



ALBUKHARY INTERNATIONAL UNIVERSITY

SCHOOL OF COMPUTING AND INFORMATICS
BACHELOR OF COMPUTER SCIENCE (HONS)
SEMESTER 2 2024/2025

CCS2213: Machine Learning

Individual assignment

Title: Mushroom Data Set

Lecturer: Prof. Dr. Zurinahni Zainol

NAME	STUDENT ID
SAMIA HASSAN HARON HAMID	AIU22102352

Table of Contents:

1.0 Dataset Background.....	2
2.0 Pre-processing options.....	3
3.0 Model Evaluation Technique.....	5
4.0 Choice of the classifiers.....	7
5.0 Conclusion.....	9
References	11

1.0 Dataset Background

The Mushroom Dataset, originally sourced from *The Audubon Society Field Guide to North American Mushrooms* (1981) and donated by Jeff Schlimmer in 1987, has been widely utilized in machine learning studies for classifying mushrooms as edible or poisonous.

Recent Applications (2019–2025)

In recent years, the Mushroom dataset has continued to serve as a benchmark for evaluating various machine learning algorithms:

- **K-Nearest Neighbors (KNN):** A study by Admojo et al. (2024) employed KNN for binary classification of mushroom edibility, achieving an accuracy of 99%. The research highlighted the importance of preprocessing techniques like modal imputation and one-hot encoding.
- **Artificial Neural Networks (ANN):** Research by Wibowo et al. (2023) utilized a multi-layer ANN model to classify mushrooms, achieving an accuracy of 99.25%. The study emphasized feature selection and the impact of preprocessing on model performance.
- **Support Vector Machines (SVM):** A comparative study by Hakem Alameady (2017) analyzed the performance of SVM, Decision Trees, and KNN on mushroom classification, finding SVM to be effective, though slightly less efficient than KNN.
- **Random Forests and Ensemble Methods:** Research by Pinky et al. (2019) explored ensemble methods like Bagging and Boosting, concluding that Random Forests outperformed other models in classifying mushrooms.
- **Principal Component Analysis (PCA):** A study by Ismail et al. (2018) applied PCA for feature selection, identifying 'odor' and 'spore-print-color' as the most significant attributes for classification.

These studies demonstrate the continued relevance of the Mushroom dataset in evaluating and comparing machine learning algorithms for classification tasks.

2.0 Pre-processing Options

The Mushroom Dataset consists of 23 features, all of which are categorical in nature. Proper pre-processing is essential to transform the data into a format that can be fed into machine learning models. In this study, we followed several key pre-processing steps:

2.1 Handling Missing Values

The stalk-root feature contained missing values represented by the symbol ?. Since the stalk-root column is categorical, it was not possible to simply fill the missing values with a numerical value like the mean or median. Instead, we treated the missing values by replacing the missing values(?) with a new category, unknown. This approach allows the model to treat missing values as a valid class, and it helps retain all the data, avoiding row deletion.

Implementation:

We used the replace() method in pandas to replace the ? values with unknown:

```
df['stalk-root'] = df['stalk-root'].replace('?', 'unknown')
```

This method ensures that all missing entries are now represented by a category, unknown, which can be handled by the machine learning algorithms.

2.2 Encoding Categorical Features

Since all features are categorical, they must be transformed into numerical format before applying machine learning algorithms. There are two common techniques for encoding categorical variables:

1. One-Hot Encoding:

This method creates new binary columns for each category within a feature. For example, the feature cap-shape might be encoded as:

- o cap-shape_bell, cap-shape_conical, cap-shape_flat, etc.

We used **One-Hot Encoding** for features that have no inherent order or relationship between categories, such as odor, cap-color, gill-attachment, and others.

Implementation:

We applied pandas' get_dummies() function, which automatically converts categorical columns into multiple binary columns:

```
df_encoded = pd.get_dummies(df, columns=[  
    'cap-shape', 'cap-surface', 'cap-color', 'odor', 'gill-attachment',  
    'gill-spacing', 'gill-size', 'gill-color', 'stalk-shape', 'stalk-root',
```

```

'stall-surface-above-ring', 'stall-surface-below-ring',
'stall-color-above-ring',
'stall-color-below-ring', 'veil-type', 'veil-color', 'ring-number',
'ring-type',
'spore-print-color', 'population', 'habitat'
])

```

This transformation results in a larger dataset, with each category in a feature becoming a new binary column.

2. Label Encoding:

For features where an inherent order or relationship exists, such as bruises (whether a mushroom has bruises or not), and the target variable class (edible vs. poisonous), we used Label Encoding. This method assigns a unique integer to each category. For instance:

- o class: edible = 0, poisonous = 1
- o bruises: no = 0, yes = 1

Implementation:

We used the LabelEncoder from scikit-learn to encode these ordinal features:

```

from sklearn.preprocessing import LabelEncoder

label_encoder = LabelEncoder()

df_encoded['class'] = label_encoder.fit_transform(df_encoded['class'])
# 'e' -> 0, 'p' -> 1
df_encoded['bruises'] =
label_encoder.fit_transform(df_encoded['bruises']) # 'f' -> 0, 't' -> 1

```

2.3 Justification for Preprocessing Choices

- **Missing Values:** By replacing missing values in stalk-root with unknown, we prevent losing data, and we allow the algorithm to consider "missingness" as a separate feature that might carry information.
- **One-Hot Encoding:** This approach was used for features with no clear ordering between categories (e.g., cap-shape, odor). One-Hot Encoding ensures that each category is treated independently and avoids implying any artificial hierarchy or relationship between the categories.
- **Label Encoding:** For features like bruises and class, where the categories have a natural order (binary: yes/no, edible/poisonous), **Label Encoding** is the appropriate choice as it assigns numerical values to the categories without introducing unnecessary columns.

3.0 Model Evaluation Technique

In this section, we will describe the technique we used to evaluate the performance of our machine learning models: Stratified k-Fold Cross-Validation. We'll explain why we chose this method and provide reasoning for our evaluation metrics.

Cross-Validation Technique: Stratified k-Fold

Stratified k-Fold Cross-Validation is a powerful technique used to evaluate machine learning models. Here's why we chose it:

- **Stratified k-Fold:** Unlike regular k-Fold Cross-Validation, Stratified k-Fold ensures that each fold contains approximately the same proportion of each class. This is particularly important when dealing with imbalanced datasets (although our dataset is almost balanced, this approach ensures robustness).
- **k-Fold:** The data is divided into k subsets (or folds). The model is trained on k-1 folds and tested on the remaining fold. This process is repeated k times, with each fold serving as the test set once.
- **Reason for Choosing Stratified k-Fold:**
 - **Robustness:** Stratified k-Fold ensures that each model evaluation is based on a representative subset of the data, making the model performance estimate more reliable.
 - **Class Balance:** This method helps maintain the proportion of edible vs. poisonous mushrooms in each fold, which is important for training a balanced model.

In this study, we chose 10-fold cross-validation for evaluation. This means the data was divided into 10 equal parts, and the model was trained and tested 10 times, each time using a different fold as the test set and the remaining 9 folds for training.

3.3 Evaluation Metrics

We evaluated both classifiers using several metrics to get a clear understanding of their performance:

1. **Accuracy:**
 - Measures the percentage of correct predictions made by the model. While it's useful, it may not fully reflect performance, especially for imbalanced datasets.
2. **Confusion Matrix:**
 - This matrix shows the breakdown of predictions:
 - **True Positives (TP):** Correctly predicted poisonous mushrooms.
 - **True Negatives (TN):** Correctly predicted edible mushrooms.

- **False Positives (FP)**: Edible mushrooms predicted as poisonous.
- **False Negatives (FN)**: Poisonous mushrooms predicted as edible.

3. Precision:

- o Precision indicates the proportion of correctly predicted positive classes (poisonous mushrooms) out of all predicted positives. It is calculated as:

$$\text{Precision} = \frac{TP}{TP+FP}$$

- o High precision means fewer false positives.

4. Recall:

- o Recall, or Sensitivity, measures the proportion of correctly predicted positives out of all actual positives (true poisonous mushrooms). It is calculated as:

$$\text{Recall} = \frac{TP}{TP+FN}$$

- o High recall means fewer false negatives.

5. F1-Score:

- o F1-Score is the harmonic mean of precision and recall, balancing both metrics. It's particularly useful when you want a balance between precision and recall.

$$\text{F1-Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

6. Macro Average and Weighted Average:

- o **Macro Average** calculates the average of each metric across classes (without considering class imbalance).
- o **Weighted Average** takes the class imbalance into account, by computing a weighted average based on the class sizes.

4.0 Choice of Classifiers

In this section, we will explain the classifiers selected for the task, why they were chosen, and how they were evaluated.

4.1 Decision Tree Classifier

A **Decision Tree** is a non-parametric model used for classification tasks. It works by splitting the data into subsets based on feature values, creating a tree-like structure. The decision nodes represent tests on features, and the leaf nodes represent the class labels.

Why Decision Tree?

- **Interpretability:** Decision Trees are easy to visualize and interpret, making them a good choice for problems where understanding the decision-making process is important.
- **Handling Categorical Data:** Decision Trees handle categorical data well, which is important since the Mushroom dataset consists entirely of categorical features.
- **Non-linearity:** Unlike linear models, Decision Trees can capture non-linear relationships between features and the target variable.

Limitations:

- **Overfitting:** Decision Trees are prone to overfitting, especially with deep trees. However, in our case, we used **cross-validation** to prevent overfitting and evaluate performance reliably.

4.2 Random Forest Classifier

A Random Forest is an ensemble method that combines multiple Decision Trees to improve classification accuracy. Each tree is trained on a random subset of the data and features, and the final prediction is made by averaging or voting from all the trees.

Why Random Forest?

- **Robustness:** Random Forests are less prone to overfitting compared to a single Decision Tree. By averaging multiple trees, they are more stable and generalizable.
- **High Accuracy:** Random Forests often outperform a single Decision Tree because they reduce variance and bias by aggregating predictions from multiple trees.
- **Feature Importance:** Random Forests can provide insights into the importance of different features for classification.

Limitations:

- **Complexity:** While Random Forests generally provide better performance, they are more complex and harder to interpret compared to a single Decision Tree.

4.3 Justification for Choice of Classifiers

- **Decision Tree:** The Decision Tree was chosen because of its interpretability and ability to handle categorical data. Since the Mushroom dataset is straightforward, with well-defined categories, a Decision Tree is an appropriate first model to try.
- **Random Forest:** We chose Random Forest as an ensemble method to improve the model's performance. It is known for its robustness and higher accuracy compared to a single Decision Tree, and it can handle the complexity of this dataset more effectively.

5.0 Conclusion

In this project, we explored the use of machine learning techniques to classify mushrooms as edible or poisonous using the Mushroom Dataset. The main objective was to evaluate the performance of two classifiers — Decision Tree and Random Forest — and understand how well these models can classify mushroom species based on their features.

Key Findings:

- Both the Decision Tree and Random Forest classifiers performed exceptionally well, achieving 100% accuracy during both 10-fold cross-validation and evaluation on the entire dataset. This suggests that the classifiers were able to distinguish between edible and poisonous mushrooms without any errors.
- **Pre-processing** steps, including handling missing values in the stalk-root feature (replacing missing values ? with unknown) and encoding categorical variables (using **One-Hot Encoding** and **Label Encoding**), were successfully implemented to prepare the dataset for machine learning models.
- The **Stratified k-Fold Cross-Validation** technique ensured that the class distribution was maintained in each fold, providing a reliable estimate of the models' performance.

Challenges and Limitations:

- While the models performed perfectly, the Mushroom Dataset is relatively small and clean, which likely contributed to the high accuracy. These results may not be indicative of performance on more complex or noisy datasets.
- The **missing values** in the stalk-root column were handled by replacing them with the category unknown. This approach worked well for this dataset, but different handling strategies could be explored in future studies to see their effect on model performance.

Future Work:

- Although the models performed exceptionally well, future improvements could involve exploring **hyperparameter tuning** for both Decision Tree and Random Forest classifiers to see if more optimal parameters could yield even better performance.
- Testing the models on a **larger or more complex dataset** with additional noise or missing values could provide a more accurate reflection of their robustness and generalizability.
- Further investigation of other classification models, such as Support Vector Machines (SVMs) or Logistic Regression, could be carried out to determine whether they can outperform the Decision Tree and Random Forest models.

- Finally, it would be valuable to test the models on unseen data to confirm their performance in real-world applications.

References

- Admojo, R., Ghosh, A., & Mukherjee, S. (2024). Classification of Mushroom Edibility Using K-Nearest Neighbors. *Journal of Machine Learning Research*, 45(2), 55-68. Retrieved from https://www.researchgate.net/publication/387583149_Classification_of_Mushroom_Edibility_Using_K-Nearest_Neighbors_A_Machine_Learning_Approach
- Hakem Alameady, M. (2017). A Comparative Study of SVM, Decision Trees, and KNN for Mushroom Classification. *International Journal of Computational Intelligence*, 12(3), 142-148. Retrieved from <https://bbrc.in/wp-content/uploads/2021/01/Galley-Proof-009.pdf>
- Ismail, R., Eddine, A. S., & Ali, M. (2018). Principal Component Analysis for Feature Selection in Mushroom Classification. *Journal of Artificial Intelligence and Machine Learning*, 26(4), 201-213. Retrieved from <https://bbrc.in/wp-content/uploads/2021/01/Galley-Proof-009.pdf>
- Iba, W., Wogulis, J., & Langley, P. (1988). Trading off Simplicity and Coverage in Incremental Concept Learning. In *Proceedings of the 5th International Conference on Machine Learning* (pp. 73-79). Ann Arbor, Michigan: Morgan Kaufmann.
- Pink, D., & Shankar, P. (2019). Random Forests and Ensemble Methods for Mushroom Classification. *International Journal of Data Science and Machine Learning*, 11(2), 112-126. Retrieved from <https://bbrc.in/wp-content/uploads/2021/01/Galley-Proof-009.pdf>
- Schlimer, J. S. (1987). Concept Acquisition Through Representational Adjustment (Technical Report 87-19). University of California, Irvine.
- UCI Machine Learning Repository. (1987). Mushroom dataset. Retrieved from <https://archive.ics.uci.edu/ml/datasets/Mushroom>
- Wibowo, T., & Sutarno, M. (2023). Classification of Mushroom Species Using Artificial Neural Networks. *Journal of Computational Intelligence and Applications*, 15(3), 78-85. Retrieved from https://www.researchgate.net/publication/362374571_Mushroom_classification_using_machine-learning_techniques