**\*\*\*** ==Gradient Boosting==

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

D.T-2

Ex:-

using D.T-1

| $X_1$ Age | $X_2$ City | $y$ Income | | Predction $\hat{y}$ | Residuals Error | Predi ction | Residu als error | | $y$-Err |
|---|---|---|---|---|---|---|---|---|---|
| 32 | A | 51000 | | 53500 | – 2500 | -5500 | 3000 | | 48000 |
| 30 | B | 78000 | model: | 61000 | 17000 | 8000 | 9000 | model | 69000 |
| 21 | A | 20000 | 1→ | 28500 | -8500 | -5500 | -3000 | 2→ | 23000 |
| 27 | B | 44000 | | 61000 | –17000 | -4300 | -12700 | | 56700 |
| 36 | B | 89000 | | 90500 | –15000 | 8000 | -9500 | | 98500 |
| 25 | A | 37000 | | 28500 | 8500 | | | | 36500 |
| | | | | | | 8000 | 5800 | | |

X

3000-5100
= 48000

↓Target
(error)$_{min}$

Decision Tree – 2

| MODEL 2 Income | = | MODEL 1 income | + | Predicted Errors |
|---|---|---|---|---|

(Model 0) ..... (M₁) ..... (M₂) ..... (Mₙ)

| Feature | X | X | X | X |
|---|---|---|---|---|
| Target | Y   y-ý→ | e₀   y-y⁰→ | e₁   ... | eₙ₋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".

Errors calculated

D.T$_1$

Errors

*Minimize "loss function"

D.T.2

D.T.n

Sequential Order

→ 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)

Que :- What is difference between cost function & loss function?

Total = 100

Ans :- Ex :-

| | Exam | Error |
|-----------|------|-------|
| Student 1 | 70 | 30 |
| Stundent 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

Que :- What is difference b/w "Adaboost" and "gradient boost$^{-n}$"

A:- * Adaboost :- We Consider Only Stump depth = 1

* Gradient boost :- We Consider Some what grown
   The Tree, Max. Leaf nodes = 8 to 32

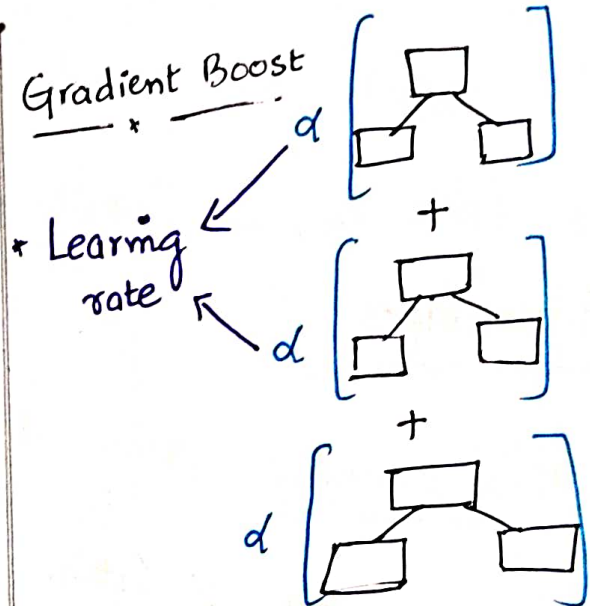How it works *

* $\hat{y}$ hat

* $[y\_actual - y\_hat]^2$

* $y = m(x) + error\ 1$

* $error\ 1 = G[x] + error\ 2$

* $error\ 2 = H[x] + error\ 3$

* $y = m(x) + G[x] + H[x] + error\ 3$

* $y = alpha * M(x) + beta * G(x) + gamma * H[x] + error\ 4$

Gradient Boost *

* Learning rate



1st Tree . Calculate residuals /Errors

2nd Tree . Tree after adjusting for Residuals / Errors

Nth Tree . 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 in indentified as /definetly edible/, /definitely poisonous / or of unknown edibility and not recommended. this latter class was Was combined with the poisinous one. the Guide Clearly states that there is no simple rule for determing the edibility of a mushroom; no rule like "leaflets three, Let it be " for poisinous Oak and lvy.

### Attribute Information →

 → Mushroom

1. 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 - Colour :   brown = n, buff = b, Cinnamon = c, gray = g, green = r, pink = p, purple = u, red = e, white = w, Yellow = y.

ruises ? : bruises = t, no _F

dor : almond = a, anise = l, Creoste = 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 = r, 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 = n

1. gill-color : black = k, brown = n, buff = b, chocolate = h, grave = 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 - 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 = 0, two = t

19. ring - type :- cobwebby = C, evanescent = e, flaring = f, large = l, none = D, Pendent = P, sheathing = S, Zone = Z,

20. Spore - print Colour :- black = K, brown = n, buff = b, chocolate = h, green = r, Orange = 0, Purple = u, white = w, Yellowsy

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)

## Bussiness Problem

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

# df = Pd. redd_csv ( "mushrooms. csv")

# df. head ()

| Class | Cap Shape | Cap Surface | Cap Color | bruise | Odor | gill attachment | gill spacing | gill size | gill color | ... | stalk surface above Ring | stalk color above ring | stalk color below ring | veil type | veil color | ring number | ring type | Spore Print Color | Population | habitat |
|-------|-----------|-------------|-----------|--------|------|-----------------|--------------|-----------|------------|-----|--------------------------|------------------------|------------------------|-----------|------------|-------------|-----------|-------------------|-----------|---------|
| | | | | | | | | | | | | | | | | | | w | 0 | P | K | S | u |
| P | X | s | n | t | P | f | c | n | K | | S | w | w | P | w | 0 | P | n | n | g |
| e | X | s | y | t | a | f | c | b | K | | S | W | W | P | w | 0 | P | n | n | m |
| e | b | s | w | t | l | f | c | b | n | | S | w | W | P | w | 0 | P | K | s | u |
| P | X | y | w | t | P | f | c | n | n | | S | w | W | P | w | 0 | P | n | s | u |
| P | X | s | g | f | n | f | w | b | K | | S | w | W | P | w | 0 | e | n | a | g |
| e | X | s | g | | | f | | | | | | | | | | | | | | |

# df. shape

Out : [8124, 23]

# df. info ()

Out : Range Index : 8124 Entries, 0 to 8123

| column | Non null Count | Dtype. |
|---|---|---|
| * class | 8124 Nonull | object |
| * Cap-Shape | " | " |
| * Cap Surface | " | " |
| . | | |
| * habitat | " | |

⊦ df. ismull (). Sum ()

ut : 0.

EDA

: sns. Count plot ( data = df, X = "class", palette = "Dark 2")
                                    ‾‾‾‾‾
                                    output

plt. show ().

# 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 ring | 8124 | 4 | s | 5176 |
| Stalk - Surface - below ring | 8124 | 4 | s | 4936 |
| Stalk - Color above ring | 8124 | 9 | S | 4464 |
| Stalk - color below ring | 8124 | 9 | w | 4384 |
| Veil - type | 8124 | 1 | w | 8124 |
| Veil - color | 8124 | 4 | P | 7924 |
| ring - number | 8124 | 3 | w | 7488 |
| ring - type | 8124 | 5 | o | 3968 |
| Spore - Print color | 8124 | 9 | P | 2388 |
| Population | 8124 | 6 | w | 4040 |
| habitant | 8124 | 7 | v | 3148 |
| | | | d | |

\* Posimous
\* non posinous

CODE :

Same Data. and Same Procedure from import to

X and Y   (Mushroom DataSet)

X & Y

```
#  x = pd.get-dummies (df.drop("class") axis = 1 drop_first = True)

#  y = df ["class"]


#  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.
                                      random_state = 29)
```

modelling

```
from  sklearn.ensemble import Gradient Boosting Classifier

#  gradmodel = Gradient Boosting classifier ()

#  gradmodel. fit (x_train, y_train)
```
Out :  Gradient Boosting classifier ()

## Prediction

```
# ypred_train = grad model . predict (x_train)
# ypred_test = grad model . predict (x_test)
```

## Evaluation

* **Accuracy**

```
from sklearn.metrics import accuracy_score
# print ("train accuray", accuracy_score (y_train, ypred_train)
# print ("test accuray", accuracy_score (y_test, ypred_test)
```

Out: train accuray : 1.0
     Test accuray : 1.0

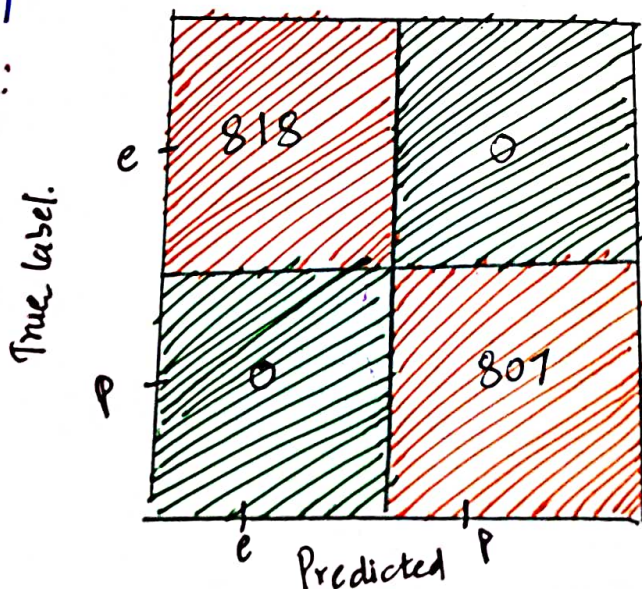* **Confusion Matrix**

```
from sklearn.metrics import plot_confusion_matrix
# plot_confusion_matrix (grad model , 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 | f_score | Support |
|---|---|---|---|---|
| | | | | 818 |
| e | 1.00 | 1.00 | 1.00 | 807 |
| P | 1.00 | 1.00 | 1.00 | |
| | | | | 1625 |
| 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 ([[ 0.00 e+00 , 1.0465 , 1.93018 , 0.000e+0,
6. 19492 , 1.36 263 , 4.49731 , -3.5778 , ...]]

## Hyper Parameter Tunning

from sklearn.model_selection import GridSearchCV

\# estimator = GradientBoostingClassifier()

\# param_grid = { "n_estimators" : [1,5,10,20,40,100], "learning rate" :    *Important Parameter for Gradient Boost*

                                                  [0.1, 0.2, 0.5, 0.9] }

\# grid = GridSearchCV (estimator, param_grid, Scoring = "accuracy")
                                , cv = 5

\# grid. fit (x_train, y_train)

\# grid. best_params_

[out] : { "learning-rate" : 0.2 , "n_estimators" : 100)

## Final Model

\# final model = GradientBoostingClassifier ( n_estimators = 100, learning_rate = 0.2 )

\# 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, 1.27008, 4.735
1.137228        1.485440        1.923168 , 3.20965)
                                0.000ber]]

# important = pd. Data Frame (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")

## importance_features

Out :

| | importance_features |
|---|---|
| gill_size_n | 0.010084 |
| Spore print color_h | 0.014901 |
| Odor_l | 0.019874 |
| Spore print color_y | 0.032914 |
| Stalk root_r | 0.052216 |
| bruises_t | 0.060678 . . . . |

```python
# plt.figure (figsize=(14,6), dpi=200)
# sns.barplot (data = Jack.sort_values ("importance_features",
                    x = Jack.index, y = "importance_feature")

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

Out: