# CSE 422: ARTIFICIAL INTELLIGENCE

Fall 2024

# Project Report

**MUSHROOM CLASS DETECTOR**

Student Information:

B M Rauf - 22201782

Mahdi Hasan - 22201760

Section: 16, Group : 16

Date : 06 January 2025

Submitted to:

Hasnat Jamil Bhuiyan

Maazin Munawar

Department of Computer Science and Engineering (CSE)

# TABLE OF CONTENTS

# Introduction

Mushroom classification is a critical task in the field of mycology and food safety. This project aims to classify mushrooms into categories based on their characteristics using machine learning models. The dataset contains various attributes of mushrooms, such as cap shape, surface, and color, which serve as predictors for the classification task.

The goal is to preprocess the data, scale the features, split the dataset, and train multiple models to identify the most effective one for the classification task. With the growing popularity of mushroom hunting, this model is highly relevant and can help prevent the consumption of toxic mushrooms, which can be life-threatening.

# Dataset Description

- Source

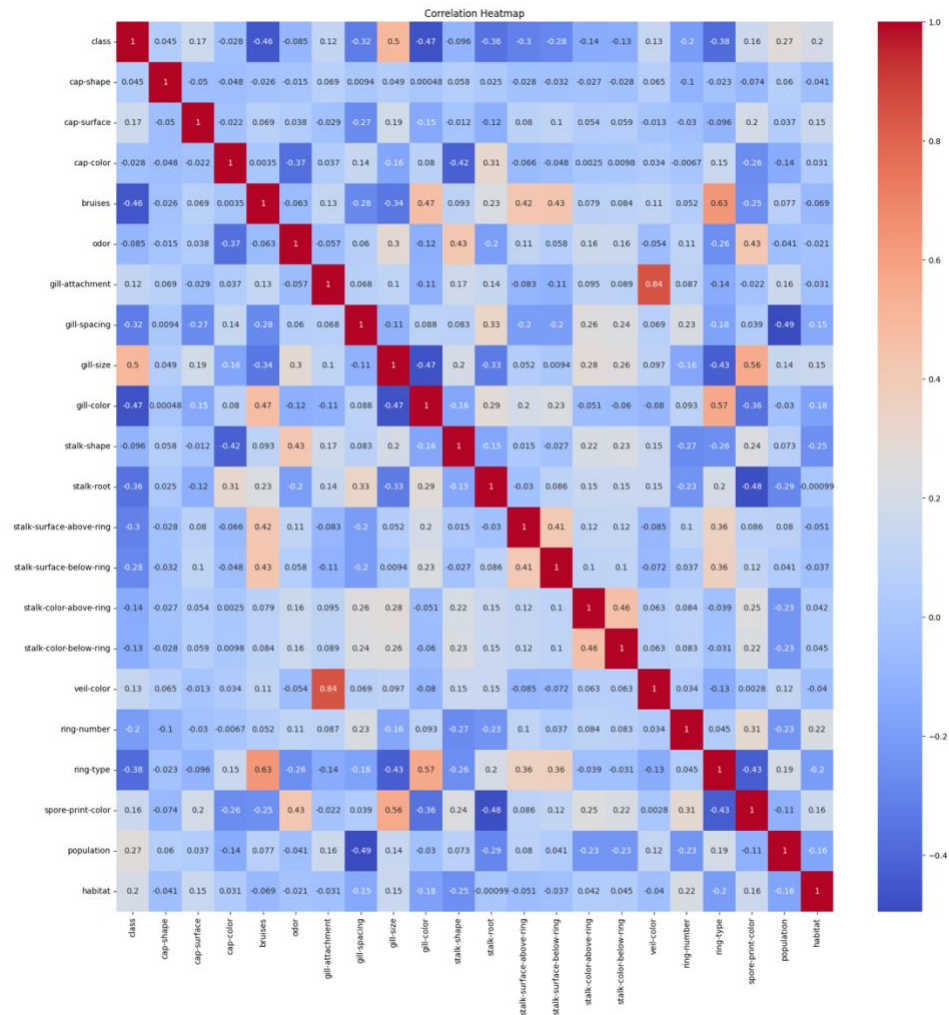  The dataset was originally contributed to the UCI Machine Learning Repository nearly 30 years ago.

  - Link:

    HTTPS://WWW.KAGGLE.COM/DATASETS/UCIML/MUSHROOM-CLASSIFICATION

  - Reference : The dataset is based on descriptions from "The Audubon Society Field Guide to North American Mushrooms" (1981).
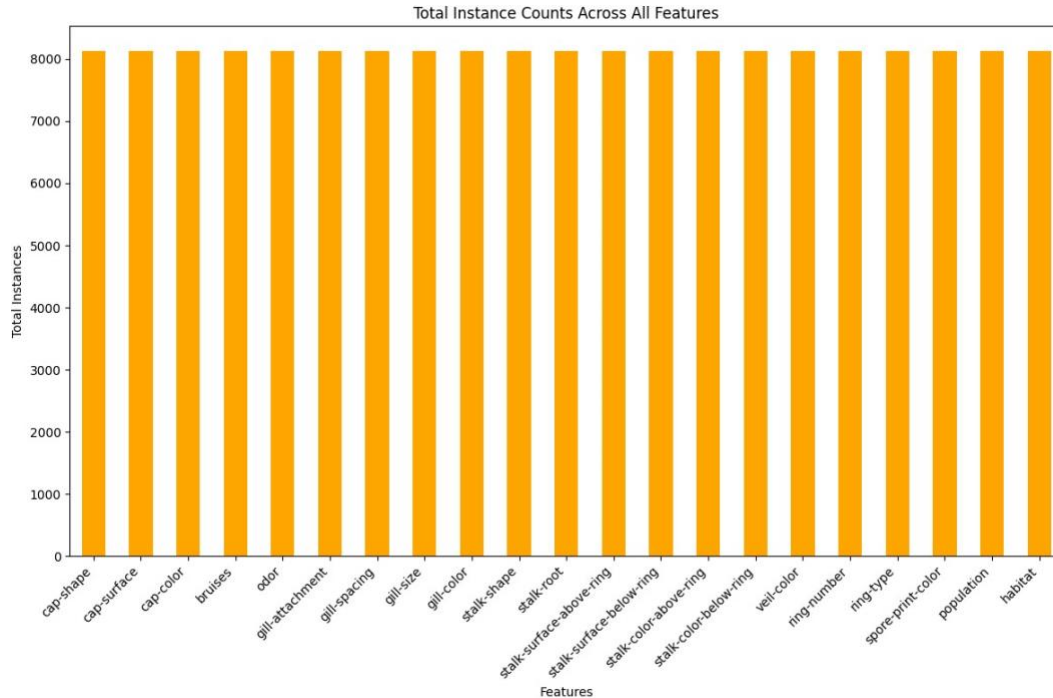
- DATASET DESCRIPTION

  o NUMBER OF FEATURES: 22 FEATURES, INCLUDING CAP SHAPE, ODOR, AND HABITAT.

  o PROBLEM TYPE: CLASSIFICATION PROBLEM, AS THE TARGET IS CATEGORICAL (EDIBLE VS. POISONOUS).

  o NUMBER OF DATA POINTS: 8,124 SAMPLES CORRESPONDING TO 23 SPECIES OF GILLED MUSHROOMS.

  o FEATURE TYPES: ALL FEATURES ARE CATEGORICAL AND WERE ENCODED INTO NUMERICAL VALUES USING LABELENCODER.

  o CORRELATION ANALYSIS: A HEATMAP WAS GENERATED TO VISUALIZE CORRELATIONS, HIGHLIGHTING HIGHLY CORRELATED FEATURES THAT MIGHT IMPACT THE MODEL.



Correlation Heatmap

- **BALANCED DATASET**

  FOR THE OUTPUT FEATURE, ALL THE UNIQUE CLASSES HAVE AN EQUAL NUMBER OF INSTANCES.



Total Instance Counts Across All Features

# DATASET PREPROCESSING

- **FAULTS**
  - NULL VALUES: IN THE ACTUAL DATASET THERE WAS NO NULL VALUES, SO WE RANDOMLY INSERTED 5 – 10% NULL VALUES.
  - CATEGORICAL VALUES: ALL FEATURES WERE CATEGORICAL.

- **SOLUTIONS**
  - DROPPED COLUMNS: THE 'VEIL-TYPE' COLUMN WAS REMOVED DUE TO IRRELEVANCE, AS IT CONTAINED ONLY ONE UNIQUE VALUE. HANDLED NULL VALUES AND FILL CATEGORICAL NULL VALUES WITH MODE . STILL IF A ROW CONTAINS ANY NULL VALUE, THAT ROW WILL BE DROPPED.

- Encoding: Categorical variables were encoded using LabelEncoder to convert categorical values into numeric representations.

# Feature Scaling

Feature scaling is essential for ensuring that the machine learning algorithms treat all features equally. The dataset was scaled using the MinMaxScaler from scikit-learn, which normalizes the values between 0 and 1. This step ensures the dataset is suitable for algorithms sensitive to magnitude differences in features.

# Dataset Splitting

The dataset was split into training and testing subsets to evaluate model performance:

- Splitting Technique: Stratified Random Sampling was used to ensure that the distribution of the target variable in the training and testing subsets matches the original dataset.
- Training Set: 70% of the data.
- Testing Set: 30% of the data.

The splitting ensured an unbiased evaluation of the model's performance on unseen data, while maintaining the representativeness of the target variable distribution. The random_state parameter was used to ensure reproducibility of the splits.

# Model Training & Testing

Multiple machine learning models were trained to classify mushrooms, including:

1. ## Decision Tree
   - **Description:**
     A decision tree is a supervised learning algorithm used for classification and regression tasks. It works by splitting the dataset into subsets based on the feature that provides the maximum information gain or minimizes impurity (e.g., Gini index or entropy). Each internal node represents a decision based on a feature, branches represent the outcomes, and leaf nodes represent the final prediction.
   - **How It Works:**
     - The algorithm starts with the root node and recursively splits the dataset based on features.
     - It continues splitting until a stopping criterion (e.g., maximum depth, minimum samples per leaf) is met.
     - The final tree structure is used to classify new data by traversing the nodes based on feature values.

2. ## Random Forest Classifier
   - **Description:**
     Random Forest is an ensemble learning method that builds multiple decision trees during training and aggregates their predictions (e.g., majority voting for

classification) to improve accuracy and reduce overfitting. It is robust to noise and works well on large datasets.

- How It Works:
  - Multiple decision trees are created using random subsets of the data and features (bagging).
  - Each tree provides a classification result, and the most frequent class is chosen as the final prediction.
  - The randomness in data and feature selection ensures diversity among trees, reducing overfitting.
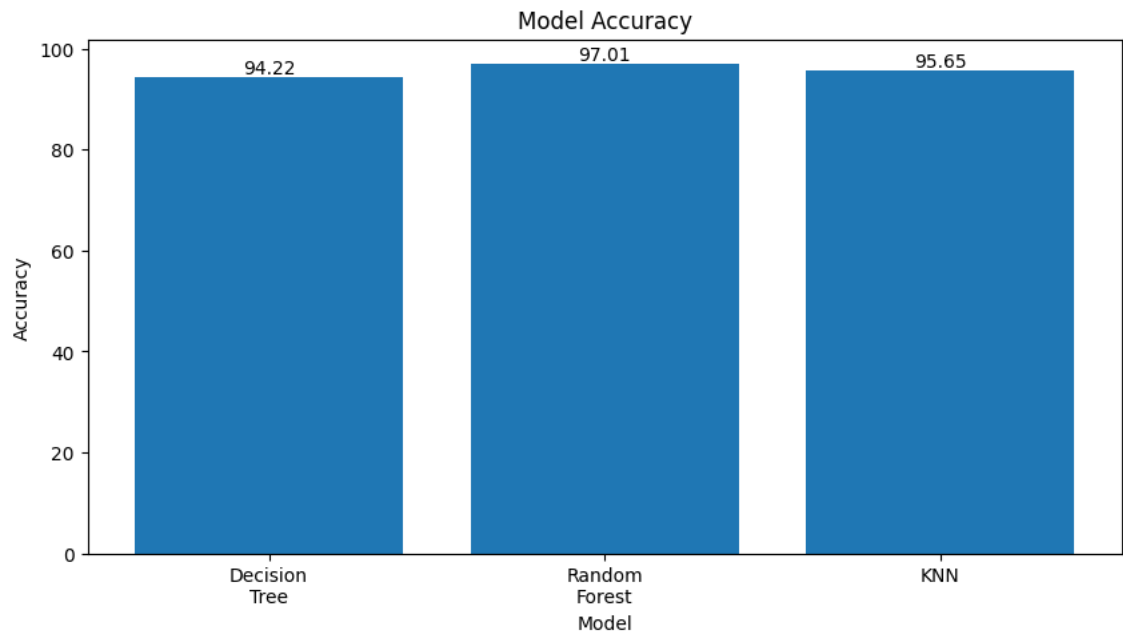
3. K-Nearest Neighbors (KNN) Classifier

- Description:

  KNN is a simple, non-parametric algorithm used for classification and regression. It makes predictions by finding the k-th nearest data points in the training set to a given test point, based on a distance metric (e.g., Euclidean distance).

- How It Works:
  - The algorithm calculates the distance between the test point and all training points.
  - It selects the k closest points and assigns the test point to the class most common among these neighbors.
  - The choice of k and the distance metric significantly affect the algorithm's performance.

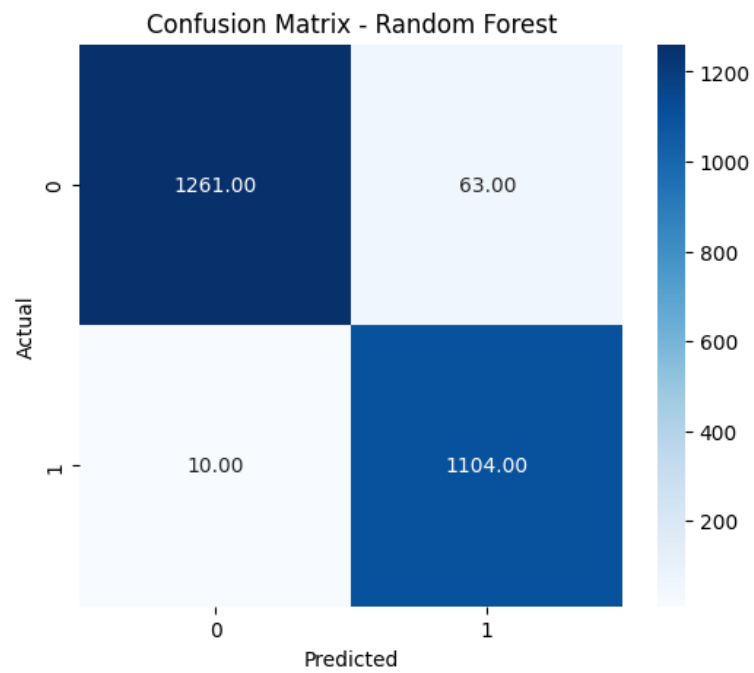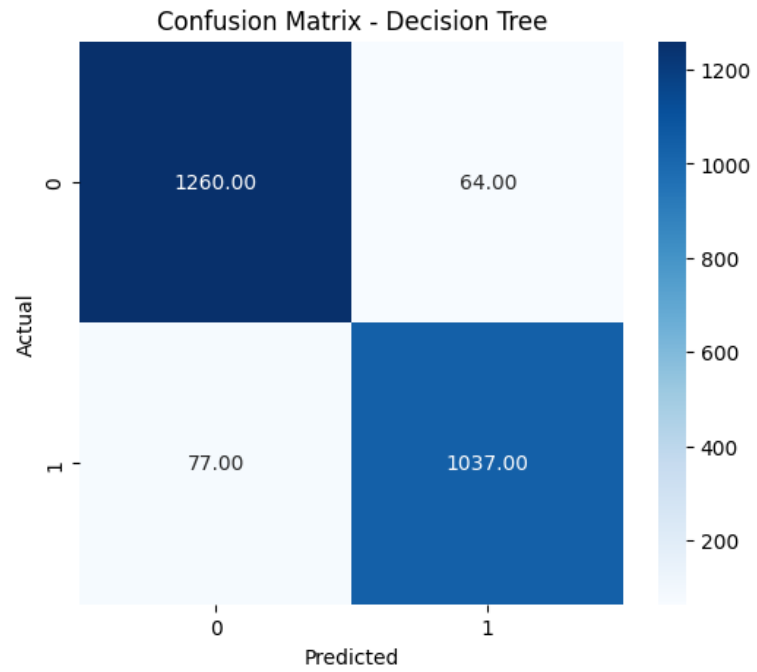# Model Selection/Comparison Analysis

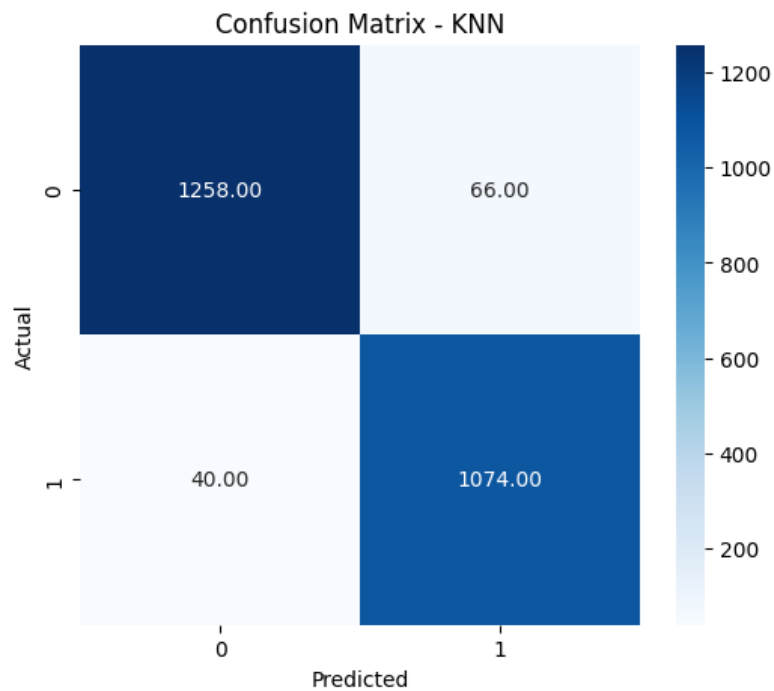- BAR CHART: SHOWCASED THE PREDICTION ACCURACY OF ALL MODELS.



- PRECISION AND RECALL:

  - PRECISION: THE PROPORTION OF CORRECTLY PREDICTED POSITIVE INSTANCES OUT OF ALL PREDICTED POSITIVES. IT MEASURES ACCURACY IN POSITIVE PREDICTIONS.

  - RECALL: THE PROPORTION OF CORRECTLY PREDICTED POSITIVE INSTANCES OUT OF ALL ACTUAL POSITIVES. IT MEASURES THE ABILITY TO find ALL POSITIVE INSTANCES.

| | Model | Precision | Recall |
|---|---|---|---|
| 0 | Decision Tree | 0.942163 | 0.942166 |
| 1 | Random Forest | 0.971060 | 0.970057 |
| 2 | KNN | 0.956811 | 0.956522 |

- CONFUSION MATRIX: ANALYZED FOR EACH MODEL TO EVALUATE classification PERFORMANCE IN DETAIL.

Confusion Matrix - Decision Tree



Confusion Matrix - Random Forest

## Confusion Matrix - KNN

|  | Predicted 0 | Predicted 1 |
|---|---|---|
| Actual 0 | 1258.00 | 66.00 |
| Actual 1 | 40.00 | 1074.00 |

# CONCLUSION

The Mushroom Class Detector project successfully used machine learning to classify mushrooms as edible or poisonous. Through proper preprocessing, scaling, and splitting of the dataset, the models—Decision Tree, Random Forest, and K-Nearest Neighbors—achieved 94.22%, 97.01%, and 95.65% accuracy respectively. Metrics like precision, recall, and confusion matrix confirmed their effectiveness. This project shows how machine learning can be a reliable tool to improve food safety and prevent the consumption of harmful mushrooms.