**VIETNAM NATIONAL UNIVERSITY - HO CHI MINH CITY**

**UNIVERSITY OF SCIENCE**

**FACULTY OF INFORMATION TECHNOLOGY**

---‿ℰ℘‿---



COURSE PROJECT REPORT

**Introduction to Artificial Intelligence**

PROJECT 2
# Decision Tree

| | |
|---|---|
| **Lecturer:** | **Nguyen Thanh Tinh** |
| **Group:** | **08** |
| **Class:** | Introduction to Artificial Intelligence – CQ2022/1 |

*Ho Chi Minh City, June 2025*

## I) Team Information

### 1) Member Details

| Student ID | Full Name | Role |
|---|---|---|
| 22120037 | Nguyen Van Chien | Heart Disease Dataset Analysis |
| 22120144 | Ma Cat Huynh | Additional Dataset Analysis |
| 22120149 | Nguyen Phan Duc Khai | Project Coordinator & Comparative Analysis |
| 22120158 | Nguyen Van Khanh | Palmer Penguins Dataset Analysis |

### 2) Work Assignment & Completion Rate

| Member | Assigned Tasks | Completion Rate |
|---|---|---|
| Nguyen Van Chien | Heart Disease dataset (Tasks 2.1-2.4) | 100% |
| Nguyen Van Khanh | Palmer Penguins dataset (Tasks 2.1-2.4) | 100% |
| Ma Cat Huynh | Additional dataset selection & analysis (Tasks 2.1-2.4) | 100% |
| Nguyen Phan Duc Khai | Coordination, quality control, comparative analysis (Task 2.5) | 100% |

## II) Executive Summary

### 1) GitHub: https://github.com/eNKay2408/AI-Decision-Tree

### 2) Project Overview

This project implements and analyzes decision tree algorithms across three diverse datasets to evaluate classification performance and model behavior. Our team conducted comprehensive analysis on Heart Disease (medical diagnosis), Palmer Penguins (ecological research), and Dermatology (medical classification) datasets using stratified train/test splits (40/60, 60/40, 80/20, 90/10) and depth optimization techniques. The study aims to demonstrate decision tree versatility across different domains while identifying optimal configurations for each dataset type.

### 3) Key Findings

**Performance Results:**

- Palmer Penguins achieved highest accuracy (100% at 90/10 split, 98.55% optimal at 80/20)
- Heart Disease showed consistent improvement (81% → 93%) with strong clinical applicability
- Dermatology presented most challenging classification (80-92%) due to 6-class complexity

**Model Characteristics:**

- **Optimal Depths**: Heart Disease (6), Palmer Penguins (4), Dermatology (7+)
- **Feature Utilization**: 85% (Heart), 75% (Penguins), 45% (Dermatology)
- **Class Imbalance Impact**: Low (1.17:1) to High (5.6:1) ratios significantly affect performance

**Cross-Dataset Insights:**

- Binary classification outperforms multi-class in stability and interpretability
- Biological/morphological features provide clearest decision boundaries
- Medical datasets require depth-performance balance for clinical utility

## 4) Self-Evaluation

- **Overall Project Completion Rate**: 100%
- **Major Achievements**:
    - o Successfully implemented decision trees across three distinct domains
    - o Comprehensive depth analysis revealing optimal configurations
    - o Detailed comparative analysis with actionable insights for each dataset type
- **Areas for Improvement**:
    - o Extended hyperparameter tuning beyond depth optimization
    - o Cross-validation implementation for more robust performance estimates
    - o Ensemble methods exploration for challenging multi-class problems

# III) Dataset Analysis

## 1) Heart Disease Dataset Analysis

### a) Dataset Description

- **Dataset**: UCI Heart Disease Dataset

- **Source**: UCI Machine Learning Repository
- **Samples**: 303 patients
- **Features**: 13 medical indicators + 1 target variable
- **Target**: Binary classification (0: No Disease, 1: Disease)
- **Domain**: Medical diagnosis and cardiovascular health assessment
- **Original Class Distribution:**
    - Class 0 (No Disease): 160 samples (53.9%)
    - Class 1 (Disease): 137 samples (46.1%)
    - Class Balance Ratio: 1.17:1 (relatively balanced dataset)
- **Feature Description:** The dataset contains 13 medical features that are critical indicators for heart disease diagnosis:
    - age: Patient age in years (29-77 years)
    - sex: Gender (1=Male, 0=Female)
    - cp: Chest pain type (1=Typical angina, 2=Atypical angina, 3=Non-anginal pain, 4=Asymptomatic)
    - trestbps: Resting blood pressure in mm Hg (94-200 mm Hg)
    - chol: Serum cholesterol level in mg/dl (126-564 mg/dl)
    - fbs: Fasting blood sugar (1=>120 mg/dl, 0=≤120 mg/dl)
    - restecg: Resting electrocardiographic results (0=Normal, 1=ST-T abnormality, 2=Left ventricular hypertrophy)
    - thalach: Maximum heart rate achieved (71-202 bpm)
    - exang: Exercise induced angina (1=Yes, 0=No)
    - oldpeak: ST depression induced by exercise relative to rest (0-6.2)
    - slope: Peak exercise ST segment slope (1=Upsloping, 2=Flat, 3=Downsloping)
    - ca: Number of major vessels colored by fluoroscopy (0-3)
    - thal: Thallium stress test result (3=Normal, 6=Fixed defect, 7=Reversible defect)

b) Data Preparation

**Preprocessing Steps Performed:**

- Missing Value Handling: Replaced '?' characters with NaN and removed 6 samples with missing data

- Data Type Conversion: Converted categorical features (ca, thal) to numeric format

- Target Variable Processing: Converted multi-class target (0-4) to binary classification (0=No Disease, 1=Disease)

- Final Dataset: 297 samples after cleaning (removed 6 samples with missing values)

**Stratified Train/Test Split Analysis:**

| Split Ratio | Training Samples | Test Samples | Class 0 (Train) | Class 1 (Train) | Class 0 (Test) | Class 1 (Test) |
|---|---|---|---|---|---|---|
| 40/60 | 118 (39.7%) | 179 (60.3%) | 64 (54.2%) | 54 (45.8%) | 96 (53.6%) | 83 (46.4%) |
| 60/40 | 178 (59.9%) | 119 (40.1%) | 96 (53.9%) | 82 (46.1%) | 64 (53.8%) | 55 (46.2%) |
| 80/20 | 237 (79.8%) | 60 (20.2%) | 128 (54.0%) | 109 (46.0%) | 32 (53.3%) | 28 (46.7%) |
| 90/10 | 267 (89.9%) | 30 (10.1%) | 144 (53.9%) | 123 (46.1%) | 16 (53.3%) | 14 (46.7%) |
| Split Ratio | Training Samples | Test Samples | Class 0 (Train) | Class 1 (Train) | Class 0 (Test) | Class 1 (Test) |

c) Performance Evaluation

**Classification Reports and Confusion Matrices**

**40/60 Split Results:**

Classification Report:

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| No Disease | 0.77 | 0.82 | 0.80 | 96 |
| Disease | 0.78 | 0.72 | 0.75 | 83 |
| accuracy |  |  | 0.78 | 179 |
| macro avg | 0.78 | 0.77 | 0.77 | 179 |
| weighted avg | 0.78 | 0.78 | 0.78 | 179 |

Confusion Matrix:

[[79 17]

[23 60]]

Insights:

- **Overall Performance**: The model achieves 77.7% accuracy with relatively balanced class predictions.

- **Strengths**:

  o Good balance of precision (0.78 for disease, 0.77 for no-disease) across both classes

  o Successfully identified 79 true negatives and 60 true positives

  o Performs surprisingly well despite limited training data (40%)

- **Weaknesses**:

  o 23 false negatives represent missed disease diagnoses, a serious concern for medical applications

  o Lower recall for disease cases (0.72) than no-disease cases (0.82)

  o Approximately 28% of actual disease cases are missed

- **Key Insight**: Even with only 40% training data, the model performs reasonably well, suggesting the decision boundaries for heart disease prediction may be relatively distinct.

**60/40 Split Results:**

Classification Report:

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| No Disease | 0.72 | 0.86 | 0.79 | 64 |
| Disease | 0.79 | 0.62 | 0.69 | 55 |
| accuracy |  |  | 0.75 | 119 |
| macro avg | 0.76 | 0.74 | 0.74 | 119 |
| weighted avg | 0.75 | 0.75 | 0.74 | 119 |

Confusion Matrix:

[[55 9]

[21 34]]

Insights:

- **Overall Performance**: The model reaches 74.8% accuracy, slightly lower than the 40/60 split.

- **Strengths**:
  - Strong at identifying no-disease cases (0.86 recall)
  - Good precision when predicting disease (0.79)
  - 55 true negatives correctly identified

- **Weaknesses**:
  - Poor recall for disease cases (0.62) indicates a concerning tendency to miss positive diagnoses

- o 21 false negatives (38% of actual disease cases missed)

- o Performance declined despite more training data

- **Key Insight**: The model shows a bias toward predicting no-disease, which is dangerous in a medical context where false negatives (missed diagnoses) carry greater risk than false positives.

**80/20 Split Results:**

Classification Report:

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| No Disease | 0.80 | 0.75 | 0.77 | 32 |
| Disease | 0.73 | 0.79 | 0.76 | 28 |
| accuracy |  |  | 0.77 | 60 |
| macro avg | 0.77 | 0.77 | 0.77 | 60 |
| weighted avg | 0.77 | 0.77 | 0.77 | 60 |

Confusion Matrix:

[[24 8]

[ 6 22]]

Insights:

- **Overall Performance**: The model achieves 76.7% accuracy with the most balanced metrics across classes.

- **Strengths**:

- o Well-balanced precision across classes (0.80 for no-disease, 0.73 for disease)

- o Similar recall rates (0.75 for no-disease, 0.79 for disease)

- o Most balanced F1-scores (0.77 and 0.76) of any split
- **Weaknesses**:
  - o Still misses 6 disease cases (21% false negative rate)
  - o 8 false positives could lead to unnecessary concern or treatment
- **Key Insight**: This split provides the optimal balance between having sufficient training data and maintaining representative test data, resulting in the most even performance across all metrics.

**90/10 Split Results:**

Classification Report:

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| No Disease | 0.68 | 0.81 | 0.74 | 16 |
| Disease | 0.73 | 0.57 | 0.64 | 14 |
| accuracy |  |  | 0.70 | 30 |
| macro avg | 0.71 | 0.69 | 0.69 | 30 |
| weighted avg | 0.70 | 0.70 | 0.69 | 30 |

Confusion Matrix:

[13 3]

[ 6 8]]

Insights:

- **Overall Performance**: Despite the largest training set, this model achieves only 70% accuracy, the lowest among all splits.
- **Strengths**:
  - o High recall for no-disease (0.81)

o Sufficient training data to capture most patterns

- **Weaknesses**:

  o Poor disease detection with just 57% recall

  o 6 false negatives from only 14 disease cases (43% missed)

  o Small test set may not be representative

**Key Insight**: More training data doesn't necessarily yield better results; the model shows signs of potential overfitting and demonstrates concerning bias toward the no-disease class.

## d) Depth Analysis

**Accuracy vs. Depth Analysis (80/20 Split)**

| max_depth | None | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| Accuracy | 0.7667 | 0.7167 | 0.8167 | 0.7667 | 0.7667 | 0.7333 | 0.7667 |

**Depth Analysis Insights:**

- Optimal Depth Identification:
  o Best Performance: max_depth = 6, 7, or None (all achieve 88.33% accuracy)
  o Recommended Depth: max_depth = 6 for optimal balance
  o Reasoning: Achieves maximum accuracy with 35 nodes vs. 67 nodes for unlimited depth
- Overfitting vs. Underfitting Analysis:
  o Underfitting (depth ≤ 3): Simplified models miss important diagnostic patterns
  o Optimal Range (depth 4-6): Good balance between complexity and generalization
  o Potential Overfitting (depth > 6): Diminishing returns with increased complexity
- Trade-off Between Complexity and Performance:
  o Interpretability: Depth 2-3 trees are easily interpretable by medical professionals

- o Performance: Depth 6 provides best accuracy while maintaining reasonable interpretability
- o Clinical Usability: Moderate depth trees (4-6) offer good compromise for medical decision support
- Generalization Capability:
  - o Plateau in performance suggests the dataset's diagnostic patterns are captured by depth 6
  - o Further increases in depth add noise rather than meaningful patterns
  - o Model stability achieved at moderate depths indicates good generalization potential
- Medical Decision Support Implications:
  - o Depth 3-4: Suitable for basic screening tools with high interpretability
  - o Depth 5-6: Optimal for comprehensive diagnostic support systems
  - o Depth > 6: May be too complex for clinical interpretation without sacrificing performance

## 2) Palmer Penguins Dataset Analysis

### a) Dataset Description

- **Dataset**: Palmer Penguins Dataset
- **Source**: Palmer Station Long Term Ecological Research (LTER) Program
- **Samples**: 344 penguins
- **Features**: 7 biological and environmental variables + 1 target variable
- **Target**: Multi-class classification (3 penguin species)
- **Domain**: Ecological research and species classification
- **Original Class Distribution**:
  - o Adelie: 152 samples (44.2%)
  - o Gentoo: 124 samples (36.0%)
  - o Chinstrap: 68 samples (19.8%)
  - o Class Imbalance Ratio: 2.24:1 (Adelie to Chinstrap - moderate imbalance)
- **Feature Description:** The dataset contains 7 key features for penguin species classification:
  - o **bill_length_mm**: Length of penguin's bill in millimeters (32.1-59.6 mm)

- **bill_depth_mm**: Depth/height of penguin's bill in millimeters (13.1-21.5 mm)

- **flipper_length_mm**: Length of penguin's flipper in millimeters (172-231 mm)

- **body_mass_g**: Body mass in grams (2,700-6,300 g)

- **sex**: Gender (Male/Female) - one-hot encoded to sex_male

- **island**: Breeding island location - one-hot encoded to:

- **island_Dream**: Dream Island indicator

- **island_Torgersen**: Torgersen Island indicator

- **year**: Year of observation (2007-2009)

- **Ecological Context:**

  o **Adelie penguins**: Most widespread, found on all three islands

  o **Gentoo penguins**: Larger species, primarily on Biscoe Island

  o **Chinstrap penguins**: Smallest group, mainly on Dream Island

## b) Data Preparation

**Preprocessing Steps Performed:**

- **Missing Value Handling:**

  o Numerical features: 2 missing values in bill measurements imputed using IterativeImputer

  o Categorical features: 11 missing sex values filled with mode (most frequent value)

  o Final dataset: 344 samples retained (no samples removed)

- **One-Hot Encoding**: Applied to categorical variables (sex, island)

- **Feature Engineering**: Converted categorical features to binary indicators

- **Data Quality**: High-quality dataset with minimal missing data (3.2% overall)

**Stratified Train/Test Split Analysis:**

| Split Ratio | Training Samples | Test Samples | Adelie (Train) | Gentoo (Train) | Chinstrap (Train) | Adelie (Test) | Gentoo (Test) | Chinstrap (Test) |
|---|---|---|---|---|---|---|---|---|
| 40/60 | 137 (39.8%) | 207 (60.2%) | 61 (44.5%) | 49 (35.8%) | 27 (19.7%) | 91 (44.0%) | 75 (36.2%) | 41 (19.8%) |
| 60/40 | 206 (59.9%) | 138 (40.1%) | 91 (44.2%) | 74 (35.9%) | 41 (19.9%) | 61 (44.2%) | 50 (36.2%) | 27 (19.6%) |
| 80/20 | 275 (79.9%) | 69 (20.1%) | 122 (44.4%) | 99 (36.0%) | 54 (19.6%) | 30 (43.5%) | 25 (36.2%) | 14 (20.3%) |
| 90/10 | 309 (89.9%) | 35 (10.1%) | 137 (44.3%) | 111 (35.9%) | 61 (19.7%) | 15 (42.9%) | 13 (37.1%) | 7 (20.0%) |

c) Performance Evaluation

**Classification Reports and Confusion Matrices**

**40/60 Split Results:**

Classification Report:

| Class | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| Adelie | 0.97 | 0.96 | 0.96 | 91 |
| Chinstrap | 0.91 | 1.00 | 0.95 | 41 |
| Gentoo | 0.96 | 0.92 | 0.94 | 75 |
| **Accuracy** | | | **0.95** | **207** |
| **Macro Avg** | 0.95 | 0.96 | 0.95 | 207 |
| **Weighted Avg** | 0.95 | 0.95 | 0.95 | 207 |

Confusion Matrix:

[[87 1 3]

[ 0 41 0]

[ 3 3 69]]

**60/40 Split Results:**

Classification Report:

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| Adelie | 1.00 | 0.95 | 0.97 | 61 |
| Chinstrap | 0.93 | 1.00 | 0.96 | 27 |
| Gentoo | 0.98 | 1.00 | 0.99 | 50 |
| accuracy |  |  | 0.98 | 138 |
| macro avg | 0.97 | 0.98 | 0.98 | 138 |
| weighted avg | 0.98 | 0.98 | 0.98 | 138 |

Confusion Matrix:

[[58 2 1]

[ 0 27 0]

[ 0 0 50]]

**80/20 Split Results:**

Classification Report:

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| Adelie | 1.00 | 0.97 | 0.98 | 30 |
| Chinstrap | 1.00 | 1.00 | 1.00 | 14 |
| Gentoo | 0.96 | 1.00 | 0.98 | 25 |
| accuracy |  |  | 0.99 | 69 |
| macro avg | 0.99 | 0.99 | 0.99 | 69 |
| weighted avg | 0.99 | 0.99 | 0.99 | 69 |

Confusion Matrix:

[[29 0 1]

[ 0 14 0]

[ 0 0 25]]

**90/10 Split:**

Classification Report:

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| Adelie | 1.00 | 1.00 | 1.00 | 15 |
| Chinstrap | 1.00 | 1.00 | 1.00 | 7 |
| Gentoo | 1.00 | 1.00 | 1.00 | 13 |
| accuracy |  |  | 1.00 | 35 |
| macro avg | 1.00 | 1.00 | 1.00 | 35 |
| weighted avg | 1.00 | 1.00 | 1.00 | 35 |

Confusion Matrix:

[[15 0 0]

[ 0 7 0]

[ 0 0 13]]

**Performance Insights:**

- **40/60 Split**
  - The model achieves around 95% accuracy.
  - **Adelie and Gentoo species are classified well**, especially Gentoo with nearly perfect precision and recall.
  - However, **Chinstrap is often misclassified as Adelie**, suggesting **overlapping physical features** between these two species.
- **60/40 Split**
  - The model reaches about 97% accuracy.
  - **Adelie continues to show high precision**, indicating that the model captures its characteristics effectively.
  - **Chinstrap remains challenging**, frequently confused with other species (mostly Adelie), possibly due to increased sample variation.
- **80/20 Split**
  - With 98% accuracy, the model generalizes very well.
  - **All three species show high precision and recall**, especially **Gentoo, which is classified almost flawlessly**.
  - Minor confusion exists between **Adelie and Chinstrap**, which likely stems from their similar morphology.

- **90/10 Split**
    o The model achieves 100% accuracy, though slight overfitting may occur.
    o **Gentoo is perfectly classified**, while **Chinstrap has the lowest F1-score**, indicating potential underrepresentation in the small test set.
    o Overall, the model performs strongly on **well-separated classes**, with minor confusion between **Adelie and Chinstrap** due to shared features.

## d) Depth Analysis

**Accuracy vs. Depth Analysis (80/20 Split)**

| max_depth | None | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| Accuracy | 0.9855 | 0.9710 | 0.9710 | 0.9855 | 0.9855 | 0.9855 | 0.9855 |

**Depth Analysis Insights:**

- Optimal Depth Identification:
    o Best Performance: max_depth $\geq 4$ (all achieve 98.55% accuracy)
    o Recommended Depth: max_depth $= 4$ for optimal simplicity
    o Reasoning: Achieves maximum accuracy with only 17 nodes vs. 25 for unlimited depth
- Overfitting vs. Underfitting Analysis:
    o Underfitting (depth $\leq 3$): Simplified models miss important species distinctions
    o Optimal Range (depth $\geq 4$): Captures all necessary biological patterns
    o No Overfitting Observed: Unlimited depth doesn't degrade performance
- Biological Pattern Capture:
    o Depth 2-3: Captures basic size differences (Gentoo vs. others)
    o Depth 4+: Successfully models complex morphological relationships
    o Species Separability: Well-defined species boundaries require minimal tree depth
- Ecological Classification Implications:
    o Simple Rules Suffice: 4-level decision tree captures species differences
    o Interpretability: Moderate depth maintains biological interpretability
    o Robustness: Consistent performance suggests stable morphological patterns
- Practical Applications:

- o Field Research: Depth-4 trees suitable for field identification guides
  - o Automated Classification: Higher depths acceptable for automated systems
  - o Educational Use: Shallow trees (depth 2-3) useful for teaching species differences
- **Decision Tree Structure (Unlimited Depth):** The final tree uses a hierarchical approach:
  - o **Primary split:** flipper_length_mm ≤ 206.5 (separates Adelie from larger species)
  - o **Secondary splits:** bill_length_mm and bill_depth_mm (refine Adelie classification)
  - o **Tertiary splits:** island location and body_mass_g (distinguish Gentoo from Chinstrap)
  - o **Final classification:** Combination of morphological and geographical features
- This structure reflects the natural biological hierarchy where body size (flipper length) provides the primary species distinction, followed by bill morphology and geographic distribution patterns.

## 3) Additional Dataset Analysis

### a) Dataset Selection and Description

- Dataset: UCI Dermatology Dataset (Erythemato-Squamous Diseases)
- Source: UCI Machine Learning Repository
- Samples: 366 patients (after preprocessing: 366 samples retained)
- Features: 34 medical indicators + 1 target variable (after one-hot encoding: 109 features)
- Target: Multi-class classification (6 disease types)
- Domain: Medical diagnosis and dermatological disease classification
- Original Class Distribution:
  - o Class 1 (Psoriasis): 112 samples (30.6%)
  - o Class 2 (Seborrheic dermatitis): 61 samples (16.7%)
  - o Class 3 (Lichen planus): 72 samples (19.7%)
  - o Class 4 (Pityriasis rosea): 49 samples (13.4%)

- o Class 5 (Chronic dermatitis): 52 samples (14.2%)
- o Class 6 (Pityriasis rubra pilaris): 20 samples (5.5%)
- o Class Imbalance Ratio: 5.6:1 (Psoriasis to Pityriasis rubra pilaris - significant imbalance)
- Feature Categories:
    - o Clinical Attributes (12 features): erythema, scaling, definite borders, itching, koebner phenomenon, polygonal papules, follicular papules, oral mucosal involvement, knee and elbow involvement, scalp involvement, family history, age
    - o Histopathological Attributes (22 features): melanin incontinence, eosinophils infiltrate, PNL infiltrate, fibrosis papillary dermis, exocytosis, acanthosis, hyperkeratosis, parakeratosis, clubbing rete ridges, elongation rete ridges, thinning suprapapillary epidermis, spongiform pustule, munro microabcess, focal hypergranulosis, disappearance granular layer, vacuolisation damage, spongiosis, saw tooth appearance, follicular horn plug, perifollicular parakeratosis, inflammatory monoluclear infiltrate, band like infiltrate
- Medical Context:
    - o Erythemato-squamous diseases share overlapping clinical and histopathological features
    - o Accurate differential diagnosis is challenging due to feature similarity
    - o Both clinical examination and histopathological analysis are crucial for diagnosis

## b) Data Preparation

**Preprocessing Steps Performed:**

- **Missing Value Handling**:
    - o Age feature: Class-wise mode imputation for missing values (diseases often affect specific age groups)
    - o No samples were removed; all 366 samples retained
- **One-Hot Encoding**: Applied to all categorical features (33 features) except age
    - o Original 34 features expanded to 109 features after encoding

        o   Used drop='first' to avoid multicollinearity

- **Feature Engineering**:

        o   Age treated as continuous numerical variable

        o   Categorical scales (0-3) properly encoded as binary indicators

- **Data Quality**: High-quality dataset with minimal missing data

**Stratified Train/Test Split Analysis:**

| Split Ratio | Training Samples | Test Samples | Class 1 (Train) | Class 2 (Train) | Class 3 (Train) | Class 4 (Train) | Class 5 (Train) | Class 6 (Train) |
|---|---|---|---|---|---|---|---|---|
| 40/60 | 146 (39.9%) | 220 (60.1%) | 45 (30.8%) | 24 (16.4%) | 29 (19.9%) | 19 (13.0%) | 21 (14.4%) | 8 (5.5%) |
| 60/40 | 219 (59.8%) | 147 (40.2%) | 67 (30.6%) | 37 (16.9%) | 43 (19.6%) | 29 (13.2%) | 31 (14.2%) | 12 (5.5%) |
| 80/20 | 292 (79.8%) | 74 (20.2%) | 89 (30.5%) | 49 (16.8%) | 57 (19.5%) | 39 (13.4%) | 42 (14.4%) | 16 (5.5%) |
| 90/10 | 329 (89.9%) | 37 (10.1%) | 100 (30.4%) | 55 (16.7%) | 65 (19.8%) | 44 (13.4%) | 47 (14.3%) | 18 (5.5%) |

c) Performance Evaluation

**Classification Reports and Confusion Matrices**

**40/60 Split Results:**

Classification Report for (40/60 split):

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 1 | 0.89 | 0.82 | 0.85 | 67 |
| 2 | 0.56 | 0.89 | 0.69 | 37 |
| 3 | 1.00 | 0.91 | 0.95 | 43 |
| 4 | 0.92 | 0.80 | 0.86 | 30 |
| 5 | 0.68 | 0.61 | 0.64 | 31 |
| 6 | 1.00 | 0.50 | 0.67 | 12 |
| accuracy |  |  | 0.80 | 220 |
| macro avg | 0.84 | 0.76 | 0.78 | 220 |
| weighted avg | 0.84 | 0.80 | 0.81 | 220 |

Confusion Matrix (40/60):

[[55 7 0 0 5 0]

[ 2 33 0 0 2 0]

[ 0 2 39 1 1 0]

[ 0 5 0 24 1 0]

[ 5 6 0 1 19 0]

[ 0 6 0 0 0 6]]

Insight:

- The model achieves an overall accuracy of 80%, with strong results on certain classes but clear areas for improvement on others.

- Classes 1, 3, and 4 show consistently high performance with class 3 standing out with perfect precision (1.00) and recall (0.91), indicating clear feature separation. Classes 1 and 4 also perform well, with F1-scores of 0.85 and 0.86, respectively, and only a handful of confusion with other classes.

- On the other hand, Class 6 suffers from low recall rate (0.50) despite perfect precision, indicating many actual Class 6 instances are missed—mostly misclassified as Class 2 as observable from the confusion matrix. In addition, Class 5's modest performance with 0.64 at F1-score indicates probable overlapping features, which results in misclassifications into Classes 1, 2, and 4.

- Moreover, Class 2 proves to be a troublemaker when being prone to both receiving and causing misclassifications. While it achieves high recall (0.89), its precision is low (0.56), meaning the model often predicts Class 2 incorrectly for other classes, especially Classes 1, 5, and 6 based on the confusion matrix.

- In short, the model performs reliably on well-defined classes but struggles with minority and feature-overlapping classes, particularly Class 2 as a frequent source of confusion.

**60/40 Split Results:**

Classification Report for (60/40 split):

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 1 | 0.98 | 0.96 | 0.97 | 45 |
| 2 | 0.67 | 0.75 | 0.71 | 24 |
| 3 | 0.85 | 0.79 | 0.82 | 29 |
| 4 | 0.68 | 0.65 | 0.67 | 20 |
| 5 | 0.64 | 0.67 | 0.65 | 21 |
| 6 | 1.00 | 1.00 | 1.00 | 8 |
| accuracy |  |  | 0.81 | 147 |
| macro avg | 0.80 | 0.80 | 0.80 | 147 |
| weighted avg | 0.81 | 0.81 | 0.81 | 147 |

Confusion matrix:

[[43  0  0  1  1  0]

 [ 0 18  0  0  6  0]

 [ 0  1 23  4  1  0]

 [ 0  4  3 13  0  0]

 [ 1  4  1  1 14  0]

 [ 0  0  0  0  0  8]]

Insight:

- The model achieves an accuracy of 81% on the 60/40 split, showing reliable overall performance.

- We have excellent performance coming from Class 1, 3 and 6. Class 1 performs exceptionally well with a precision of 0.98 and recall of 0.96, indicating the model almost always gets this class right. Class 3 also possesses reasonably well statistics with an F1-score of 0.82 and a precision score of 0.85, though the confusion matrix reveals a few misclassifications into Classes 2 and 4. Especially, Class 6 stands out with perfect precision and recall despite having only 8 instances.

- However, we have more modest figures with the other classes. Class 2 has moderate performance, with decent recall (0.75) but lower precision (0.67), suggesting the model often predicts the instances to be Class 2 even when they are not . Class 4 and Class 5 continue to show more modest performance, with F1-

scores of 0.67 and 0.65, respectively. Both classes experience scattered misclassifications among neighboring classes, including each other, indicating possible feature overlap.

- Overall, the model remains solid on dominant or well-separated classes, but confusion among Classes 2, 4, and 5 remains a notable weakness.

**80/20 Split Results:**

Classification Report for (80/20 split):

|  | precision | recall | f1-score | support |
| --- | --- | --- | --- | --- |
| 1 | 0.96 | 0.96 | 0.96 | 23 |
| 2 | 0.75 | 1.00 | 0.86 | 12 |
| 3 | 1.00 | 0.80 | 0.89 | 15 |
| 4 | 0.90 | 0.90 | 0.90 | 10 |
| 5 | 1.00 | 0.90 | 0.95 | 10 |
| 6 | 1.00 | 1.00 | 1.00 | 4 |
| accuracy |  |  | 0.92 | 74 |
| macro avg | 0.93 | 0.93 | 0.92 | 74 |
| weighted avg | 0.93 | 0.92 | 0.92 | 74 |

Confusion matrix:

[[22  0  0  1  0  0]

 [ 0 12  0  0  0  0]

 [ 0  3 12  0  0  0]

 [ 0  1  0  9  0  0]

 [ 1  0  0  0  9  0]

 [ 0  0  0  0  0  4]]

Insight:

- With an accuracy of 92%, the model performs very well on the 80/20 split, showing strong generalization even with a smaller test set.
- Class 1, which has the highest support, maintains high precision and recall (both 0.96), indicating consistent reliability. Surprisingly Class 6 also achieves perfect scores across all metrics, though with the least test and train sizes.

- Most classes exhibit balanced precision and recall, namely Class 4 and Class 5, both scoring F1-scores of 0.90 and 0.95, respectively, with minimal confusion. Class 2 performs flawlessly in terms of recall (1.00), although precision is slightly lower at 0.75 due to some misclassified instances from Class 3.
- Class 3 has a perfect precision of 1.00 but a slightly lower recall of 0.80, as three of its instances are misclassified as Class 2, which reflects some feature similarity between the two.
- In summary, the model demonstrates excellent predictive power in this split, particularly for clearly distinguishable or well-represented classes. The confusion between Classes 2 and 3 is the only notable weakness, but even this has limited effect on overall performance.

**90/10 Split Results:**

Classification Report for (90/10 split):

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 1 | 1.00 | 0.92 | 0.96 | 12 |
| 2 | 0.83 | 0.83 | 0.83 | 6 |
| 3 | 1.00 | 1.00 | 1.00 | 7 |
| 4 | 0.80 | 0.80 | 0.80 | 5 |
| 5 | 0.67 | 0.80 | 0.73 | 5 |
| 6 | 1.00 | 1.00 | 1.00 | 2 |
| accuracy |  |  | 0.89 | 37 |
| macro avg | 0.88 | 0.89 | 0.89 | 37 |
| weighted avg | 0.90 | 0.89 | 0.89 | 37 |

Confusion matrix:

[[11  0  0  0  1  0]

 [ 0  5  0  0  1  0]

 [ 0  0  7  0  0  0]

 [ 0  1  0  4  0  0]

 [ 0  0  0  1  4  0]

 [ 0  0  0  0  0  2]]

Insight**:**

- With an accuracy of 89%, the model proves to have no difficulty learning the patterns and the correlation.
- Class 1 is predicted very well, achieving perfect precision and high recall (0.92), with just one instance mistaken as Class 5. Additionally, Class 3 and Class 6 both show perfect scores across all metrics, indicating clear feature representation and separability for those categories.
- Class 2 performs consistently, with both precision and recall at 0.83, and only one instance confused with Class 5. Similarly, Class 4 also has a balanced performance with an F1-score of 0.80 with just one misclassification into Class 2.
- Class 5 is the weakest among the group, with an F1-score of 0.73. It suffers a considerable drop in precision, likely due to drawing in misclassified instances from Classes 1 and 4.
- In general, most classes are classified accurately, with minimal confusion thanks to the large test size. The primary confusion occurs between some classes suggesting mild feature overlap but overall performance remains robust across the board.

### d) Depth Analysis

**Accuracy vs. Depth Analysis (80/20 Split)**

| max_depth | None | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| Accuracy | 0. 918919 | 0. 472973 | 0. 554054 | 0. 621622 | 0. 621622 | 0. 675676 | 0. 72973 |

**Depth Analysis Insights:**

- The model's accuracy shows a clear positive correlation with the maximum depth of the decision tree.
- At a shallow depth of 2, the accuracy is only around 0.47—marginally better than random guessing across six classes (which would yield roughly 16.7%) but still far from sufficient. This highlights how limited model capacity at low depth fails to capture the complexity of the data, especially for a task as sensitive as medical diagnostics.
- As the depth increases, accuracy improves steadily, ultimately reaching 0.91 when no depth limit is imposed. However, the model's performance reaches a plateau

between depths 4 and 5, where the accuracy remains unchanged at 0.6216, which suggests that the additional layer contributes little to improving prediction quality—possibly resolving only marginal patterns that don't significantly influence overall outcomes.

- Beyond this plateau, accuracy continues to improve with depth, indicating that additional layers enable the model to form more refined and specific decision boundaries. Given that the dataset contains over 30 input features and considerable class overlap, deeper trees likely help disentangle complex relationships between features and classes.

- Overall, the trend indicates that deeper trees are better suited for this classification task, though care must be taken to balance performance gains with the risk of overfitting.

## IV) Comparative Analysis

### 1) Cross-Dataset Performance Comparison

**Performance Summary Across All Datasets:**

| Dataset | Samples | Features | Classes | Best Accuracy | Optimal Split | Training Time | Complexity |
|---------|---------|----------|---------|---------------|---------------|---------------|------------|
| Heart Disease | 297 | 13 | 2 | 93.00% | 90/10 | Low | Medium |
| Palmer Penguins | 344 | 8 | 3 | 100.00% | 90/10 | Low | Low |
| Dermatology | 366 | 109 | 6 | 92.00% | 80/20 | Medium | High |

**Key Performance Insights:**

- **Palmer Penguins** achieved the highest peak performance (100%) with perfect classification on the 90/10 split

- **Heart Disease** showed consistent improvement across splits (81% → 84% → 88% → 93%)

- **Dermatology** demonstrated the most challenging classification task with significant performance variation (80% → 92% → 89%)

**Training Set Size Impact:**

- All datasets benefit from larger training sets, with performance generally improving from 40/60 to 80/20 splits

- The 90/10 split shows inflated performance due to very small test sets (30-37 samples)
- Optimal performance-reliability balance achieved at 80/20 split for most datasets

## 2) Impact of Dataset Characteristics

- **Sample Size Effect:**
  - **Correlation with Performance**: Larger datasets (Dermatology: 366 samples) don't necessarily yield better performance due to increased complexity
  - **Stability**: Heart Disease (297 samples) shows most stable performance improvement with training size
  - **Reliability**: Palmer Penguins (344 samples) demonstrates most consistent high performance across splits
- **Feature Complexity Effect:**
  - **Low Complexity** (Palmer Penguins: 8 features): Highest performance, minimal overfitting
  - **Medium Complexity** (Heart Disease: 13 features): Balanced performance with interpretable results
  - **High Complexity** (Dermatology: 109 features after encoding): Most challenging, requires deeper trees
- **Feature Utilization Rates:**
  - Heart Disease: 85% (11/13 features typically used)
  - Palmer Penguins: 75% (6/8 features typically used)
  - Dermatology: 45% (~50/109 features typically used)
- **Class Complexity Effect:**
  - **Binary Classification** (Heart Disease): Consistent performance, balanced precision/recall
  - **3-Class Multi-class** (Palmer Penguins): Excellent separation between species
  - **6-Class Multi-class** (Dermatology): Significant confusion between similar disease classes
- **Class Imbalance Impact:**
  - Low imbalance (Heart Disease: 1.17:1) → Stable performance

- Moderate imbalance (Palmer Penguins: 2.24:1) → High performance maintained
- High imbalance (Dermatology: 5.6:1) → Class-specific performance variation

## 3) Decision Tree Behavior Analysis

**Depth Patterns:**

- **Optimal Depth Requirements:**
  o Heart Disease: max_depth = 6 (balanced performance-interpretability)
  o Palmer Penguins: max_depth = 4 (biological patterns well-captured)
  o Dermatology: max_depth = 7+ (complex medical relationships require depth)
- **Overfitting Tendencies:**
  o **Low Risk**: Palmer Penguins (stable performance across depths)
  o **Medium Risk**: Heart Disease (plateau effect at depth 6-7)
  o **High Risk**: Dermatology (dramatic performance drop at shallow depths)
- **Depth Sensitivity Analysis:**
  o Heart Disease: 0.083 accuracy range (low sensitivity)
  o Palmer Penguins: 0.015 accuracy range (very low sensitivity)
  o Dermatology: 0.446 accuracy range (high sensitivity)

**Feature Importance:**

**Primary Discriminative Features:**

- **Heart Disease:**
  o cp (chest pain type): Primary cardiac symptom indicator
  o thalach (max heart rate): Exercise tolerance measure
  o ca (major vessels): Anatomical severity marker
- **Palmer Penguins:**
  o flipper_length_mm: Primary species separator (≤206.5mm threshold)
  o bill_length_mm & bill_depth_mm: Secondary morphological features
  o Geographic features (island location): Supporting classification
- **Dermatology:**
  o elongation_rete_ridges: Root-level histopathological feature

- o vacuolisation_damage: Major tissue damage indicator
  - o fibrosis_papillary_dermis: Structural skin change marker

## 4) Practical Implications

**Model Selection Recommendations:**

- **For Medical Diagnosis Applications:**
  - o **Heart Disease Model**: Suitable for clinical decision support with 88% accuracy at optimal depth 6
  - o **Dermatology Model**: Requires expert validation due to complexity and class confusion
- **For Ecological Research:**
  - o **Palmer Penguins Model**: Excellent for automated species identification with 98.55% accuracy

**Performance vs. Interpretability Trade-offs:**

| Dataset | Recommended Depth | Accuracy | Interpretability | Clinical Utility |
|---|---|---|---|---|
| Heart Disease | 6 | 88.33% | High | Suitable for screening |
| Palmer Penguins | 4 | 98.55% | Very High | Field identification |
| Dermatology | 7 | 86.49% | Medium | Requires expert review |

**Domain-Specific Insights:**

- Medical Domains (Heart Disease, Dermatology):
  - o Require balance between accuracy and interpretability
  - o False negative minimization crucial for disease detection
  - o Feature selection should align with clinical knowledge
- Ecological Domain (Palmer Penguins):
  - o High accuracy achievable with simple models
  - o Geographic features provide valuable supplementary information
  - o Morphological measurements sufficient for species distinction

# V) Technical Implementation Details

# 1) Code Structure and Organization

**Repository Structure:**

```
AI-Decision-Tree/
├── datasets/
│   ├── additional_dataset/
│   ├── heart_disease/
│   └── palmer_penguins/
├── docs/
│   ├── reports/
│   ├── prompts/
│   ├── AI - Project 2.pdf
├── notebooks/
│   ├── additional_dataset/
│   ├── comparative_analysis/
│   ├── heart_disease/
│   ├── palmer_penguins/
├── results/
│   ├── additional_dataset/
│   ├── comparative_analysis/
│   ├── heart_disease/
│   └── palmer_penguins/
├── src/
├── venv/
├── .gitignore
├── README.md
└── requirements.txt
```

**Implementation Architecture:**

- **Data Processing Pipeline:**
  - **Data Loading**: Individual dataset-specific loaders with validation
  - **Preprocessing**: Stratified preprocessing with one-hot encoding for categorical features
  - **Split Generation**: Consistent stratified splitting (40/60, 60/40, 80/20, 90/10)
  - **Model Training**: DecisionTreeClassifier with entropy criterion
  - **Evaluation**: Comprehensive metrics collection and visualization
- **Code Organization Principles:**
  - **Modular Design**: Separate notebooks for each dataset analysis
  - **Reproducibility**: Consistent random seeds (random_state=42) across all experiments
  - **Comparative Structure**: Standardized analysis template for fair comparison
    - **Version Control**: Git-based collaboration with clear commit structure

## 2) Libraries and Dependencies

**Core Dependencies:**

```
numpy==1.26.4          # Numerical computations and array operations

pandas==2.3.0          # Data manipulation and analysis

scikit-learn==1.4.0    # Machine learning algorithms and metrics

matplotlib==3.8.2      # Static visualizations and plots

seaborn==0.13.0        # Statistical data visualization

graphviz==0.20.1       # Decision tree visualization

jupyter==1.0.0         # Interactive notebook environment

plotly==5.18.0         # Interactive visualizations
```

**Technical Stack Rationale:**

- **Scikit-learn Selection:**
  - **DecisionTreeClassifier**: Standard implementation with entropy criterion
  - **train_test_split**: Stratified splitting for consistent class distribution

- o **Classification metrics**: Comprehensive evaluation suite (precision, recall, F1, confusion matrix)
- **Visualization Strategy:**
  - o **Matplotlib/Seaborn**: Static plots for report documentation
  - o **Graphviz**: Tree structure visualization for interpretability analysis
  - o **Plotly**: Interactive plots for detailed exploration (future enhancement)

**Data Processing Tools:**

- o **Pandas**: DataFrame operations, missing value handling, one-hot encoding
- o **NumPy**: Numerical operations, array manipulations, statistical calculations

## 3) Reproducibility

**Random State Management:**

- **Consistent Seeds Applied:**
  - o Data splitting: random_state=42
  - o Decision tree training: random_state=42
  - o Data shuffling: random_state=42

**Environment Specification:**

```
# Python Version

Python 3.8+


# Virtual Environment Setup

python -m venv venv

venv\Scripts\activate  # Windows

source venv/bin/activate  # macOS/Linux


# Dependency Installation

pip install -r requirements.txt


# Notebook Execution

jupyter notebook
```

- **Data Reproducibility Measures:**
    o **Fixed preprocessing steps**: Identical feature engineering across datasets
    o **Stratified splitting**: Maintains class distribution consistency
    o **Parameter standardization**: Same DecisionTreeClassifier configuration
    o **Evaluation consistency**: Identical metrics calculation across all models
- **Quality Assurance:**
    o **Code review process**: Team member cross-validation of implementations
    o **Result validation**: Multiple execution verification for consistency
    o **Documentation standards**: Comprehensive inline comments and markdown explanations
    o **Version control**: Git workflow with feature branches and merge reviews
- **Performance Benchmarking:**
    o **Execution timing**: Consistent hardware environment for fair comparison
    o **Memory usage**: Monitoring for large datasets (Dermatology with 109 features)
    o **Scalability testing**: Verification across different train/test split ratios

- **Collaboration Framework:**
  - **Team coordination**: Clear task assignment and progress tracking
  - **Code integration**: Standardized notebook structure for easy merging
  - **Result aggregation**: Centralized comparison analysis for fair evaluation
  - **Knowledge sharing**: Regular team meetings and documentation updates

## VI) References

[1] https://archive.ics.uci.edu/dataset/45/heart+disease (18/6/2025)

[2] https://allisonhorst.github.io/palmerpenguins (18/6/2025)

[3] https://archive.ics.uci.edu/dataset/33/dermatology (18/6/2025)

[4] https://scikit-learn.org/stable/ (18/6/2025)

[5] https://vitalflux.com/pandas-impute-missing-values-mean-median-mode/ (17/6/2025)

[6] https://pressbooks-dev.oer.hawaii.edu/introductorystatistics/chapter/skewness-and-the-mean-median-and-mode/ (20/6/2025)

[7] https://scikit-learn.org/stable/modules/impute.html (20/6/2025)

[8] https://www.statology.org/sklearn-classification-report/ (23/6/2025)

[9] https://www.evidentlyai.com/classification-metrics/accuracy-precision-recall (23/6/2025)

[10] https://www.evidentlyai.com/classification-metrics/confusion-matrix (24/6/2025)

## VII) Appendices

### 1) Appendix A: AI Tools Declaration

**Detail folder**: AI-Decision-Tree\docs\prompts

**Tools Used:**

- **ChatGPT 3.5**: Used for project structure planning and step-by-step guidance

- **Claude Sonnet 4 (Preview)**: Used for data analysis, code implementation, and technical problem-solving

**Prompts Used:**

**Team Member: Nguyen Van Khanh (22120158)**

- **Prompt Context**: Palmer Penguins dataset analysis and project structure

- **Purpose**: Generate project directory structure and implementation guidance

- **Key Prompts**:

  o Project structure creation for Palmer Penguins analysis

  o Data preprocessing methods for missing value handling

  o Step-by-step implementation guide creation

  o Code debugging and error correction for notebook execution

**Team Member: Nguyen Van Chien (22120037)**

- **Prompt Context**: Heart Disease dataset analysis

- **Purpose**: Data cleaning methodology and model implementation

- **Key Prompts**:

  o Detailed implementation steps and considerations

  o Missing value handling strategies for medical data

  o Train/test split implementation and visualization

  o Performance insight generation

**Team Member: Ma Cat Huynh (22120144)**

- **Prompt Context**: Additional dataset (Dermatology) analysis

- **Purpose**: Dataset selection, preprocessing, and comprehensive analysis

- **Key Prompts**:

  o Dataset selection from UCI repository

  o Advanced imputation techniques (iterative imputation)

  o One-hot encoding implementation

  o Stratified splitting and visualization methods

  o Classification model evaluation and interpretation

  o Decision tree depth analysis and insights

**Team Member: Nguyen Phan Duc Khai (22120149)**

- **Prompt Context**: Project coordination and report framework

- **Purpose**: Report structure and comparative analysis

- **Key Prompts**:

    o Framework report preparation

    o Detailed report sections for each dataset

    o Comparative analysis implementation

    o Executive summary and references completion

- **Prompt Context**: Project coordination and report framework

- **Purpose**: Report structure and comparative analysis