

~~09.10.2023
Wednesday~~

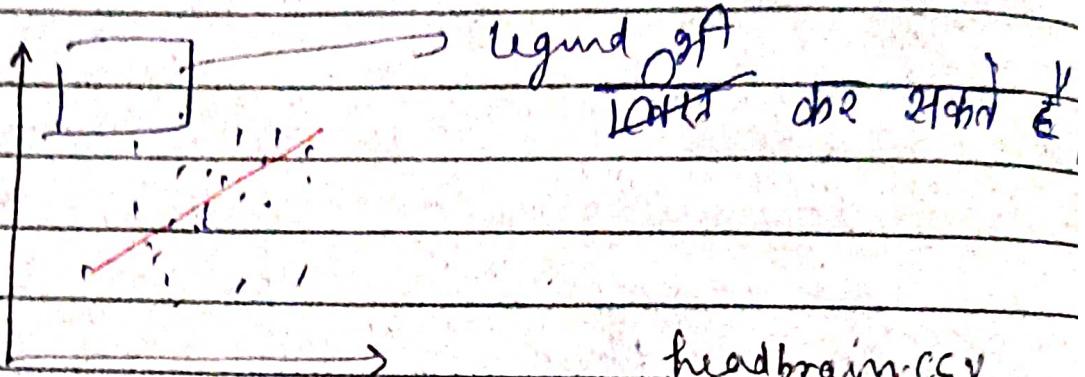
`.ravel()` → contiguous array

`import matplotlib.pyplot as plt`

`plt.scatter(x, y, color='r')` → Data point

`plt.plot(x, pred_y, color='g')` → Regression line

`plt.show()`



M.2: Scikit Learn Method :-

install sklearn.

Age	Gender	headsiz	brainw

```
from sklearn.linear_model import LinearRegression  
import numpy as np  
import pandas as pd
```

`data = pd.read_csv('headbrain.csv')`

`x = data['headsiz']`

`y = data['brainweight']`

`reg = LinearRegression()`

`reg.fit(x, y)`

`y_pred = reg.predict(x)`

Logistic Regression

↳ supervised
M.L. Algo.

import matplotlib.pyplot as plt

per = scatter(x, y) # points = original

plt. plot(x, y-pred) # Regression = Predicted values

LOGISTIC REGRESSION :-

from sklearn.linear_model import LogisticRegression

import pandas as pd each has

import numpy as np so corresponding samples

from sklearn.model_selection

import train-test-split

from sklearn.metrics import confusion-
matrix

data = pd.read_csv('iris.csv')

x = data[['sepal length', 'sepal width', 'petal length',
'petal width']]

y = data['variety']

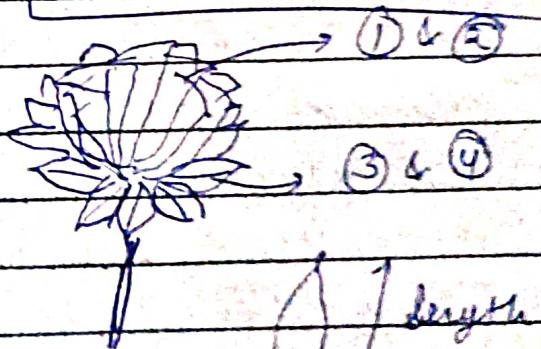
x-train, x-test, y-train, y-test = train-test-split(x, y,
test_size=0.3, random_state=1)

70% Training data set

30% Testing data set

(maximum samples $\frac{42}{70}$ Train $\frac{18}{30}$ Test)

150 rows
iris.csv
iris flower has
3 categories.
setosa, versicolor,
virginica



① & ② length
③ & ④ width

↔ width

We can also use cross validation.
random state = 1 \Rightarrow dividing randomly
↳ if we again execute the program then we will
get same samples in training & testing set.
for multiple iteration & we will get same performance
- wise. (until we have same random state)

Constructing the model :-

logReg = LogisticRegression()

Training :-

logReg.fit(x-train, y-train)

Prediction :-

y-pred = logReg.predict(x-test)

print (confusion_matrix(y-test, y-pred))
predicted

prec. & recall will

be calculated

individually &
then take avg.

	setosa	versicolor	virginica
Actual	8	1	1
versicolor	0	9	1
virginica	2	1	7

$$\text{Accuracy} = \frac{8+9+7}{30} = 0.8$$

Actual	Setosa	Versicolor	Virginica	Avg.	Accuracy
Predicted	$\frac{8}{10}$	$\frac{9}{10}$	$\frac{7}{10}$	$\frac{0+9+1}{3} = 0.8$	0.8
Accuracy Precision	$\frac{8}{10}$	$\frac{9}{11}$	$\frac{7}{9}$	2.396	100%

out of total predicted satara many are correct.

FPI

Date: _____ Page: _____

1/19/23

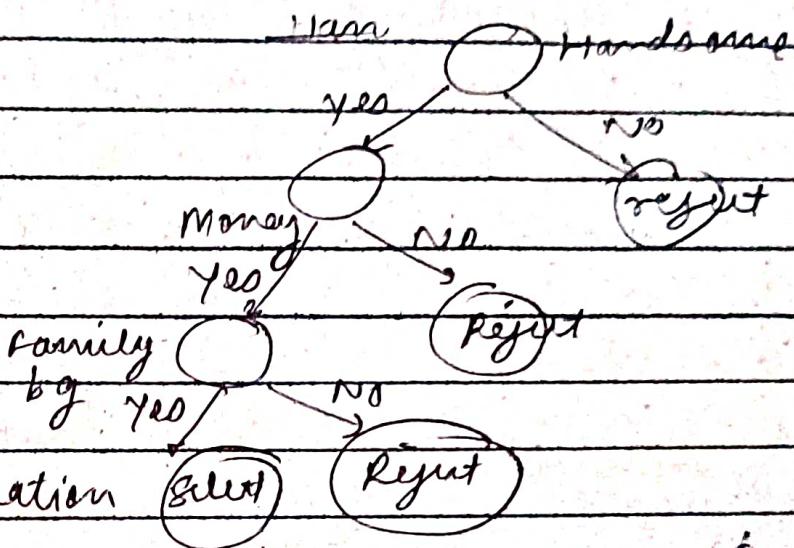
$$f_1\text{-score} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} = 0.500$$

DECISION TREE

- ① root node
- ② leaf node
- ③ parent node
- ④ child node
- ⑤ subtree

supervised ml

algo. used for classification & regression



ASM Attribute Selection metric :-

(Evaluation on basis of)

① Entropy & Information Gain (IG)

randomness

due to feature (Impurity)

→ purity

of a sample
(dataset)

IG

Evaluation
on
basis

② Gini Index

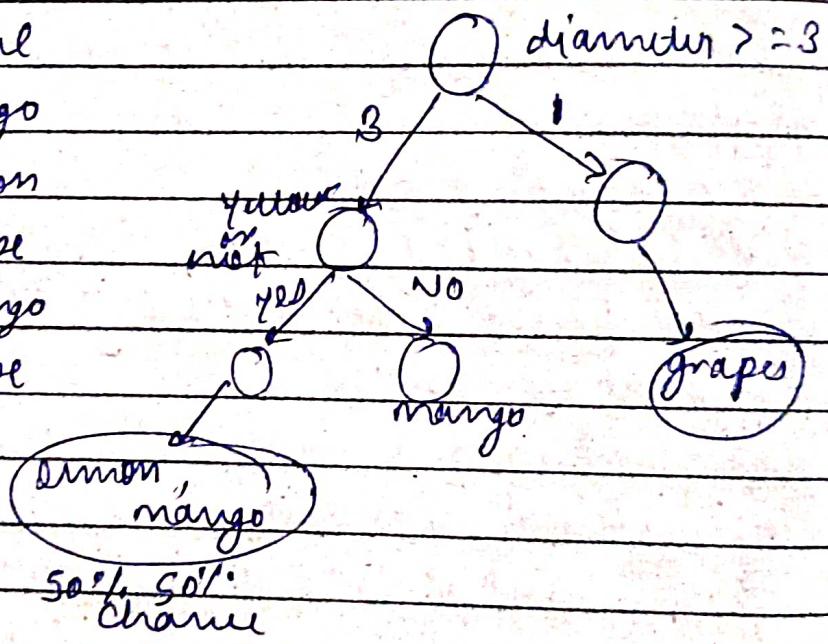
Decision Tree :-

- ① Start with the root node, having complete dataset

as an input.

- (i) Selection of most significant attributes based on attribute selection metric (IG)
- (ii) Evaluate the attribute & divide the tree into subtree based on possible values of attribute.
- (iii) Repeat step (ii) & step (iii) for each subtree until all the attributes are evaluated.

Color	Diameter	Class
Green	3	Mango
Yellow	3	Lemon
Red	1	Grape
Yellow	3	Mango
Red	1	Grape



Entropy of dataset, S $E(S) = - \sum p_i \log_2 p_i$

$$= - \frac{6}{10} \log_2 \frac{6}{10} - \frac{4}{10} \log_2 \frac{4}{10}$$

$$= 0.97$$

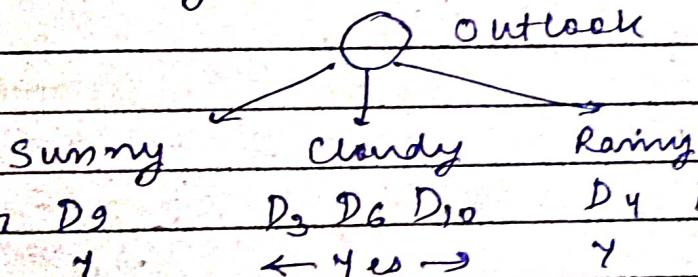
$$IG(\text{Outlook}) = 0.37$$

$$IG(\text{Humidity}) = 0.045$$

$$IG(\text{Wind}) = 0.02$$

10.10.2023
Tuesday

Day	Outlook	Humidity	Wind	Play
1	Sunny	High	Weak	No
2	Sunny	High	Strong	No
3	Cloudy	High	Weak	Yes
4	Rainy	High	Weak	Yes
5	Rainy	Normal	Strong	No
6	Cloudy	Normal	Strong	Yes
7	Sunny	High	Weak	No
8	Rainy	Normal	Weak	Yes
9	Sunny	Normal	Strong	Yes
10	Cloudy	High	Weak	Yes



Wind . Wind

D1	High	Weak	No
D2	High	Strong	No
D3	High	Strong	No
D4	Normal	Strong	Yes

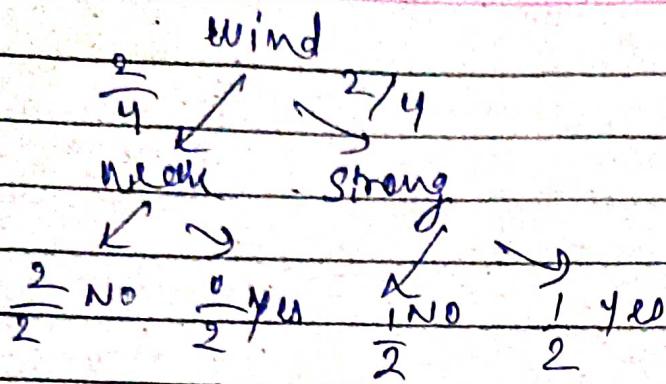
	Humidity	Wind	Play
Yes	High	weak	y
Dy	Normal	Strong	N
Ds	Normal	Weak	y
Dg			

$$IG(\text{CH-midy}) = \frac{3}{4} \left[-3 \log_2 \left(\frac{3}{3} \right) \right] = 0.81$$

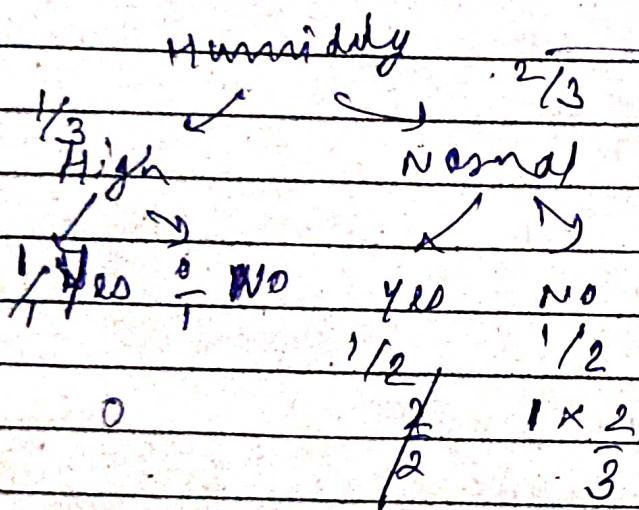
~~Humidity~~ ~~Wind~~

High Normal $\xrightarrow{\text{?}}$ No

$$\text{Entropy} = -\frac{1}{4} \log_2 \frac{1}{4} - \frac{3}{4} \log_2 \frac{3}{4} = 0.81$$



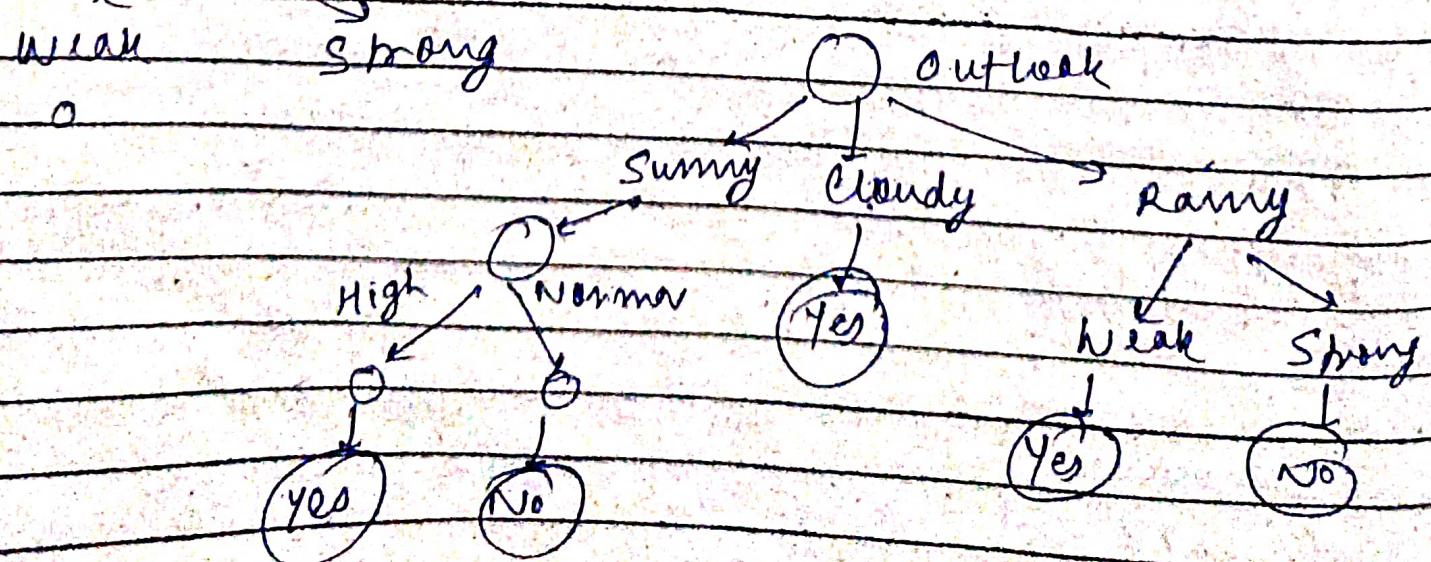
$$\begin{aligned} & \frac{2}{4} \left(-\frac{2}{3} \log_2 \frac{2}{3} \right) + \frac{2}{4} \left(-\frac{2}{3} \log_2 \frac{1}{3} \right) \\ &= 0.81 - \frac{1}{2} \\ &= 0.31 \end{aligned}$$



$$\begin{aligned} H(S) &= -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \\ &= 0.918 \end{aligned}$$

$$IG(\text{Humidity}) = 0.25$$

$$IG(\text{Wind}) = 0.91 \checkmark$$



16.10.2023
Mandavi

Random forest or 10 cross validation
→ To overcome disadvantages

Date: _____
Page: _____

Decision Tree

'iris' Data Set

```
import pandas as pd
```

```
from sklearn.tree import DecisionTreeClassifier  
from sklearn.model_selection import train_test_split  
from sklearn.metrics import confusion_matrix
```

```
Data = pd.read_csv('iris.csv')
```

```
x = Data [ ['sepal length', 'sepal width', 'petal length',  
           'petal width']]
```

```
y = Data ['variety']
```

~~x-test~~
~~x-train, y-train, y-test = train-test-split (x, y, test-size=0.3, random-state=1)~~

```
DTmodel = DecisionTreeClassifier()
```

```
DTmodel.fit(x-train, y-train)
```

```
y-pred = DTmodel.predict(x-test)
```

Advantages :-

→ Human

+ type method
to make

- ① Simple & Easy to understand.
- ② It is useful for decision making problem. decision
- ③ less requirement of data cleaning

Disadvantages:-

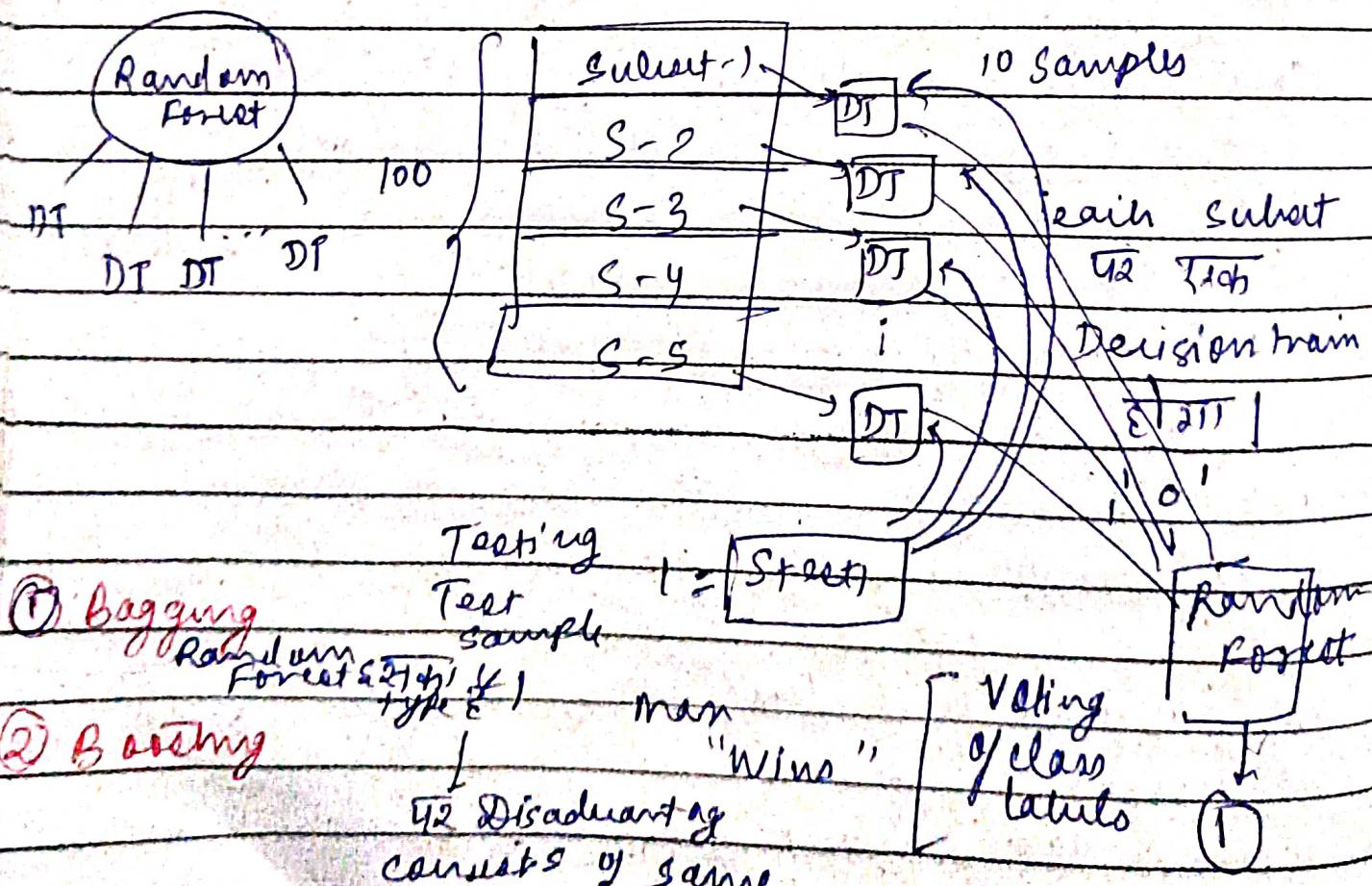
- ① Increase in layers & tree results in increase in complexity.
- ② Overfitting problem may occur.
- ③ more the number of categories class labels, complexity of DT may increase

Overfitting :- When we train the model too much, such that dataset is trained also on noise, which will create conflict.

Underfitting :- Not completely trained on data, so classification ^{will} not be correct.

ENSEMBLE LEARNING

Random Forest consists of multiple trees (Decision (Supervised Classification Algorithm) trees)



① Bagging

Random Forest type 1

Testing

Test sample

Street

② Boosting

Type 2 man

42 Disadvantage
consists of same

"Wins"

Voting

of class

labels

①

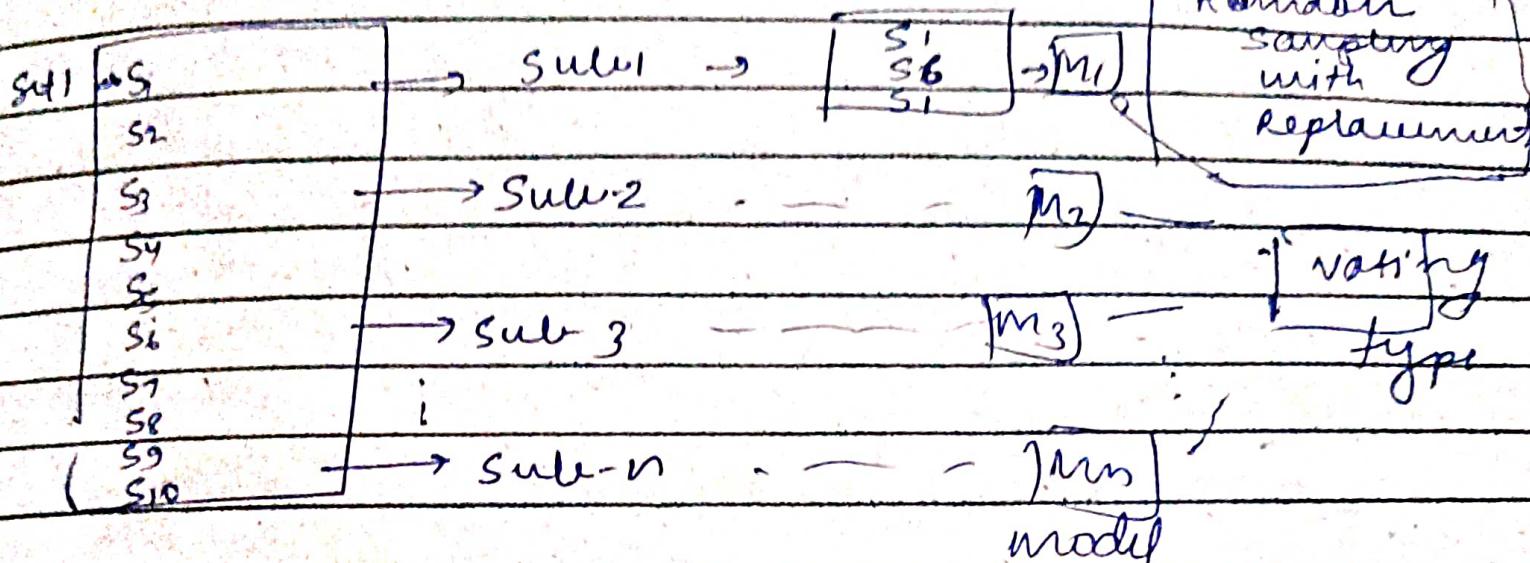
DT, but in bagging multiple different classification algs can be used

Trained
Parallelly models

Also known as

Data Concept used

① Bagging:- Bootstrap Aggregation



Bagging uses

trained

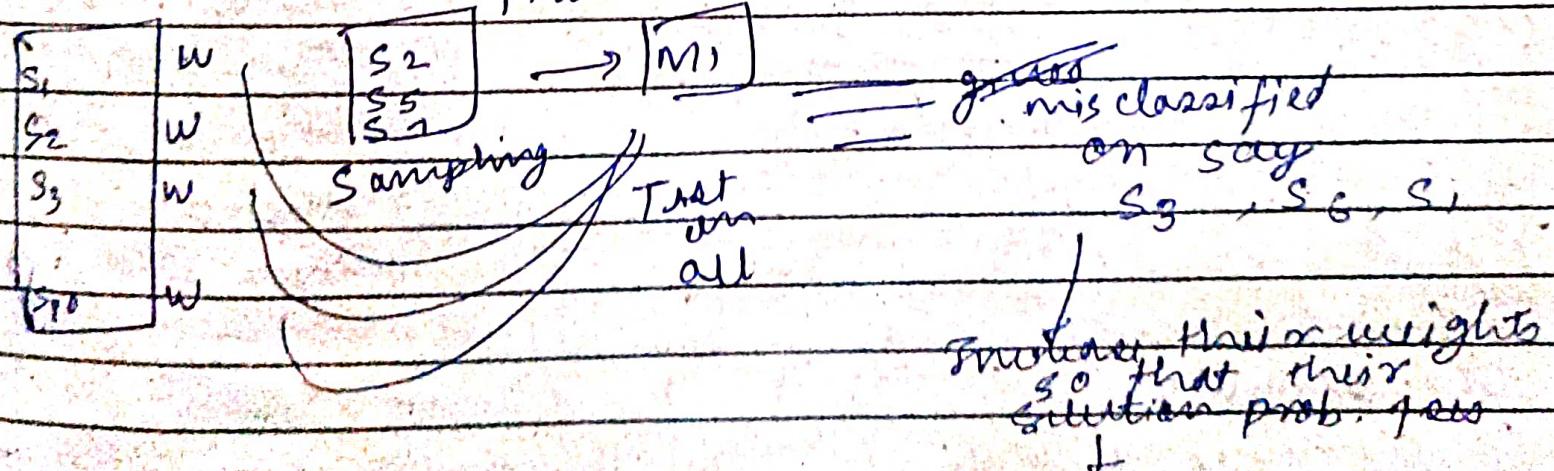
can be

→ Model creation: Sequentially
different

② Boosting : Selection on the basis of wt.

$$wt. \text{ } S_2 \text{ } T_1 \Rightarrow \text{prob. } w_2 \text{ } T_1$$

Train



Decide your
iterations

(models)

Fill your
stack
so do!

Again

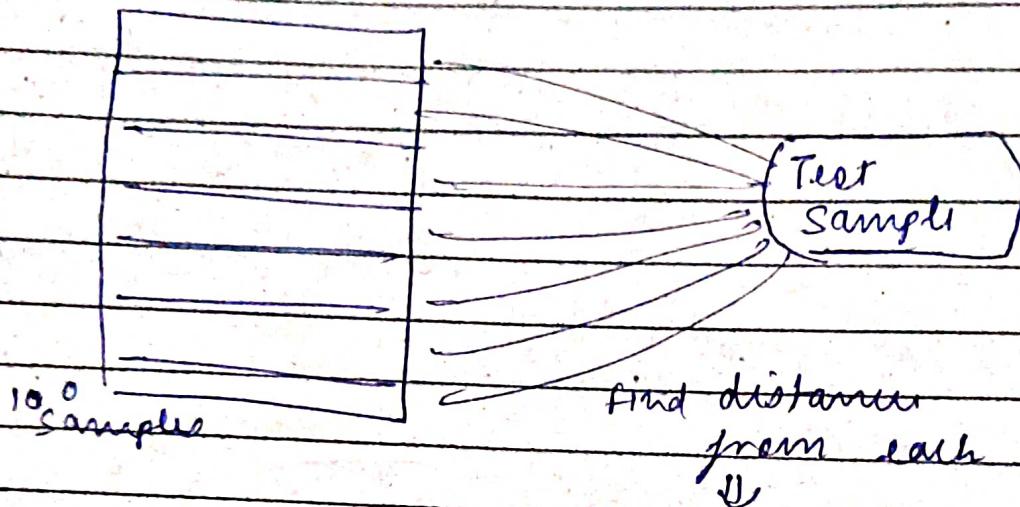
also
the
same
method

s_1	w_1
s_2	w_2
s_3	w_3
\vdots	\vdots
s_{10}	w_{10}

12.10.2023
Tuesday

Date _____ Page _____

K-NN (K-nearest Neighbours) (Lazy Learner)



Then take K nearest neighbours
from among them (Say K=3)

Distance is calculated using Euclidean distance

$$E\text{-distance} = \sqrt{\sum_{i=1}^N (x_i - y_i)^2}$$

Training sample's feature Testing sample's feature

Two things we need to decide :-

① Features to be used in E-distance.

② K value decides

- (i) For binary classification, take $|K=3|$
- (ii) For more than that,

counter = class labels = 3
ex:

But not fully
correct

Date _____
Page _____

$$K = \frac{\sqrt{N}}{2}$$

Generally K odd

(iii) Trial & Error method (Algorithmic approach)

Vary $K = 1$ to N

& find best suited value of K for Test sample

& use that further.

Example

Data Score	OS	Result	New Student	
			E-dist	T ₁ = (Data Score = 6, OS = 8)
S(1)	4	3	fail	5.385
S(2)	6	7	pass	1
S(3)	7	8	pass	1
S(4)	5	5	fail	3.162
S(5)	8	8	pass	2

50% above (strictly) to
Pass

↳ Don't use this.

$K=3$

3 Nearest Neighbours			Pass	Fail
S(2)	S(3)	S(5)	3/3	2/3
Pass	Pass	Pass	Pass	Fail

T₁

Supervised Classification
 (i) KNN (ii) Naive Bayes classifier (iii) Decision Tree
 (iv) Logistic Regression

Unsupervised Classification (Clustering)

Date _____
 Page _____
 Data set : 'iris.csv'

Python implementation

```
import pandas as pd
from sklearn.neighbors import KNeighborsClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import confusion_matrix
```

```
df = pd.read_csv('iris.csv')
```

```
x = df[['sepal length', 'sepal width', 'petal length',  

         'petal width']]
```

```
y = df['varieties']
```

```
x-train, x-test, y-train, y-test = train_test_split(x, y,  

                                                 test_size=0.3, random_state=1)
```

KNN model = KNeighborsClassifier(n_neighbors=7)

 || ↳ Value of k

 || Default, k=5

KNN model.fit(x-train, y-train)

y-pred = KNNmodel.predict(x-test)

print(confusion_matrix(y-test, y-pred))

 ↳ || there are other libraries

 || which calculate Accuracy, precision directly

Accuracy = (Sum of diagonal values) / (Sum of all)

Precision = (Single column) / (Column sum)

Recall = (Single value) / (Row sum)