

## 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 - 0.5 - 0.5
$\leq 50k$	Good	1	0.5
$\leq 50k$	Good	1	0.5
$> 50k$	Bad	0	- 0.5 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

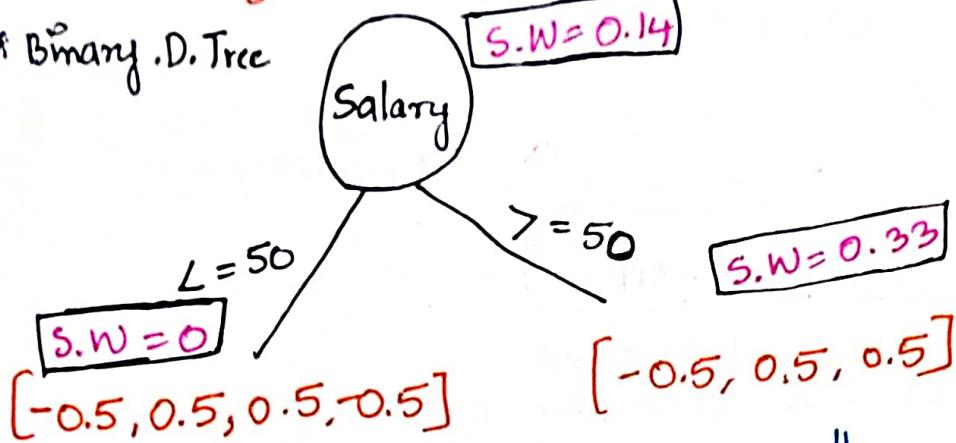
Similarity weight

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

3. Information Gain.

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

# Binary D. Tree



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

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

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

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

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

↓ = [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

$$* S.W = \frac{E(\text{res})^r}{\overline{EP(1-P)}\lambda}$$

$$* S.W = \frac{E(\text{res})^r}{\overline{EP(1-P)}\lambda}$$

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

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

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

= 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  $\rightarrow 0/1 [0-1]$

Base model = 0.5

Real probability =

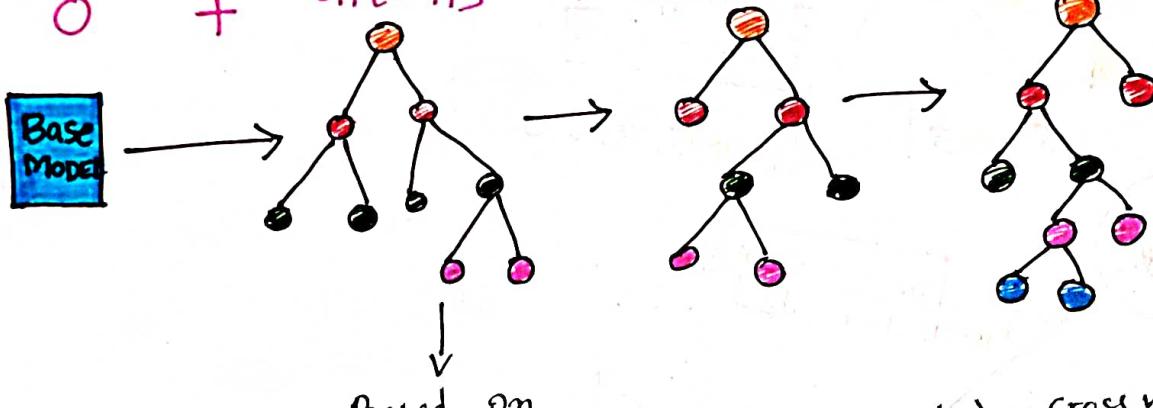
$$\log \frac{(P)}{(1-P)} = \log \left( \frac{0.5}{1-0.5} \right) = 0$$

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 - 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

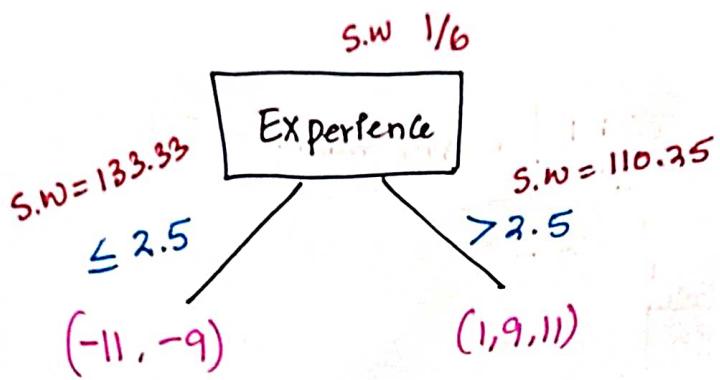
$$\text{Similarity weight} = \frac{\sum [\text{Residual}]^2}{\text{no. of residuals} + \lambda}$$

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

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

$$\therefore \lambda = 1$$

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



$$S.W = \frac{[-11-9]^2}{2+1} = \frac{[1+9+11]^2}{3+1}$$

$$= \frac{[-20]^2}{3} = \frac{[21]^2}{4}$$

$$= \frac{400}{3} \Rightarrow 133.33$$

$$= 110.25$$

any records goes from base model = 51

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

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

↳ 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 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()

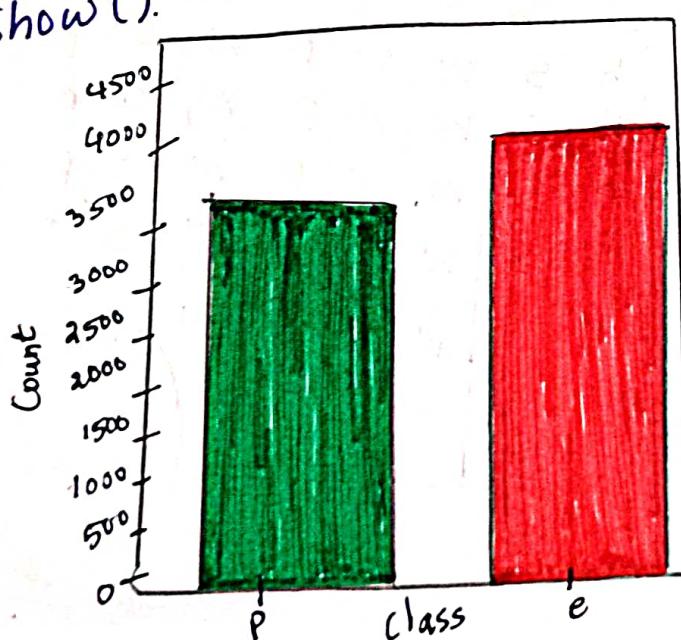
[out]: 0.

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 ring	8124	4	s	5176
Stalk - Surface - below ring	8124	4	s	4936
Stalk - color above ring	8124	9	w	4464
Stalk - color below 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

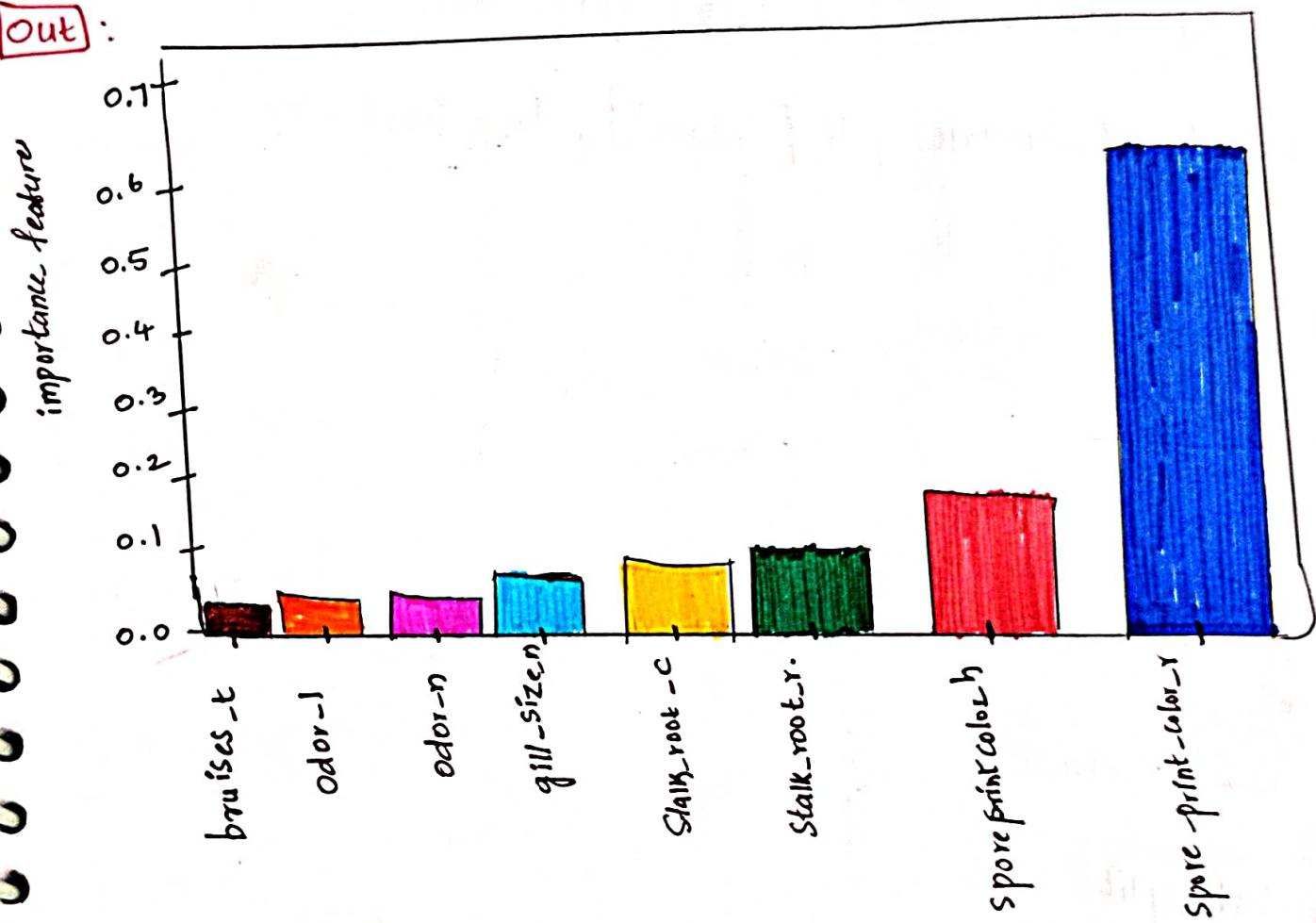
\* posinouss  
\* non posinouss

```
# 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))
```

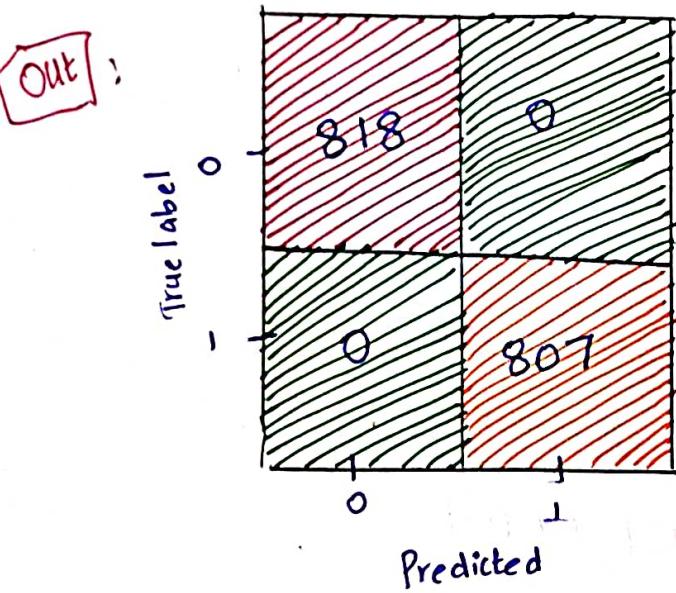
Train accuracy : 1.0

Test accuracy : 1.0

```
[Out]:
```

## \* Confusion Matrix

```
from sklearn.metrics import plot_confusion_matrix  
# plot_confusion_matrix (nbg-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

```
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}
```

Important parameters in XGBoost

## 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:**

