

26/09/22
10:06pm

Ensemble Learning

Ex :- Avengers.

Boosting Technique :- it is methodology, not a m.L
~~~~~\*~~~~~  
Algorithm, it is applied to an Existing m.L, mostly  
applied on Decision Tree.

Ex :- Student

| * Exam | * Attempt again | * Again Exam |
|--------|-----------------|--------------|
| ✓      |                 |              |
| ✗      |                 |              |
| ✗      |                 |              |
| ✓✓     |                 |              |
| ✓      |                 |              |

⇒ not prepared well, focus on wrong answers only

⇒ rectified but missed another section

⇒ Total correctly done.

Boosting :-

⇒ Decision Tree 1 :- misclassified records

⇒ Decision Tree 2 :- misclassified records

⇒ Decision Tree 3 :-

- \* 2<sup>nd</sup> weak learner is dependent on 1<sup>st</sup> weak learner  
 ↴  
 algorithm doesn't give more accuracy.

\* Formula for Boosting :-

$$F_T(x) = \sum_{t=1}^T f_t(x)$$

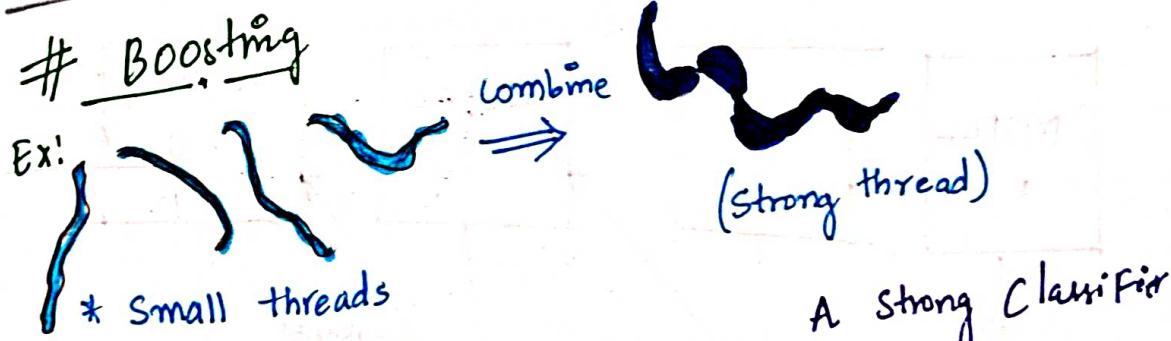
↑ Summation of  $f_t(x)$  to  $F_T(x)$

↑ Learning rate  
↓ [Each model]

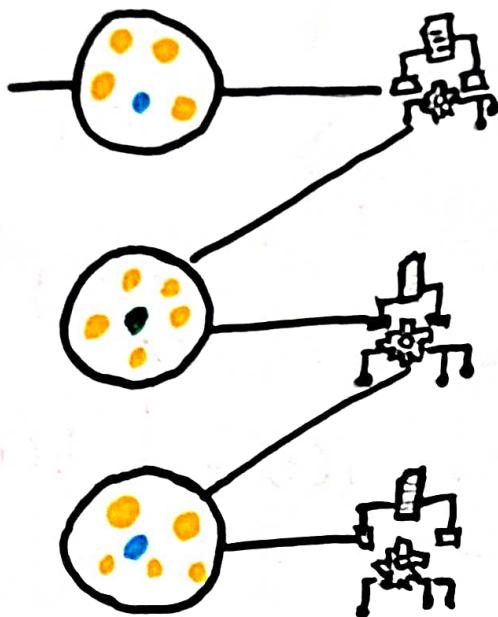
$$f_t(x) = \alpha_t h(x)$$

- \* boosting is the Combination of weak learners:

\* Implies that Combination of Estimators (models) with an applied Coefficient could act as an effective Ensemble Estimator.



- \* Multiple Weak classifier (stumps)



"Sequential learning:

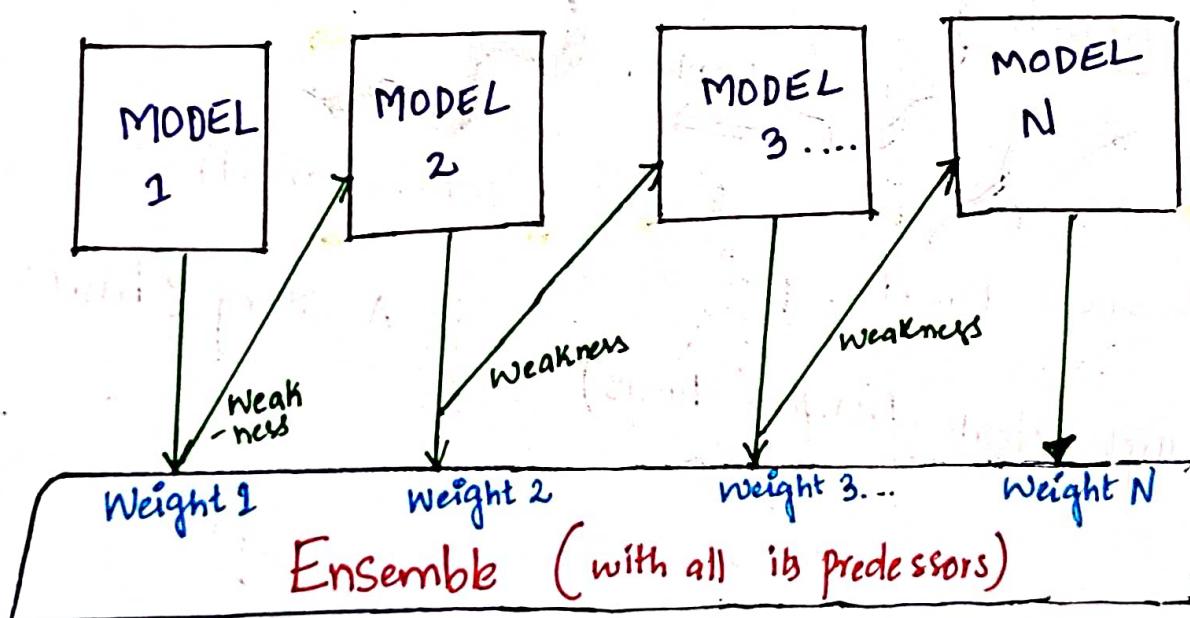
- \* From the first algorithm whatever mistakes was done.
- \* we share to the second algorithm.
- \* again it shared to third algorithm. . . .
- \* it is in "sequential order"

# Boosting - "Sequential" Learning.

## I. Ada Boost → Boosting.

Adaptive

Ex:- Model, 1, 2, . . . N are individual models (e.g. Decision Tree)



covid (positive/negative)

Ex:- ✓

Sample DataSet

|   |   |   |
|---|---|---|
| + | + | - |
| - | - | + |
| + | - | + |
| + | - | - |

Weak Learner 1

Weak learner 2

Weak learner 3

iteration 1

|                |   |   |   |   |
|----------------|---|---|---|---|
| x <sub>0</sub> | + | + | + | - |
| x <sub>1</sub> | - | - | - | - |

miss classified

iteration 2

|                |   |   |   |   |
|----------------|---|---|---|---|
| x <sub>0</sub> | + | + | + | - |
| x <sub>1</sub> | - | - | - | - |

miss classified

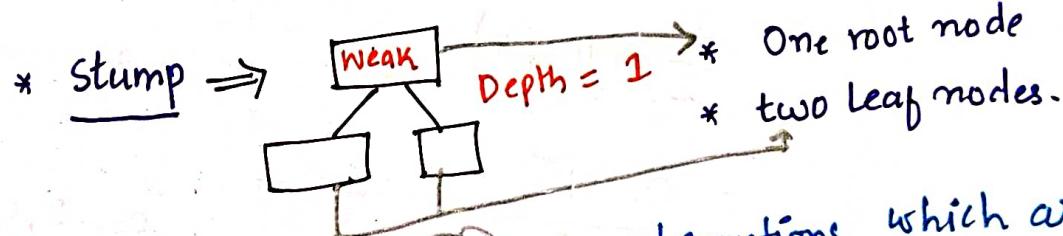
iteration 3.

|                |   |   |   |   |
|----------------|---|---|---|---|
| x <sub>0</sub> | + | + | + | - |
| x <sub>1</sub> | - | - | - | + |

#Final classifier / strong classifier.

|   |    |   |
|---|----|---|
| + | ++ | - |
| + | -  | - |
| + | -  | - |

\* In AdaBoost : Each weak learner is called "stump"



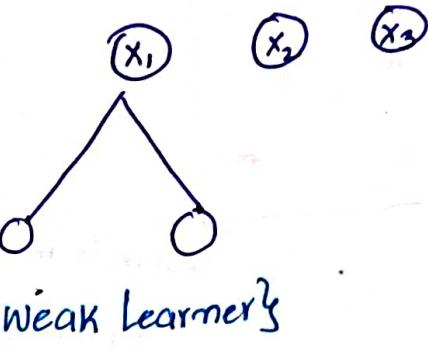
\* Step 1 : AdaBoost assigns weights to observations which are incorrectly predicted to predict these values

Step 2 : Next model works which are wrongly predicted before.

Ex:-

|   | $X_1$ | $X_2$ | $X_3$ | $X_4$ | O/p | weights |
|---|-------|-------|-------|-------|-----|---------|
| 1 | -     | -     | -     | -     | Yes | $y_1$   |
| 2 | -     | -     | -     | -     | No  | $y_1$   |
| 3 | -     | -     | -     | -     |     | $y_1$   |
| 4 | -     | -     | -     | -     |     | $y_1$   |
| 5 | -     | -     | -     | -     |     | $y_1$   |
| 6 | -     | -     | -     | -     |     | $y_1$   |
| 7 | -     | -     | -     | -     |     | $y_1$   |

# Taken by Entropy/Inform. gain



\*1. Total Error [TE] :-  $\frac{1}{7}$

$\therefore \epsilon_t = \epsilon_w$   $\therefore \epsilon_t = \text{"Episimol"} T$

2. Performance of Stump :-

$\frac{1}{2} \log_e \left[ \frac{1-TE}{TE} \right]$   $\therefore TE = \text{Total Error}$

$$= \frac{1}{2} \log_e \left[ \frac{1 - \frac{1}{7}}{\frac{1}{7}} \right]$$

$$= \frac{1}{2} \log_e \left[ \frac{6/7}{1/7} \right]$$

$$= \frac{1}{2} \log_e (6)$$

$$= 0.895$$

$\therefore PS = \text{Performance Stump}$

New Sample weight

$$= \text{Weight} * e^{-PS}$$

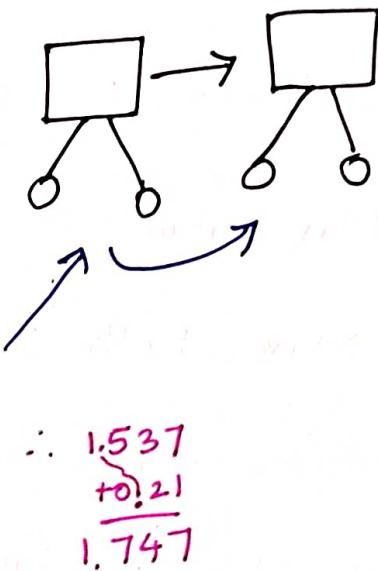
$$= \frac{1}{7} * e^{-0.895} = 0.05$$

+ incorrect record =  $\text{Weight} * e^{PS}$   $\rightarrow$  incorrect records.

$$= \frac{1}{7} * e^{0.895} = 0.349$$

is it 1??

| Weights       | New weights       | Normalized weights | Buckets           |
|---------------|-------------------|--------------------|-------------------|
| $\frac{1}{7}$ | $0.05 \div 0.649$ | 0.077              | (0 - 0.07)        |
| $\frac{1}{7}$ | $0.05 \div 0.649$ | 0.077              | (0.07 - 0.14)     |
| $\frac{1}{7}$ | 0.05              | 0.077              | (0.14 - 0.21)     |
| $\frac{1}{7}$ | 0.349             | 0.537              | (0.21 - 0.747)    |
| $\frac{1}{7}$ | 0.05              | 0.077              | [1.747 - 0.751] ✓ |
| $\frac{1}{7}$ | 0.05              | 0.077              | -                 |
| $\frac{1}{7}$ | 0.05              | 0.077              | -                 |
| <hr/>         |                   | $\sim 1$           |                   |



\* STEPS :-

$$\Rightarrow \text{Weight} = \frac{1}{n}$$

$$\Rightarrow \text{Total Error}_t = \frac{[\text{Correct Prediction} - n]}{n}$$

$$\Rightarrow \text{Total weighted error}_t = \frac{\sum (\text{Weight}(i) * \text{Error}(i))}{\sum(\text{Weight})}$$

$$\Rightarrow \alpha_t = \ln \left[ \frac{1 - TE}{TE} \right]$$

$$\Rightarrow \text{Weight} = \text{weight} * \exp(\alpha_t * \text{Total error})$$

$$\Rightarrow \text{total error} = 0 \text{ if } (y == \hat{y}), \text{ else } 1.$$

## \* Parameters

- \* n\_estimators :- no. of weak learners
- \* Criterion :- Gini / Entropy
- \* Max\_Feature :- all feature / selected ones
- \* Max\_depth :
- \* min\_samples\_split
- \* n-jobs :- System is quadcore :- no. of cores = 4  
                  dual core :- no. of cores = 2  
                  octo core : no. of cores = 8  
                  -1 (all cores)
- \* random\_state :
- \* learning\_rate :

Dr. W. H. H.

11:20 AM

CODE: AdaBoost

Data :-

This Data Set includes descriptions of hypothetical Samples corresponding to 23 species of gilled mushrooms in 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 → Mushroom
2. Cap-Surface :- fibrous = f, grooves = g, scaly = y, smooth = -
3. Cap-colour :- brown = n, buff = b, cinnamon = c, gray = g, green = o, 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 : black = k, brown = n, buff = b, chocolate = h,  
gray = g, green = t, orange = o, pink = p,  
purple = u, red = e, white = w, yellow = y

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

8. gill-size : broad = b, narrow = m

9. gill-color : black = k, brown = n, buff = b, chocolate = h, gray = g,  
green = t, 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 - colour - above - ring : brown = n, buff = b, cinnamon = c,  
gray = g, orange = o, pink = p, red = e,  
white = w, yellow = y.

stalk - colour - below - ring :

veil - type : partial = p, universal = u

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

18. Ring - number :- none = n, One = o, two = t, large = l
19. Ring - type :- cobwebby = c, evanescent = e, flaring = f, large = l, none = n, Pendent = p, Sheathing = s, Zone = z
20. Spore - Print :- black = b, brown = n, buff = b, chocolate = h, green = g, 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.

### 23. Class (Output)

#### 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 general Guideline for features people should look out for when picking mushrooms.

# df = pd.read\_csv("mushrooms.csv")  
# df.head()

| class | cap shape | cap surface | cap color | bruises | odor | gill attachment | gill spacing | gill size | gill color | 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 |
|-------|-----------|-------------|-----------|---------|------|-----------------|--------------|-----------|------------|--------------------------|--------------------------|------------------------|------------------------|-----------|------------|-------------|-----------|-------------------|------------|---------|
| P     | x         | s           | n         | t       | p    | f               | c            | n         | b          | s                        | w                        | w                      | w                      | p         | w          | o           | p         | b                 | s          | u       |
| E     | x         | s           | y         | t       | a    | f               | c            | b         | n          | s                        | w                        | w                      | w                      | p         | w          | o           | p         | n                 | n          | g       |
| e     | b         | s           | w         | t       | p    | f               | c            | n         | n          | s                        | w                        | w                      | w                      | p         | w          | o           | p         | k                 | s          | u       |
| P     | x         | y           | w         | f       | n    | f               | w            | b         | b          | s                        | w                        | w                      | w                      | p         | w          | o           | e         | n                 | a          | g       |
| e     | x         | s           | g         | f       | n    | f               | w            | b         | b          | s                        | w                        | w                      | w                      | p         | w          | o           | p         | n                 | n          | g       |

# df.shape

[out]: [8124, 23]

# df.info()

[out]: RangeIndex: 8124 Entries, 0 to 8123

| column        | Nonnull Count | Dtype.  |
|---------------|---------------|---------|
| *             | 8124          | Nonnull |
| class         | "             | Object  |
| * cap-shape   | "             | "       |
| * cap-surface | "             | "       |
| :             | "             | "       |
| * habitat     | "             | "       |

F df.isnull().sum()

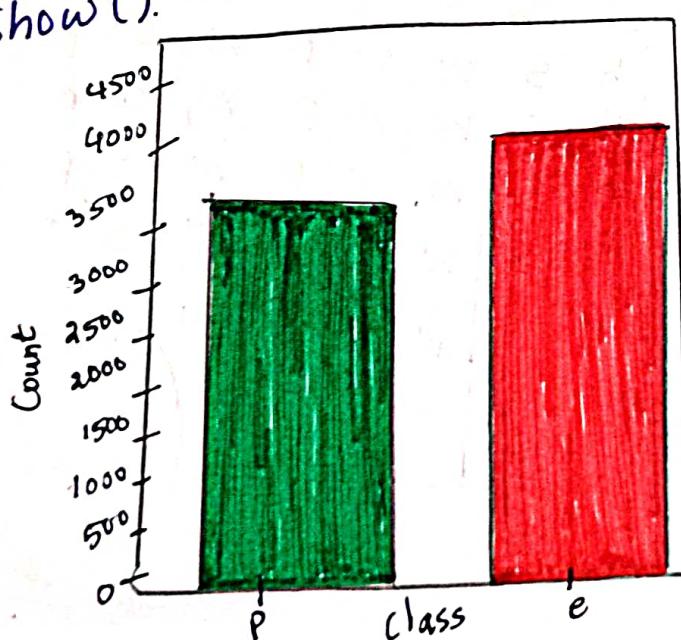
[out]: 0.

EDA

: sns.countplot(data=df, x="class", palette="Dark2")

plt.show()

AE



# Transpose. gives total no. columns, without ...

# df.describe().transpose()

Out

|                                 | Count | Unique | top | freq |
|---------------------------------|-------|--------|-----|------|
| (o/p) Class                     | 8124  | 2      | e   | 4208 |
| Cap Shape                       | 8124  | 6      | x   | 3656 |
| Cap - Surface                   | 8124  | 4      | y   | 3244 |
| cap - color                     | 8124  | 10     | n   | 2284 |
| bruises                         | 8124  | 2      | f   | 4748 |
| odor                            | 8124  | 9      | n   | 3528 |
| gill attachment                 | 8124  | 2      | f   | 7914 |
| gill spacing                    | 8124  | 2      | c   | 6812 |
| gill size                       | 8124  | 2      | b   | 5612 |
| gill - color                    | 8124  | 12     | b   | 1728 |
| Stalk - shape                   | 8124  | 2      | t   | 4608 |
| Stalk - root                    | 8124  | 5      | b   | 3776 |
| Stalk - Surface - above<br>ring | 8124  | 4      | s   | 5176 |
| Stalk - Surface - below<br>ring | 8124  | 4      | s   | 4936 |
| Stalk - color above<br>ring     | 8124  | 9      | w   | 4464 |
| Stalk - color below<br>ring     | 8124  | 9      | w   | 4384 |
| Veil - type                     | 8124  | 1      | p   | 8124 |
| Veil - color                    | 8124  | 4      | w   | 7924 |
| ring - number                   | 8124  | 3      | o   | 7488 |
| ring - type                     | 8124  | 5      | p   | 3968 |
| Spore - Print color             | 8124  | 9      | w   | 2388 |
| Population                      | 8124  | 6      | v   | 4040 |
| habitant                        | 8124  | 7      | d   | 3148 |

\* posinouss  
\* non posinouss

$x \& y$

Name ...  
(nominal data)  $\Rightarrow$  One hot Encoding  
 $\uparrow$

#  $X = \text{pd.get_dummies}(\text{df.drop}(\text{"class", axis=1}), \text{drop\_first} = \text{True})$

# y = df [ "class" ]

廿  $x$

8124 rows x 95 columns

# Y

|            |   |   |   |
|------------|---|---|---|
| <b>Out</b> | : | 0 | P |
|            |   | 1 | e |
|            |   | 2 | e |
|            |   | 3 | P |
|            |   | 4 | e |

# x.shape, y.shape

out:  $(8124, 95), (8124, 1))$

## train / test split

```
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=29)
```

## MODELLING

```
from sklearn.ensemble import AdaBoostClassifier  
# model = AdaBoostClassifier ()  
# model.fit (x-train, y-train)  
out : AdaBoostClassifier()
```

## PREDICTION

```
# ypred-train = model.predict (x-train)  
# ypred-test = model.predict (x-test)
```

## Evaluation

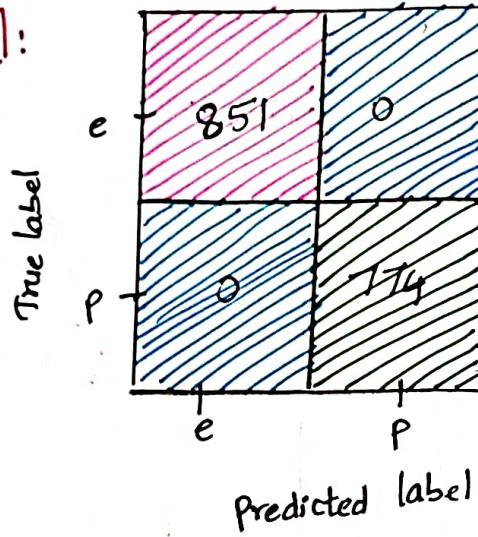
x Accuracy

```
from sklearn.metrics import accuracy_score  
# print ("Train accuracy":, accuracy_score (y-train, ypred-train)  
# print ("Test accuracy":, accuracy_score (y-test, ypred-test)  
out : Train Accuracy : 1.0  
      Test Accuracy : 1.0
```

## \* Confusion Matrix

```
# Confusion Matrix  
from sklearn.metrics import plot_confusion_matrix  
# plot_confusion_matrix(model, X-test, y-test)  
# plt.show()
```

Out:



## \* classification Report

```
from sklearn.metrics import classification_report  
print(classification_report(y-test, y-pred-test))
```

Out:

|              | Precision | recall | f1 Score | Support |
|--------------|-----------|--------|----------|---------|
| e            | 1.00      | 1.00   | 1.00     | 851     |
| P            | 1.00      | 1.00   | 1.00     | 774     |
| Accuracy     | 1.00      | 1.00   | 1.00     | 1625    |
| macro avg    | 1.00      | 1.00   | 1.00     | 1625    |
| Weighted Avg | 1.00      | 1.00   | 1.00     | 1625    |

## \* Cross-validation-Score

```
from sklearn.model_selection import cross_val_score.  
# Scores = cross_val_score(model, X, y, cv=5)  
# print("Cross Validation Score:", scores.mean())  
Out: Cross validation Score : 0.925
```

## \* model.features.importances

```
# model.feature_importances -
```

```
Out array([0., 0., 0., 0., 0., 0., 0.,  
          0., 0., 0., 0.02, 0., 0.04,  
          0.06, 0.02, 0.0, 0., 0.1, 0.12,  
          0.]
```

```
# f_imp = pd.DataFrame(index=X.columns, data=model.feature_importances_, columns=[["Feature Importances"]])
```

```
# f_imp
```

```
Out:
```

|                     | feature importances |
|---------------------|---------------------|
| cap-shape=C         | 0.00                |
| cap-shape=F         | 0.00                |
| cap-shape=K         | 0.00                |
| cap-shape=S         | 0.00                |
| cap-shape=X         | 0.00                |
| :                   | :                   |
| 95 rows × 1 columns |                     |

## # Data Extraction (Pandas)

```
# f-imp [f-imp ["feature importance"] > 0]
```

Out :

|                | Feature Importance |
|----------------|--------------------|
| cap-color w    | 0.02               |
| odor-c         | 0.04               |
| odor-f         | 0.04               |
| odor-n         | 0.06               |
| odor-p         | 0.02               |
| gill-spacing w | 0.10               |
| :              | :                  |
| :              | :                  |

## Hyper Parameter Tuning

```
from sklearn.model_selection import GridSearchCV
```

```
# estimator = AdaBoostClassifier()
```

```
# estimator = { "n_estimators": list(range(1, 101)) }
```

```
# param_grid = { "n_estimators": list(range(1, 101)), cv=5, scoring="Accuracy" }
```

```
# grid = GridSearchCV(estimator, param_grid, cv=5, scoring="Accuracy")
```

```
# grid . fit (x-train, y-train)
```

```
# grid . best_params_
```

```
# grid . best_params_
```

Out : { "n\_estimators": 20 }

## final model

```
# final_model = AdaBoost Classifier (n_estimators = 20)
```

```
# final_model . fit (x-train, y-train)
```

```
# pred_train = final_model . predict (x_train)
```

```
# pred_test = final_model . predict (x-test)
```

```
# print ("train accuracy:", accuracy_score (y_train, pred_train))
```

```
# print ("test accuracy:", accuracy_score (y-test, pred-test))
```

[out]: train accuracy : 1.0

Test accuracy : 1.0

28/04/22  
10:00pm

```
# final_model . feature_importances_
```

[out]: array ([ 0. , 0. , 0. , 0. , 0.  
0.0434 0. , 0. , 0. , 0.  
0. , 0. , 0. , 0. , 0.434... ])

anjali  
27/4/22  
4:30pm.

```
# f_imp2 = pd. DataFrame (index = x. columns, data = final_model  
. feature_importances_, columns = ["importance-  
feature"])
```

```
# f_imp2
```

Out:

importance\_Feature

cap\_shape\_c 0.0

cap\_shape\_f 0.0

habitat\_p 0.0

habitat\_u 0.0

# important = f\_imp2 [ f\_imp2 [ "importance\_Feature" ] > 0 ]

# important = Sort\_values ( "importance\_Feature" )

Out:

importance\_Feature

cap\_color\_w 0.043478

odor\_c 0.043478

odor\_f 0.043478

odor\_p 0.043478

:  
order\_n 0.130435

gill\_size\_n 0.173913

# plt.figure ( fig\_size = (14, 6), dpi = 200 )

# sns.barplot ( data = important, sort\_values ( "importance\_Feature" ), x = important.index,

y = "importance\_Feature" )

# plt.xticks ( rotation = 90 )

# plt. show()

