

18

TUESDAY
DAY (04-03-17)
8th WeekMLFA
KNN & DTFebruary
Feb

1)	\rightarrow	x_1	x_2	y
1		-1	-1	2
2		-1	2	1
3		-1	4	1
4		0	-1	1
5		0	0	1
6		0	3	2
7		1	-1	2
8		1	0	2
9		2	0	3
10		2	3	3
11		-1	0	?

Feb	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T
20	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20

Manhattan distance =

WEDNESDAY
DAY (056-314)
8th Week

19

$$|x_1 - x_2| + |y_1 - y_2|$$

Distance matrix : \rightarrow (To calculate distance amongst samples.)

20

THURSDAY
DAY (031-315)
8th WeekFeb 15
February~~For sample 91 ↳~~ $k=1 \Rightarrow$ Nearest Neighbours = Samples 1 & 5. $\Rightarrow Y_1 = 2$ and $Y_5 = 1$.
(It's a tie)⇒ Increment k , $k=2 \Rightarrow$ Nearest neighbours = Samples 1 & 5.
(It's a tie)⇒ Increment k , $k=3 \Rightarrow$ Nearest neighbours = Samples 1, 5 and
2, 8 $\Rightarrow Y_1 = 2$ $Y_8 = 2$ $Y_5 = 1$ $Y_2 = 1$

(It's a tie)

S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16

Increment k , $k=4 \Rightarrow$

Nearest neighbours = Samples 1, 5, 2, 8.
(It's a tie)

Increment k , $k=5 \Rightarrow$

Nearest neighbours = samples 1, 5, 2, 8
7

$$\Rightarrow Y_1 = 2, Y_8 = 2$$

$$Y_5 = 1, Y_7 = 1$$

$$Y_2 = 1$$

\Rightarrow Final Y for
sample 11 = 1

For sample 12 :
→

$k=1$, Nearest neighbours = samples 2, 6 and 10

$$\Rightarrow Y_2 = 1$$

$$Y_6 = 2$$

\Rightarrow (It's a tie)

$$Y_{10} = 3$$

S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T							
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31

22

SATURDAY
DAY (083-313)
8th Week

February
February

Increment k , $k = 2 \Rightarrow$

Nearest neighbours = Samples 2, 6, 10, 5
and 9

$$Y_2 = 1, Y_5 = 1.$$

$$Y_6 = 2, Y_9 = 3$$

$$Y_{10} = 3$$

(to a tie)

Increment k , $k = 3 \Rightarrow$

Nearest neighbours = samples 2, 6, 10, 5, 9,
3, 4 and 11

$$Y_2 = 1, Y_3 = 1$$

DAY

$$Y_6 = 2, Y_9 = 1$$

$$Y_{10} = 3, Y_{11} = 1$$

$$Y_5 = 1$$

$$Y_9 = 3$$

\Rightarrow Final classification
for sample 12 = 1

S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M												
2	3	4	5	6	7	8	9	10	11	12	13	4	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29

Q2

WEDNESDAY

DAY (064-302)

10th Week

04

Pinal label /

Classification Notes in weighted k-NN: →

$$y' = \operatorname{argmax}_v \sum_{(x_i, y_i) \in N_k} w_i \times I(v = y_i)$$

, where N_k is the k -nearest neighbor
-hood,

x_i is the feature vector,

y_i is the label of sample x_i ,

and w_i is the weight of sample i .

In our case w_i will be the inverse
distance from the query point.

For sample 1.2: \rightarrow

$k=1 \Rightarrow$ Nearest neighbours = samples 6 and 10

<u>labels</u>	<u>weightage</u>
$y_6 = C$	$y_6 = 1$
$y_{10} = B$	$y_{10} = 1$

\Rightarrow (It's a tie)

29

SATURDAY

DAY (060-306)

9th Week

February

$k=2$ \Rightarrow Nearest neighbours = samples 6, 10, 13, 7 and 9.

Labels

$$Y_6 = C$$

$$Y_{10} = B$$

$$Y_3 = B$$

$$Y_7 = C$$

$$Y_9 = C$$

Weightage

$$Y_6 = 1$$

$$Y_{10} = 1$$

$$Y_3 = \frac{1}{2} = 0.5$$

$$Y_7 = \frac{1}{2} = 0.5$$

$$Y_9 = \frac{1}{2} = 0.5$$

$$\Rightarrow Y = \text{Argmax} \left\{ \begin{array}{l} ((1) \cdot Y_6 + (0.5) Y_7 + (0.5) Y_9), \\ ((1) \cdot Y_{10} + (0.5) Y_3) \end{array} \right\}$$

$$\Rightarrow \text{Argmax}_Y \left(\left[2C \right], \left[1.5B \right] \right)$$

$$\Rightarrow Y_{11} = C$$

\Rightarrow Final classification of sample 11 = C

FEB '20	S	S	M	T	W	F	S	S	M	T	W	F	S	S	M	T	W	F	S
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20

for sample 12 : \rightarrow

Work to do

$k = 1 \Rightarrow$ Nearest neighbours = samples 1, 3, 5 and 9

labels

weights

$$Y_1 = A$$

$$Y_1 = \frac{1}{2} = 0.5$$

$$Y_3 = B$$

$$Y_3 = 0.5$$

$$Y_5 = B$$

$$Y_5 = 0.5$$

$$Y_9 = C$$

$$Y_9 = 0.5$$

$$\Rightarrow Y_{12} = \underset{y}{\operatorname{argmax}} \left([(0.5)Y_1], [(0.5)Y_3 + (0.5)Y_5], [(0.5)Y_9] \right)$$

$$= \underset{y}{\operatorname{argmax}} \left([0.5 A], [1. B], [0.5 C] \right)$$

W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F															
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	.	APR '20

03

TUESDAY
DAY (063-303)
10th Week

March
March

Appointment

$$\Rightarrow Y_{12} = B$$

Notes

Work to do

\Rightarrow Final classification for sample 12 = B

$$Z \cdot D = 1 = X$$

$$A = 1$$

$$Z \cdot D = \Sigma X$$

$$B = \Sigma$$

$$Z \cdot D = \Sigma X$$

$$C = 2$$

$$Z \cdot D = pX$$

$$D = p$$

Q. In case of the following dataset, you want to build a decision stump. Which splitting criteria will you choose?

ID	x_1	x_2	y
1	1	5	A
2	6	5	B
3	2	4	A
4	6	4	B
5	1	2.5	B
6	8	2.5	A
7	5	2	A
8	3	1.5	B
9	1	1	B
10	6	0	A

- (a) $x_1 > 4$ (b) $x_1 > 5.5$ (c) $x_2 > 3$ (d) $x_2 > 0.5$

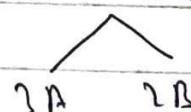
→ Solution:-

Here, we have 2 classes, A & B; out of which 5 are labelled as A & 5 as B.

$$\text{Entropy}([5A, 5B]) = -\left(\frac{5}{10}\right)\log_2\left(\frac{5}{10}\right) - \left(\frac{5}{10}\right)\log_2\left(\frac{5}{10}\right)$$

$$= 1.$$

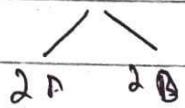
For $x_1 > 4 \Rightarrow 5$ labels



$$\text{Entropy}([3A, 2B]) = -\frac{3}{5}\log_2\frac{3}{5} - \frac{2}{5}\log_2\frac{2}{5} = 0.9709$$

$$\text{Gain} = (1 - 0.9709) = \underline{0.0291}$$

For $x_1 > 5.5 \Rightarrow 4$ labels

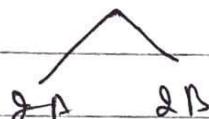


$$\text{Entropy } ([2A, 2B]) = 1.$$

$$\text{Gain} = 1 - 1$$

$$= 0$$

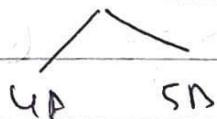
for $x_1 > 3 \Rightarrow 4$ labels



$$\text{Entropy } ([2A, 2B]) = 0.$$

$$\text{Gain} = 1 - 1 = 0.$$

for $x_1 > 0.5 \Rightarrow 9$ labels



$$\text{Entropy } ([4A, 5B]) = 0.9910.$$

$$\text{Gain} = 1 - 0.9910$$

$$= 0.009$$

Highest gain for $x_1 > 4$

$$\therefore \underline{x_1 > 4}.$$

Q. Consider the following dataset.

ID	1	2	3	4	5	6	7	8
X_1	1	1	1	1	2	2	2	2
X_2	1	1	2	2	1	1	2	2
X_3	1	2	1	2	1	2	1	2
Y	3	10	4	12	2	9	2	?

We want to build a regression stump, using variance reduction as the splitting criteria
 (original variance - split1 fraction * split1 variance -
 split2 fraction * split2 variance).

On which attribute will you split?

- (a) X_1 (b) X_2 (c) X_3 (d) All equivalent.

→ Solution:-

By Observation,

When $X_3 = 1 \Rightarrow Y$ values are low.

When $X_3 = 2 \Rightarrow Y$ values are high.

do not
do
this
way!

However, this is not the case with X_1 & X_2 .

∴ Splitting should be on X_3 .

$$\text{Original Variance} = \text{Var}(Y) = \frac{\sum (Y - \bar{Y})^2}{n}$$

$$\bar{Y} = 6$$

$$\text{Var}(Y) = \frac{\sum (Y - \bar{Y})^2}{n} = 15.14$$

We have 3 attributes x_1, x_2 & x_3 .
for x_1

$$x_1 [1, 2]$$

$$1 \rightarrow [3, 10, 4, 12]$$

$$2 \rightarrow [7, 9, 2]$$

Split 1 variance

$$\bar{x}_1(1) = \frac{29}{4}$$

$$\text{Var}(x_1[1]) = \frac{235}{16}$$

Split 2 variance

$$\bar{x}_1(2) = \frac{13}{3}$$

$$\text{Var}(x_1[2]) = \frac{98}{9}$$

$$\begin{aligned} \text{Variance Reduction } (x_1) &= \frac{106}{7} - \frac{4}{7} \times \frac{235}{16} - \frac{3}{7} \times \frac{98}{9} \\ &= \frac{25}{12} = 2.083 \end{aligned}$$

For x_2

$$x_2 [1, 2]$$

$$1 \rightarrow [3, 10, 2, 9]$$

$$2 \rightarrow [4, 12, 7]$$

Split 1 variance

$$\bar{x}_2(1) = 6$$

Split 2 variance

$$\bar{x}_2(2) = 6$$

$$\text{Var}(\bar{x}_2[1]) = \frac{25}{2}$$

$$\text{Var}(\bar{x}_2[2]) = \frac{56}{3}$$

$$\text{Variance Reduction} = \frac{106}{7} - \frac{4}{7} \times \frac{25}{2} - \frac{3}{7} \times \frac{56}{3} = 0$$

for X_3

$$X_3 [1, 2]$$

$$1 \rightarrow [3, 4, 2, 1]$$

$$2 \rightarrow [10, 12, 9]$$

Split 1 Variance

$$\bar{X}_3(1) = \frac{11}{4}$$

$$\text{Var}(X_3(1)) = \frac{11}{16}$$

Split 2 Variance

$$\bar{X}_3(2) = \frac{31}{3}$$

$$\text{Var}(X_3(2)) = \frac{14}{9}$$

$$\begin{aligned} \text{Variance Reduction} &= \frac{106}{7} - \frac{4}{7} \times \frac{11}{16} - \frac{3}{7} \times \frac{14}{9} \\ &= \frac{169}{12} = 14.083 \end{aligned}$$

We choose the attribute with minimum variance reduction.

\therefore splitting will be done on X_3 .