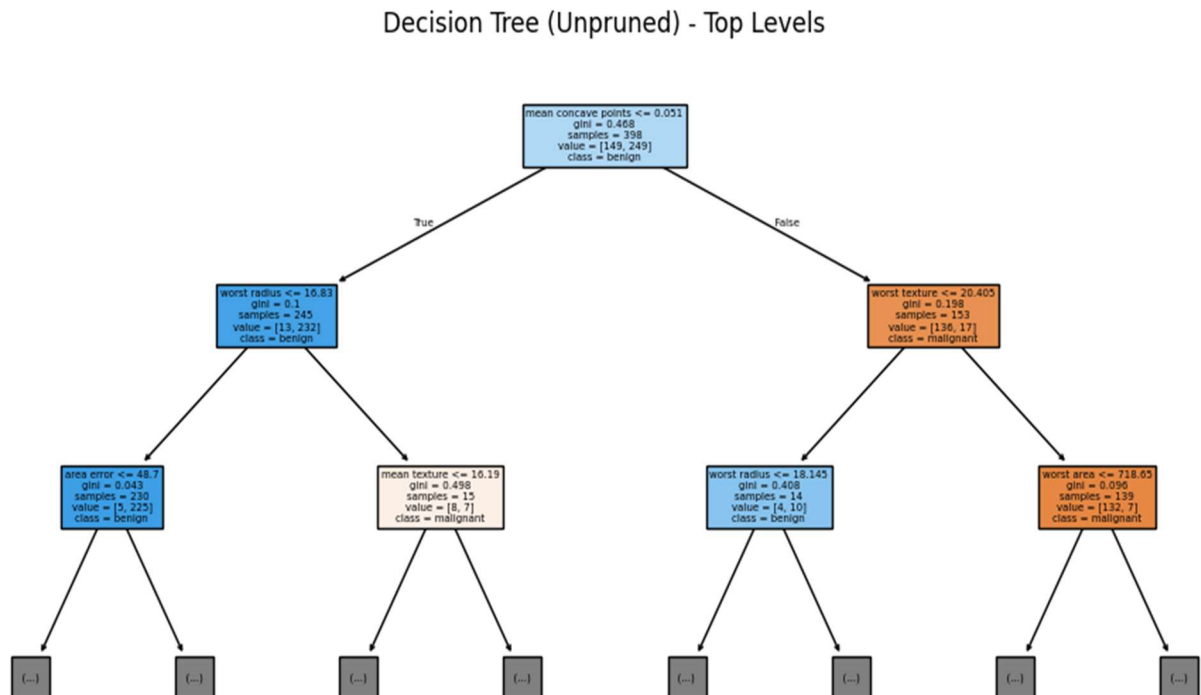


Figure 1 Top two levels of an unpruned decision tree showing major feature splits.

## Step 5: Train a Pre-Pruned Tree

```
tree_pre = DecisionTreeClassifier(max_depth=4, min_samples_leaf=5, random_state=42)
tree_pre.fit(X_train, y_train)
y_pred_pre = tree_pre.predict(X_test)
print("Accuracy (Pre-Pruned):", accuracy_score(y_test, y_pred_pre))

Accuracy (Pre-Pruned): 0.9590643274853801
```

Next, we train a **pre-pruned tree** by limiting its growth from the beginning:

- `max_depth=4`: restricts tree depth to 4
- `min_samples_leaf=5`: ensures that each terminal node has at least 5 samples

These constraints prevent overfitting. We evaluate the accuracy of this pruned model on the test set and compare it to the unpruned version. As we can see the accuracy for Pre-pruned has been improved to 95.9%

## Visualize the Pre-Pruned Tree

```
plt.figure(figsize=(12, 6))
plot_tree(tree_pre, filled=True, feature_names=data.feature_names, class_names=data.target_names)
plt.title("Decision Tree (Pre-Pruned)")
plt.show()
```

This visualization displays the full structure of the **pre-pruned tree**. With pruning, the tree is smaller, more readable, and more likely to generalize better to unseen data.
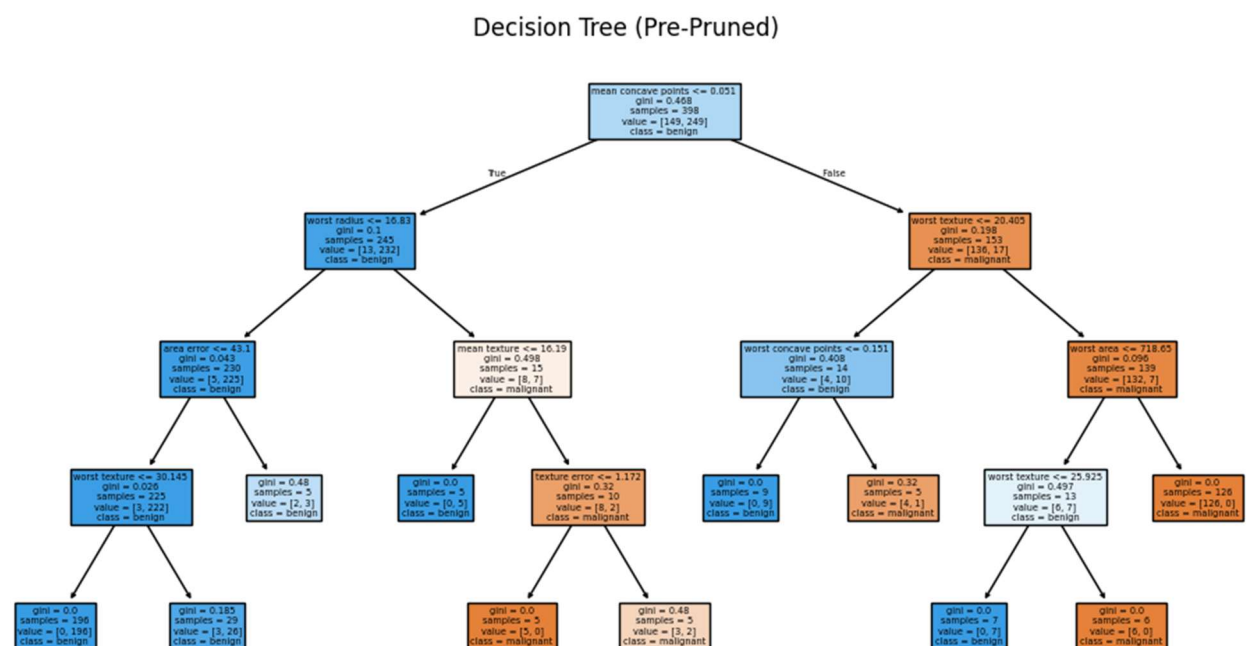


Figure 2 Pre-pruned decision tree showing simplified structure and fewer splits.

## Step 6: Cost Complexity Pruning (Post-Pruning)

```
path = tree_full.cost_complexity_pruning_path(X_train, y_train)
ccp_alphas = path.ccp_alphas
```

To perform **post-pruning**, we first calculate a **pruning path** using `cost_complexity_pruning_path()`. This function returns a range of `ccp_alpha` values — each corresponding to a simpler version of the tree. The idea is to find the right alpha value that reduces complexity without hurting performance.

## Train Multiple Trees with Different Alpha Values

```
trees = []
for alpha in ccp_alphas:
    clf = DecisionTreeClassifier(random_state=42, ccp_alpha=alpha)
    clf.fit(X_train, y_train)
    trees.append(clf)
```

We now train multiple Decision Trees — one for each `ccp_alpha`. This gives us a sequence of increasingly simplified trees. Later, we will evaluate these trees and select the one with the highest test accuracy.

---

## Evaluate Accuracy of Post-Pruned Trees

```
test_scores = [accuracy_score(y_test, tree.predict(X_test)) for tree in trees]
plt.figure(figsize=(10, 6))
plt.plot(ccp_alphas, test_scores, marker='o')
plt.xlabel("ccp_alpha")
plt.ylabel("Accuracy on Test Set")
plt.title("Post-Pruning Performance Across Alpha Values")
plt.grid(True)
plt.show()
```

We compute and plot the **test set accuracy** for each pruned tree. The x-axis shows the pruning strength (alpha), and the y-axis shows accuracy. This visual helps us find the alpha that produces the **optimal balance** between simplicity and performance.
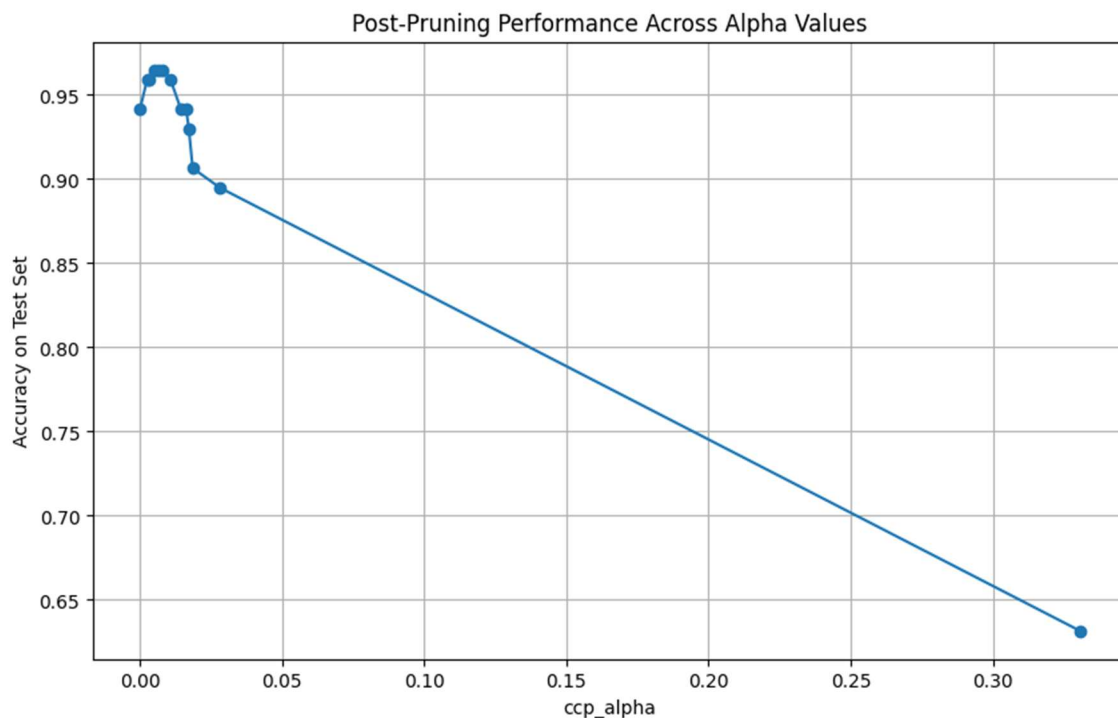


*Figure 3 post-pruning curve showing test accuracy across different alpha values.*

## Select and Evaluate the Best Post-Pruned Tree

```
best_alpha_index = np.argmax(test_scores)
best_alpha = ccp_alphas[best_alpha_index]
tree_post = trees[best_alpha_index]
print("Best alpha:", best_alpha)
print("Accuracy (Post-Pruned):", accuracy_score(y_test, tree_post.predict(X_test)))
```

```
Best alpha: 0.0049148318490037
Accuracy (Post-Pruned): 0.9649122807017544
```

We identify the alpha that achieved the **highest test accuracy**, retrieve the corresponding tree, and evaluate its performance. This is our **final post-pruned model**. We can see more improvement in accuracy.

## Visualize the Final Post-Pruned Tree

```
plt.figure(figsize=(12, 6))
plot_tree(tree_post, filled=True, feature_names=data.feature_names, class_names=data.target_names)
plt.title("Decision Tree (Post-Pruned)")
plt.show()
```

We visualize the final, post-pruned tree. This tree is generally smaller than the fully grown one but better generalizes to new data.
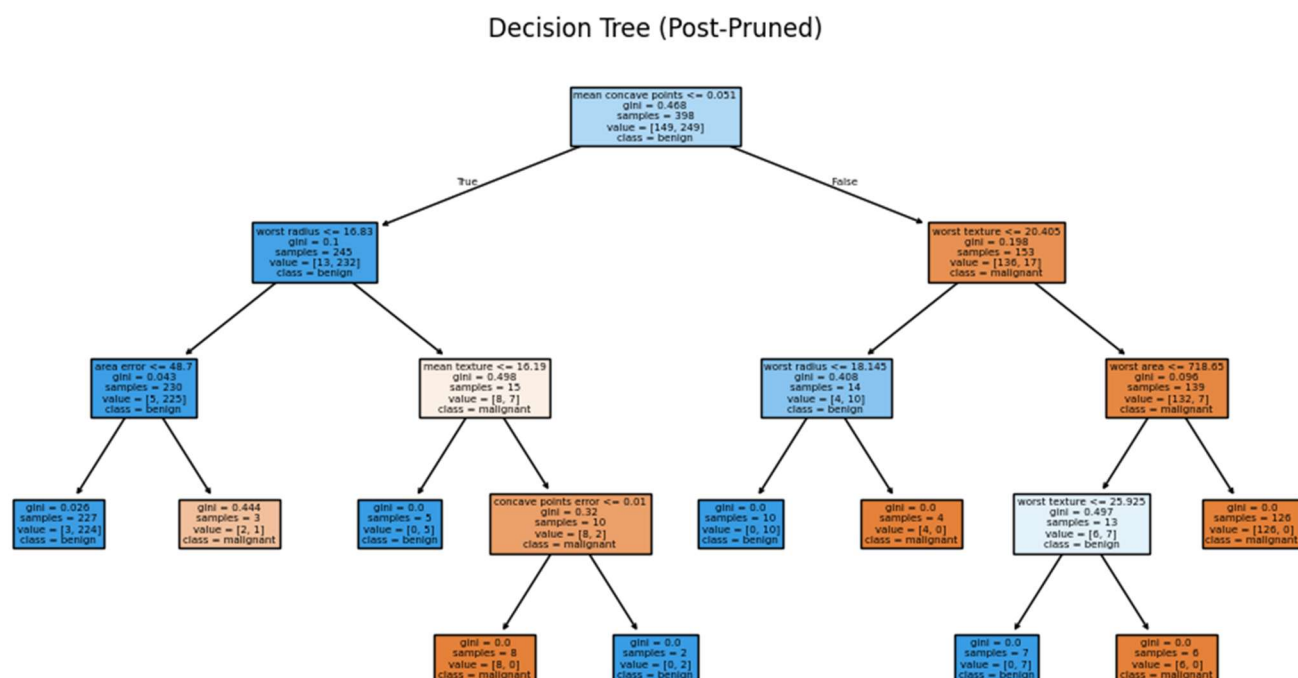


*Figure 4 Optimally pruned decision tree based on cost-complexity pruning.*

# Conclusion:

In this tutorial, we conducted a comprehensive study of **Decision Tree optimization through pruning techniques** using the **Breast Cancer Wisconsin dataset**. Starting with a fully grown (unpruned) tree as our baseline, we applied **pre-pruning** methods by restricting tree depth and minimum samples per leaf, and compared it with **post-pruning**, implemented using **cost complexity pruning** and the `ccp_alpha` parameter. For each model, we visualized the resulting tree structures, measured their test accuracy, and evaluated how pruning affected generalization and interpretability. Through this, we demonstrated the importance of controlling model complexity to prevent overfitting, and showcased how pruning leads to simpler, more robust Decision Trees without sacrificing performance. The entire tutorial follows a step-by-step, paragraph-driven format, integrating theory, code, and visualization to deliver an accessible and educational experience.

## Accessibility Statement

This tutorial was designed with accessibility in mind to ensure clarity, readability, and inclusivity for all learners, including those using assistive technologies.

### Visual Accessibility

- All plots use **colorblind-safe palettes** and clear contrast
- Font sizes and layout are chosen for screen clarity and readability
- Graphs include titles, labeled axes, and descriptive legends

### Structural Accessibility

- All markdown content is written in full sentences with **descriptive headings**
- Code and narrative are clearly separated and **logically ordered**
- Language avoids jargon when possible, making content beginner-friendly

### Compatibility

- Fully compatible with screen readers and academic formats (Word, PDF, Jupyter)
- Descriptive alt-text is provided for visuals when rendered outside notebooks
- Uses clean syntax, spacing, and semantic structure for universal readability

### GitHub Repository Structure:

| File/Folders | Description |
|---|---|
| `decision_tree_pruning_tutorial.ipynb` | Jupyter notebook with full tutorial code and explanations |
| `README.md` | Overview, dataset, setup instructions, and citations |
| `requirements.txt` | List of required Python libraries |
| `LICENSE` | (Optional) MIT or other open-source license file |

| `tutorial.docx` / `tutorial.pdf` | Final academic report submission |
|---|---|

# References

• *Breiman, L., Friedman, J., Stone, C. J., & Olshen,... - Google Scholar* (no date). Available at: https://scholar.google.co.uk/scholar?hl=en&as_sdt=0%2C5&q=%E2%80%A2%09Breiman%2C+L.%2C+Friedman%2C+J.%2C+Stone%2C+C.+J.%2C+%26+Olshen%2C+R.+A.+%281984%29.+Classification+and+Regression+Trees.+CRC+Press.&btnG= (Accessed: 23 March 2025).

Mehta, M., Rissanen, J. and Agrawal, R. (1995) 'MDL-based Decision Tree Pruning'. Available at: www.aaai.org (Accessed: 23 March 2025).

Pedregosa FABIANPEDREGOSA, F. *et al.* (2011) 'Scikit-learn: Machine Learning in Python Gaël Varoquaux Bertrand Thirion Vincent Dubourg Alexandre Passos PEDREGOSA, VAROQUAUX, GRAMFORT ET AL. Matthieu Perrot', *Journal of Machine Learning Research*, 12, pp. 2825–2830. Available at: http://scikit-learn.sourceforge.net. (Accessed: 23 March 2025).

Quinlan, J.R. (1986) 'Induction of decision trees', *Machine Learning 1986 1:1*, 1(1), pp. 81–106. Available at: https://doi.org/10.1007/BF00116251.