

# **Binary Classification of Edible and Poisonous Mushrooms using Classical ML techniques**



**SUBMITTED BY:**

Priyanka Borwanker(BT19CSE018)

**UNDER THE GUIDANCE OF:**

Dr. Ravindra Keskar

Department of Computer Science and Engineering, VNIT Nagpur

1. **Problem Statement:** Using Machine Learning to determine whether a mushroom is edible to eat or poisonous based on certain features.

2.

Input Parameters(23)	Output Parameters(1)
cap-shape cap-surface cap-color bruises odor gill-attachment gill-spacing gill-color stalk-shape stalk-root stalk-surface-above-ring stalk-color-above-ring stalk-color-below-ring veil-type veil-color ring-number ring-type spore-print-color population habitat	class (edible/poisonous)

3. **Dataset Used:** [University of California, Irvine - Machine Learning - Mushroom Classification | Kaggle](#)

**a. Insights about data:**

- Dataset contains 8124 data samples with 22 features each, classified into edible and poisonous.
- Each data sample consists of class type and feature set (nominal values).

	count	unique	top	freq
<b>class</b>	8124	2	e	4208
<b>cap-shape</b>	8124	6	x	3656
<b>cap-surface</b>	8124	4	y	3244
<b>cap-color</b>	8124	10	n	2284
<b>bruises</b>	8124	2	f	4748
<b>odor</b>	8124	9	n	3528
<b>gill-attachment</b>	8124	2	f	7914
<b>gill-spacing</b>	8124	2	c	6812
<b>gill-size</b>	8124	2	b	5612
<b>gill-color</b>	8124	12	b	1728
<b>stalk-shape</b>	8124	2	t	4608
<b>stalk-root</b>	8124	5	b	3776
<b>stalk-surface-above-ring</b>	8124	4	s	5176
<b>stalk-surface-below-ring</b>	8124	4	s	4936
<b>stalk-color-above-ring</b>	8124	9	w	4464
<b>stalk-color-below-ring</b>	8124	9	w	4384
<b>veil-type</b>	8124	1	p	8124
<b>veil-color</b>	8124	4	w	7924
<b>ring-number</b>	8124	3	o	7488
<b>ring-type</b>	8124	5	p	3968
<b>spore-print-color</b>	8124	9	w	2388
<b>population</b>	8124	6	v	4040
<b>habitat</b>	8124	7	d	3148

#### 4. Proposed Solution:

##### a. Data Preprocessing:

- i. Remove any missing or null attributes from the dataset.
- ii. Dimensionality Reduction: Remove attributes with low variance for e.g. veil-type has only one value.
- iii. Encode labels to digits

##### b. Possible Algorithms for Classification:

- i. **SVMs:**
  1. Since this is a problem of binary classification, SVMs are an obvious choice.
  2. It is not sure whether the data is linearly separable or not, so we can try SVMs with different types of kernels (linear, rbf, etc)
- ii. **Decision Tree Learning:**
  1. Under the assumption that all the attributes are independent of each other, we can apply decision tree learning.
  2. Simple model which can use the whole dataset and all of the features
  3. Highly intuitive solution which can be easily visualized.
- iii. **Multi Layer Perceptrons:**
  1. Multilayer perceptrons are a good option for supervised learning when the relationship between the attributes and output classes is not apparent.
  2. A multilayer perceptron can model even non linear classification problems quite well.

##### c. Libraries Used:

- i. **Pandas:**
  1. Open source, easy to use, data management framework
- ii. **Scikit-learn**
  1. Simple and efficient tools for predictive data analysis
  2. Built on NumPy, SciPy, and matplotlib
  3. Open source, commercially usable

##### d. Implementation:

<https://colab.research.google.com/drive/1I5a6Dsk1G0Nop8ZdWbczVkWX4Q0pt3bX?usp=sharing>

- i. SVMs implemented with two types of kernels (rbf and linear)
  1. Linear - standard linear svm
  2. Rbf - radial basis function, can model non-linearity in data as well
- ii. MLP Classifier:
  1. ReLu activation function
  2. 1 hidden layer of 70 neurons(decreased to avoid overfitting)
  3. sgd optimizer
  4. Learning\_rate = 0.0005 (decreased to avoid overfitting)
- iii. Decision Tree:
  1. Pre Pruning of decision trees with maximum depth set to 4

Before Pruning:

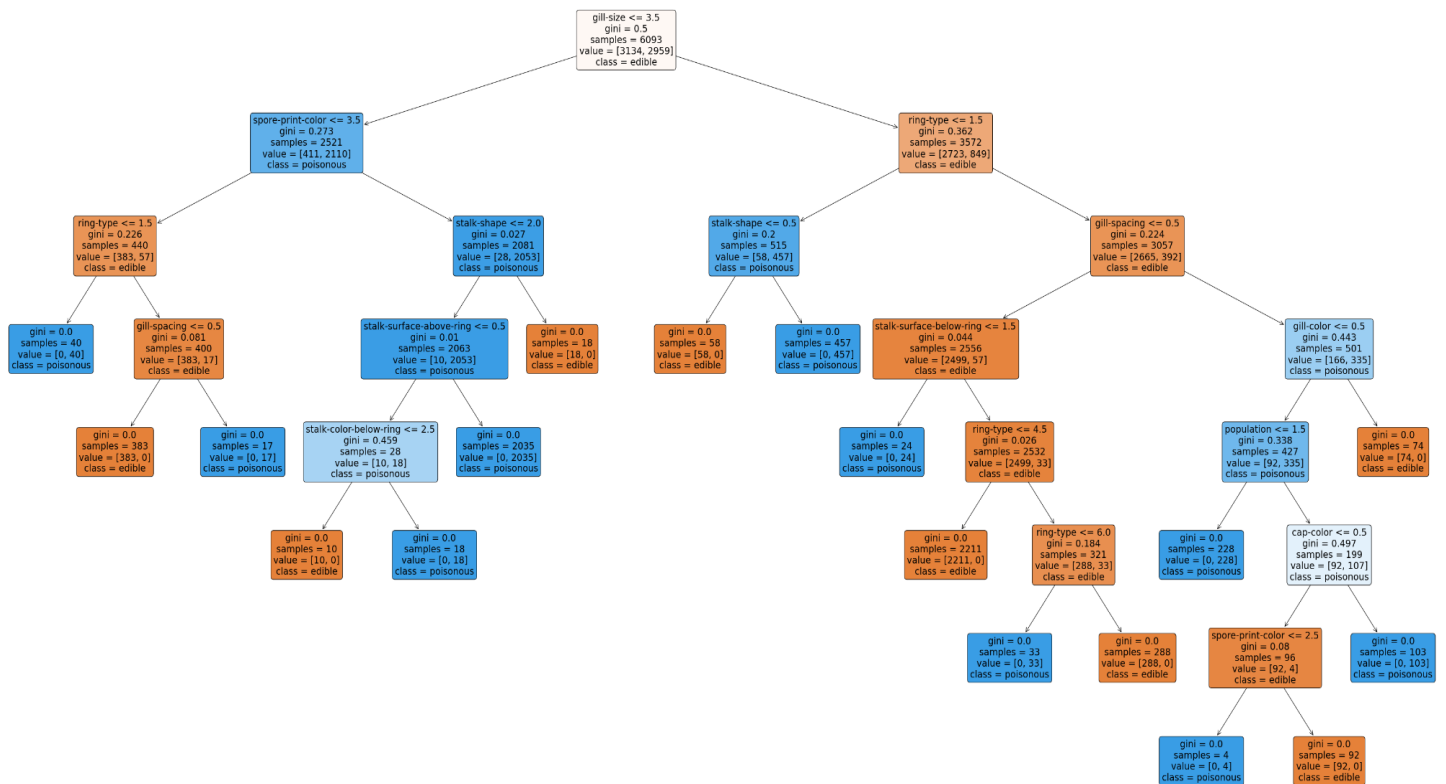


Fig: Highly irregular and too specific decision tree

After Pruning:

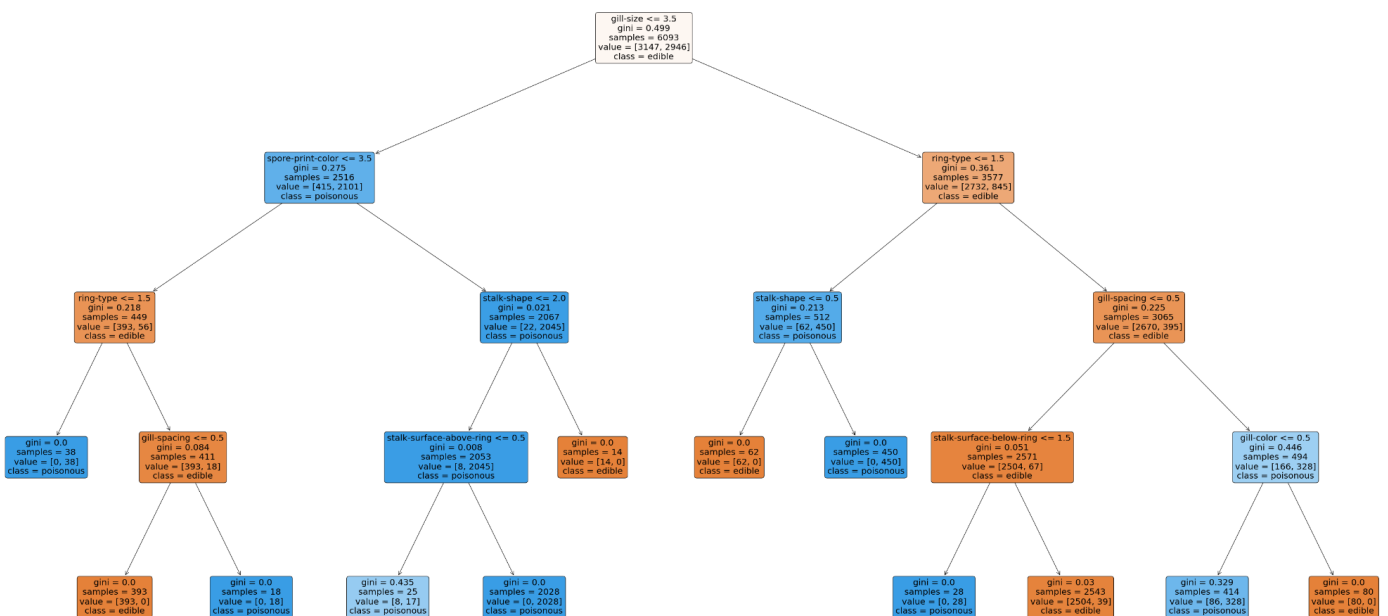


Fig: Balanced decision tree

**e. Accuracy Measure:**

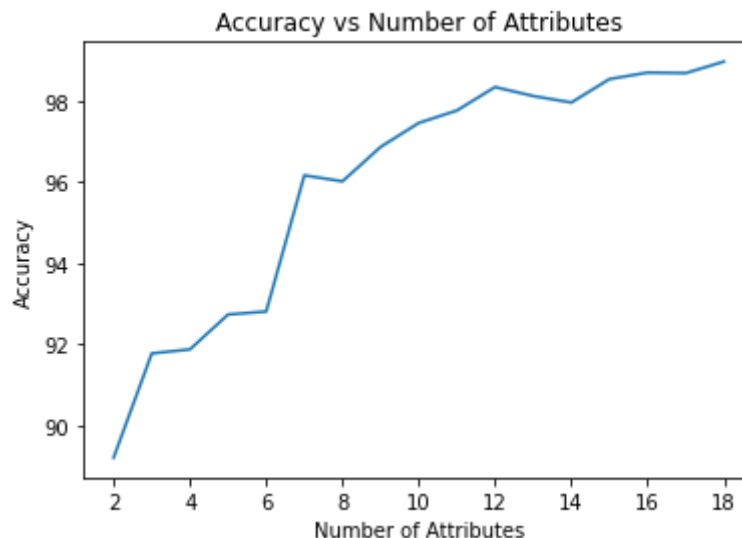
- i. We track when the predicted class and actual class is the same and include that in our score, else the score does not increase.
- ii. **K-Fold Cross Validation:** The dataset is divided into k parts and then k-1 parts are used for training and the kth part is used for validation. The model accuracy is then defined as the average of the validation accuracy of each part.
- iii. **Results:**

	Fold 1	Fold 2	Fold 3	Fold 4	Average Acc.
<b>SVM (linear kernel)</b>	95%	92%	94%	95%	94.04%
<b>SVM (rbf kernel)</b>	99%	99%	99%	99%	98.93%
<b>MLP Classifier</b>	97%	96%	98%	96%	96.80%
<b>DecisionTrees</b>	98%	98%	97%	98%	97.78%

→ *SVMs with the rbf kernel are giving the best results for binary classification, which is in accordance with the fact that SVMs are still considered the gold standard for binary classification.*

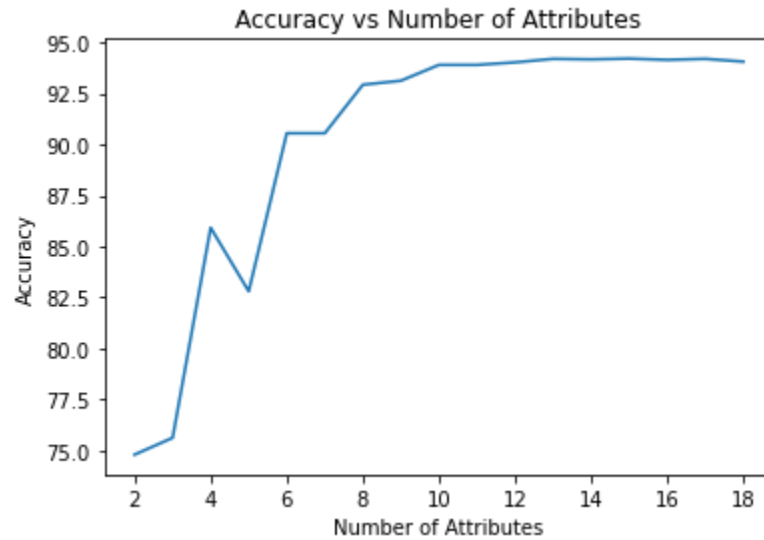
**f. Feature Selection: Can we achieve similar accuracy with less features?**

- i. **Chi-Square Analysis** - Using chi-square analysis we can pick the top most independent attributes and only use them for classification and check the accuracy. Here are some plots which show the variation of accuracy when we chose 2-18 top features.
  1. SVC (kernel = 'rbf')



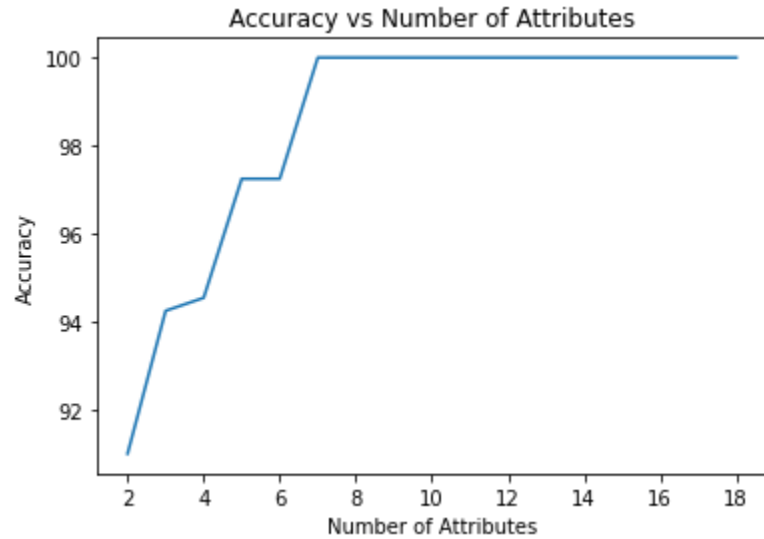
Here we can see that accuracy is increasing almost linearly with the number of features.

## 2. SVC (kernel = 'linear')



Here we can observe that accuracy is approximately constant after 10 topmost features.

## 3. Decision Trees



In this, accuracy is the same whether we use the topmost 7 features or all the 18 features. So it is better to use only top 7 features as others are just increasing the computational complexity.

**g. Advantages and Disadvantages**

	<b>Advantages</b>	<b>Disadvantages</b>
<b>SVMs</b>	<ul style="list-style-type: none"><li>• Computationally inexpensive</li><li>• More productive in high dimensional spaces</li><li>• Better when data is linearly separable with a clear margin</li></ul>	<ul style="list-style-type: none"><li>• Does not work well when relationship between attributes may not be linear (Our case)</li><li>• Does not work well in cases of noise or large dataset</li></ul>
<b>MLP</b>	<ul style="list-style-type: none"><li>• Can handle non-linearity in data well</li><li>• Quick predictions</li><li>• Can handle large amount of data easily</li></ul>	<ul style="list-style-type: none"><li>• Computationally more expensive than DecisionTrees and SVMs</li></ul>
<b>DecisionTrees</b>	<ul style="list-style-type: none"><li>• Less preprocessing of data is required</li><li>• Normalization/scaling not required</li><li>• Missing values has minimal impact on decision trees</li></ul>	<ul style="list-style-type: none"><li>• Noise can lead to unbalanced decision trees</li><li>• Training time is expensive but testing is fast</li><li>• Overfitting</li></ul>