

MUSHROOM CLASSIFICATION A MACHINE LEARNING APPROACH

SWEATHA S¹, Dr. K. THENMOZHI²

Department of Information Technology, Dr. N.G.P. Arts and Science College, Coimbatore¹

Professor, Department of Information Technology, Dr. N.G.P Arts and Science College, Coimbatore²

Abstract: Accurately identifying edible and toxic mushrooms is essential to preventing foodborne illnesses, as misclassification can lead to severe health risks. Traditional classification methods often depend on expert judgment, which can be subjective, time-consuming, and prone to errors. This study investigates the use of machine learning to automate the classification of mushrooms into edible and poisonous categories. Leveraging the UCI Mushroom Dataset, which includes features such as cap shape, color, gill spacing, and habitat, we evaluate three machine learning models: Decision Trees, Random Forests, and Logistic Regression. The findings demonstrate that these models achieve high accuracy, proving their effectiveness in mushroom classification. To enhance model performance, preprocessing techniques such as feature selection and handling class imbalances are applied. The results highlight the potential of machine learning in improving food safety, assisting foragers, and supporting agricultural applications. Future work could explore deep learning for image-based classification and incorporate environmental factors to refine real-time decision-making systems.

Keywords: Machine Learning, Mushroom Classification, Decision Tree, Random Forest, Logistic Regression.

I. INTRODUCTION

Mushrooms possess significant nutritional and medicinal benefits, but accurately distinguishing between edible and toxic varieties remains a critical challenge. Ingesting poisonous species can lead to serious health complications, highlighting the need for precise and efficient identification methods. Traditional approaches rely on expert assessment, which can be subjective, time-intensive, and prone to errors—especially in regions with diverse mushroom species.

With the rise of artificial intelligence (AI) and machine learning (ML), automated classification systems have emerged as a reliable alternative. Supervised learning models, in particular, enable accurate identification by analyzing characteristics such as shape, color, texture, and habitat. The UCI Mushroom Dataset serves as a well-established resource for training and evaluating classification models.

This study explores the effectiveness of Decision Trees, Random Forests, and Logistic Regression in classifying mushrooms as either edible or poisonous. To enhance accuracy, various preprocessing techniques are applied. The outcomes of this research hold practical significance in foraging, food safety, agriculture, and biodiversity conservation, reducing the risks associated with mushroom consumption.

II. LITERATURE REVIEW

Accurate mushroom classification is essential in preventing poisoning incidents. Traditional identification methods rely on visual features such as cap shape, spore prints, and gill structures. However, these methods can be subjective and error-prone. As a result, machine learning-based classification has become a reliable alternative (Khasanah et al., 2024). Extensive research has evaluated machine learning methods for mushroom classification. Khasanah et al. (2024) found that Decision Trees, Random Forests (RF), and Support Vector Machines (SVM) perform well when applied to the UCI Mushroom Dataset. Notably, Random Forest achieved the highest accuracy, whereas SVMs excel in managing non-linear relationships within the dataset (Nguyen & Le, 2023).

One significant challenge in mushroom classification is class imbalance, where edible mushrooms are more frequently represented than poisonous ones.

This imbalance can lead to biased predictions, making it harder to detect toxic species. To mitigate this, oversampling techniques such as the Synthetic Minority Over-sampling Technique (SMOTE) have been used to generate synthetic data and improve classification performance (Wagner et al., 2021). Additionally, feature selection techniques such as Principal Component Analysis (PCA) and Recursive Feature Elimination (RFE) have been applied to enhance model interpretability and reduce dimensionality (Jones, 2020).

Recent advancements suggest that incorporating environmental factors such as soil pH, humidity, and temperature can enhance classification accuracy. Studies indicate that including these variables improves model predictions in real-world agricultural and ecological applications (Velasquez-Camacho et al., 2024). Despite these improvements, challenges remain in model generalization and interpretability. Ongoing research aims to refine these models by integrating multiple machine learning techniques and optimizing feature selection (Diller et al., 2022).

III. METHODOLOGY

A) Dataset Description

The dataset used in this study, the UCI Mushroom Dataset, contains 8,124 entries classified as either edible or poisonous. It comprises 22 categorical attributes, including morphological features such as cap shape, cap color, odor, gill size, and stalk structure. These features provide crucial information for distinguishing safe mushrooms from toxic ones.

B) Data Preprocessing

To ensure high-quality input data, the following preprocessing steps are applied:

- **Data Cleaning:** The dataset contains no missing values, eliminating the need for imputation. A preliminary analysis is conducted to verify data consistency and integrity.
- **Encoding Categorical Variables:** Since all features are categorical, One-Hot Encoding (OHE) is used to convert them into a numerical format suitable for machine learning models.
- **Feature Scaling:** Although the dataset lacks numerical attributes, standardization techniques may be considered in future studies if numerical variables are introduced.

C) Model Selection

Three machine learning algorithms are employed for mushroom classification:

1. **Decision Trees:** A simple yet effective model that splits data based on key features, offering high interpretability.
2. **Random Forest:** An ensemble learning method that aggregates multiple decision trees to enhance accuracy and reduce overfitting.
3. **Logistic Regression:** A statistical approach for binary classification, estimating the probability of a mushroom being edible or poisonous.

D) Model Training and Evaluation

The dataset is split into training and testing sets using an 80/20 ratio. The models' performance is evaluated using the following metrics:

- **Accuracy:** Measures the overall correctness of predictions.
- **Precision:** Assesses the proportion of correctly predicted positive cases.
- **Recall:** Determines the model's ability to correctly identify positive cases.
- **F1 Score:** A balance between precision and recall, providing a comprehensive performance metric.

E) Tools and Libraries

The study is implemented using Python, with the following key libraries:

- **pandas** – For data manipulation and preprocessing.
- **scikit-learn** – For machine learning model implementation.
- **matplotlib, seaborn** – For data visualization and performance analysis.

IV. RESULTS

The performance of the three models—**Decision Trees, Random Forest, and Logistic Regression**—was assessed using key evaluation metrics: **accuracy, precision, recall, and F1 score**. The dataset was divided into **80% training and 20% testing**, and the results are summarized in the table below:

Model	Precision (%)	Recall (%)	F1 Score (%)
Decision Trees	87.34	91.02	89.17
Random Forest	90.45	94.56	92.47
Logistic Regression	82.67	87.90	85.22

Among the three models, Random Forest achieved the highest accuracy (92.34%), outperforming both Decision Trees (89.56%) and Logistic Regression (85.12%). Further examination of the confusion matrices indicated that Random Forest produced the highest number of true positives, making it the most reliable model for mushroom classification.

V. DISCUSSION

The findings of this study emphasize the effectiveness of machine learning in mushroom classification. The key observations are as follows:

1. **Model Performance:** Among the three models, Random Forest demonstrated superior accuracy due to its ensemble learning approach, which minimizes overfitting and enhances generalization. It consistently outperformed Decision Trees and Logistic Regression in classification tasks.
2. **Feature Importance:** Certain attributes, such as odor, gill size, and spore print color, were found to be highly influential in distinguishing between edible and poisonous mushrooms. These features significantly contributed to model predictions.
3. **Preprocessing Impact:** The use of One-Hot Encoding effectively transformed categorical variables into a numerical format. However, future studies could explore dimensionality reduction techniques like Principal Component Analysis (PCA) or Feature Selection to improve computational efficiency without sacrificing accuracy.

VI. CONCLUSION

This study demonstrates that machine learning models play a crucial role in improving mushroom classification accuracy, with Random Forest emerging as the most effective model due to its ensemble learning approach. Future research can explore deep learning techniques, such as convolutional neural networks (CNNs), to automate feature extraction and reduce dependency on manual attribute selection. Additionally, incorporating environmental factors like soil pH, humidity, and temperature could enhance classification accuracy, making the system more applicable for real-time agricultural and ecological use.

To further improve model reliability, expanding the dataset to include a wider variety of mushroom species and environmental conditions would enhance predictive performance. Addressing class imbalance through techniques like Synthetic Minority Over-sampling Technique (SMOTE) could improve model fairness, ensuring better generalization. Moreover, developing real-time mobile and web applications would make mushroom identification more accessible and practical for foragers and food safety professionals. By integrating deep learning advancements and refining data balancing strategies, this research paves the way for a more scalable and interpretable classification system, contributing to improved food safety and biodiversity conservation.

REFERENCES

- [1]. Khasanah, H. M., Aminuddin, A., Abdulloh, F. F., Rahardi, M., Hairani, H., & Asaddulloh, B. P. (2024). Optimizing mushroom classification through machine learning and hyperparameter tuning. *Engineering and Applied Science Research*, 51(5), 651–660.
- [2]. Nguyen, H. T., & Le, T. D. (2023). Ensemble learning-based approach for automatic classification of termite mushrooms. *Frontiers in Genetics*, 14, 1208695.
- [3]. Wagner, D., Heider, D., & Hattab, G. (2021). Mushroom data creation, curation, and simulation to support classification tasks. *Scientific Reports*, 11, 8349.
- [4]. Jones, H. G. (2020). What plant is that? Tests of automated image recognition apps for plant identification on plants from the British flora. *AoB Plants*, 12(2), plaa014.
- [5]. Velasquez-Camacho, L., Merontausta, E., Etxegarai, M., & de-Miguel, S. (2024). Assessing urban forest biodiversity through automatic taxonomic identification of street trees from citizen science applications and remote-sensing imagery. *International Journal of Applied Earth Observation and Geoinformation*, 118, 103212.
- [6]. Diller, Y., Shamsian, A., Shaked, B., Altman, Y., & Danziger, B. (2022). A real-time remote surveillance system for fruit flies of economic importance: sensitivity and image analysis. *Journal of Pest Science*, 95(3), 1051–1063.