

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 |
|--------|-----------------|--------------|
| =✓     |                 |              |
| -x     |                 |              |
| -x     |                 |              |
| =✓     |                 |              |
| =✓     |                 |              |
| -      |                 |              |

⇒ 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 :-

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

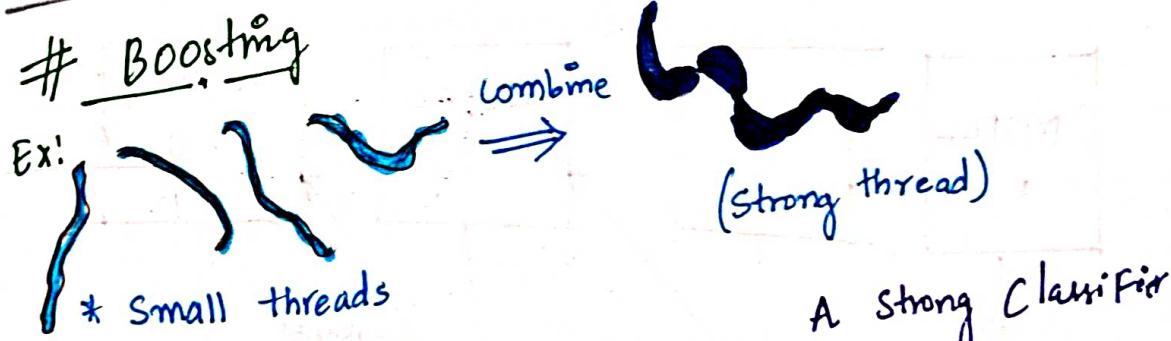
↑ Summation of  $f_t(x)$  to  $F_T(x)$

$$f_t(x) = \alpha_t \frac{h(x)}{\text{Learning rate}}$$

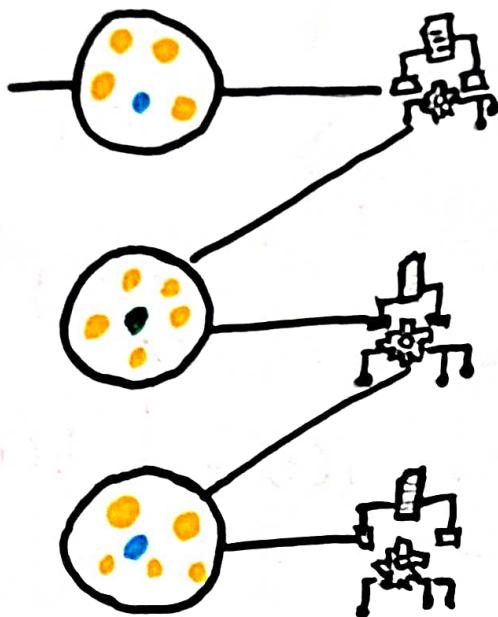
↑ Learning rate  
↓ [Each model]

- \* 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:

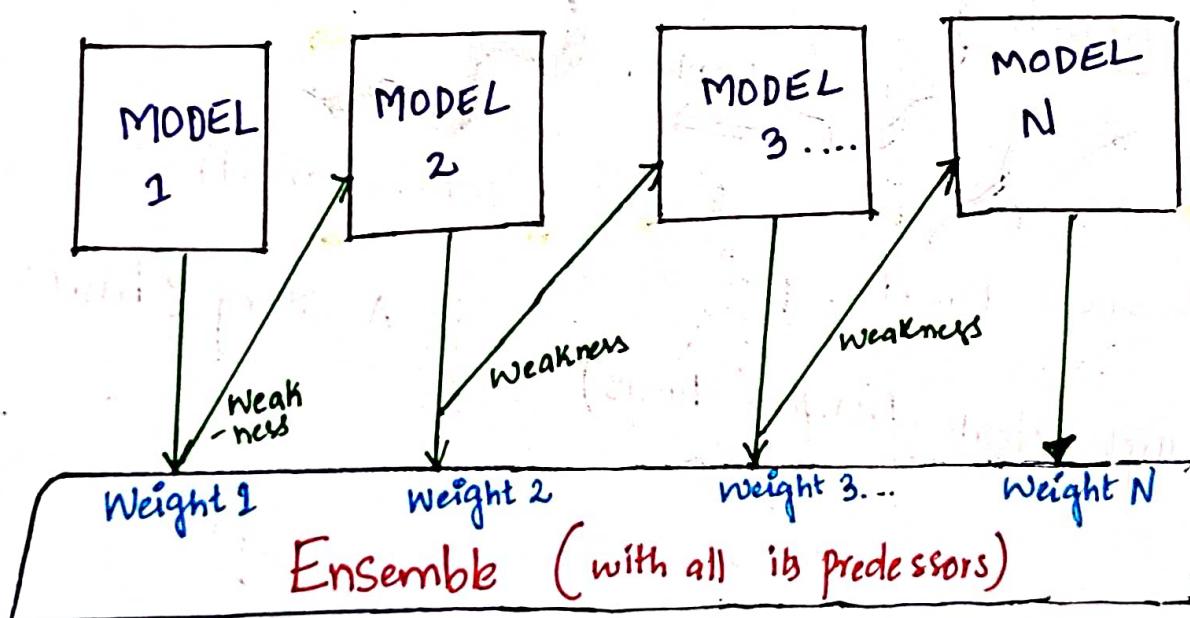
- \* From the first algorithm whatever mistakes was done.
- \* we share to the second algorithm.
- \* again it shared to third algorithm. . . .
- \* it is in "sequential order"

# 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

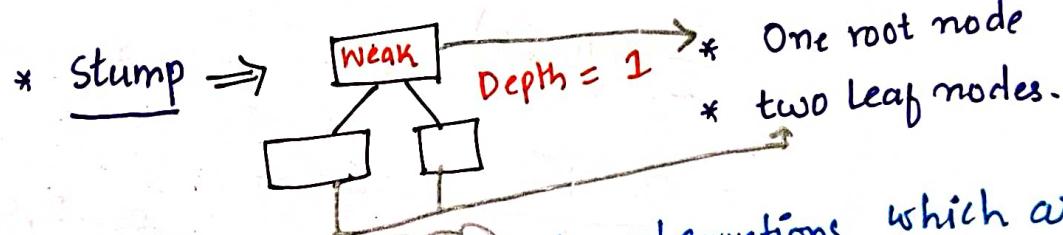
iteration 3.

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

#Final classifier / strong classifier.

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

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



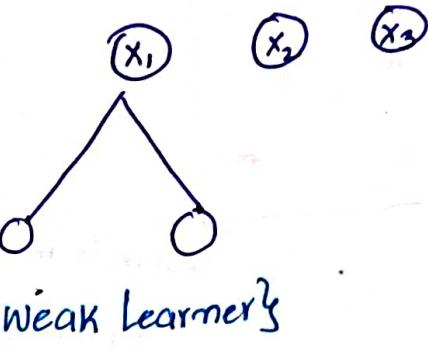
\* Step 1 : AdaBoost assigns weights to observations which are incorrectly predicted to predict these values.

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

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

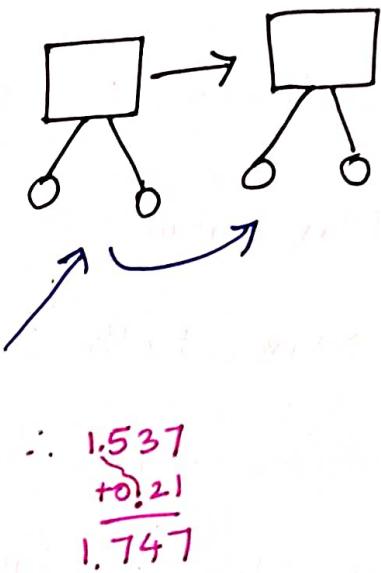
$$= \frac{1}{7} * e^{-0.895} = 0.05$$

+ 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 color below ring | stalk color above 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                      | p         | w          | o           | p         | b                 | s          | u       |
| E     | x         | s           | y         | t       | a    | f               | c            | b         | n          | s                        | w                      | w                      | p         | w          | o           | p         | n                 | n          | g       |
| e     | b         | s           | w         | t       | p    | f               | c            | n         | n          | s                        | w                      | w                      | p         | w          | o           | p         | k                 | s          | u       |
| P     | x         | y           | w         | f       | n    | f               | w            | b         | b          | s                        | w                      | w                      | p         | w          | o           | e         | n                 | a          | g       |
| e     | x         | s           | g         | f       | n    | f               | w            | b         | b          | s                        | w                      | w                      | p         | w          | o           |           |                   |            |         |

# 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()

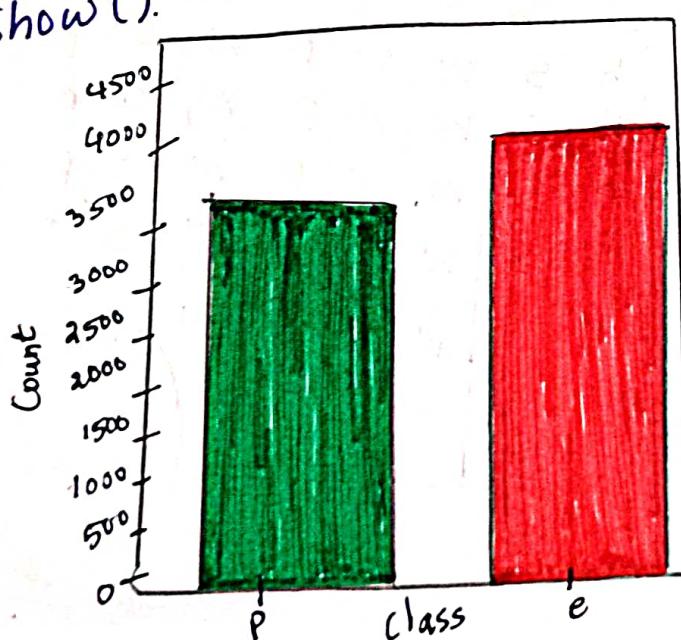
[out]: 0.

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

$x \in Y$

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

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

# y = df [ "class" ]

廿二

|     | Cap | Popult | Popult | Popult | Popult | habit | habit | habit | habit |    | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| out | -c  | -s  | -k  | -s  | -x  | -g  | -s  | -4  | -c  | -e  | -   | -   | -   | -   | -   | -   | -n     | -s     | -v     | -y     | -g    | -i    | -m    | -p    | -u |
| o   | o   | o   | o   | 1   | 0   | 1   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 1      | 0      | 0      | 0      | 0     | 0     | 0     | 0     | 1  |
| o   | o   | o   | o   | 1   | 0   | 1   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 1      | 0      | 0      | 0      | 1     | 0     | 0     | 0     | 0  |
| o   | o   | o   | o   | 0   | 0   | 1   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 1      | 0      | 0      | 0      | 0     | 0     | 1     | 0     | 0  |
| o   | o   | o   | o   | 0   | 0   | 0   | 1   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0      | 1      | 0      | 0      | 0     | 0     | 0     | 0     | 1  |
| o   | o   | o   | o   | 1   | 0   | 0   | 1   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0      | 1      | 0      | 0      | 0     | 1     | 0     | 0     | 0  |
| o   | o   | o   | o   | 1   | 1   | 0   | 1   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0      | 0      | 0      | 0      | 1     | 0     | 0     | 0     | 0  |
| o   | o   | o   | o   | 1   | 1   | 0   | 1   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0      | 0      | 0      | 0      | 1     | 0     | 0     | 0     | 0  |

8124 rows x 95 columns

# 4

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

# x.shape, y.shape

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

## 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

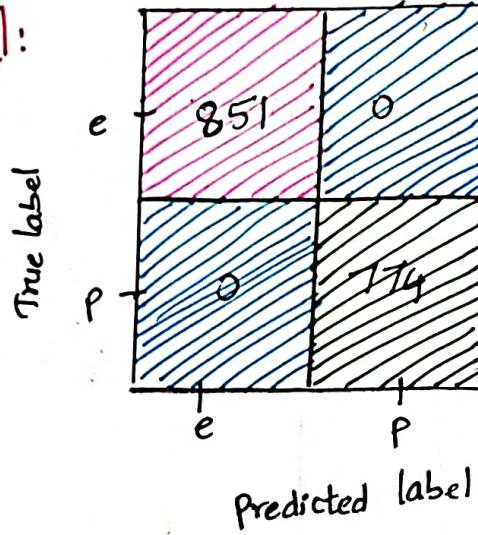
x Accuracy

```
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

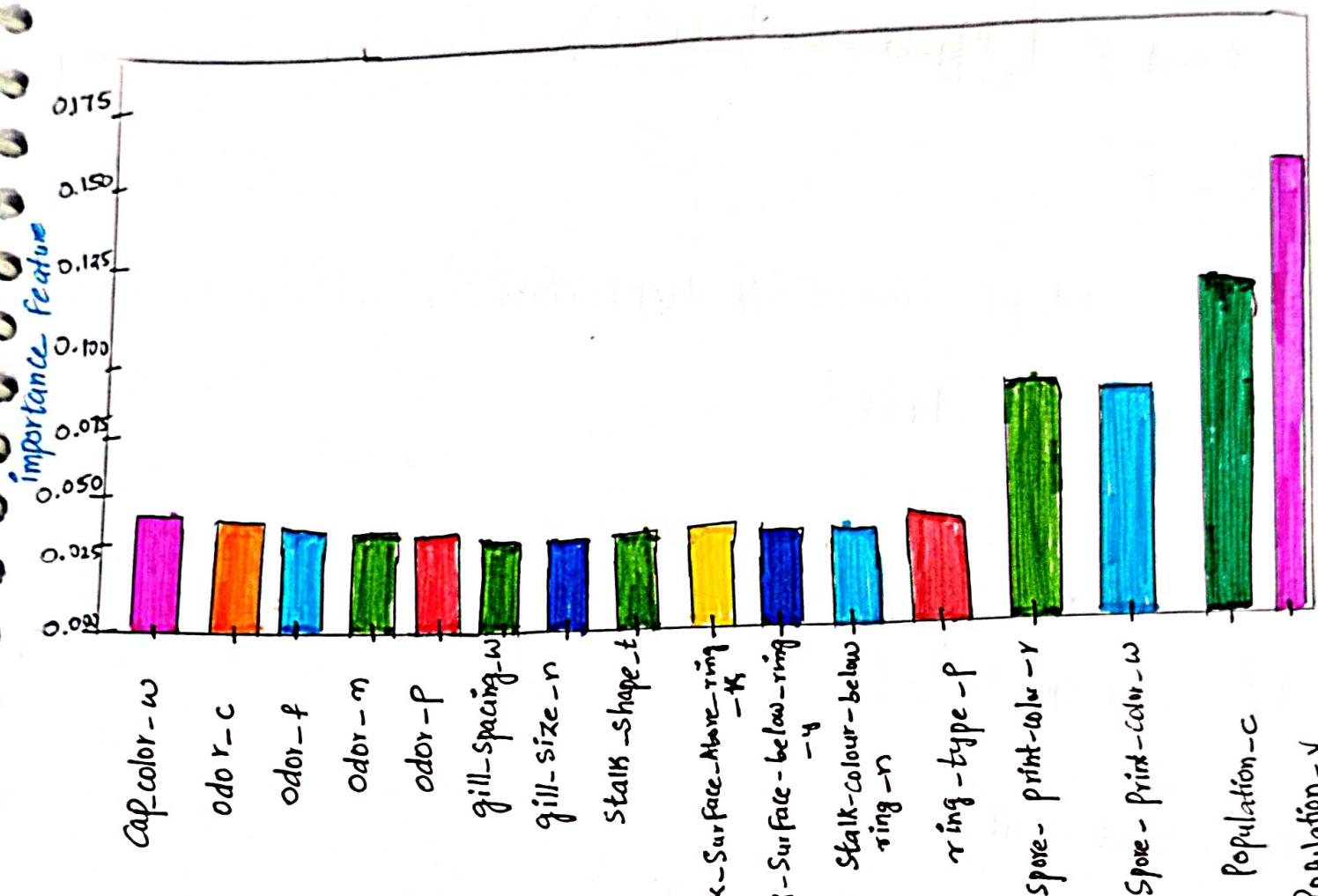
gill\_size\_n 0.173913

# 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()



## \*\*\* Gradient Boosting

\* Within Gradient Boost, The Decision Trees are Pruned To Maximum Leaf nodes from 8 to 32.

Ex:-

| $X_1$<br>Age | $X_2$<br>City | $y$<br>Income |        | Prediction<br>$\hat{y}$ | Residuals<br>Error | D.T-2 |        | $y - \hat{y}$ |
|--------------|---------------|---------------|--------|-------------------------|--------------------|-------|--------|---------------|
| 32           | A             | 51000         |        | 53500                   | -2500              | -5500 | 3000   | 48000         |
| 30           | B             | 78000         | model: | 61000                   | 17000              | 8000  | 9000   | 69000         |
| 21           | A             | 20000         | 1 →    | 28500                   | -8500              | -5500 | -3000  | 23000         |
| 27           | B             | 44000         |        | 61000                   | -17000             | -4300 | -12700 | 56700         |
| 36           | B             | 89000         |        | 90500                   | -15000             | 8000  | -9500  | 98500         |
| 25           | A             | 37000         |        | 28500                   | 8500               | 8000  | 5800   | 36500         |

x

Decision Tree - 2

$(\text{error})_{\min}$

$3000 - 51000 = 48000$

$$\text{MODEL 2 Income} = \text{MODEL 1 income} + \text{Predicted Errors}$$

Model

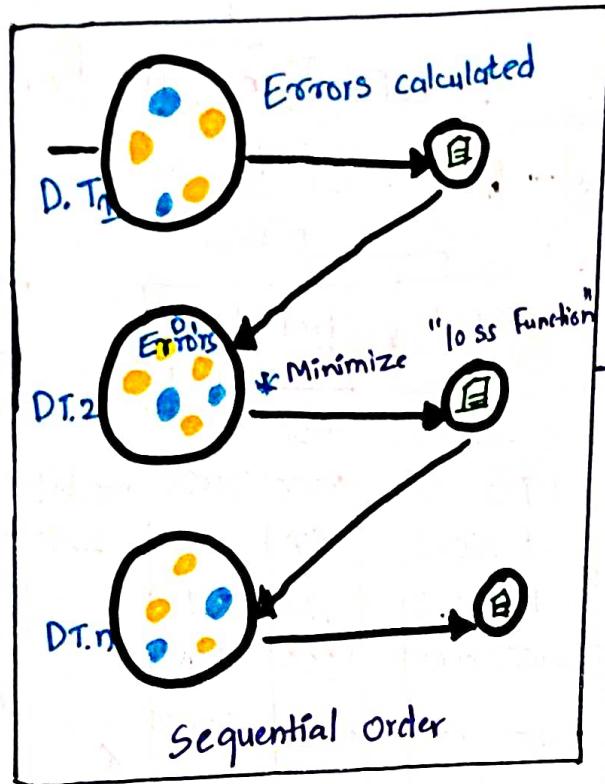
$M_1$

$M_2$

$M_n$

| Feature     | X | X             | X             | X               | X             |
|-------------|---|---------------|---------------|-----------------|---------------|
| Target      | Y | $y - \hat{y}$ | $e_0$         | $y - \hat{y}_0$ | $e_1$         |
| $H_0(x, y)$ |   |               | $H_1(x, e_0)$ |                 | $H_2(x, e_1)$ |
|             |   |               |               |                 | $H_n(x, e_n)$ |

\* it Continues, till it gets "Minimum Error".



- GBM  
"Gradient Boosting method"
- ↓ How it works
1. A loss function to be optimized
  2. A weak learner to make predictions
  3. An additive mode to add weak learners to minimum  
(Or) minimize the loss function.

\* Loss Function :-  $(\text{Sum of Error})_{\min}$   
(Or)

(Overall Error should be Minimum)

Ques : What is difference between cost function & loss function ?

Ans :

Ex :-

|           | Exam | Error | Total = 100 |
|-----------|------|-------|-------------|
| Student 1 | 70   | 30    |             |
| Student 2 | 85   | 15    |             |
| Student 3 | 60   | 40    |             |

cost function.

(Focusing on only one student)

loss function

: Cost function : Error of individual record

: Loss function : Overall Error

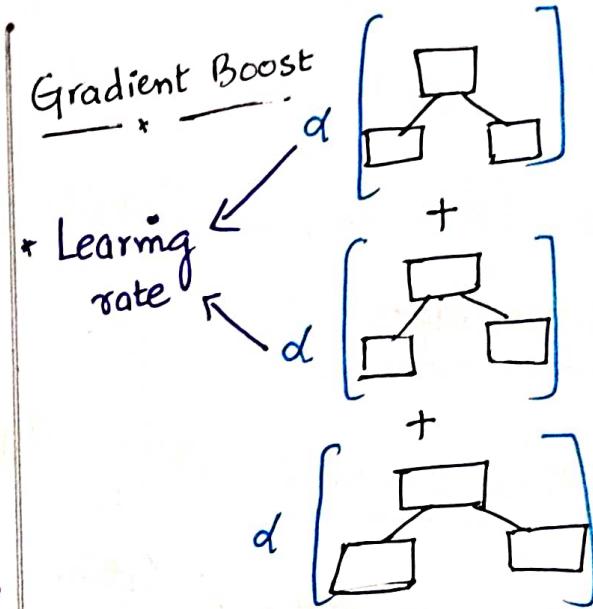
Cost function

Ques:- What is difference b/w "Adaboost" and "gradient boost";

- A:-
- \* Adaboost :- We consider only stump. depth = 1
  - \* Gradient boost :- We consider some what grown The Tree, Max. leaf nodes = 8 to 32

How it works

- \*  $\hat{y}$
- \*  $(y_{\text{actual}} - \hat{y})^2$
- \*  $y = m(x) + \text{error 1}$
- \*  $y = G[x] + \text{error 2}$
- \*  $\text{error 1} = H[x] + \text{error 3}$
- \*  $\text{error 2} = H[x] + \text{error 3}$
- \*  $y = m(x) + G[x] + H[x] + \text{error 3}$
- \*  $y = \alpha * M(x) + \beta * G(x) + \gamma * H(x) + \text{error 4}$
- \*  $y = \text{alpha} * M(x) + \text{beta} * G(x) + \text{gamma} * H(x) + \text{error 4}$



calculate residuals / Errors

Tree after adjusting for Residuals / Errors

Final Tree after adjusting  
Errors multiple times.

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
- .. Cap-Surface : fibrous = f, grooves = g, scaly = y, smooth = -
- Cap-colour : brown = n, buff = b, cinnamon = c, gray = g, green = o, 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, punget = 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()

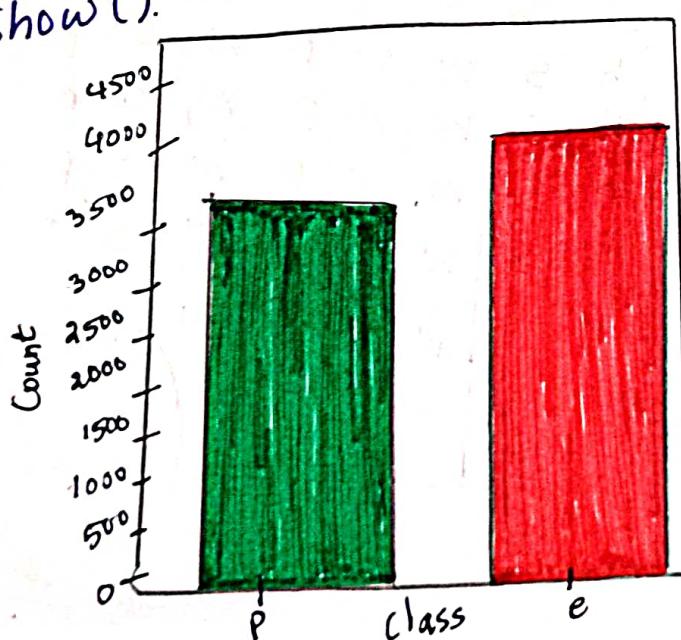
[out]: 0.

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

CODE :

Gradient Boost :-

Same Data. and Same Procedure from import to  
X and Y (Mushroom DataSet)

X & Y

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

#  $y = \text{df}["\text{class}"]$

#  $x.\text{shape}, y.\text{shape}$

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

train Test split

from sklearn.model\_selection import train\_test\_split

#  $x.\text{train}, x.\text{test}, y.\text{train}, y.\text{test} = \text{train\_test\_split}(x, y, \text{test\_size=0.2, random\_state=29})$

modelling

from sklearn.ensemble import GradientBoostingClassifier

# gradmodel = GradientBoostingClassifier()

# gradmodel.fit(x.train, y.train)

Out: GradientBoostingClassifier()

## Prediction

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

## Evaluation

### 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))
```

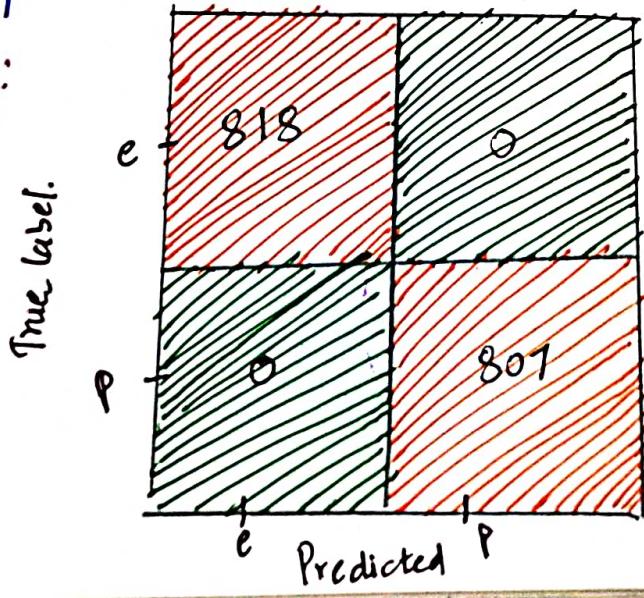
**Out**: train accuracy : 1.0  
Test accuracy : 1.0

### Confusion Matrix

```
* from sklearn.metrics import plot_confusion_matrix  
# plot_confusion_matrix(gradmodel, x-test, y-test)
```

```
# plt.show()
```

**Out**:



## \* Classification report

```
from sklearn.metrics import classification_report
# Print (classification_report(y-test, ypred-test)))
```

**Out** :

|              | Precision | Recall | f1-score | Support |
|--------------|-----------|--------|----------|---------|
| e            | 1.00      | 1.00   | 1.00     | 818     |
| P            | 1.00      | 1.00   | 1.00     | 807     |
| 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

```
from sklearn.model_selection import cross_val_score
# scores = cross_val_score (gradmodel, X, y, cv=5)
# print ("cross-val-score : ", scores.mean())
```

**Out** : cross-val-score : 0.9192

## \* feature importance

```
# gradmodel.feature_importances_
```

**Out** : array ([ 0.00e+00 , 1.0465 , 1.93018 , 0.000e+00 ,  
               6.19492 , 1.36263 , 4.49731 -3.5778 , ...])

## Hyper Parameter Tuning

```
from sklearn.model_selection import GridSearchCV  
# estimator = Gradient Boosting Classifier()  
# param_grid = {"n_estimators": [1, 5, 10, 20, 40, 100], "learning_rate":  
# grid = GridSearchCV(estimator, param_grid, scoring="accuracy")  
# grid.fit(x_train, y_train)  
# grid.best_params_  
out: {"learning_rate": 0.2, "n_estimators": 100}
```

Important Parameter for GradientBoost

## Final Model

```
# final model = Gradient Boosting Classifier(n_estimators=100,  
# final model.fit(x_train, y_train)  
# ypred_train = final model.predict(x_train)  
# ypred_test = final model.predict(x_test)  
# 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
```

```
# final model . feature_importances_
```

```
[out]: array([ 0.000e+00, 2.45745e-16, 1.27008e-08, 4.735  
           1.237228, 1.485440, 1.923168, 3.20965],  
           [0.000bed])
```

```
# important = pd.DataFrame(index=x.columns, data=final model.  
                           feature_importances_, columns=[ "importance_feature" ])
```

```
# important
```

```
[out]: importance_Feature
```

| Cap_shape-c | 0.0000e+00   |
|-------------|--------------|
| Cap_shape-f | 2.45745e-16  |
| Cap_shape-K | 1.270085e-23 |
| Cap_shape-S | 4.735679e-08 |
| ;           | ;            |
| habitat-u   | 9.502012e-05 |
| habitat-w   | 0.00000e+00  |

```
# Jack = important [important [ "importance_Feature" ] > 0.01]
```

```
# Jack. sort_Values ( "importance_Feature" )
```

```
[out]: importance_Feature
```

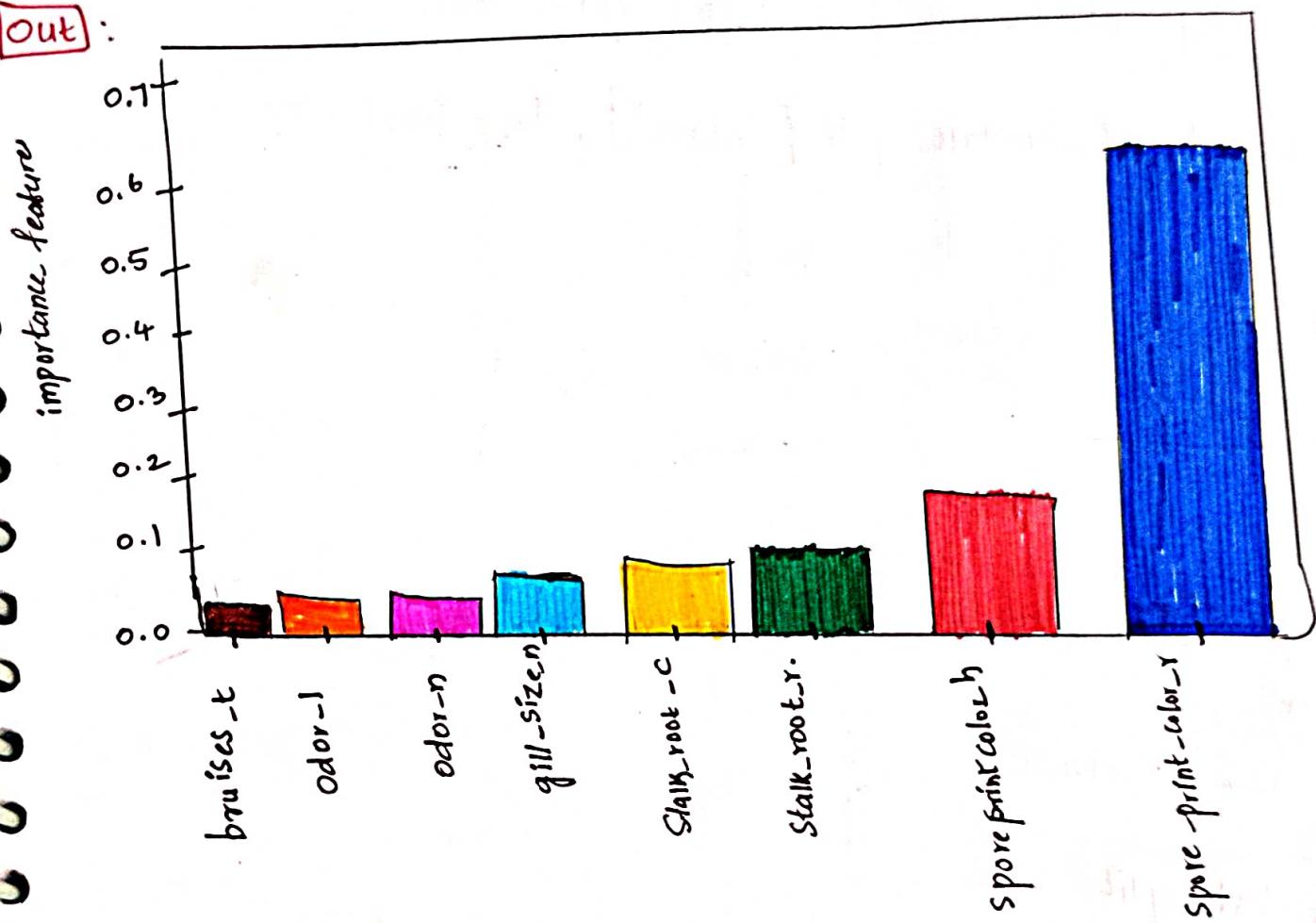
|                       |                |
|-----------------------|----------------|
| gill-size-n           | 0.010084       |
| Spore print color - h | 0.014901       |
| Odor - l              | 0.019874       |
| Spore print color - r | 0.032914       |
| Stalk root - r        | 0.052216       |
| bruises - t           | 0.060678 . . . |

```
# plt.figure(figsize=(14,6), dpi=200)  
# sns.barplot(data=Jacks.sort_values("importance-feature"),  
# x=Jacks.index, y="importance-feature")
```

```
# plt.xticks(rotation=90)
```

```
# plt.show()
```

[Out]:



## XG BOOST

Gradient

Xtreme

\* Adding "penalty" term Reduce overall Error

\* "10" times Faster than Gradient Boosting.

\* XG Boost also called as "Regularized Boosting".

\* Like L<sub>1</sub>, L<sub>2</sub> Regularization

\* Advantages

\* Tree Pruning

Using depth first approach

\* High-flexibility

- Catche awareness and out of core computing. (-) it uses all cores in system

\* Reduce Overfitting:

"Regularization" for avoiding overfitting

Efficient Handling

missing data

Inbuilt

"cross-validation" Capability

XG-boost

error + Penalty Term

$$L_2: \frac{d}{2} (\text{slope})^2$$

$$L_1: \frac{d}{2} |\text{slope}|$$

with in XGBoost :-  
 $\alpha = \frac{1}{\gamma}$

## \* XG boost classifier      Black box Model.

Base model = Probability = 0.5

Ex:- DataSet

(Pr - Approval)

| Salary     | Credit | Approval | Residuals                 |
|------------|--------|----------|---------------------------|
| $\leq 50k$ | Bad    | 0        | 0 - 0.5<br>- 0.5<br>- 0.5 |
| $\leq 50k$ | Good   | 1        | 0.5                       |
| $\leq 50k$ | Good   | 1        | 0.5                       |
| $> 50k$    | Bad    | 0        | - 0.5<br>0.5              |
| $> 50k$    | Good   | 1        | 0.5                       |
| $> 50k$    | Normal | 1        | 0.5                       |
| $\leq 50k$ | Normal | 0        | - 0.5                     |

## \* Steps :-

1. Create Binary Decision Tree using Features

2. Calculate

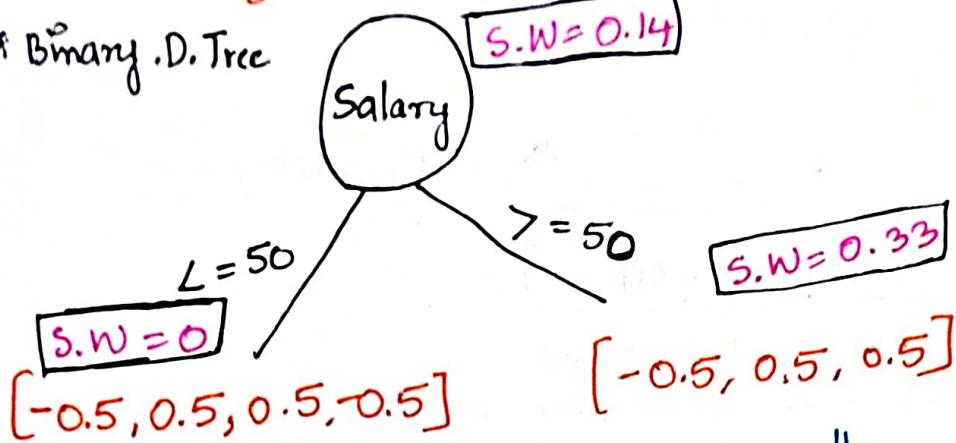
Similarity weight

$$= \frac{E [ \text{Residual} ]^2}{E [ \text{Probability} [1 - \text{probability}] ]} \rightarrow + \lambda (\text{hyper parameter})$$

3. Information Gain.

$$\text{Errors} := [-0.5, 0.5, 0.5, -0.5, 0.5, 0.5, 0.5, -0.5]$$

# Binary D. Tree



$$\begin{aligned}
 &\Downarrow \\
 \text{Similarity weight} &\Rightarrow \frac{E(\text{res})}{EP(1-P)\lambda} \\
 &= \frac{[-0.5 + 0.5 + 0.5 - 0.5]}{0.5[1-0.5] + 0.5[1-0.5]} + \underset{\because \lambda=0}{\cancel{0}} \\
 &= \frac{0.5[1-0.5] + 0.5[1-0.5]}{0.5[1-0.5] + 0.5[1-0.5]} \\
 &= 0
 \end{aligned}$$

$$\begin{aligned}
 \text{Similarity weight} &= \frac{E(\text{res})^2}{EP(1-P)\lambda} \\
 &= \frac{[-0.5 + 0.5 + 0.5]^2}{0.5[1-0.5] + 0.5[1-0.5]} \\
 &= \frac{0.5[1-0.5] + 0.5[1-0.5]}{0.5[1-0.5] + 0.5[1-0.5]} \\
 &= \frac{0.25}{0.75} = \frac{1}{3} \Rightarrow 0.33
 \end{aligned}$$

$$\# \text{ for Root node} \quad \text{Similarity weight} = \frac{E(\text{res})^2}{EP(1-P)\lambda}$$

$$\begin{aligned}
 &= \frac{[-0.5, 0.5, 0.5, -0.5, 0.5, 0.5, -0.5]^2}{0.5[1-0.5] + 0.5[1-0.5] + 0.5[1-0.5] + 0.5[1-0.5] \\
 &\quad + 0.5[1-0.5] + 0.5[1-0.5]} \\
 &= \frac{0.25}{1.75} = \frac{1}{7} = 0.14
 \end{aligned}$$

$$* \text{ information Gain} = 0 + 0.33 - 0.14$$

$\hookrightarrow = 0.19$

Error: [0.5, 0.5, 0.5, -0.5, 0.5, 0.5, -0.5]

S.W = 0.14

Salary



E = 0.5

$$* S.W = \frac{E(\text{res})^r}{E(p(1-p))}$$

$$* S.W = \frac{E(\text{res})^r}{E(p(1-p))}$$

$$= [0.5, 0.5, -0.5]^r$$

$$= \frac{0.5[1-0.5] + 0.5[1-0.5] + 0.5[-0.5]}{0.5[1-0.5] + 0.5[1-0.5] + 0.5[-0.5]}$$

$$= \frac{0.25}{0.75} = \frac{1}{3} = 0.33$$

= 1

$$* \text{ information Gain} = 1 + 0.33 - 0 \quad ?$$

$$= 1.33$$

| Salary | credit | Approval |
|--------|--------|----------|
| L = 50 | B      | 0        |

activation function  
sigmoid  $\sigma$   $[0 + \alpha(1)] \Rightarrow$  Based on "d" value  $\rightarrow 0/p [0-1]$

Base model = 0.5

Real probability =

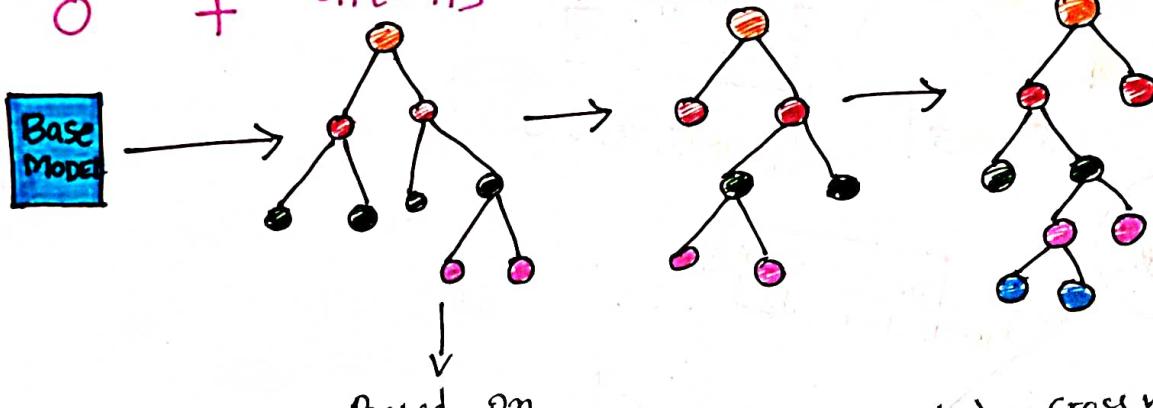
$$\log \frac{p}{1-p} = \log \left( \frac{0.5}{1-0.5} \right) = 0$$

Output :-

$$\sigma [o + \alpha_1 [DT_1] + \alpha_2 [DT_2] + \alpha_3 [DT_3] + \dots + \alpha_n [DT_n]]$$

For any new record  $\uparrow$

$$o + \alpha_1 [DT_1] + \alpha_2 [DT_2] + \alpha_3 [DT_3] + \dots$$



Based on  
independent features.

$\therefore \lambda = \text{cross validation.}$

## \* XG boost regressor

\* Base model = 51K

Ex:-

| $x_1$ | $x_2$ | $y$    | Resid<br>-<br>base |
|-------|-------|--------|--------------------|
| Exp   | Gap   | Salary |                    |
| 2     | Yes   | 40 K   | -11 K              |
| 2.5   | Yes   | 42.5   | 9                  |
| 3     | No    | 52 K   | 1                  |
| 4     | No    | 60 K   | 9                  |
| 4.5   | Yes   | 62 K   | 11                 |

$$\text{Avg} = 51.$$

$$S.W = \frac{1}{6}$$

$$[-11, -9, 1, 9, 11]$$

$$\frac{[-11, -9, 1, 9, 11]}{5+1} = \frac{1}{6}$$

Experience

$$S.W = \frac{1}{6} \cdot 2$$

$$[-11]$$

$$S.W = 28.8$$

$$7.2$$

$$[-9, 1, 9, 11]$$

\* Similarity weight

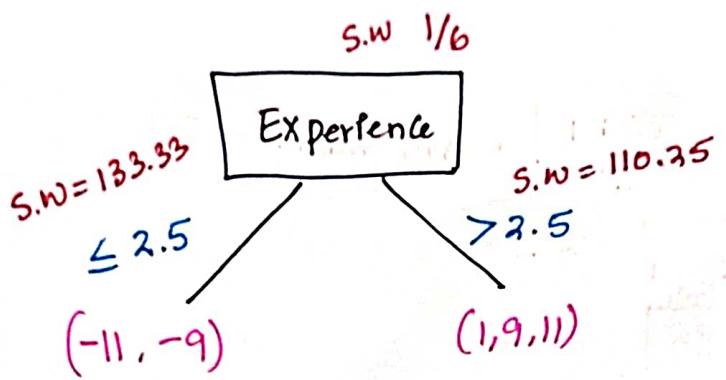
$$\text{Similarity weight} = \frac{\sum [\text{Residual}]^2}{\text{no. of residuals} + \lambda}$$

$$\frac{[-9, 1, 9, 11]^2}{4 + 1} = \frac{12^2}{5} = \frac{144}{5} = 28.8 \quad \lambda = 1$$

$$= \frac{[-11]^2}{1+1} = \frac{121}{2} = 60.5$$

$$\therefore \lambda = 1$$

$$\text{Information Gain} = 60.5 + 28.8 - \frac{1}{6} = 89.13$$



$$\begin{aligned}
 S.W &= \frac{[-11-9]^2}{2+1} = \frac{[1+9+11]^2}{3+1} \\
 &= \frac{[-20]^2}{3} = \frac{[21]^2}{4} \\
 &= \frac{400}{3} \Rightarrow 133.33
 \end{aligned}$$

any records goes from base model = 51

|                                                                                                                  |                                                                                                                      |
|------------------------------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------------------------|
| $\leq 2.5$<br>Arg. of S.W<br>$51 + \alpha_1 \left[ \frac{[-11-9]}{2} - (-10) \right]$<br>$* 51 + \alpha_1 [-10]$ | Avg. of S.W<br>$51 + \alpha_1 \left[ \frac{1+11+9}{3} \right] \Rightarrow \frac{21}{3} = 7$<br>$* 51 + \alpha_1 [7]$ |
|------------------------------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------------------------|

$$* 51 + \alpha_1 (-10) + \alpha_2 [DT_2] + \alpha_3 [DT_3] \dots \alpha_n [DT_n]$$

↳ Output:

21/04/22  
6:00 AM.

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
- .. Cap-Surface : fibrous = f, grooves = g, scaly = y, smooth = -
- Cap-colour : brown = n, buff = b, cinnamon = c, gray = g, green = o, 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, punget = 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 color below ring | stalk color above 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                      | p         | w          | o           | p         | b                 | s          | u       |
| E     | x         | s           | y         | t       | a    | f               | c            | b         | n          | s                        | w                      | w                      | p         | w          | o           | p         | n                 | n          | g       |
| e     | b         | s           | w         | t       | p    | f               | c            | n         | n          | s                        | w                      | w                      | p         | w          | o           | p         | k                 | s          | u       |
| P     | x         | y           | w         | f       | n    | f               | w            | b         | b          | s                        | w                      | w                      | p         | w          | o           | e         | n                 | a          | g       |
| e     | x         | s           | g         | f       | n    | f               | w            | b         | b          | s                        | w                      | w                      | p         | w          | o           |           |                   |            |         |

# 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()

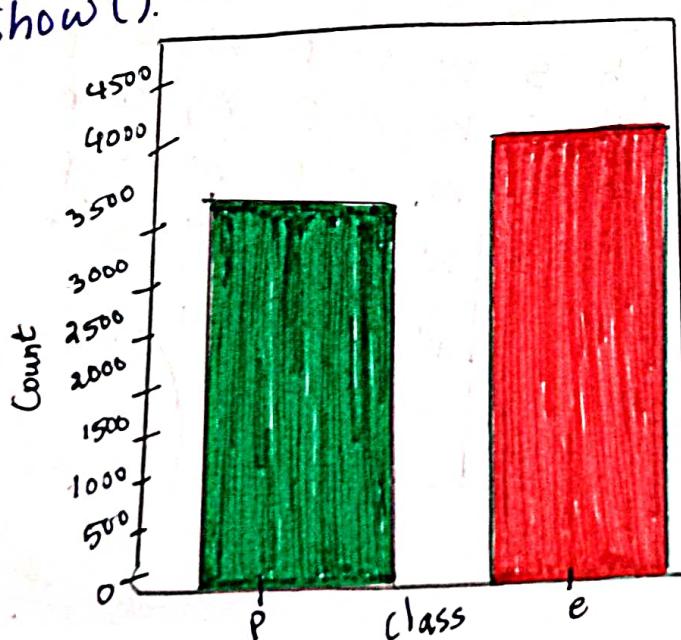
[out]: 0.

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 |

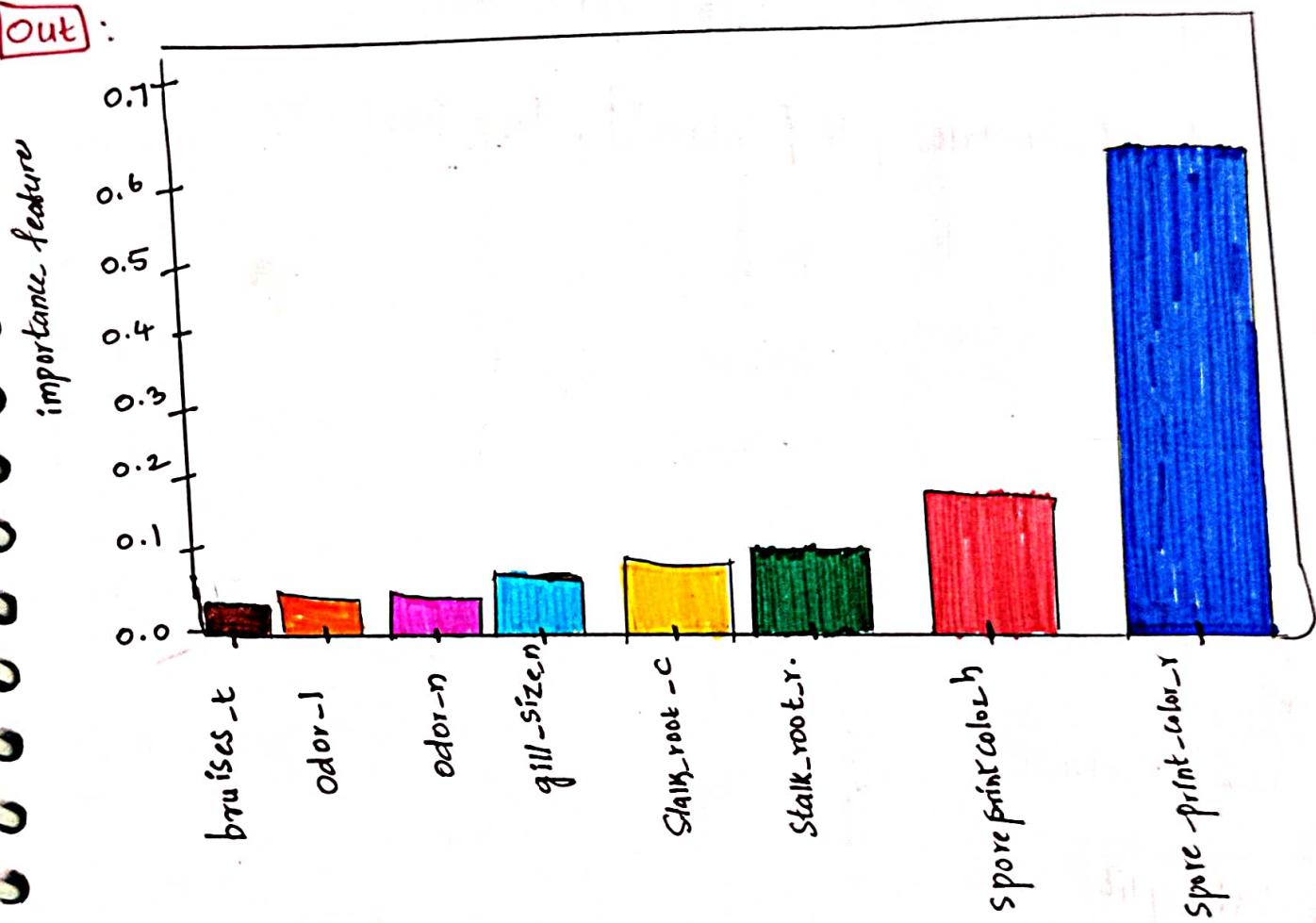
\* posinouss  
\* non posinouss

```
# plt.figure(figsize=(14,6), dpi=200)  
# sns.barplot(data=Jacks.sort_values("importance-feature"),  
# x=Jacks.index, y="importance-feature")
```

```
# plt.xticks(rotation=90)
```

```
# plt.show()
```

[Out]:



Code :

## XG boost Classifier

Same Example, Same Question of

Mushroom dataset.

import to X & Y.

X & Y

```
# X = pd.get_dummies(df.drop("class", axis=1), drop_first=True)  
# y = pd.get_dummies(df[["class"]], drop_first=True)
```

# y

out : P

|      | P          |
|------|------------|
| 0    | 1          |
| 1    | 0          |
| 2    | 0          |
| 3    | 1          |
| 4    | 0          |
| :    | :          |
| 8124 | x 1 column |

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

! Pip install Xgboost

```
from xgboost import XGBClassifier  
# xgb-model = XGBClassifier()  
# xgb-model.fit(x-train, y-train)  
[Out]: XGBClassifier(base_score=0.5, booster='gbtree', ...)
```

## Prediction

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

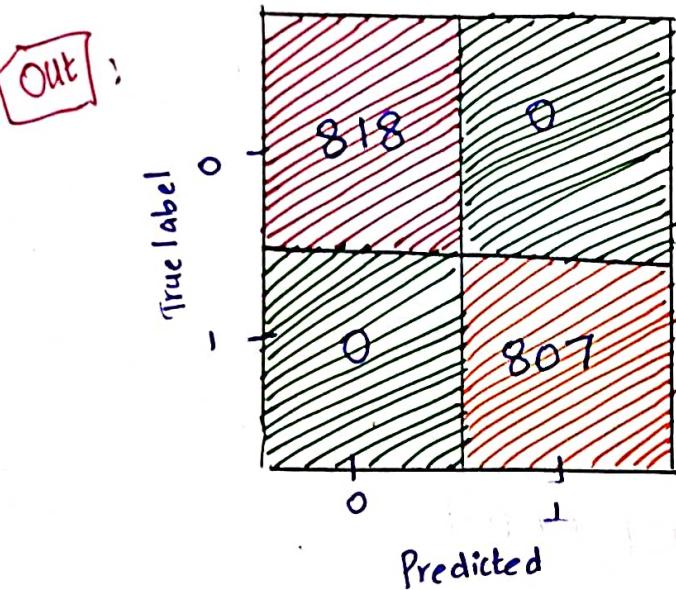
## Evaluation

### \* 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

```
from sklearn.metrics import plot_confusion_matrix  
# plot_confusion_matrix (ngb-model, x-test, y-test)  
# plt.show()
```



## \* classification report

```
from sklearn.metrics import classification_report  
# Print (classification_report (y-test, ypred-test))
```

Out:

|              | Precision | recall | f1 score | Support |
|--------------|-----------|--------|----------|---------|
| 0            | 1.00      | 1.00   | 1.00     | 818     |
| 1            | 1.00      | 1.00   | 1.00     | 807     |
| 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

```
from sklearn.model_selection import cross_val_score
# scores = cross_val_score(xgb_model, x, y, cv=5)
# print("cross-val-score:", scores.mean())
Out: cross-val-score: 0.935
```

## \* Feature importance

```
# xgb_model.feature_importances_
Out: array([ 0.0000e+0,  0.0000e+0,  0.0000e+0,  0.0000e+0,
            3.961512e-5,  0.00000e+0, 9.23225e-0])
```

## Hyper Parameter tuning

Important parameters in XGBoost

```
from sklearn.model_selection import GridSearchCV
# estimator = XGBClassifier()
# param_grid = {"n_estimators": [1, 5, 10, 20, 40, 100], "gamma": [0, 1, 0.5, 0.9], "max_depth": [2, 6]}
# grid = GridSearchCV(estimator, param_grid, scoring="accuracy")
# grid.fit(x_train, y_train) CV=5
# grid.best_params_
Out: {"gamma": 0.1, "max_depth": 2, "n_estimators": 20}
```

## final Model

```
# final_model = XGBClassifier(n_estimators=20, gamma=0.1, max_depth=2)  
# final_model . fit (x_train, y_train)  
  
# y_predictions_test = final_model. Predict (x-test)  
# y_predictions_train = final_model. Predict (x-train)  
  
# print ("train accuracy:", accuracy_score (y_prediction_train, y_train))  
# print ("test accuracy:", accuracy_score (y_prediction_test, y-test))
```

**out**: Test accuracy : 0.99815  
Train accuracy : 0.9979

```
# final_model. feature_importances_
```

**out**: array ([0., 0., 0., 0.,  
0., 0., 0., 0.])

```
# important = pd. DataFrame (index=x.columns, data=final-  
model. feature_importances_, columns=[importance-  
Features])
```

```
# Jack = important[important["importance-feature"] > 0.01]
```

```
# Jack . sort_values("importance-feature")
```

### important\_features

**Out:**

|                     |          |
|---------------------|----------|
| odor-p              | 0.015831 |
| spore-print-Color-W | 0.020152 |
| :                   | :        |
| odor-n              | 0.290767 |

```
# plt.figure(figsize=(14,6), dpi=200)
```

```
# sns.barplot(data=Jack.sort_values("importance-feature"),  
#               x=Jack.index, y="importance-feature")
```

```
# plt.xticks(rotation=90)
```

```
# plt.show()
```

**Out:**

