

Business Problem

Goal here is to see if we can harness the power of machine learning and boosting to help create not just a predictive model, but a general guideline for features people should look out for when picking mushrooms.

```
In [1]: import numpy as np  
import pandas as pd  
import matplotlib.pyplot as plt  
import seaborn as sns
```

```
In [2]: df = pd.read_csv("mushrooms.csv")  
df.head()
```

Out[2]:

	class	cap-shape	cap-surface	cap-color	bruises	odor	gill-attachment	gill-spacing	gill-size	gill-color	...	stalk-surface-below-ring	?
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

Data

Mushroom: Edible or Poisonous?

Data Source: <https://archive.ics.uci.edu/ml/datasets/Mushroom>
[\(https://archive.ics.uci.edu/ml/datasets/Mushroom\)](https://archive.ics.uci.edu/ml/datasets/Mushroom)

This data set includes descriptions of hypothetical samples corresponding to 23 species of gilled mushrooms in the Agaricus and Lepiota Family (pp. 500-525). 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.

Attribute Information:

1. cap-shape: bell=b,conical=c,convex=x,flat=f, knobbed=k,sunken=s
2. cap-surface: fibrous=f,grooves=g,scaly=y,smooth=s
3. cap-color: brown=n,buff=b,cinnamon=c,gray=g,green=r, pink=p,purple=u,red=e,white=w,yellow=y
4. bruises?: bruises=t,no=f
5. odor: almond=a,anise=l,creosote=c,fishy=y,foul=f, musty=m,none=n,pungent=p,spicy=s
6. gill-attachment: attached=a,descending=d,free=f,notched=n
7. gill-spacing: close=c,crowded=w,distant=d
8. gill-size: broad=b,narrow=n
9. 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
10. stalk-shape: enlarging=e,tapering=t
11. stalk-root: bulbous=b,club=c,cup=u,equal=e, rhizomorphs=z,rooted=r,missing=?
12. stalk-surface-above-ring: fibrous=f,scaly=y,silky=k,smooth=s
13. stalk-surface-below-ring: fibrous=f,scaly=y,silky=k,smooth=s
14. stalk-color-above-ring: brown=n,buff=b,cinnamon=c,gray=g,orange=o, pink=p,red=e,white=w,yellow=y
15. stalk-color-below-ring: brown=n,buff=b,cinnamon=c,gray=g,orange=o, pink=p,red=e,white=w,yellow=y
16. veil-type: partial=p,universal=u
17. veil-color: brown=n,orange=o,white=w,yellow=y
18. ring-number: none=n,one=o,two=
19. ring-type: cobwebby=c,evanescent=e,flaring=f,large=l, none=n,pendant=p,sheathing=s,zone=z
20. spore-print-color: black=k,brown=n,buff=b,chocolate=h,green=r, orange=o,purple=u,white=w,yellow=y
21. population: abundant=a,clustered=c,numerous=n, scattered=s,several=v,solitary=y
22. habitat: grasses=g,leaves=l,meadows=m,paths=p, urban=u,waste=w,woods=d

```
In [3]: 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
```

EDA

```
In [ ]:
```

```
In [ ]:
```

X & y

```
In [4]: X = pd.get_dummies(df.drop('class',axis=1),drop_first=True)
y = df['class']
```

```
In [5]: X.shape
```

```
Out[5]: (8124, 95)
```

Train Test Split

```
In [6]: from sklearn.model_selection import train_test_split  
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=101)
```

Modelling

Gradient Boosting with default parameters

```
In [7]: from sklearn.ensemble import GradientBoostingClassifier  
  
gb_model = GradientBoostingClassifier()  
  
gb_model.fit(X_train,y_train)  
  
Out[7]: GradientBoostingClassifier()
```

prediction

```
In [8]: y_pred_train = gb_model.predict(X_train)  
  
y_pred_test = gb_model.predict(X_test)
```

Evaluation

Accuracy

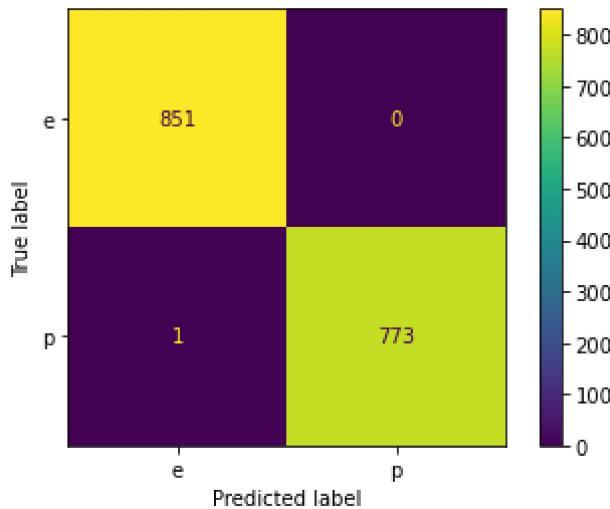
```
In [9]: from sklearn.metrics import accuracy_score  
print(accuracy_score(y_train,y_pred_train)) # train accuracy  
print(accuracy_score(y_test,y_pred_test)) # test accuracy
```

```
0.9996922603477458  
0.9993846153846154
```

Confusion Matrix

```
In [10]: from sklearn.metrics import plot_confusion_matrix  
print(plot_confusion_matrix(gb_model,X_test,y_test))  
plt.show()
```

```
<sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay object at 0x00  
0001F602A48370>
```



Classification Report

```
In [11]: from sklearn.metrics import classification_report  
print(classification_report(y_test,y_pred_test))
```

	precision	recall	f1-score	support
e	1.00	1.00	1.00	851
p	1.00	1.00	1.00	774
accuracy			1.00	1625
macro avg	1.00	1.00	1.00	1625
weighted avg	1.00	1.00	1.00	1625

Cross Validation Score

```
In [12]: from sklearn.model_selection import cross_val_score  
scores = cross_val_score(gb_model,X,y,cv=5)  
print("Cross Validation Score:",scores.mean())
```

```
Cross Validation Score: 0.9192312239484653
```

Hyperparameter Tuning

```
In [13]: from sklearn.model_selection import GridSearchCV  
  
estimator= GradientBoostingClassifier()  
  
param_grid = {"n_estimators": [1,5,10,20,40,100], "learning_rate": [0.1,0.2,0.3, 0.5,0.8,1]}  
  
grid = GridSearchCV(estimator, param_grid, cv=5, scoring='accuracy')  
  
grid.fit(X_train,y_train)  
  
grid.best_params_
```

```
Out[13]: {'learning_rate': 0.5, 'n_estimators': 40}
```

Final Model

```
In [14]: final_model = GradientBoostingClassifier(n_estimators=40,learning_rate=0.5)  
final_model.fit(X_train,y_train)  
  
preds_train = final_model.predict(X_train)  
preds_test = final_model.predict(X_test)  
  
print("Train Accuracy Score: ", accuracy_score(y_train,preds_train))  
print("Test Accuracy Score: ",accuracy_score(y_test,preds_test))
```

Train Accuracy Score: 1.0
Test Accuracy Score: 1.0

Feature Importance

```
In [15]: final_model.feature_importances_
```

```
Out[15]: array([ 1.40010963e-04, -2.27670640e-19,  3.99382234e-16,  0.00000000e+00,
   2.54477173e-18,  1.27243027e-03,  3.88388138e-17,  6.64765904e-05,
   2.86767688e-07,  4.95221088e-18,  0.00000000e+00,  4.15268210e-19,
   3.63154330e-18,  0.00000000e+00,  0.00000000e+00,  3.50864898e-11,
   3.92928487e-03,  4.90732370e-02,  1.66731484e-04,  3.52657438e-03,
   1.85945411e-02,  1.13118827e-02,  6.29333275e-01,  1.13542514e-03,
   0.00000000e+00,  0.00000000e+00, -1.87229145e-22,  2.10079187e-03,
   2.67813742e-03,  0.00000000e+00,  6.62984253e-19,  0.00000000e+00,
   0.00000000e+00,  0.00000000e+00,  0.00000000e+00,  0.00000000e+00,
   0.00000000e+00,  0.00000000e+00,  2.34167492e-17,  0.00000000e+00,
   2.61917856e-03,  6.65336405e-05,  1.35971789e-01,  9.82006596e-03,
   5.77031764e-02,  1.02335109e-03,  0.00000000e+00,  2.39024180e-05,
   1.24870781e-06,  1.79784641e-05,  1.34641438e-02,  0.00000000e+00,
   0.00000000e+00,  0.00000000e+00,  6.57691430e-04,  0.00000000e+00,
   0.00000000e+00,  0.00000000e+00,  7.85386395e-07,  2.11936791e-04,
   0.00000000e+00,  0.00000000e+00,  1.60357125e-04,  0.00000000e+00,
   0.00000000e+00,  3.41340855e-08,  2.16751132e-03,  0.00000000e+00,
   1.13712783e-05,  3.36698865e-04,  5.77099185e-06,  1.15416733e-04,
   1.90839921e-06,  0.00000000e+00,  1.32175358e-05,  8.29146281e-04,
   1.65412375e-02,  1.79600532e-05,  3.24100055e-04,  0.00000000e+00,
   2.89889685e-02,  4.37437969e-03,  1.38099928e-04,  0.00000000e+00,
   1.66239015e-04,  0.00000000e+00,  0.00000000e+00,  4.85825265e-07,
   0.00000000e+00,  1.77420093e-18,  8.96173293e-04, -9.28656559e-20,
   0.00000000e+00,  5.60513440e-08,  0.00000000e+00])
```

```
In [16]: imp_feats = pd.DataFrame(index=X.columns,data=grid.best_estimator_.feature_importances_,columns=[ 'Importance'])
imp_feats
```

```
Out[16]:
```

	Importance
cap-shape_c	8.138390e-05
cap-shape_f	1.127187e-16
cap-shape_k	0.000000e+00
cap-shape_s	0.000000e+00
cap-shape_x	-1.462527e-23
...	...
habitat_l	8.960505e-04
habitat_m	9.520533e-17
habitat_p	0.000000e+00
habitat_u	3.372033e-06
habitat_w	0.000000e+00

95 rows × 1 columns

```
In [17]: imp_feats = imp_feats[imp_feats[ 'Importance'] > 0.01]
```

```
In [18]: plt.figure(figsize=(14,6),dpi=200)
sns.barplot(data=imp_feats.sort_values('Importance'),x=imp_feats.index,y='Importance')
plt.xticks(rotation=90);
```

