## Mushroom Classification

Safe to eat or deadly poison?

**Project Title** : Mushroom Classification

**Technologies** : Machine Learning Technology

**Domain** : Agriculture

**Project Difficulties level** : Intermediate

## Problem Statement:

The Audubon Society Field Guide to North American Mushrooms contains descriptions of hypothetical samples corresponding to 23 species of gilled mushrooms in the Agaricus and Lepiota Family Mushroom (1981). Each species is labelled as either definitely edible, definitely poisonous, or maybe edible but not recommended. This last category was merged with the toxic category. The Guide asserts unequivocally that there is no simple rule for judging a mushroom's edibility, such as "leaflets three, leave it be" for Poisonous Oak and Ivy. The main goal is to predict which mushroom is poisonous & which is edible.

Approach: The classical machine learning tasks like Data Exploration, Data Cleaning, Feature Engineering, Model Building and Model Testing. Try out different machine learning algorithms that's best fit for the above case.

Results: You have to build a solution that should able to predict which mushroom is poisonous & which is edible.

## About Dataset

Context

Although this dataset was originally contributed to the UCI Machine Learning repository nearly 30 years ago, mushroom hunting (otherwise known as "shrooming") is enjoying new peaks in popularity. Learn which features spell certain death and which are most palatable in this dataset of mushroom characteristics. And how certain can your model be?

Content

This dataset includes descriptions of hypothetical samples corresponding to 23 species of gilled mushrooms in the Agaricus and Lepiota Family Mushroom drawn from The Audubon Society Field Guide to North American Mushrooms (1981). Each species is identified as definitely edible, definitely poisonous, or of unknown edibility and not recommended. This latter class was combined with the poisonous one. The Guide clearly states that there is no simple rule for determining the edibility of a mushroom; no rule like "leaflets three, let it be" for Poisonous Oak and Ivy.

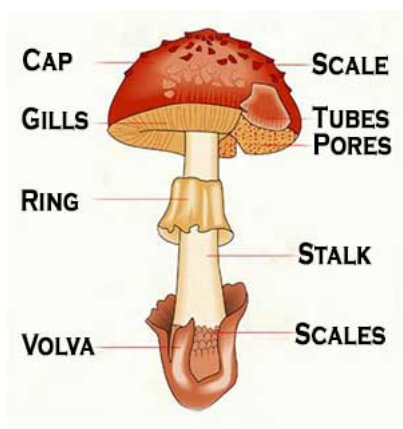Time period: Donated to UCI ML 27 April 1987

Data Link : https://www.kaggle.com/datasets/uciml/mushroom-classification

Detail about data set : https://www.nature.com/articles/s41598-021-87602-3
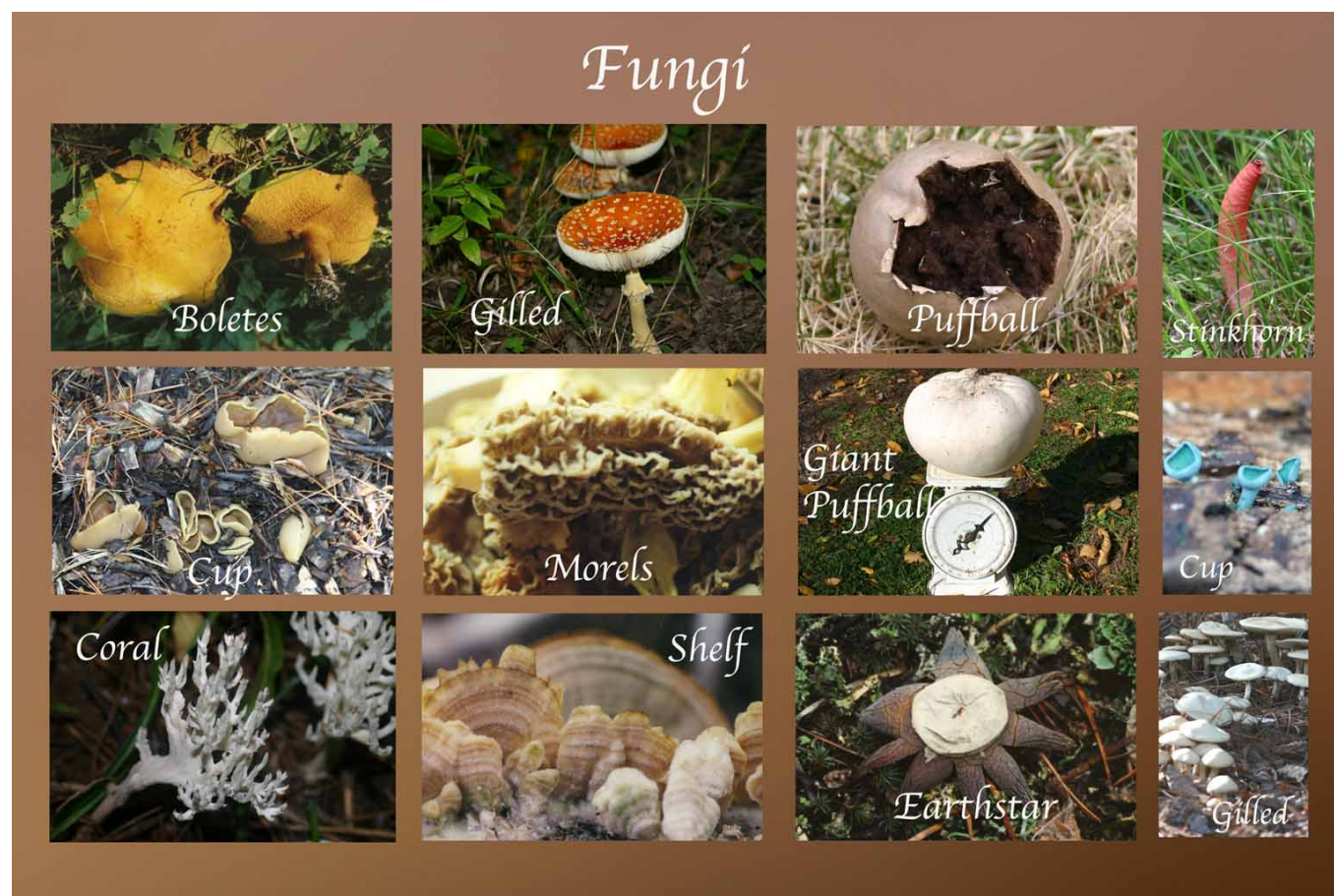
## Lets understand about data set and mushroom

Lets disscuss about the mushroom :

The term "gilled mushroom" refers to the mushrooms that have gills on the underside of their cap. These gills are thin, papery structures that radiate out from the stem and produce spores. The order of fungi that includes gilled mushrooms is called Agaricales

Mushroom Structure for understanding data set :

Mushrooms are a type of fungus that belong to the kingdom Fungi. They are characterized by their umbrella-shaped fruiting body, which is called a sporophore [1]. The most common classification of mushrooms is based on the structure of their fruiting body. The fruiting body of mushrooms can be classified into two main types: those with gills and those with pores [1].



The data set is given in csv file. Column and about data set as follows :

About this file you can also refer above mushroom stucture for better understanding data set.

The dataset used in this project contains 8124 instances of mushrooms with 23 features like cap-shape, cap-surface, cap-color, bruises, odor, etc.

Attribute Information:

(**classes**: edible=e, poisonous=p)

**cap-shape**: bell=b,conical=c,convex=x,flat=f, knobbed=k,sunken=s

**cap-surface**: fibrous=f,grooves=g,scaly=y,smooth=s

**cap-color**: brown=n,buff=b,cinnamon=c,gray=g,green=r,pink=p,purple=u,red=e,white=w,yellow=y

**bruises**: bruises=t,no=f

**odor**: almond=a,anise=l,creosote=c,fishy=y,foul=f,musty=m,none=n,pungent=p,spicy=s

**gill-attachment**: attached=a,descending=d,free=f,notched=n

**gill-spacing**: close=c,crowded=w,distant=d

**gill-size**: broad=b,narrow=n

**gill-color**: black=k,brown=n,buff=b,chocolate=h,gray=g, green=r,orange=o,pink=p,purple=u,red=e,white=w,yellow=y

**stalk-shape**: enlarging=e,tapering=t

**stalk-root**: bulbous=b,club=c,cup=u,equal=e,rhizomorphs=z,rooted=r,missing=?

**stalk-surface-above-ring**: fibrous=f,scaly=y,silky=k,smooth=s

**stalk-surface-below-ring**: fibrous=f,scaly=y,silky=k,smooth=s

**stalk-color-above-ring**: brown=n,buff=b,cinnamon=c,gray=g,orange=o,pink=p,red=e,white=w,yellow=y

**stalk-color-below-ring**: brown=n,buff=b,cinnamon=c,gray=g,orange=o,pink=p,red=e,white=w,yellow=y

**veil-type**: partial=p,universal=u

**veil-color**: brown=n,orange=o,white=w,yellow=y

**ring-number** : none=n,one=o,two=t

**ring-type** : cobwebby=c,evanescent=e,flaring=f,large=l,none=n,pendant=p,sheathing=s,zone=z

**spore-print-color** : black=k,brown=n,buff=b,chocolate=h,green=r,orange=o,purple=u,white=w,yellow=y

**population**: abundant=a,clustered=c,numerous=n,scattered=s,several=v,solitary=y

**habitat**: grasses=g,leaves=l,meadows=m,paths=p,urban=u,waste=w,woods=d

```
#Import libraries
import pandas as pd
```

**Loding the dataset**

```
# read data
df = pd.read_csv('/content/mushrooms.csv')
```

**Examin the data**

```
df.head()
```

| | class | cap-shape | cap-surface | cap-color | bruises | odor | gill-attachment | gill-spacing | gill-size | gill-color | ... | stalk-surface-below-ring | stalk-color-above-rin |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | p | x | s | n | t | p | f | c | n | k | ... | s | |
| 1 | e | x | s | y | t | a | f | c | b | k | ... | s | |
| 2 | e | b | s | w | t | l | f | c | b | n | ... | s | |
| 3 | p | x | y | w | t | p | f | c | n | n | ... | s | |
| 4 | e | x | s | g | f | n | f | w | b | k | ... | s | |

5 rows × 23 columns

```
df.all()
```

```
    class                     True
    cap-shape                 True
    cap-surface               True
    cap-color                 True
    bruises                   True
    odor                      True
    gill-attachment           True
    gill-spacing              True
    gill-size                 True
    gill-color                True
    stalk-shape               True
    stalk-root                True
    stalk-surface-above-ring  True
    stalk-surface-below-ring  True
```

```
stalk-color-above-ring     True
stalk-color-below-ring     True
veil-type                  True
veil-color                 True
ring-number                True
ring-type                  True
spore-print-color          True
population                 True
habitat                    True
dtype: bool
```

df.tail()

| | class | cap-shape | cap-surface | cap-color | bruises | odor | gill-attachment | gill-spacing | gill-size | gill-color | ... | stalk-surface-below-ring | stalk-color-above-ri |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 8119 | e | k | s | n | f | n | a | c | b | y | ... | s | |
| 8120 | e | x | s | n | f | n | a | c | b | y | ... | s | |
| 8121 | e | f | s | n | f | n | a | c | b | n | ... | s | |
| 8122 | p | k | y | n | f | y | f | c | n | b | ... | k | |
| 8123 | e | x | s | n | f | n | a | c | b | y | ... | s | |

5 rows × 23 columns

df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 8124 entries, 0 to 8123
Data columns (total 23 columns):
 #   Column                    Non-Null Count  Dtype
---  ------                    --------------  -----
 0   class                     8124 non-null   object
 1   cap-shape                 8124 non-null   object
 2   cap-surface               8124 non-null   object
 3   cap-color                 8124 non-null   object
 4   bruises                   8124 non-null   object
 5   odor                      8124 non-null   object
 6   gill-attachment           8124 non-null   object
 7   gill-spacing              8124 non-null   object
 8   gill-size                 8124 non-null   object
 9   gill-color                8124 non-null   object
 10  stalk-shape               8124 non-null   object
 11  stalk-root                8124 non-null   object
 12  stalk-surface-above-ring  8124 non-null   object
 13  stalk-surface-below-ring  8124 non-null   object
 14  stalk-color-above-ring    8124 non-null   object
 15  stalk-color-below-ring    8124 non-null   object
 16  veil-type                 8124 non-null   object
 17  veil-color                8124 non-null   object
 18  ring-number               8124 non-null   object
 19  ring-type                 8124 non-null   object
 20  spore-print-color         8124 non-null   object
 21  population                8124 non-null   object
 22  habitat                   8124 non-null   object
dtypes: object(23)
memory usage: 1.4+ MB
```

df.describe()

| | class | cap-shape | cap-surface | cap-color | bruises | odor | gill-attachment | gill-spacing | gill-size | gill-color | ... | stalk-surface-below-ring | stalk-col-above-r |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| count | 8124 | 8124 | 8124 | 8124 | 8124 | 8124 | 8124 | 8124 | 8124 | 8124 | ... | 8124 | 8 |
| unique | 2 | 6 | 4 | 10 | 2 | 9 | 2 | 2 | 2 | 12 | ... | 4 | |
| top | e | x | y | n | f | n | f | c | b | b | ... | s | |
| freq | 4208 | 3656 | 3244 | 2284 | 4748 | 3528 | 7914 | 6812 | 5612 | 1728 | ... | 4936 | 4 |

4 rows × 23 columns

#define target (Y) and features (X)

```
df.columns
```

Index(['class', 'cap-shape', 'cap-surface', 'cap-color', 'bruises', 'odor',
       'gill-attachment', 'gill-spacing', 'gill-size', 'gill-color',
       'stalk-shape', 'stalk-root', 'stalk-surface-above-ring',
       'stalk-surface-below-ring', 'stalk-color-above-ring',
       'stalk-color-below-ring', 'veil-type', 'veil-color', 'ring-number',
       'ring-type', 'spore-print-color', 'population', 'habitat'],
      dtype='object')

## Preparing the data

```
y = df['class']
x = df[['cap-shape', 'cap-surface', 'cap-color', 'bruises', 'odor',
       'gill-attachment', 'gill-spacing', 'gill-size', 'gill-color',
       'stalk-shape', 'stalk-root', 'stalk-surface-above-ring',
       'stalk-surface-below-ring', 'stalk-color-above-ring',
       'stalk-color-below-ring', 'veil-type', 'veil-color', 'ring-number',
       'ring-type', 'spore-print-color', 'population', 'habitat']]
```

```
y.shape , x.shape
```

((8124,), (8124, 22))

## Data Manipulation

```
#step 4 Lincoders
from sklearn.preprocessing import LabelEncoder
le = LabelEncoder()
y = le.fit_transform(y)
```

```
y
```

array([1, 0, 0, ..., 0, 1, 0])

```
for i in x:
  x[i] = le.fit_transform(x[i])
```

<ipython-input-14-d05a489caee4>:2: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a
  x[i] = le.fit_transform(x[i])
<ipython-input-14-d05a489caee4>:2: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a
  x[i] = le.fit_transform(x[i])
<ipython-input-14-d05a489caee4>:2: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a
  x[i] = le.fit_transform(x[i])
<ipython-input-14-d05a489caee4>:2: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a
  x[i] = le.fit_transform(x[i])
<ipython-input-14-d05a489caee4>:2: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a
  x[i] = le.fit_transform(x[i])
<ipython-input-14-d05a489caee4>:2: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a
  x[i] = le.fit_transform(x[i])
<ipython-input-14-d05a489caee4>:2: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a
  x[i] = le.fit_transform(x[i])
<ipython-input-14-d05a489caee4>:2: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
```

```
      Try using .loc[row_indexer,col_indexer] = value instead

      See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a
        x[i] = le.fit_transform(x[i])
      <ipython-input-14-d05a489caee4>:2: SettingWithCopyWarning:
      A value is trying to be set on a copy of a slice from a DataFrame.
      Try using .loc[row_indexer,col_indexer] = value instead

      See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a
        x[i] = le.fit_transform(x[i])
      <ipython-input-14-d05a489caee4>:2: SettingWithCopyWarning:
      A value is trying to be set on a copy of a slice from a DataFrame.
      Try using .loc[row_indexer,col_indexer] = value instead

      See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a
        x[i] = le.fit_transform(x[i])
      <ipython-input-14-d05a489caee4>:2: SettingWithCopyWarning:
      A value is trying to be set on a copy of a slice from a DataFrame.
```

**Splitting the data into training and testing**

```
#  train and test
from sklearn.model_selection import train_test_split
x_train, x_test,y_train,y_test = train_test_split(x,y,test_size=0.3,random_state=75)
```

Classification Methods

```
#Random Forest Classification
from sklearn.ensemble import RandomForestClassifier

rf = RandomForestClassifier(n_estimators=100, random_state=42)
rf.fit(x_train,y_train)

print("Test Accuracy: {}%".format(round(rf.score(x_test,y_test)*100,2)))

      Test Accuracy: 100.0%


#Decision tree classification
from sklearn.tree import DecisionTreeClassifier

dt = DecisionTreeClassifier()
dt.fit(x_train,y_train)

print("Test Accuracy: {}%".format(round(dt.score(x_test,y_test)*100,2)))

      Test Accuracy: 100.0%


#naive bayes classfications
from sklearn.naive_bayes import GaussianNB

nb = GaussianNB()
nb.fit(x_train,y_train)

print("Test Accuracy: {}%".format(round(nb.score(x_test,y_test)*100,2)))

      Test Accuracy: 91.88%


#SVM Classification
from sklearn.svm import SVC

svm = SVC(random_state=42, gamma="auto")
svm.fit(x_train,y_train)

print("Test Accuracy: {}%".format(round(svm.score(x_test,y_test)*100,2)))

      Test Accuracy: 100.0%


#KNN
from sklearn.neighbors import KNeighborsClassifier

best_Kvalue = 0
best_score = 0
```

```
for i in range(1,10):
    knn = KNeighborsClassifier(n_neighbors=i)
    knn.fit(x_train,y_train)
    if knn.score(x_test,y_test) > best_score:
        best_score = knn.score(x_train,y_train)
        best_Kvalue = i

print("""Best KNN Value: {}
Test Accuracy: {}%""".format(best_Kvalue, round(best_score*100,2)))
```

```
    Best KNN Value: 1
    Test Accuracy: 100.0%
```

```
#Logistic Regressions
from sklearn.linear_model import LogisticRegression

## lr = LogisticRegression(solver="lbfgs")
lr = LogisticRegression(solver="liblinear")
lr.fit(x_train,y_train)

print("Test Accuracy: {}%".format(round(lr.score(x_test,y_test)*100,2)))
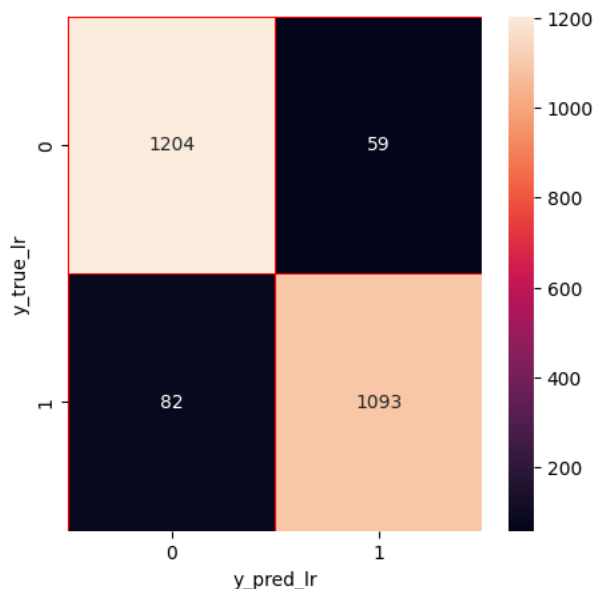```

```
    Test Accuracy: 94.22%
```

**Confusion matrix**

```
from sklearn.metrics import confusion_matrix
import matplotlib.pyplot as plt
import seaborn as sns
# Linear Regression
y_pred_lr = lr.predict(x_test)
y_true_lr = y_test
cm = confusion_matrix(y_true_lr, y_pred_lr)
f, ax = plt.subplots(figsize =(5,5))
sns.heatmap(cm,annot = True,linewidths=0.5,linecolor="red",fmt = ".0f",ax=ax)
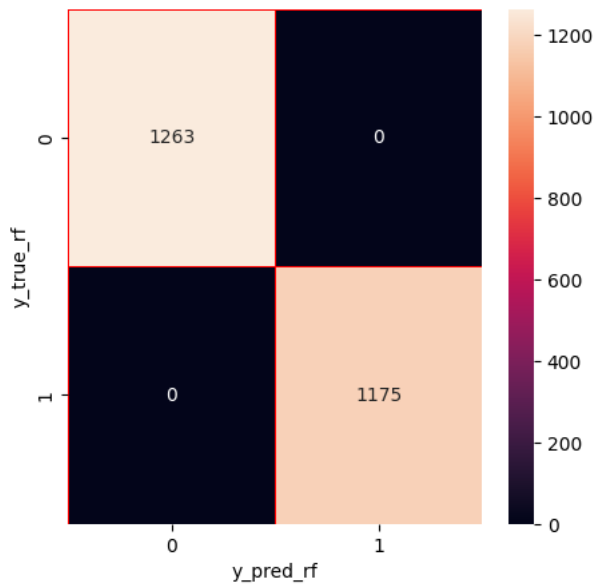plt.xlabel("y_pred_lr")
plt.ylabel("y_true_lr")
plt.show()
```



```
# Random Forest
y_pred_rf = rf.predict(x_test)
y_true_rf = y_test
cm = confusion_matrix(y_true_rf, y_pred_rf)
f, ax = plt.subplots(figsize =(5,5))
sns.heatmap(cm,annot = True,linewidths=0.5,linecolor="red",fmt = ".0f",ax=ax)
plt.xlabel("y_pred_rf")
plt.ylabel("y_true_rf")
plt.show()
```

Conclusion

From the confusion matrix, we saw that our train and test data is balanced.

Most of classfication methods hit 100% accuracy with this dataset.

In conclusion, the application of machine learning in mushroom classification has demonstrated its remarkable potential in automating and enhancing the accuracy of identifying mushroom species. Through the utilization of advanced algorithms and vast datasets, we have witnessed the development of robust models capable of distinguishing between edible and toxic mushrooms with a high degree of confidence.

As technology continues to advance and more research is conducted, we can anticipate even greater strides in the accuracy and efficiency of mushroom classification using machine learning. This, in turn, will contribute to safer mushroom foraging practices, greater understanding of fungal biodiversity, and the preservation of ecosystems.