

Clustering :-

ML :- Supervised  $\begin{matrix} \text{end label} \rightarrow \text{target} \\ x \quad y \end{matrix}$   $\rightarrow$  Lin reg, log reg, SVM, DT, RF, AB, GB, XB, KN N,  
NAIVE BAYES

ML :-  
=

Unsupervised

$x \quad y$

$\{x_1, x_2, x_3, x_4, \dots, x_n\}$

$\uparrow$   
= grouping  $\rightarrow$  clustering  $\rightarrow$  mean, kmeans, mini batch kmeans  
- hierarchical.  
DBSCAN.

Lazy  
learner  
paged  
lechner

how this  
things  
working

Supervised ml algo  
 $\begin{matrix} [X] & [Y] \end{matrix}$

sklearn

automatically

Data  $\rightarrow$  batcher

$(m \times n)$   $x_1, x_2, \dots, x_n$   
 $n \times c$

$\downarrow$   
10000  $\rightarrow$  train  $\rightarrow$  Batches  $\rightarrow$  1000, 1200, 1500  
Procedures

ML =

Data.

EDA

Preprocessing

Model / Algo

evaluation

Data

Mathematics

d.f. - d.f.

Which all math concept given;

(1) Linear Algebra = Distance

(2) Calculus

(3) matrix determinant

(4) Stats

(5) Probability theory

Data -> diff. to diff. Assumption

Linear -> Linear reg

Prob. prob (Nonlin) -> log / sum

LR reg

Optimization

Log reg

Dist

LRN

Data analysis

Naive Bayes

DT  
if prob hrs x3

tent:- Naive Bayes

Probab

Nonlinear - SVM / DT

SVM / DT

Simple idea

Sum

Linear

Non

Kernel tricks

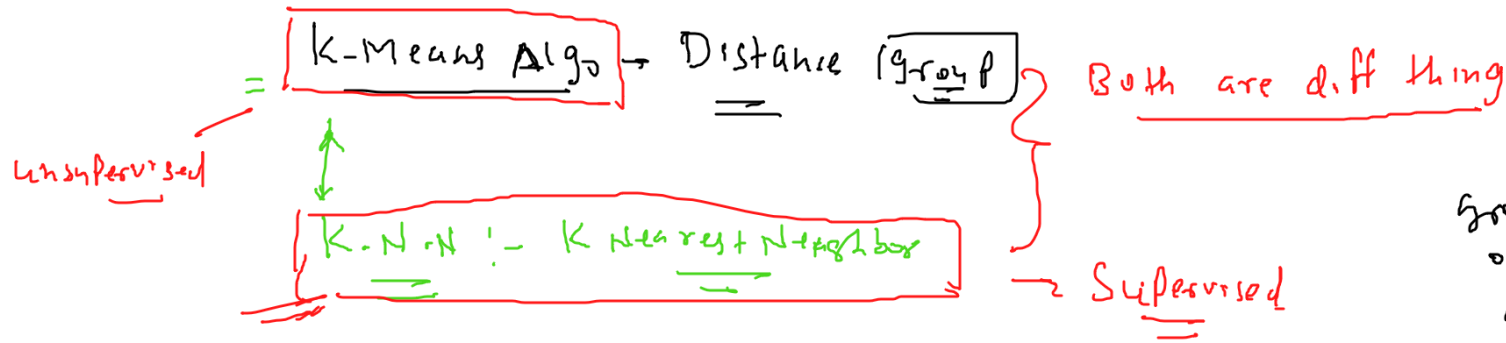
S.M.N

Math. S.A

Distance

Distance strategy

KNN



Unsupervised

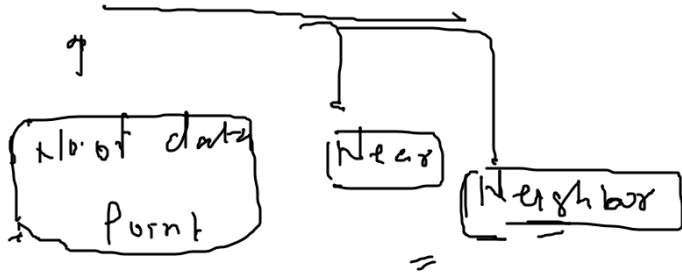
$x_1 \ x_2 \ x_3 \dots x_n$

grouping  
of the  
data =

Classification

I.F D.F → target var

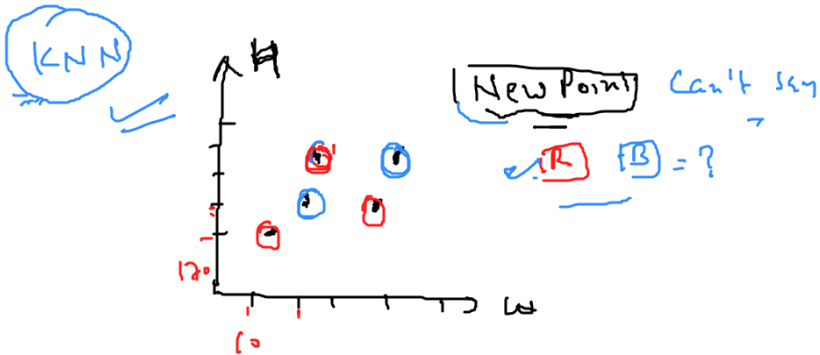
K-Nearest Neighbor :-



X	Y
(watt hr)	(Mile/hr)
60	170
70	160
80	155
90	185
100	155

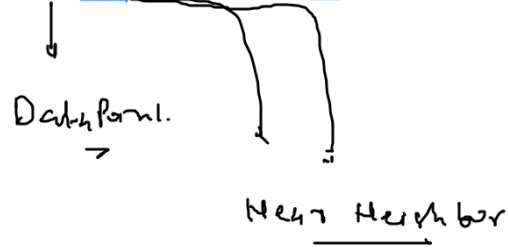
Regression

$x_1$	$x_2$	y
(w)	(h)	0/1/0
60	170	0.5
70	160	0.5
80	150	No
90	140	0.5
100	130	No

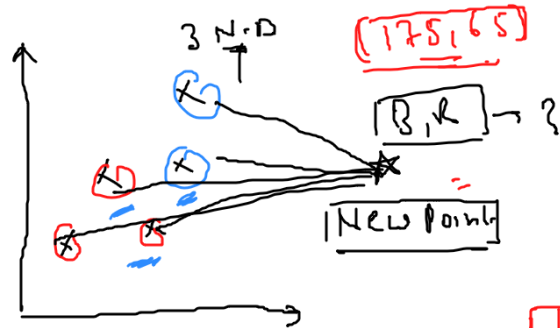


K-Nearest Neighbor → (2, 1, 3, 5, ..., N)

KNN - Dist. Classification



(Nearest data point)



	$X_1$	$X_2$	$Y$	Distance
	(W)	(H)	0/No	
(175, 65) ①	60	175	Yes	5
②	70	160	No	15.81
③	75	155	Yes	122.36
④	85	150	No	26.92
⑤	90	180	Yes	20.41

3 Nearest Neighbor

3 Nearest data point

$K=3$

$K=1, 3, 5, 7, 9, 11$  - odd

$K=2, 4, 6, 8, 10$  - even  $\times$  (Why)?

① I have to calculate distance for each data point

② we have to choose 3 Nearest Dist.

③ Prob =

④ decide to which class it will belong.

① (175, 65)

② (160, 70)

③

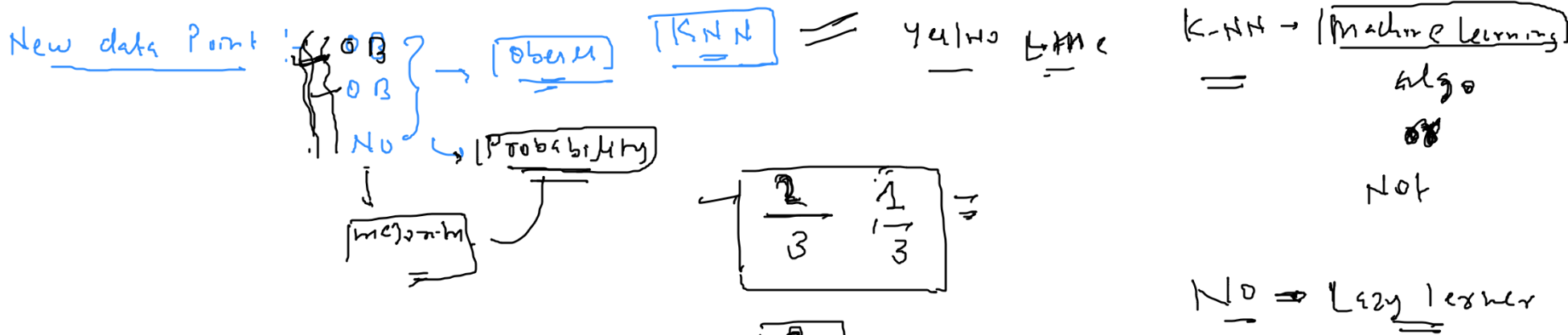
(175, 75)  
(150, 85)

(180, 95)

⑤

Eucclidean distance

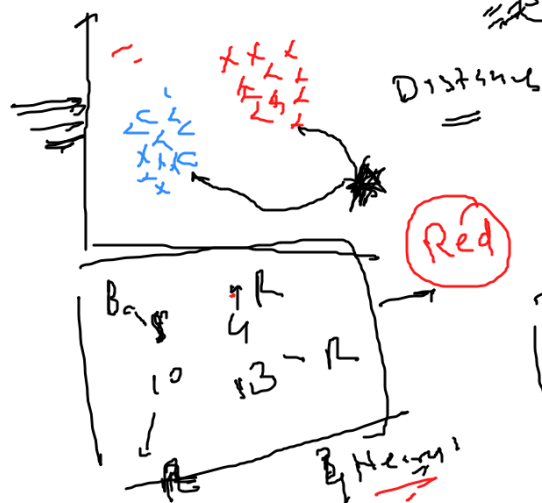
$$\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$



K-Nearest Neighbor by 4-Nearest data Point

Selected

$\begin{Bmatrix} 2 \\ 3 \end{Bmatrix}$  0B



(1) find distance with all point

(2) k  $\sim$  value  $\rightarrow$  hyperparameter

1, 3, 5, 7, 9, 11, 13, ... N

2, 4, 6, 8, 10, 12, ... N

even value

$k=4 \rightarrow \begin{Bmatrix} B & B \\ R & R \end{Bmatrix}$

$\leftarrow ?$  New data point belongs to which class

B, B, R, R, R  
 $\downarrow$   
R  
B, B, B, R, R  
 $\downarrow$   
B

B, R, R, R, R  
 $\downarrow$   
Red

Training - Calculating

Per learner

Machine learning Algo :-

1) calculation

2) loss

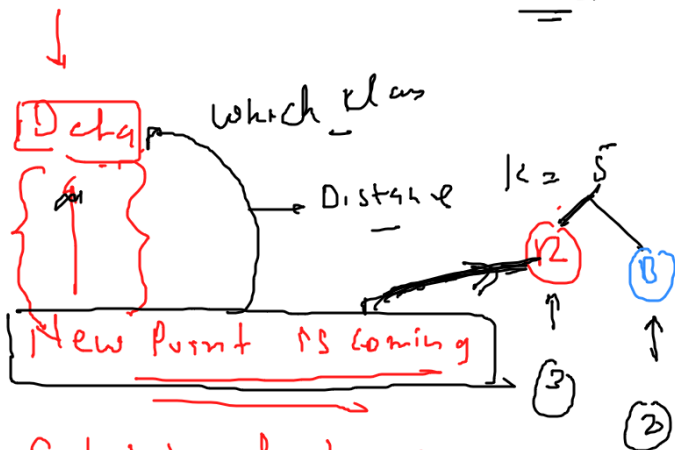
3) Optimization

KNN: IID

$x_1 x_2 x_3 x_4 \dots x_n y$

problem (C)

(training)



Calculation distance

[No LM]

[No optimization]

based on

distance

(Lazy learner)

~~✓~~

$$C_{\text{cal}} = y = wx + c$$

$$\text{loss} = (y - \hat{y})$$

optimization - gradient descent

✓

log reg - cat - Prob -> sigmoid

(log loss)

optimization => Gradient descent

Sum:  $C_{\text{cal}}$   
loss

optimization

or

Maximum

likelihood

estimation

IDT: Prob; but Node.

Data :- KNN  
           
             

Nearest data point

we don't know

Hyper Parameter

1, 3, 5, 7, 9, 11 →  
2, 4, 6, 8, 10 X

highest accuracy → plot  $\left\{ \begin{array}{l} K \text{ value} \\ \text{accuracy} \end{array} \right\}$

→ C-means algo

1, 3, 5, 7, 9, 11  
       

Machine-learning → Training → Generalize pattern

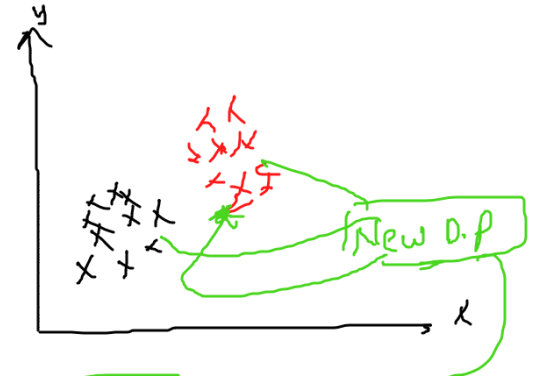
→ CNN → training  
    ↑

Rest of this  
eager learner

Lazy learner

Data if test data is coming  
    Distance N. Neighbor

KNN :- Classification - majority of the class  
Regression :- mean value



$K=3$

majority vote ?

$2R, 1B = (R)$

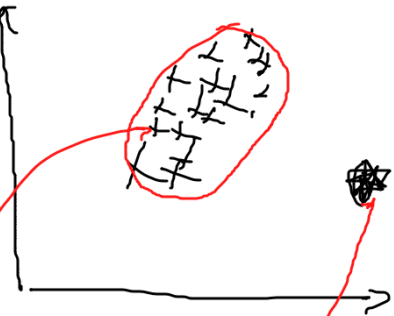
?

$$\frac{18 + 19 + 15.5 + 20}{4}$$

Final ans

K-Nearest Neighbor

	$X_1$	$X_2$	$Y$
	(W)	(H)	(BMI)
1	70	160	18
2	80	170	15
3	50	180	19.5
4	100	150	20
5	120	200	22
	130	140	2



$$\frac{18 + 19 + 19.5}{3}$$

↓ 3

18.83

Nearest distance

$K=3$

18, 19, 19.5

Final

18.83



Machine learning - Reg  
- 1 class

$X_1$ (w)	$X_2$ (h)	y obs / No
( )	( )	

→ Classification

New Point

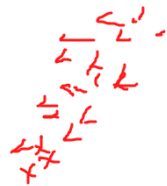
$(x_1, x_2)$   
 $(x_w, x_h)$



Distance →  
Final value with Reg

$(x_1, x_2)$   
 $(x_w, x_h)$

Final value  
with mean



$X_1$ (w)	$X_2$ (h)	y (Bmi)

→ Reg

②

New point

$(x_n, x_n)$

$x_n$   $x_n$

$(x_1, x_2) (140, 60)$

$(x_1, x_2) (60, 65)$

Hamming

$x_1 = x_2 \Rightarrow 0$   
 $x_1 \neq x_2 \Rightarrow 1$  (Categs) (Categs.)

$1 - 0$

$1 - 1$

$1 - 0$

$0 - 1$

$0 - 0$

1

0

1

1

0

③

③

1	0
0	0
0	1
1	1
0	0
1	1

$\Rightarrow ?$

$\begin{bmatrix} A & B & B & C \\ A & C & C & D \end{bmatrix} \rightarrow 3$

$\begin{bmatrix} A & B & B & C \\ A & A & C & D \end{bmatrix} \rightarrow 3$

Hamming distance

$\Rightarrow ||K=3$

$\begin{bmatrix} A & B & B & C \\ A & A & B & C \\ A & B & C & D \end{bmatrix}$

$\begin{bmatrix} A & B & B & C \\ A & A & B & C \end{bmatrix} \rightarrow 1$   
 0 1 0 1

$\begin{bmatrix} A & B & C & D \\ A & B & B & C \end{bmatrix} \Rightarrow 2$

$\begin{bmatrix} A & B & B & C \\ A & B & B & C \end{bmatrix} \rightarrow 0$

①

②

X	Y
<u>A A B C</u>	0
<u>A B C D</u>	No
<u>A B B C</u>	0
A C C D	No
A A C D	0

Final answer

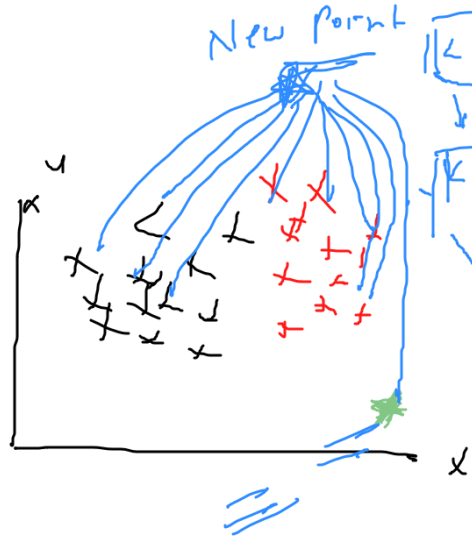
0 Line

(Hamming)

X (gender)	Y (obu/no)
(0) Male	obu
(0) Male	no
(1) Female	No
(1) Female	ob
(0) male	ob

$0 - 0 = 0$   
 $0 - 0 = 0$   
 $0 - 1 = 1$   
 $1 - 0 = 1$   
 $0 - 0 = 0$

(male)  $\rightarrow ?$   
 (0)



①  $K=NN \Rightarrow$  No. of Point

② R or L = Yes.

③ Is KNN ML algo? No/Yes

Lazy  $\Rightarrow$  no training, no loss, no optimization  
egs  $\Rightarrow$

④ Data (large)  $\Rightarrow$  Why?  $\rightarrow$  lots of memory  
(Lazy learned)

⑤ Robust to outliers? Yes.

⑥ For KNN  $\rightarrow$  do we need Scaling?

Distance + Scaling  $\rightarrow$  (standard scaling)  
(minimized)

⑦ Can we use this missing value?

⑧ Can we reduce time complexity of KNN?  
which method:- kd tree, Ball, BF  
 $\downarrow$   
Can you explain

$\begin{pmatrix} 160 & 85 \\ 70 & 190 \end{pmatrix}$   $\begin{pmatrix} 16 & 8 \\ 15 & 5 \end{pmatrix}$   
KNN = ?  $\rightarrow$  with 503

Not a null value

Missing value

Null Not Null Value

(training)

(test)

W	H	BM I
180	70	18
170	80	18.5
160	90	19.5
150	70	19.7
180	100	21.7

(W) (H) (BMI)

180	70	18
170	80	18.5
160	90	19.5
150	70	19.7

Data (training)

185	NA	20
-----	----	----

?

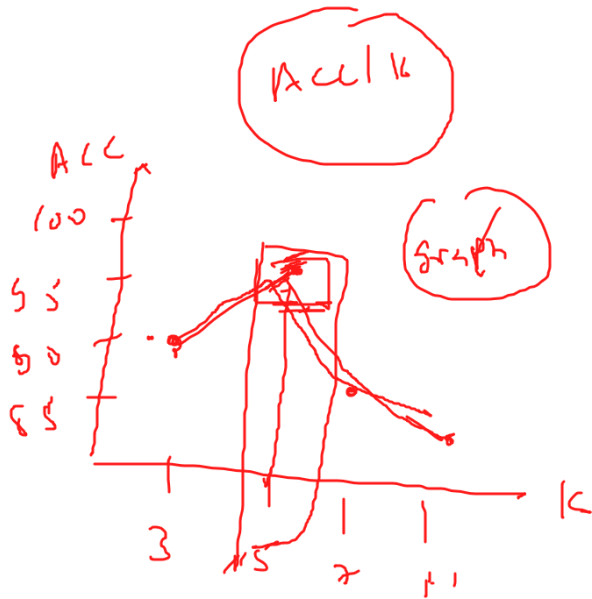
185	NA	20
NA	90	22.5
170	NA	21
NA	55	21.5
185	100	21.7

(KNN)

What null value in a test.

See 140





NAIVE BAYES →

K-value → 1, 3, 5, 7, 9, 11 = unbiasedness = Machine Learning losses loss

Data = Experiment = Data need more

learning fit → 3 test → 90% :

→ 5 = 95%

→ 7 = 89%

→ 11 = 85%

ML → (Experiment) (solution)

← 7 → loss Rest of ml

← 7 → Comput Power