

Preliminary Report - Edibility of Mushroom Species

Ishaan Saxena Nikita Rajaneesh Swaraj Bhaduri Utkarsh Jain
isaxena@purdue.edu nrajanee@purdue.edu sbhadur@purdue.edu jain192@purdue.edu

November 16, 2018

1 Introduction to the Problem

1.1 Definition of the Problem

Given a dataset \mathcal{D} with $n = 8124$ samples where each sample represents a mushroom with features being the observations about the characteristics of the mushrooms such as odor, color, etc., we aim to test and compare various supervised learning models for the problem of classifying each sample into either poisonous or edible. Further, we will optimize the Hyperparameters of the model which initially performs the best on the dataset.¹

1.2 Data Description

We are given \mathcal{D} with $n = 8124$ samples wherein each sample has the following 22 features (excluding the class label).

- | | | |
|--------------------|------------------------------|------------------|
| 1. cap-shape | 9. stalk-surface-above-ring | 17. habitat |
| 2. cap-surface | 10. stalk-surface-below-ring | 18. gill-spacing |
| 3. cap-color | 11. stalk-color-above-ring | 19. gill-size |
| 4. bruises | 12. stalk-color-below-ring | 20. stalk-shape |
| 5. odor | 13. veil-type | 21. ring-number |
| 6. gill-attachment | 14. veil-color | 22. population |
| 7. gill-color | 15. ring-type | |
| 8. stalk-root | 16. spore-print-color | |

These features have been further enumerated in Appendix A.

1.3 Encoding the Data

Note that all the features in our dataset are categorical variables. As a result, to proceed with evaluation of model performance, we must first encode these variables into numerical/binary values.

¹This dataset can be found at <https://www.kaggle.com/uciml/mushroom-classification>.

Appendix A