# UNIVERSITY OF SCIENCE, VNU-HCMC

# FACULTY INFORMATION TECHNOLOGY

## CSC14003 - Introduction to Artificial Intelligence



# Lab 2

# Decision Tree

**Instructors**: Nguyễn Trần Duy Minh, Nguyễn Ngọc Thảo,

Nguyễn Thanh Tình, Nguyễn Hải Đăng

**Class**:   22CLC06

**Student**: 22127154 – Nguyễn Gia Huy

Ho Chi Minh City – 2024

# Table of contents

# 1. Introduction

## 1.1 Description of the Dataset

The UCI Breast Cancer Wisconsin (Diagnostic) dataset is utilized in this project to classify tumors as either malignant (M) or benign (B). The dataset comprises **569 samples** and **30 numerical features** derived from imaging data. Each sample is labeled as either malignant or benign, making it a binary classification problem.

## 1.2 Objective

The primary objective of this assignment is to build a Decision Tree classifier using the scikit-learn library to predict whether a tumor is malignant or benign based on the given features. The model's performance will be evaluated across different training and testing splits, and the effect of varying the tree depth on accuracy will also be explored.

# 2. Data Preparation

## 2.1 Dataset Preparation

The dataset was first fetched using the ucimlrepo library. The features and labels were extracted, and the dataset was then divided into multiple training and testing sets with the following proportions:

- **40/60**: 40% training, 60% testing

- **60/40**: 60% training, 40% testing

- **80/20**: 80% training, 20% testing

- **90/10**: 90% training, 10% testing

Each split was performed in a stratified manner to ensure that both training and testing sets maintained the original distribution of classes.

## 2.2 Visualization of Class Distribution

To ensure that the data splits were prepared correctly, the class distributions were visualized for the original dataset as well as each of the training and testing sets. This step verifies that the stratification was successful and that each split adequately represents the overall dataset.
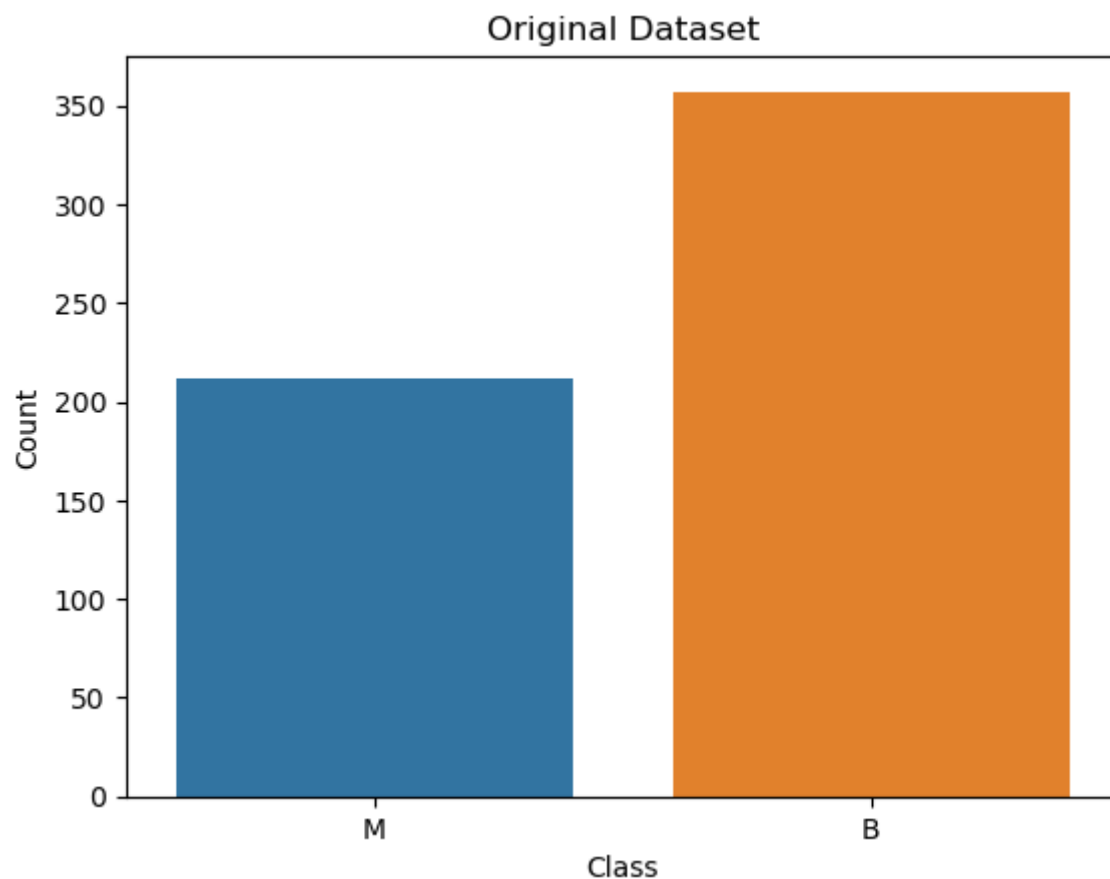
*Original dataset's class distribution*



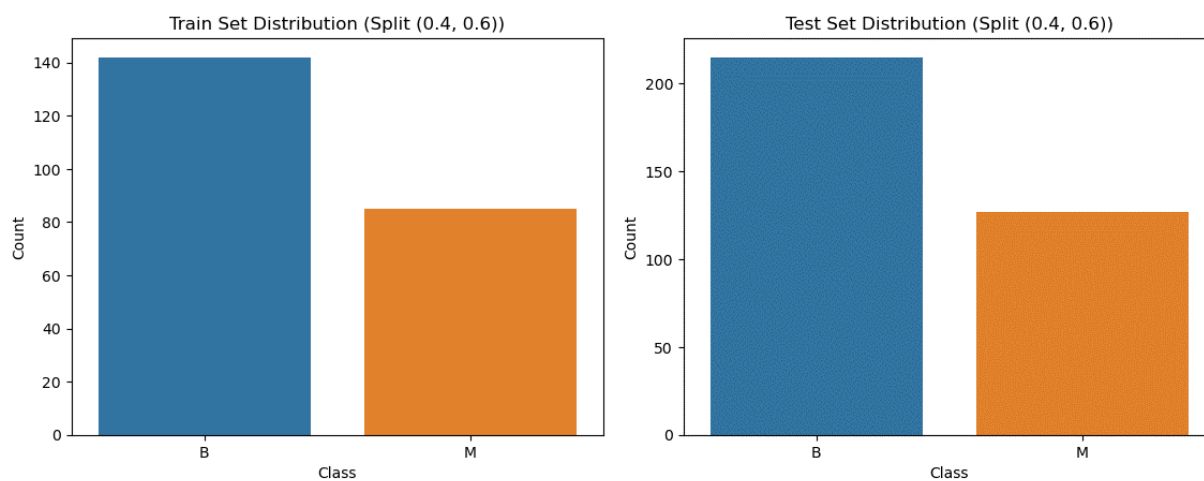Figure 2.1: Original dataset's class distribution

*Class distributions for each train/test split*



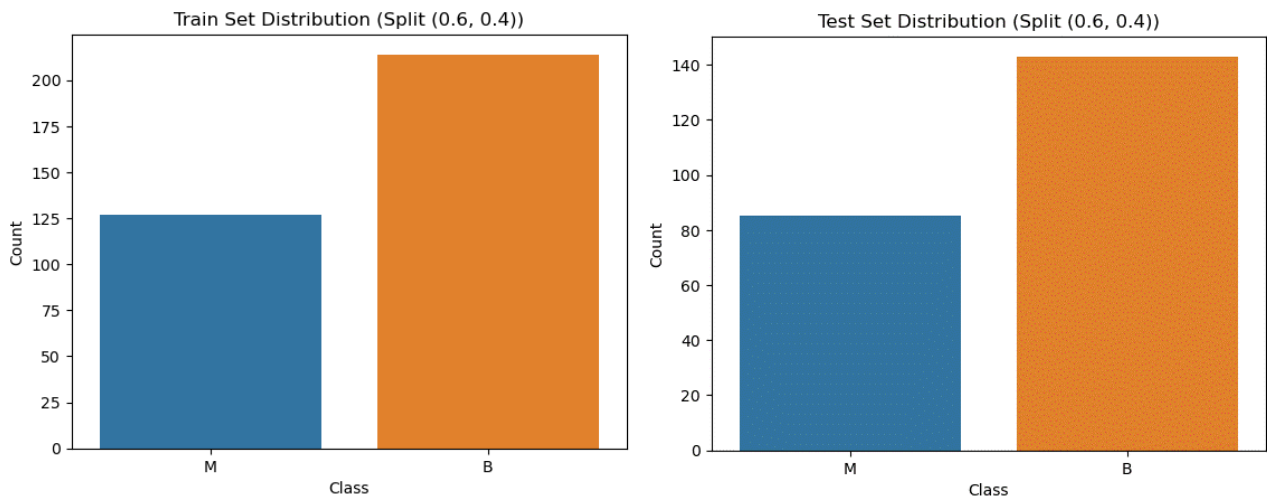Figure 2.2: Train / Test Set Distribution for Split 40/60

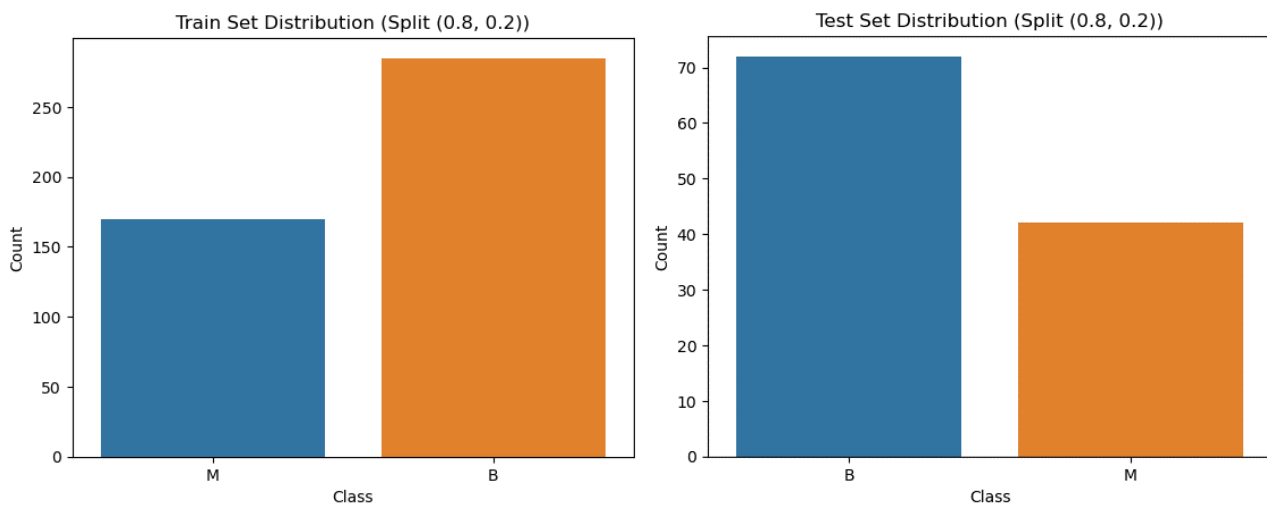Figure 2.3: Train / Test Set Distribution for Split 60/40



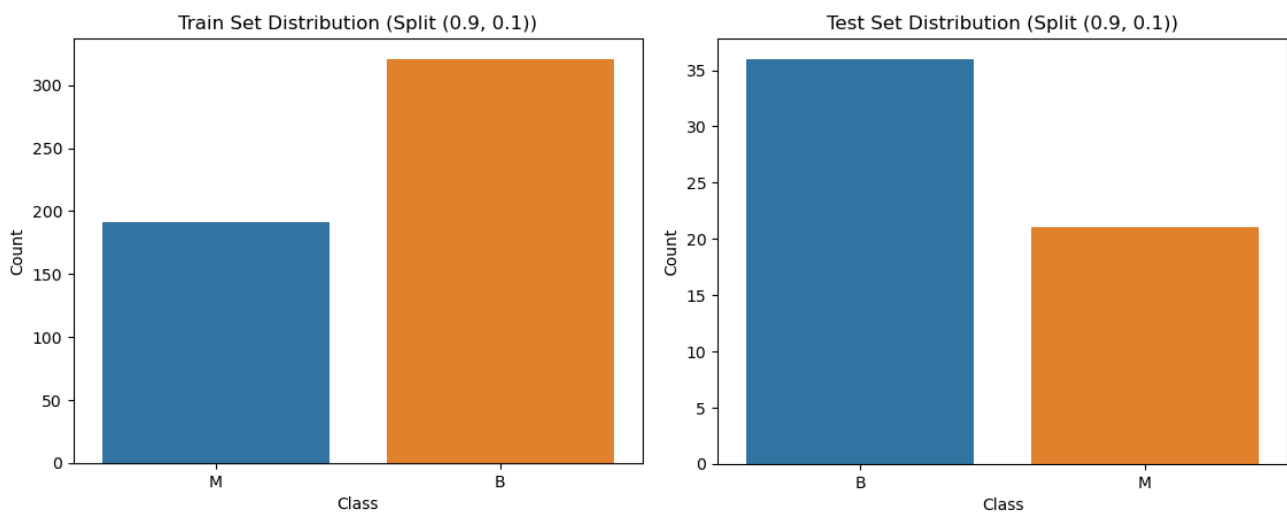Figure 2.4: Train / Test Set Distribution for Split 80/20



Figure 2.5: Train / Test Set Distribution for Split 90/10

Ho Chi Minh City – 2024

# 3. Building the Decision Tree Classifiers

## 3.1 Model Construction

For each of the train/test splits, a Decision Tree classifier was built using the DecisionTreeClassifier from the scikit-learn library. The models were trained on the respective training sets, and the resulting trees were visualized using the graphviz library.

## 3.2 Visualization of Decision Trees

The following visualizations represent the structure of the decision trees built for each of the train/test splits. These trees show the decision paths taken by the classifier and provide insight into the model's logic at different stages of the decision process.
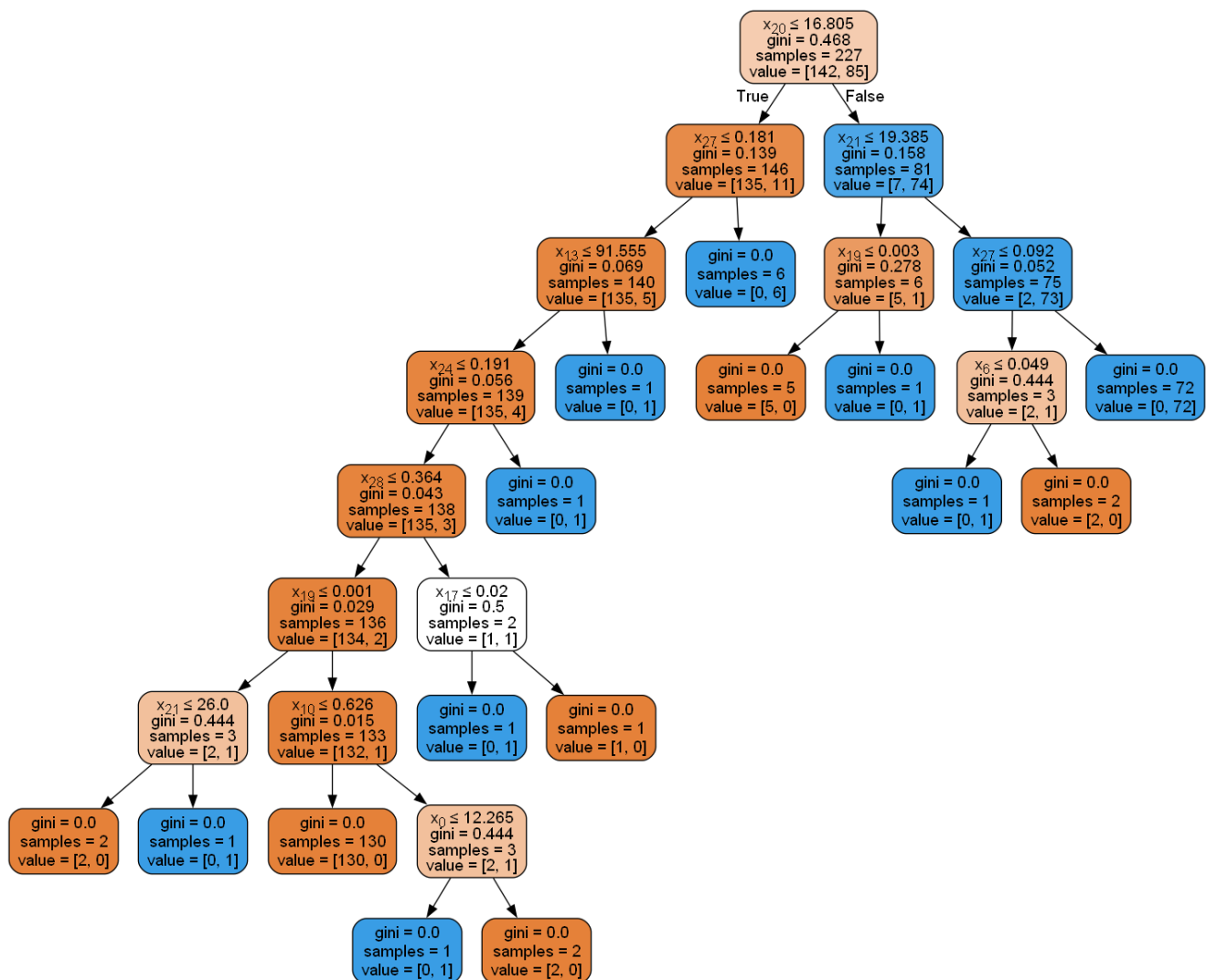


Figure 3.1: Decision tree for train/test split 40/60

Figure 3.2: Decision tree for train/test split 60/40



Figure 3.3: Decision tree for train/test split 80/20

Figure 3.4: Decision tree for train/test split 90/10

# 4. Evaluation of the Decision Tree Classifiers

## 4.1 Performance Evaluation

For each decision tree classifier, the model's performance was evaluated on the corresponding test set. The evaluation was conducted using the classification_report and confusion_matrix functions from scikit-learn.

### Classification report and confusion matrix



Figure 4.1: Classification Report and confusion matrix for train/test split 40/60



Figure 4.2: Classification Report and confusion matrix for train/test split 60/40

Classification Report

| | precision | recall | f1-score | support |
|---|---|---|---|---|
| B | 0.944 | 0.944 | 0.944 | 72.0 |
| M | 0.905 | 0.905 | 0.905 | 42.0 |
| macro avg | 0.925 | 0.925 | 0.925 | 114.0 |
| weighted avg | 0.93 | 0.93 | 0.93 | 114.0 |

Decision Tree Classifier Confusion Matrix

Overall Accuracy: 0.93

Figure 4.3: Classification Report and confusion matrix for train/test split 80/20

Classification Report

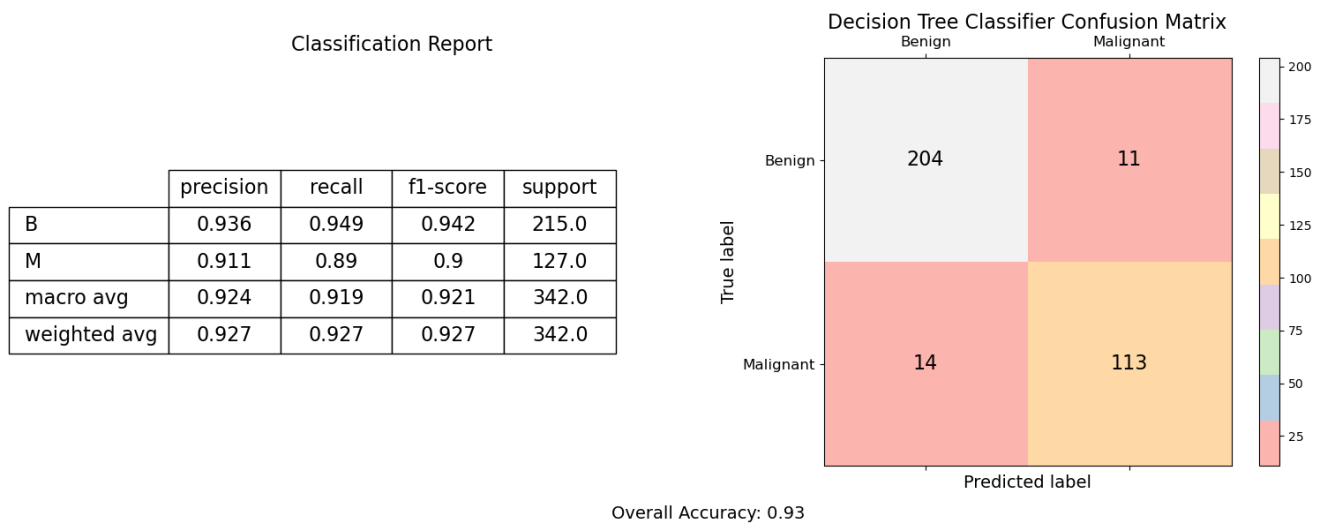| | precision | recall | f1-score | support |
|---|---|---|---|---|
| B | 0.921 | 0.972 | 0.946 | 36.0 |
| M | 0.947 | 0.857 | 0.9 | 21.0 |
| macro avg | 0.934 | 0.915 | 0.923 | 57.0 |
| weighted avg | 0.931 | 0.93 | 0.929 | 57.0 |

Decision Tree Classifier Confusion Matrix

Overall Accuracy: 0.93
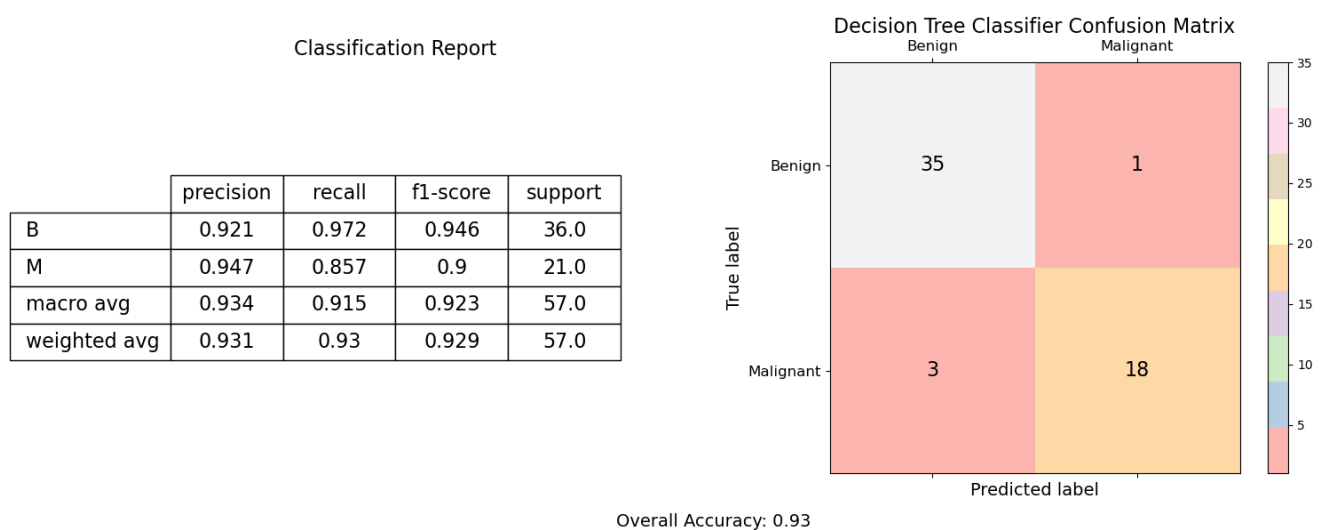
Figure 4.4: Classification Report and confusion matrix for train/test split 90/10

**Comments on the Results**

The decision tree classifiers were evaluated using four different train/test splits: 40/60, 60/40, 80/20, and 90/10. The performance of each model was assessed through classification reports and confusion matrices, as illustrated in Figures 4.1 to 4.4.

| Train /<br>Test Split | Accuracy | Precision and Recall | Confusion Matrix |
|---|---|---|---|
| 40/60 | The model achieved an overall accuracy of 0.93, which is a solid performance given that the model was trained on only 40% of the data. | The precision for the benign class (B) was 0.936, while for the malignant class (M), it was slightly lower at 0.911. This indicates that the model was slightly more likely to incorrectly classify a benign case as malignant.<br><br>The recall for benign cases was high at 0.949, showing that most benign cases were correctly identified. However, the recall for malignant cases was lower at 0.89, indicating that some malignant cases were misclassified as benign. | The model correctly classified 204 benign and 113 malignant cases. However, it misclassified 11 benign cases as malignant and 14 malignant cases as benign. These misclassifications could potentially have significant consequences in a clinical setting. |
| 60/40 | The accuracy increased to 0.952 with a larger training set, demonstrating improved performance with more data. | Both precision and recall were higher for both classes compared to the 40/60 split. The model achieved a precision of 0.946 for benign cases and 0.962 for malignant cases.<br><br>The recall was exceptionally high for benign cases at 0.979, and still strong for malignant cases at 0.906. | The confusion matrix shows that the model made fewer errors in this configuration. It correctly identified 140 benign and 77 malignant cases, while only misclassifying 3 benign cases and 8 malignant cases. This split provided the best balance between precision and recall across both classes. |
| 80/20 | The overall accuracy slightly decreased to 0.93, similar to the 40/60 split, suggesting that adding more training data does not always improve the model's accuracy. | The precision for benign cases was 0.944 and for malignant cases was 0.905. The recall values were identical at 0.944 for benign and 0.905 for malignant cases, indicating consistent performance in identifying both classes. | The model correctly classified 68 benign and 38 malignant cases, with 4 misclassifications in each class. Although the performance is still good, the smaller test set likely contributed to the slight decrease in accuracy and the balanced precision/recall. |
| 90/10 | The model maintained an accuracy of 0.93, consistent with other splits, but with much less test data. | The precision for benign cases decreased slightly to 0.921, while the precision for malignant cases increased to 0.947.<br><br>The recall for benign cases was very high at 0.972, indicating that the model identified almost all benign cases correctly. However, the recall for malignant cases was lower at 0.857, showing that some malignant cases were missed. | With a smaller test set, the model correctly identified 35 benign and 18 malignant cases, with only 1 benign case misclassified as malignant and 3 malignant cases misclassified as benign. The high recall for benign cases suggests the model is more cautious, potentially leading to more false positives (benign classified as malignant). |

**Summary**

Overall, the decision tree models performed consistently across different train/test splits, with accuracies hovering around 0.93 to 0.95. The 60/40 split provided the best performance, with the highest overall accuracy and balanced precision and recall for both classes. The model tends to have a higher recall for benign cases, suggesting it is more effective at correctly identifying benign tumors but at the risk of missing some malignant cases. The performance stability across different data splits indicates that the decision tree is a robust model for this dataset, although the slightly lower performance on the 90/10 split suggests that very small test sets might introduce some instability in the results.

The confusion matrices across all splits highlight the importance of considering both false positives and false negatives, especially in a sensitive domain like cancer diagnosis, where misclassifications can have significant implications.

# 5. The Depth and Accuracy of a Decision Tree

## 5.1 Experiment with Tree Depth

The effect of varying the tree depth on the model's accuracy was explored by adjusting the max_depth parameter of the decision tree classifier. The depths tested were: None, 2, 3, 4, 5, 6, and 7. For each depth, the tree was trained using the 80/20 train/test split, and the accuracy was recorded.

## 5.2 Results and Visualization

The results of this experiment are summarized in the table below, which shows the accuracy for each tested depth. Additionally, visualizations of the decision trees at different depths are included to illustrate how the tree structure changes as the depth increases.

*Table of accuracy vs. max depth*

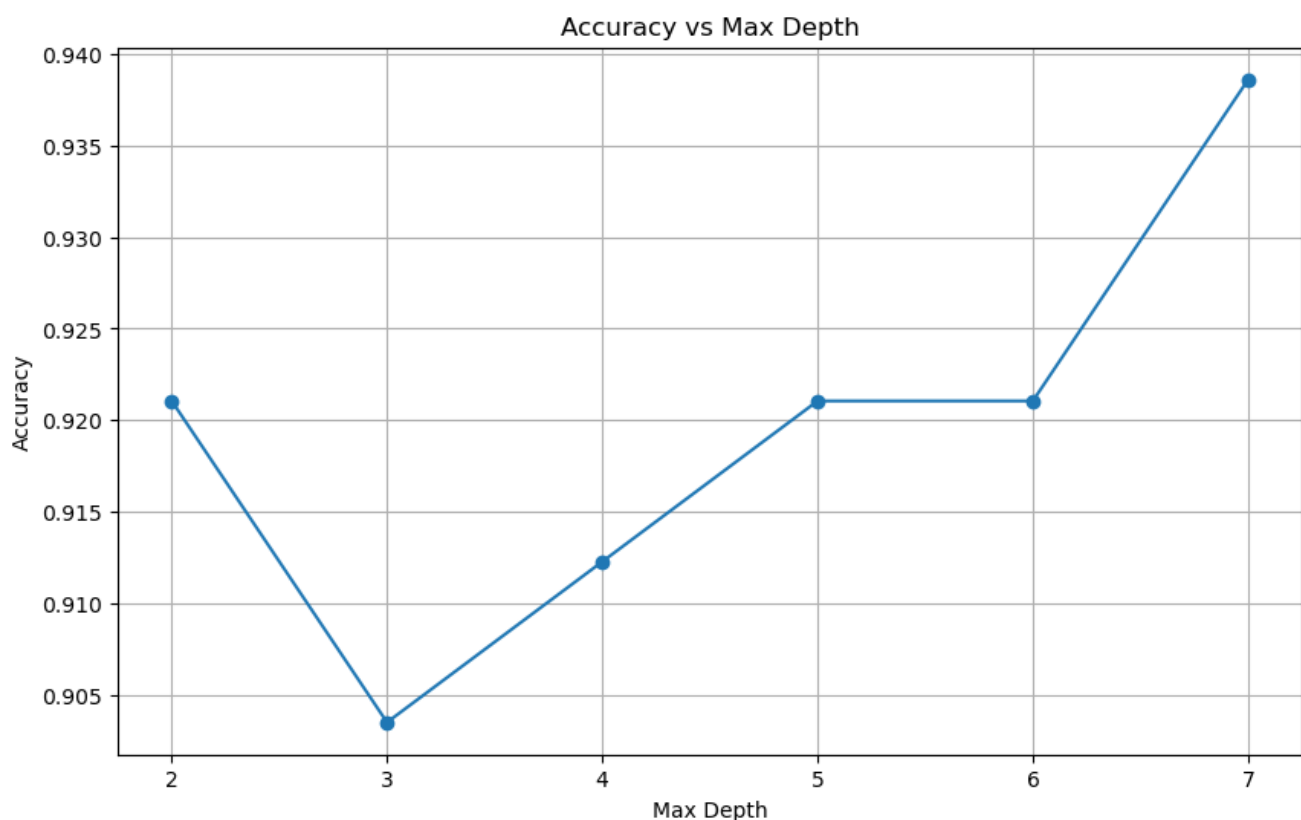| *Max Depth* | None | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| *Accuracy* | 0.93 | 0.921 | 0.904 | 0.912 | 0.921 | 0.921 | 0.939 |



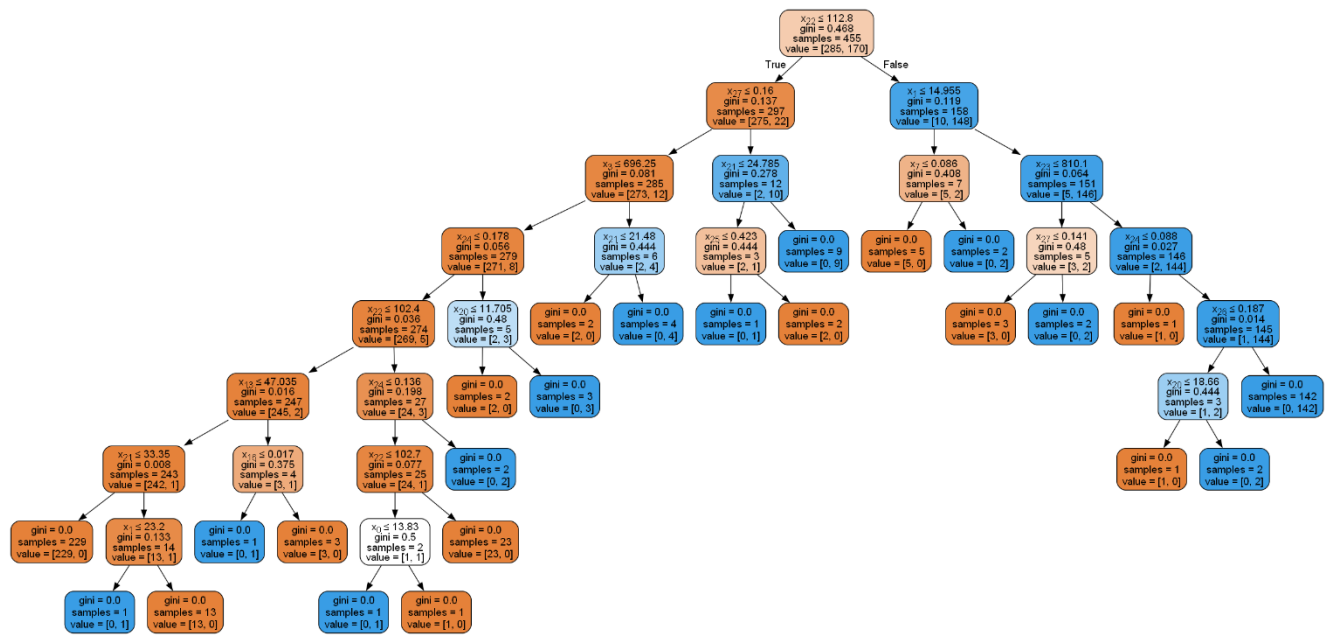Figure 5.1: Max depth vs Accuracy graph

*Visualizations of decision trees at different depths*



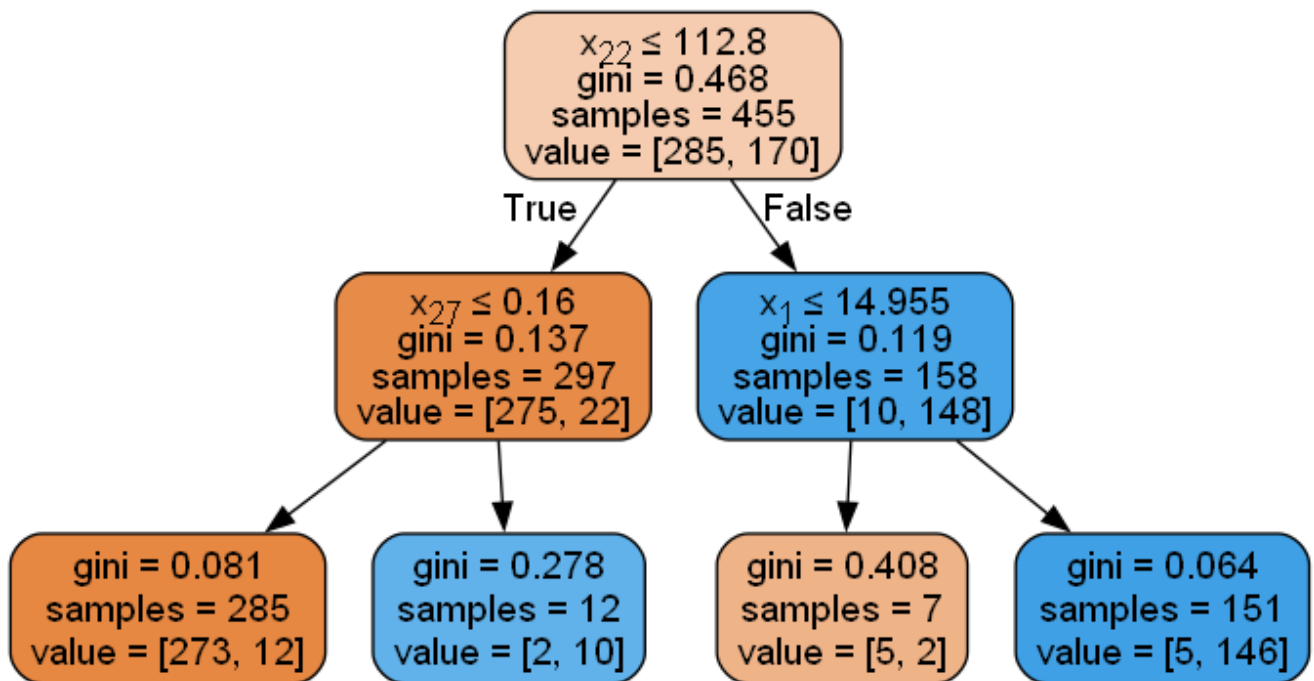Figure 5.2: Full decision tree



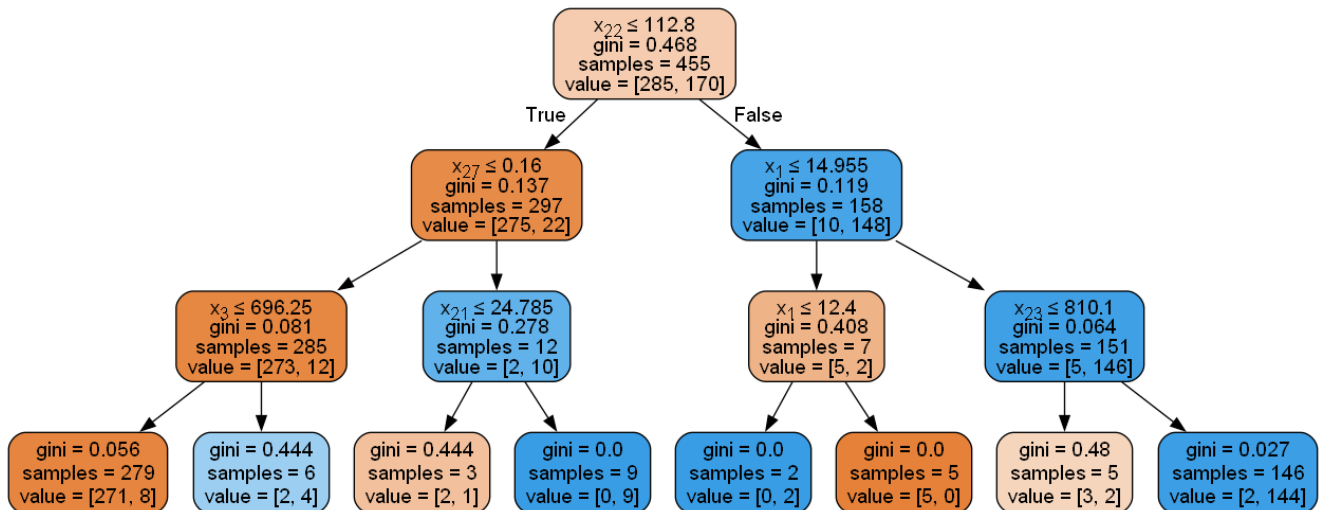Figure 5.3: Decision tree for depth 2

Figure 5.4: Decision tree for depth 3



Figure 5.5: Decision tree for depth 4



Figure 5.6: Decision tree for depth 5

Figure 5.7: Decision tree for depth 6



Figure 5.8: Decision tree for depth 7

**Comments on the Results**

The analysis of the Decision Tree classifiers with varying depths provides insights into how the complexity of the model influences its performance, particularly in terms of accuracy. The following comments reflect the observations derived from the table, graph, and visualizations of the decision trees at different depths.

## 1. Accuracy vs. Max Depth

- **General Trend:** The accuracy of the Decision Tree classifier fluctuated as the max depth increased, ranging from 0.904 at depth 3 to 0.939 at depth 7. Notably, the accuracy remained relatively stable across several depths (2, 5, and 6), where it was around 0.921, before increasing significantly at depth 7.

- **Shallow Trees (Depth 2 to 4):**

  o At a depth of 2, the tree achieved an accuracy of 0.921, demonstrating that even a shallow tree can capture the essential patterns in the data.

  o However, at depth 3, the accuracy dropped to 0.904, suggesting that the model might have been underfitting, failing to capture more complex relationships in the data.

  o At depth 4, the accuracy improved slightly to 0.912, indicating that increasing depth helped to capture more nuanced patterns without overfitting.

- **Deeper Trees (Depth 5 to 7):**

  o At depths 5 and 6, the accuracy plateaued at 0.921, similar to the performance at depth 2. This suggests that the additional splits did not contribute significantly to improving model performance.

  o Interestingly, at depth 7, the accuracy increased to 0.939, the highest observed. This indicates that the model was able to leverage the additional complexity to improve its classification performance, potentially by capturing more detailed interactions between features.

## 2. Visualizations of Decision Trees at Different Depths

- **Full Decision Tree (No Depth Limit):**
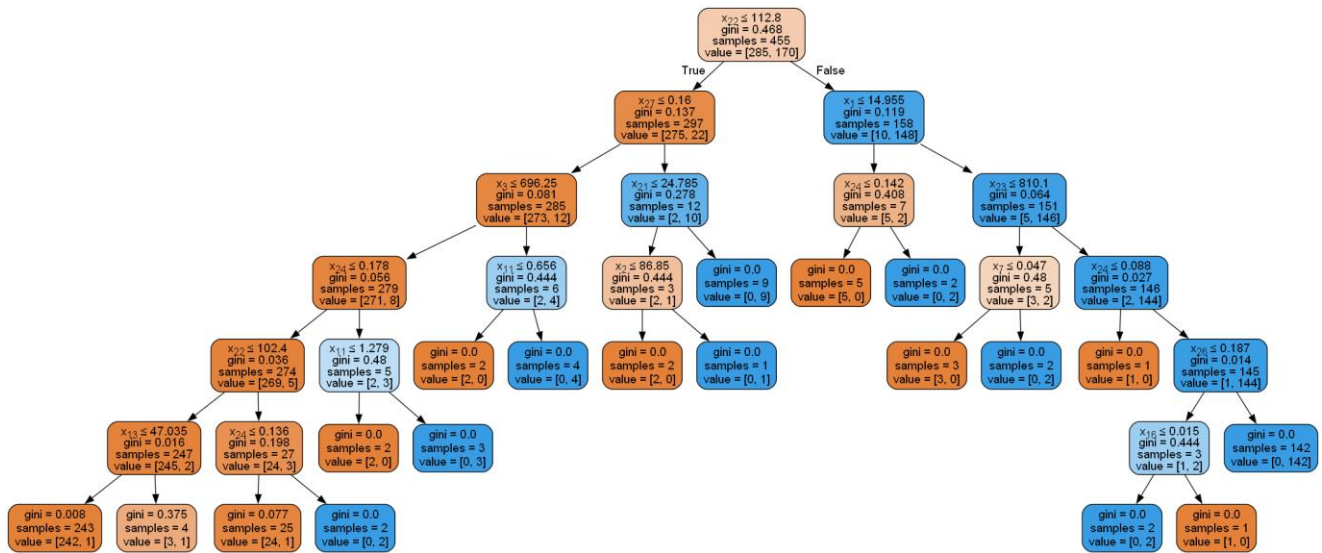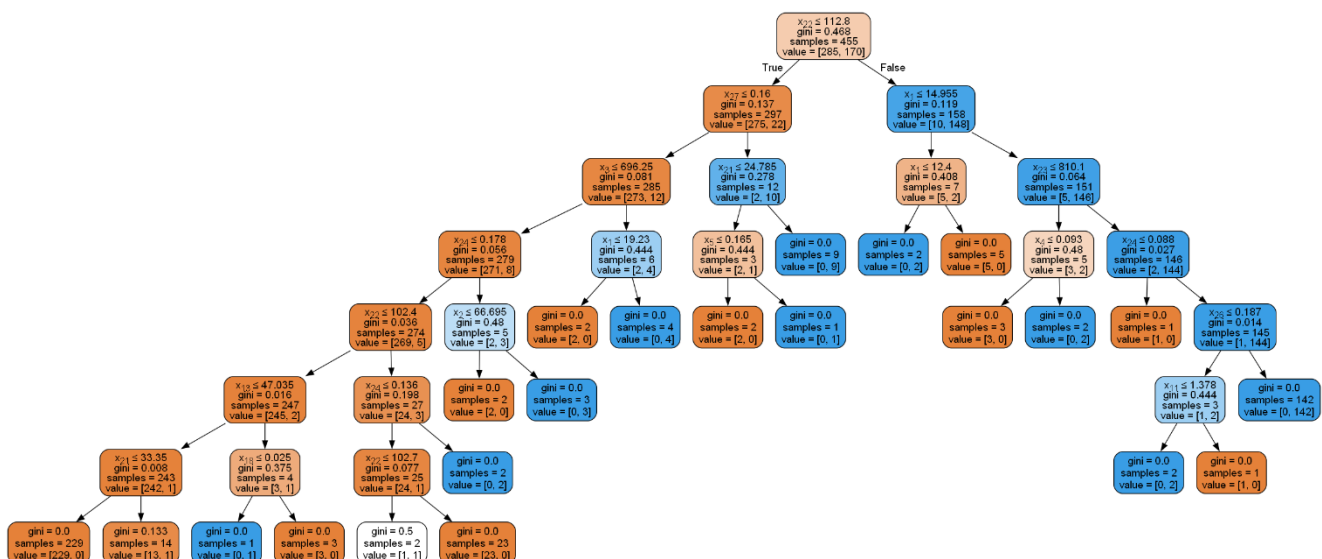
  o The full decision tree, without any depth limit, is highly complex with multiple branches. While this model likely captures intricate patterns in the training data, it risks overfitting, which is why it's essential to evaluate the effect of depth on generalization to unseen data.

- **Decision Tree at Depth 2:**

- o The tree at depth 2 is notably simple, with only a few splits. It captures the most prominent patterns but lacks the complexity to accurately classify more challenging cases. This is reflected in the accuracy, which, while decent, is not the highest.

- **Decision Trees at Depth 4 and 5:**

  - o These trees show increased complexity compared to depth 2, with more branches and decision nodes. The additional splits at these depths allow the model to capture more detailed relationships between features, which explains the slight increase in accuracy from depth 3 to 4.

  - o However, the accuracy plateau suggests that while these trees are more complex, they do not necessarily translate into better performance, likely because the added complexity does not correspond to better generalization to the test data.

- **Decision Tree at Depth 7:**

  - o The tree at depth 7 shows a significant increase in both complexity and accuracy. This depth likely strikes a balance between capturing complex patterns and maintaining generalization to the test set, as evidenced by the highest accuracy achieved.

  - o The increased depth allows the model to make more precise distinctions between classes, reducing misclassifications and leading to improved performance.

## 3. Analysis of the Trade-off Between Depth and Accuracy

- **Overfitting vs. Underfitting:**

  - o The results illustrate the classic trade-off between model complexity and accuracy. Shallow trees (e.g., depth 2 and 3) risk underfitting, as they do not capture enough detail from the data. On the other hand, extremely deep trees (beyond depth 7) could potentially overfit, memorizing training data rather than learning generalizable patterns.

  - o Depth 7 appears to be the sweet spot where the tree is complex enough to capture necessary details without overfitting, as indicated by the peak in accuracy.

- **Implications for Model Selection:**

  - o When choosing the optimal tree depth, it's essential to consider not only the accuracy but also the interpretability and computational efficiency. While a depth of 7 offers the best

accuracy, it comes at the cost of increased complexity, which may not always be desirable depending on the application.

- o For practical purposes, a tree depth that balances simplicity and performance, such as depth 4 or 5, might be preferable, especially when considering interpretability in a clinical context.

**Summary**

In conclusion, varying the depth of the Decision Tree reveals how model complexity impacts accuracy. While deeper trees generally offer better performance, there is a point beyond which additional complexity may not provide significant benefits. The tree at depth 7 provided the highest accuracy, suggesting that for this specific dataset and split, a more complex model was beneficial. However, the decision on the optimal depth should be guided by a balance between accuracy, interpretability, and the risk of overfitting.

# 6. Conclusion

## 6.1 Summary of Findings

In this report, Decision Tree classifiers were built to classify tumors as malignant or benign based on imaging features from the UCI Breast Cancer Wisconsin dataset. The models were evaluated across various train/test splits, and the impact of tree depth on accuracy was analyzed.

## 6.2 Final Remarks

The decision trees performed well across different splits, with accuracies varying depending on the depth of the tree. However, deeper trees sometimes led to overfitting, demonstrating the importance of choosing an optimal tree depth.

# 7. References

[1]. Wolberg,William, Mangasarian,Olvi, Street,Nick, and Street,W.. (1995). Breast Cancer Wisconsin (Diagnostic). UCI Machine Learning Repository. https://doi.org/10.24432/C5DW2B.

[2]. Uci-Ml-Repo. (n.d.). *GitHub - uci-ml-repo/ucimlrepo: Python package for dataset imports from UCI ML Repository*. GitHub. https://github.com/uci-ml-repo/ucimlrepo

[3]. *3.4. Metrics and scoring: quantifying the quality of predictions*. (n.d.). Scikit-learn. https://scikit-learn.org/stable/modules/model_evaluation.html

[4]. *DecisionTreeClassifier*. (n.d.). Scikit-learn. https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html

[5]. *seaborn.countplot — seaborn 0.13.2 documentation*. (n.d.). https://seaborn.pydata.org/generated/seaborn.countplot.html

[6]. *User Guide — graphviz 0.20.3 documentation*. (n.d.). https://graphviz.readthedocs.io/en/stable/manual.html

[7]. *Choosing Colormaps in Matplotlib — Matplotlib 3.9.2 documentation*. (n.d.). https://matplotlib.org/stable/users/explain/colors/colormaps.html