## Classifiers and classification.

Machine Learning

Supervised          Unsupervised          Reinforcement learning.

→ Regression (covered)

→ Classification

Q.> What is a Classifier (?)



Inputs
$n_1$ →
$n_2$ →                Classifier              label 1 →
$n_3$ →                                        label 2 →
$n_4$ →

Output
[ set of classes ].

labels.
$\{0, 1, 2\}$     1/5  2/5  :0

['tall', 'short']     only
discrete
email classifier ['spam', 'ham']     values allowed.

Output $= \{ \begin{array}{l} x_1, n_2 : label 1 \\ x_3, n_4 : label 2 \end{array} \}$

Difference b/w Classifier v/s classification model

↓                                        ↓

Algorithm itself                    Trained using a classifier algorithm.

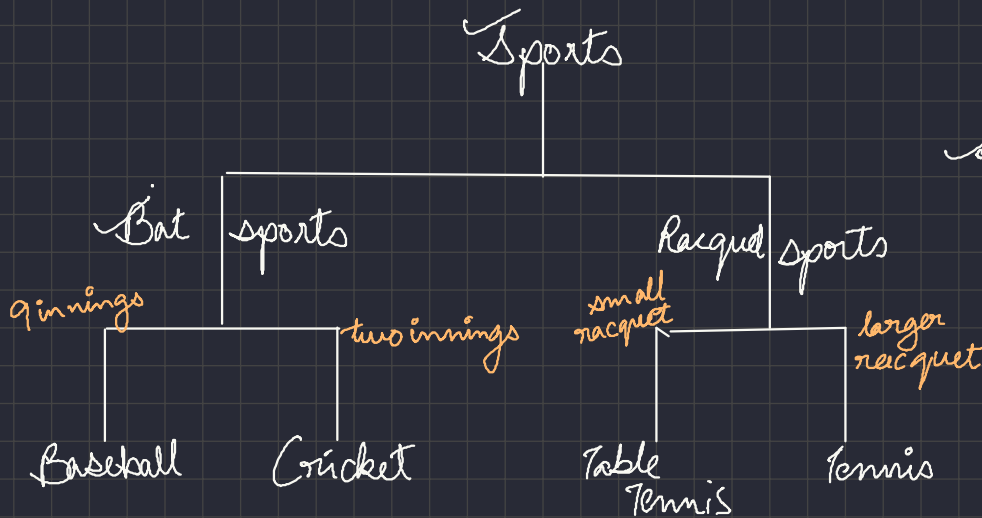methodology/set of rules to classify input data          gives output labels on input data with specific feature set.

# Types of classification algorithms

Classifiers we will cover today:
- → Decision Tree
- → Naïve Bayes Classifier
- → KNN (K-nearest neighbors)

Later portion:
- → SVMs (Support Vector Machines)
- → Artificial Neural Network (ANN).

Example:

Algo → Binary search
↓
special usecase
↓ becomes model.
Model → binary search for specific DSA problem.

## Decision Trees

Sports
- Bat sports
  - 9 innings → Baseball
  - two innings → Cricket
- Racquet sports
  - small racquet → Table Tennis
  - larger racquet → Tennis

input
↓
if-then statement
↓
if-then statement
⋮
↓
specific label given as output

output labels = [Baseball, Cricket, Table Tennis, Tennis]

## Naïve Bayes Classifier

Uses probability    [0-1]

Text classification

Comment 1: ( great )    inputs    (+ve) [0.8]   80% ✓ correct.

(Sentiment analysis) comment 2: ( (not) (so) (great) )    ✓ve [0.58]   58%

c 3: ( bad, horrible )    -ve [0.9]   90% correct.

threshold

$$P(A|B) = \frac{P(B|A) \times P(A)}{P(B)}$$

(+ve, -ve, -ve)
yes   no

Classifier that gives only two labels as output → Binary classifier = [Yes, No]

# K-nearest neighbors

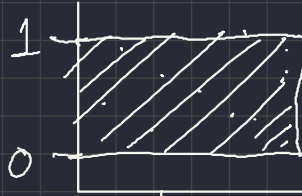$k > 0$    labels = [🟡 🔴 🟢]    ⬜ ←

newdata

Titanic = Regression dataset

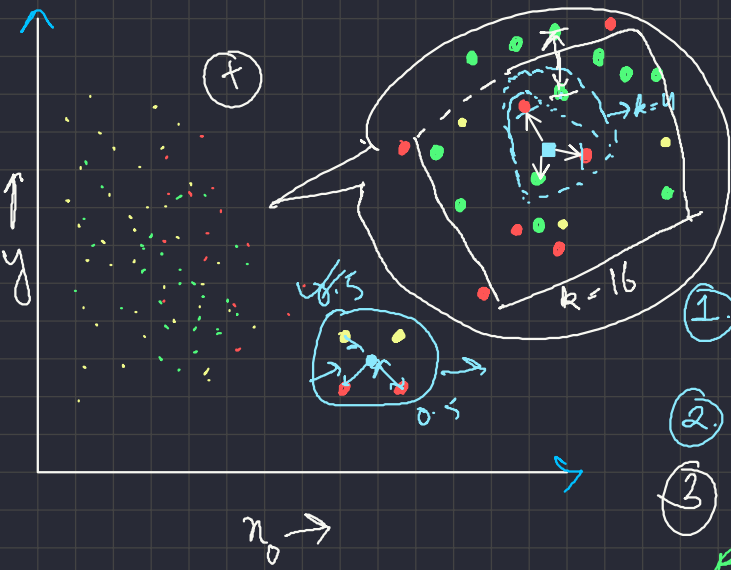new person / data → determine probability of survival.

$0.75\% — [0-1]$
$> 0.5$
∴ yes person will survive



(t)

k=4

k=16

0.5

0.5

① $k = 3$.
pred1 = red

② $k = 4$.
0.5, 0.5 inaccurate

③ $k = 16$
pred3 = green

$n'$

red more than green ∴ label = red

green more ∴ label (for k=16) = green

**regression output** [0-1]
↓
**classification output** [ yes or no based on threshold values for labels].

accuracy

0.95

0.8

acc. value = 0.97 ∴ we take k = 16

for $i$ in range [2-20]:

2

= KNN(k = i)
fit ( )
predict ( )
list.append(accuracy( ))

2    k →    16  20

For KNNs, we measure accuracy over different values of k and map those on a graph. ∴ then select the optimum value of k from the graph itself.

∴ When k = 16, we get an accuracy of 0.97 or 97%. and that k value is suitable for predicting on newer data.

X-train

$n_0$

Input $x_1$

$x_2$

Classifier model
find mathematical co-effs to correlate X and y.

$n_0 \boxed{t_1} \to y_0$  if ($y_0 > t_2$) . . .
else . . . .    more calculation

During training

labels y-train given the dataset for associated inputs $n_0, n_1, n_2$
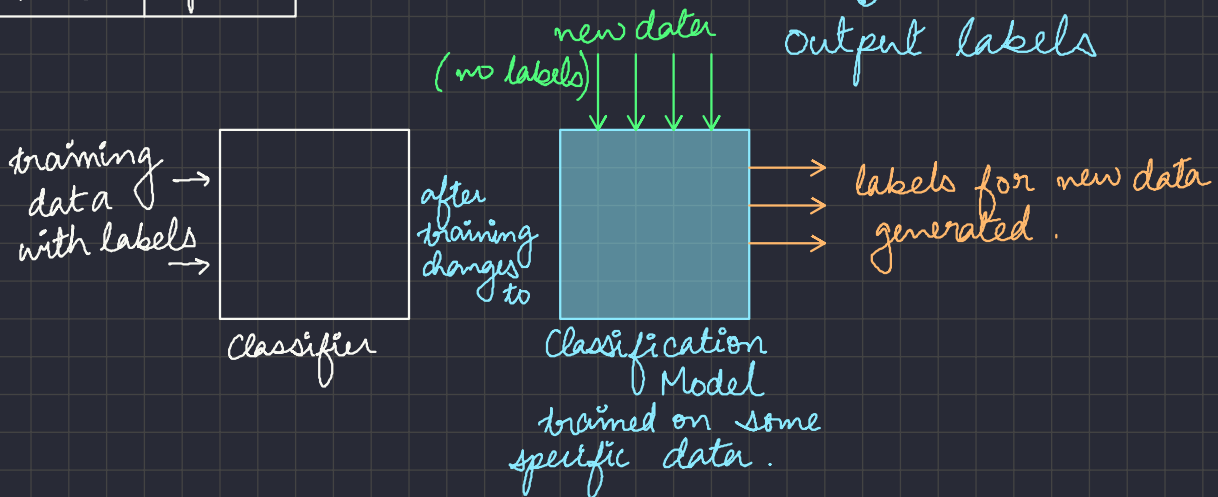⋮
$x_n$.

X-train is input data for training classifier.

y-train is input data X-train's labels.

**\*\*/** After training: (Main portion of ML) (Make predictions on newer data).

new data X′

$f_1$ → 150 →
$f_2$ → 3 →
$f_3$ → 6 →

if $f_1$:
$f_1$ → $y_1$ → if $y_1 > t_1$ : $w_1$ → if $w_1 > t_2$ : $w_1$ → $v_1$
$y_1$
$w_1$
calc 3
$f_2$ → calculation 4 → calculation 5 → $v_2$
6
$f_3$ → 7 → 8 → $v_3$

label prediction for newer data X′ which has features $(f_1, f_2, f_3)$ is {ham} (example).
∴ It is a valid email.

$f_1$ : no. of words
$f_2$ : no. of links
$f_3$ : exclamation point

(Sale !!) Classification Model

spam emails may have exclamation marks.

| $f_1$ | $f_2$ | $f_3$ | y |
|-----|-----|-----|------|
| 800 | 5 | 8 | spam |
| 80 | 1 | 0 | ham |
| 200 | 7 | 5 | spam |

already trained

new data
↓
Classification model
↓
output labels

training data with labels →

Classifier

after training changes to

new data
(no labels)
↓↓↓↓

Classification Model trained on some specific data.

→ labels for new data generated.

<u>Points to remember:</u>

1. Classification is a machine learning technique that predicts the class/category or targets of a given set of data points.

2. Output labels are discrete values.

3. Classifiers are algorithms or set of rules that are used by machines to classify input data.

4. Classification models are the end result of our classifier after it undergoes training on our dataset.