

Max
Let's see a few ex
Practice Questions

Q1.1: Use the following data and group them using K-means clustering algorithm. Show calculation of centroids.

Height	Weight
185	72
170	56
168	60
179	68
182	72
188	77
180	71
180	70
183	84
180	88
180	67
177	76

$$(185, 72) \quad | \quad (170, 56) \quad | \quad (168, 72) \quad | \quad \sqrt{(185-170)^2 + (72-56)^2} \quad | \quad \sqrt{(185-182)^2 + (72-72)^2} \quad |$$

$\uparrow \quad \uparrow \quad \uparrow \quad \uparrow \quad \uparrow$
pt c1 c2 (pt, c1) dist (pt, c2) dist

Smaller dist cluster

Iteration 1 :

Let's randomly choose two centroids.

$$\text{Centroid 1} = (170, 56) \text{ and Centroid 2} = (182, 72).$$

Now, let's calculate the distance of the data points from the chosen centroids and complete the data points table.

The data point is assigned to the cluster based on the closest centroid.

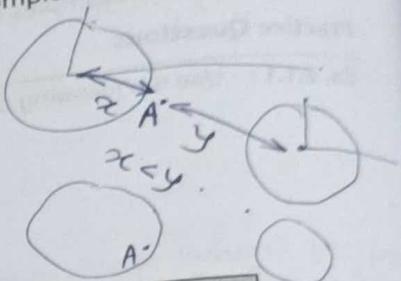
A sample distance calculation is as following for the first data point.

Distance from Centroid 1 (170, 56) for (185, 72) =

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

$$d = \sqrt{(170 - 185)^2 + (56 - 72)^2}$$

$$d = 21.93$$



Height	Weight	Distance from Centroid 1	Distance from Centroid 2	Assigned Cluster
185	72	21.93	3.00	2
170	56	0.00	20.00	1
168	60	4.47	18.44	1
179	68	15.00	5.00	2
182	72	20.00	0.00	2
188	77	27.66	7.81	2
180	71	18.03	2.24	2
180	70	17.20	2.83	2
183	84	30.87	12.04	2
180	88	33.53	16.12	2
180	67	14.87	5.39	2
177	76	21.19	6.40	2

iteration (c1)

$\sum \text{Dist from cluster centroid 1} =$
Let's re-calculate the centroids for the next iteration.

C if no much changes in next iteration
clusters stabilizing . Stopping criteria

For Centroid 1 calculation, you have two data points that fall in cluster 1. They are (170, 56) and (168, 60).

Hence, Centroid 1 is the mean of the data points which is $\left(\frac{170+168}{2}, \frac{56+60}{2} \right) = (169, 58)$.

For Centroid 2 calculation, take the remaining 10 data points that fall in cluster 2.

Hence, Centroid 2 is the mean of these 10 data points which is

$$\left(\frac{185+179+182+188+180+180+183+180+180+177}{10}, \frac{72+68+72+77+71+70+84+88+67+76}{10} \right)$$

$$= (181.4, 74.5)$$

Now, you have both the centroids ready for the next and final iteration.

Iteration 2 :

Now, let's calculate the distance of the data points from the centroids and complete the data points table. The data point is assigned to the cluster based on the closest centroid.

Height	Weight	Distance from Centroid 1	Distance from Centroid 2	Assigned Cluster
185	72	21.26	4.38	2
170	56	2.24	21.73	1
168	60	2.24	19.74	1
179	68	14.14	6.93	2
182	72	19.10	2.57	2
188	77	26.87	7.06	2
180	71	17.03	3.77	2
180	70	16.28	4.71	2
183	84	29.53	9.63	2
180	88	31.95	13.57	2
180	67	14.21	7.63	2
177	76	19.70	4.65	2

You stop here because the number of iterations that you decided are completed and also you see that the clusters assigned in the iteration 1 for the data points did not change.

Hence, you got the two clusters as shown in Fig. P. 5.1.1(a).

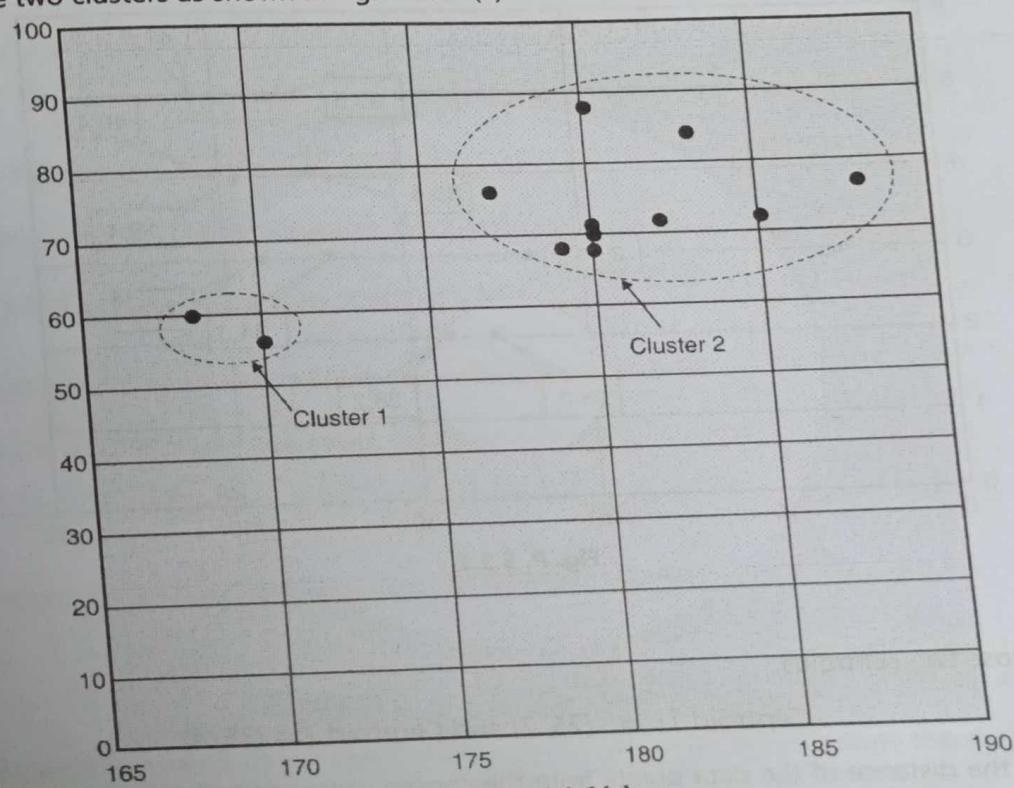


Fig. P. 5.1.1(a)

	minPts = 4 and epsilon (ϵ) = 1.9											
P1: (3, 7)	P1	P2	P3	P4	P5	P6	P7	P8	P9	P10	P11	P12
P2: (4, 6)	P1	0										
P3: (5, 5)	P2	1.41	0									
P4: (6, 4)	P3	2.83	1.41	0								
P5: (7, 3)	P4	4.24	2.83	1.41	0							
P6: (6, 2)	P5	5.66	4.24	2.83	1.41	0						
P7: (7, 2)	P6	5.83	4.47	3.16	2.00	1.41	0					
P8: (8, 4)	P7	6.40	5.00	3.61	2.24	1.00	1.00	0				
P9: (3, 3)	P8	5.83	4.47	3.16	2.00	1.41	2.83	2.24	0			
P10: (2, 6)	P9	4.00	3.16	2.83	3.16	4.00	3.16	4.12	5.10	0		
P11: (3, 5)	P10	1.41	2.00	3.16	4.47	5.83	5.66	6.40	6.32	3.16	0	
P12: (2, 4)	P11	2.00	1.41	2.00	3.16	4.47	4.24	5.00	5.10	2.00	1.41	0
	P12	3.16	2.83	3.16	4.00	5.10	4.47	5.39	6.00	1.41	2.00	1.41

DBSCAN Clustering Algorithm Solved Example – 1

	minPts = 4 and epsilon (ϵ) = 1.9											
	P1	P2	P3	P4	P5	P6	P7	P8	P9	P10	P11	P12
P1	0											
P2	1.41	0										
P3 ✓	2.83	1.41	0									
P4	4.24	2.83	1.41	0								
P5	5.66	4.24	2.83	1.41	0							
P6	5.83	4.47	3.16	2.00	1.41	0						
P7	6.40	5.00	3.61	2.24	1.00	1.00	0					
P8	5.83	4.47	3.16	2.00	1.41	2.83	2.24	0				
P9	4.00	3.16	2.83	3.16	4.00	3.16	4.12	5.10	0			
P10	1.41	2.00	3.16	4.47	5.83	5.66	6.40	6.32	3.16	0		
P11 ✓	2.00	1.41	2.00	3.16	4.47	4.24	5.00	5.10	2.00	1.41	0	
P12	3.16	2.83	3.16	4.00	5.10	4.47	5.39	6.00	1.41	2.00	1.41	0

P1: P2, P10
P2: P1, P3, P11
P3: P2, P4
P4: P3, P5
P5: P4, P6, P7, P8
P6: P5, P7
P7: P5, P6
P8: P5
P9: P12
P10: P1, P11
P11: P2, P10, P12
P12: P9, P11

Like, Share and Subscribe to Mahesh Huddar

Visit: vtupulse.com

P1 is in someone's core point set. then its going to be border

P2 is core itself

P3 is present in someone's core. hence border

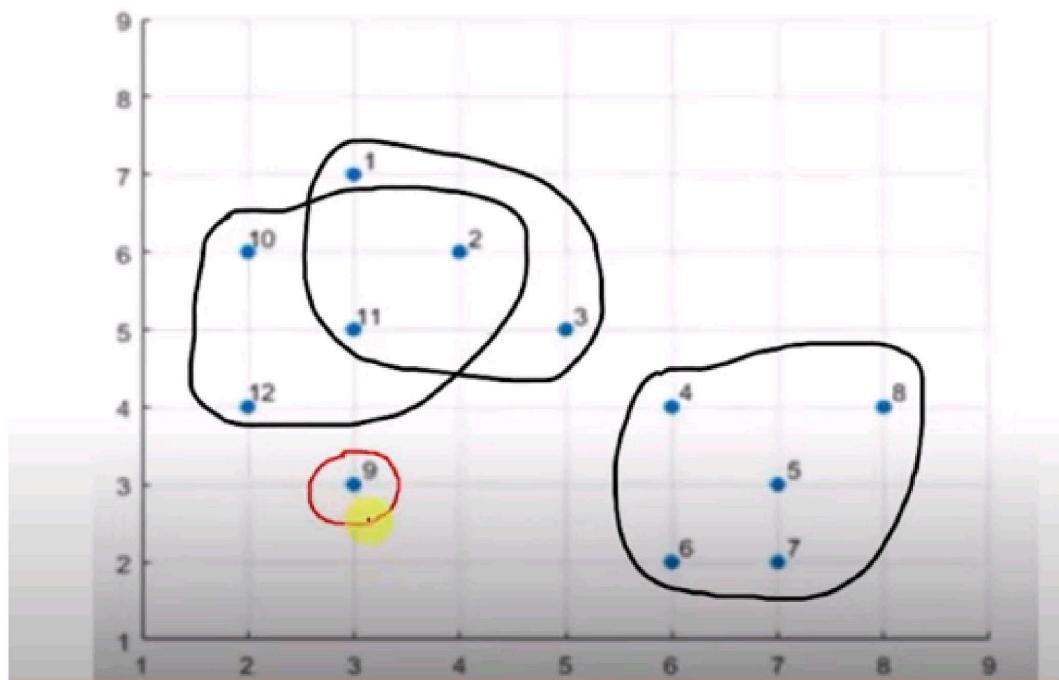
P9 is not in anyone's core. hence not a border

P1: P2, P10
P2: P1, P3, P11
P3: P2, P4
P4: P3, P5
P5: P4, P6, P7, P8
P6: P5, P7
P7: P5, P6
P8: P5
P9: P12
P10: P1, P11
P11: P2, P10, P12
P12: P9, P11

minPts = 4 and epsilon (ϵ) = 1.9

Point	Status	
P1	Noise	Border
P2	Core	
P3	Noise	Border
P4	Noise	Border
P5	Core	
P6	Noise	Border
P7	Noise	Border
P8	Noise	Border
P9	Noise	
P10	Noise	Border
P11	Core	
P12	Noise	Border

minPts = 4 and epsilon (ϵ) = 1.9



Ex. 4.3.1 : Given the following data, calculate hyperplane. Also, classify (0.6, 0.9) based on the calculated hyperplane.

A1	A2	y	α
0.38	0.47	+	65.52
0.49	0.61	-	65.52
0.92	0.41	-	0
0.74	0.89	-	0
0.18	0.58	+	0
0.41	0.35	+	0
0.93	0.81	-	0
0.21	0.1	+	0

Soln. :

As you see, the value of α is non-zero for only first two training examples. Hence, the first two training examples are support vectors.

(Copyright No. L-98904/2021)

Let's calculate the value of w and b as required to calculate the SVM hyperplanes.

$$\text{Let } w = (w_1, w_2)$$

N

$$w = \sum_{n=1}^N a_n y_n x_n$$

$$w_1 = a_1 * y_1 * A_1 x_1 + a_2 * y_2 * A_1 x_2 = \frac{y_1}{a_1} - 65.52 * 1 * 0.38 + 65.52 * -1 * 0.49 = -7.2$$

$$w_2 = a_1 * y_1 * A_2 x_1 + a_2 * y_2 * A_2 x_2 = \frac{y_1}{a_1} - 65.52 * 1 * 0.47 + 65.52 * -1 * 0.61 = -9.2$$

The parameter b can be calculated for each support vector as follows.

$$b_1 = 1 - w * x_1 = 1 - (-7.2) * 0.38 - (-9.2) * 0.47 = 8.06$$

$$b_2 = 1 - w * x_2 = 1 - (-7.2) * 0.47 - (-9.2) * 0.61 = 10$$

Averaging b_1, b_2 you get $b = 9.03$

Hence, the hyperplane line is defined as $-7.2x_1 - 9.2x_2 + 9.03 = 0$

$$w_1 x_1 + w_2 x_2 + b = 0$$

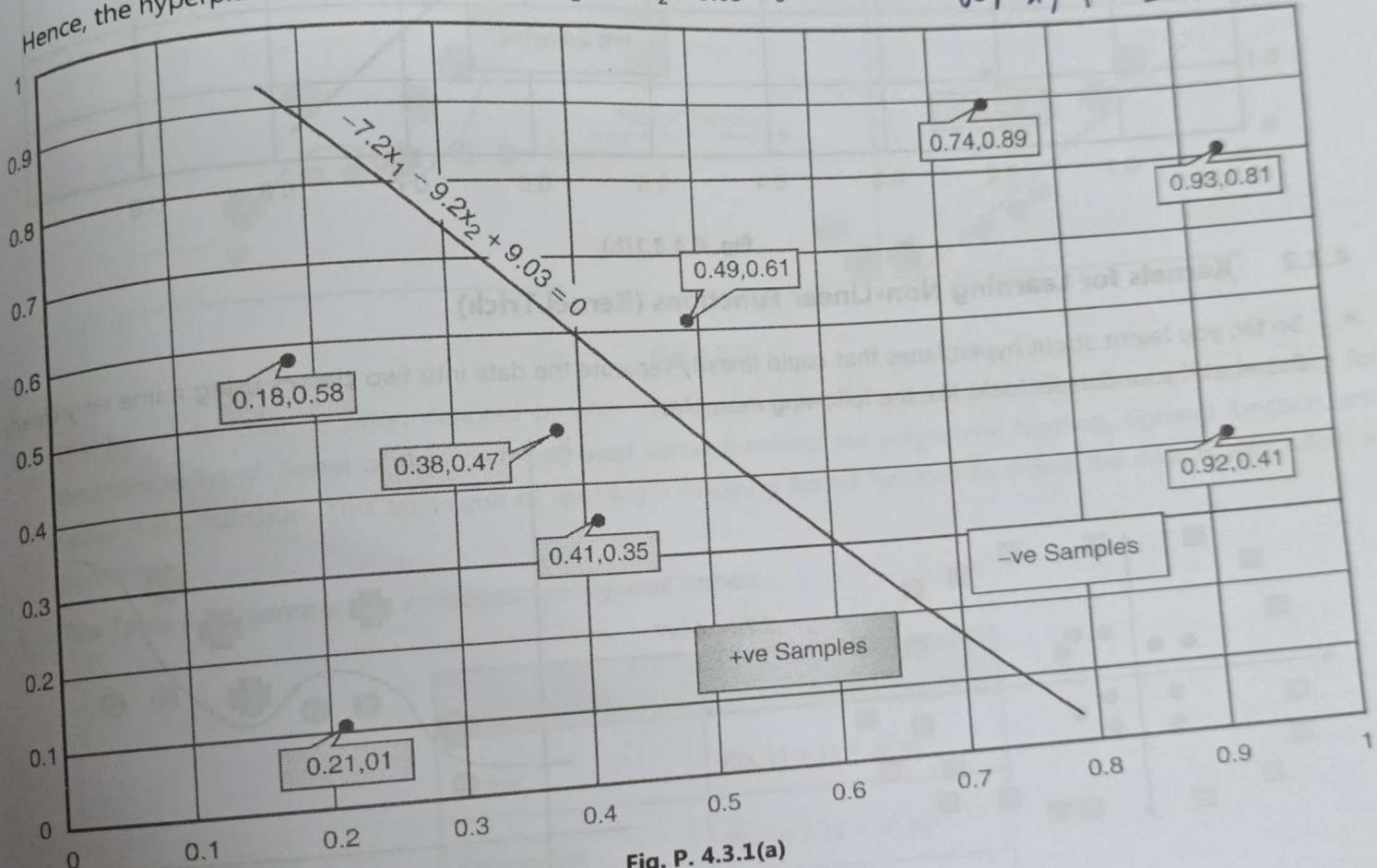


Fig. P. 4.3.1(a)

Now, suppose that there is a new data point $(0.6, 0.9)$ that you want to classify.

Putting the values in the equation $-7.2x_1 - 9.2x_2 + 9.03$ you get $-7.2 * 0.6 - 9.2 * 0.9 + 9.03 = -3.57$.

Hence, this is classified as negative sample.

Agglomerative Clustering

Q. Consider 1D Data pts. 18, 22, 25, 27, 42, 43.

18, (22, 25, 27), (42, 43)

Merge both

22, 18, (25, 27), (42, 43)

15

0

\rightarrow Distance Matrix:

18	22	25	27	42	43
18	0	4	7	9	24
22	4	0	3	5	20
25	7	3	0	2	17
27	9	5	2	0	15
42	24	20	17	15	0
43	25	21	18	16	①

(circle the min cell)

cell (42, 43) =

merge 42 & 43.



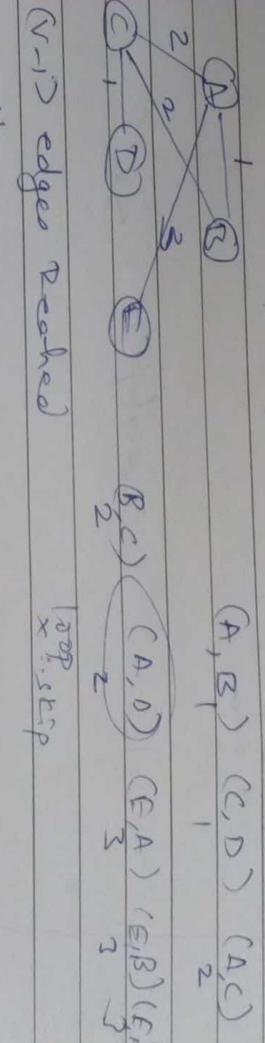
Divisive method

use MST					
A	B	C	D	E	
0	1	0			
B	1	0			
C	2	2	0		
D	0	2	1	0	
E	3	3	5	3	0

single linkage: (42, 43) \rightarrow 18 (min (42+18) & (43+18))
 complete linkage: max dist.

Merge 25, 27.

18	22	25, 27	42, 43	
18	0	4	7	24
22	4	0	3	20
25	7	3	0	17
27	9	5	②	15
42, 43	24	20	17	15



Merge (25, 27) & 22

(18)	(22, (25, 27))	(42, 43)
(22, (25, 27))	④	27
(42, 43)	24	15

Merge (25, 27) & 22

18, (22, (25, 27))

27

15

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

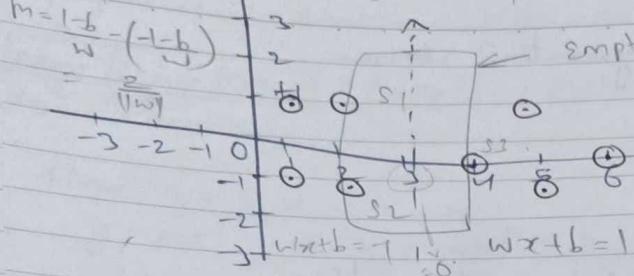
0

0

0

Q. SVM Hyperplane.

($1,1$) ($2,1$) ($1,-1$) ($2,-1$) ($4,0$) ($5,1$) ($5,-1$) ($6,0$)
 (it can be given as few pts belong to +ve class, others -ve).



$$S_1 \text{ vector} = \begin{pmatrix} 2 \\ 1 \end{pmatrix}, \quad S_2 = \begin{pmatrix} 2 \\ -1 \end{pmatrix}, \quad S_3 = \begin{pmatrix} 5 \\ 0 \end{pmatrix}$$

Bias Segment \Rightarrow

$$\frac{S_1}{S_1} = \begin{pmatrix} 2 \\ 1 \end{pmatrix}, \quad S_2 = \begin{pmatrix} 2 \\ -1 \end{pmatrix}, \quad S_3 = \begin{pmatrix} 5 \\ 0 \end{pmatrix} \quad \text{] add 1 last.}$$

$$\alpha_1 S_1 \bar{S}_1 + \alpha_2 S_1 \bar{S}_2 + \alpha_3 S_1 \bar{S}_3 = -1 \quad (\text{Eq. Constant})$$

$$\alpha_1 S_2 \bar{S}_1 + \alpha_2 S_2 \bar{S}_2 + \alpha_3 S_2 \bar{S}_3 = -1 \quad (\text{Eq. 2})$$

$$\alpha_1 S_3 \bar{S}_1 + \alpha_2 S_3 \bar{S}_2 + \alpha_3 S_3 \bar{S}_3 = +1 \quad (\text{Eq. 3})$$

assuming one side as +ve, other as -ve

$$\therefore S_1, S_2 = +\text{ve}, \quad S_3 = -\text{ve} \quad \text{or}$$

$$S_1, S_2 = -\text{ve}, \quad S_3 = +\text{ve}$$

$$\therefore \alpha_1 \begin{pmatrix} 2 \\ 1 \end{pmatrix} \begin{pmatrix} 2 \\ 1 \end{pmatrix} + \alpha_2 \begin{pmatrix} 2 \\ 1 \end{pmatrix} \begin{pmatrix} 2 \\ -1 \end{pmatrix} + \alpha_3 \begin{pmatrix} 2 \\ 1 \end{pmatrix} \begin{pmatrix} 5 \\ 0 \end{pmatrix} = -1$$

$$\therefore \alpha_1 (2 \times 2 + 1 \times 1 + 1 \times 1) + \alpha_2 (2 \times 2 - 1 \times 1 + 1 \times 1) + \alpha_3 (2 \times 4 + 0 \times 1 + 1 \times 1) = -1$$

$$\therefore \alpha_1 (6) + \alpha_2 (5) + \alpha_3 (5) = -1 \quad \text{--- (1)}$$

$$\text{Similarly, } 4\alpha_1 + 6\alpha_2 + 9\alpha_3 = -1$$

$$9\alpha_1 + 9\alpha_2 + 17\alpha_3 = 1$$

$$\therefore \alpha_1 = -3.25, \quad \alpha_2 = -3.25, \quad \alpha_3 = 3.5.$$

$$y = mx + c$$

$$0 = m(1) - 3$$

$$\bar{w} = \sum \alpha_i \bar{S}_i = -3.25 \begin{pmatrix} 2 \\ 1 \end{pmatrix} - 3.25 \begin{pmatrix} 2 \\ -1 \end{pmatrix} + 3.5 \begin{pmatrix} 5 \\ 0 \end{pmatrix} =$$

$$\left[\begin{array}{c} -3.25 \times 2 - 3.25 \times 2 + 3.5 \times 5 \\ -3.25 \times 1 + 3.25 \times 1 + 3.5 \times 0 \\ -3.25 \times 1 - 3.25 \times 1 + 3.5 \times 1 \end{array} \right] = \begin{pmatrix} 0 \\ -3 \\ 0 \end{pmatrix} \bar{w}$$

o

multiple LinReg.

y	x_1	x_2

$$\therefore \text{matrix} = \begin{pmatrix} 33 & 42 \\ 42 & 54 \end{pmatrix}$$

$$x_1 x_1 \quad x_2 x_2 \quad x_1 x_2 \quad x_2 x_1 \quad x_1 x_2$$



$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$$

$$\beta_1 = \frac{\sum x_2^2 \sum x_1 y - \sum x_1 x_2 \sum x_1 y}{\sum x_1^2 \sum x_2^2 - (\sum x_1 x_2)^2}$$

$$\beta_2 = \frac{\sum x_1^2 \sum x_2 y - \sum x_1 x_2 \sum x_2 y}{\sum x_1^2 \sum x_2^2 - (\sum x_1 x_2)^2}$$

$$\sum x_1^2 = \sum x_1 \cdot x_1 - \frac{\sum x_1 \cdot \sum x_1}{N}$$

$$\beta_0 = \bar{y} - \beta_1 \bar{x}_1 - \beta_2 \bar{x}_2$$

quad poly, $n=2$ over x

y

x^2

xy

x

x^3

x^4

x^5

x^6

x^7

x^8

x^9

x^{10}

x^{11}

x^{12}

x^{13}

x^{14}

x^{15}

x^{16}

x^{17}

x^{18}

x^{19}

x^{20}

x^{21}

x^{22}

x^{23}

x^{24}

1. Linear Reg.

$$y = mx + c + e$$

$$m = \frac{\sum y}{n} - \bar{x}\bar{y}$$

$$n \bar{x}^2 - \bar{x}^2$$

$$c = \bar{y} - m \bar{x}$$

$$\text{Avg} = \bar{y}$$

$$(\text{mean}) = \bar{x}$$

$$\bar{y}$$

When x given, find y

only 1 independent variable x

or

$$y = ax + b$$

$$\text{Cov}(x, y) = \frac{1}{n-1} \sum (x_i - \bar{x})(y_i - \bar{y})$$

$$\text{Correl. } (x, y) = \frac{\text{Cov}(x, y)}{\sqrt{\text{Var}(x) \text{Var}(y)}}$$

$$\text{Var}(x) = \frac{1}{n-1} \sum (x - \bar{x})^2$$

$$b = \frac{\text{Cov}(x, y)}{\text{Var}(x)}$$

$$a = \bar{y} - b\bar{x}$$

on linear SVM -

fully labeled -

$$\begin{pmatrix} 2 \\ 2 \end{pmatrix}$$

fully labeled data pts -

$$\begin{pmatrix} 2 \\ -2 \end{pmatrix}$$

$$\begin{pmatrix} -2 \\ -2 \end{pmatrix}$$

$$\begin{pmatrix} -2 \\ 2 \end{pmatrix}$$

fully labeled -

$$\begin{pmatrix} 1 \\ 1 \end{pmath>$$

$$\begin{pmatrix} 1 \\ -1 \end{pmath}$$

$$\begin{pmatrix} -1 \\ 1 \end{pmath}$$

$$\begin{pmatrix} -1 \\ -1 \end{pmath}$$

$$\Phi(x_1, x_2) = \begin{cases} (4 - x_2 + |x_1 - x_2|)^2 & \text{if } \sqrt{x_1^2 + x_2^2} > 2 \\ (4 - x_1 + |x_1 - x_2|)^2 & \text{otherwise} \end{cases}$$

$$\begin{pmatrix} x_1 \\ x_2 \end{pmath}$$

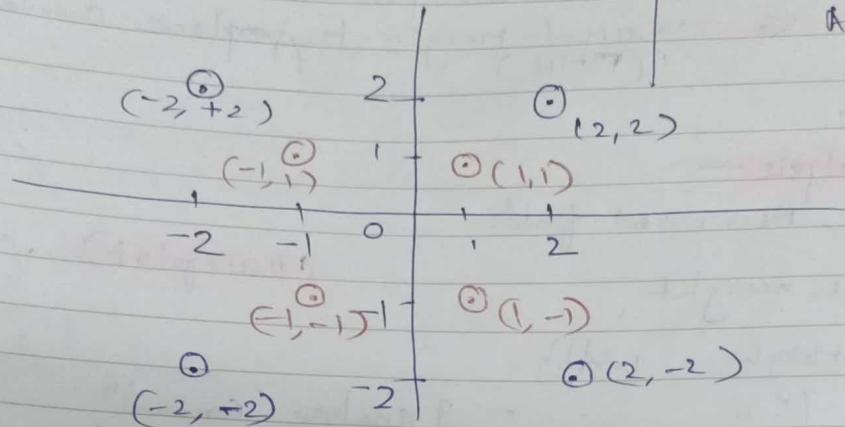
HPLA

MPHA

LPLA

LP HA

Accur.



Converting them from 1 feature space to another:

① for fully labeled Data pts :

$$\sqrt{2^2 + 2^2} > 2 \quad \text{true} \quad \therefore 4 - 2 + |2 - 2| = \binom{2}{2}$$

$$4 - 2 + |2 - 2| = \binom{2}{2}$$

$$\sqrt{2^2 + (-2)^2} > 2 \quad \text{true}$$

$$\therefore 4 + 2 + |2 + 2| = \binom{10}{6}$$

$$4 - 2 + |2 + 2| = \binom{6}{6}$$

$$\sqrt{(-2)^2 + (-2)^2} > 2 \quad \text{true}$$

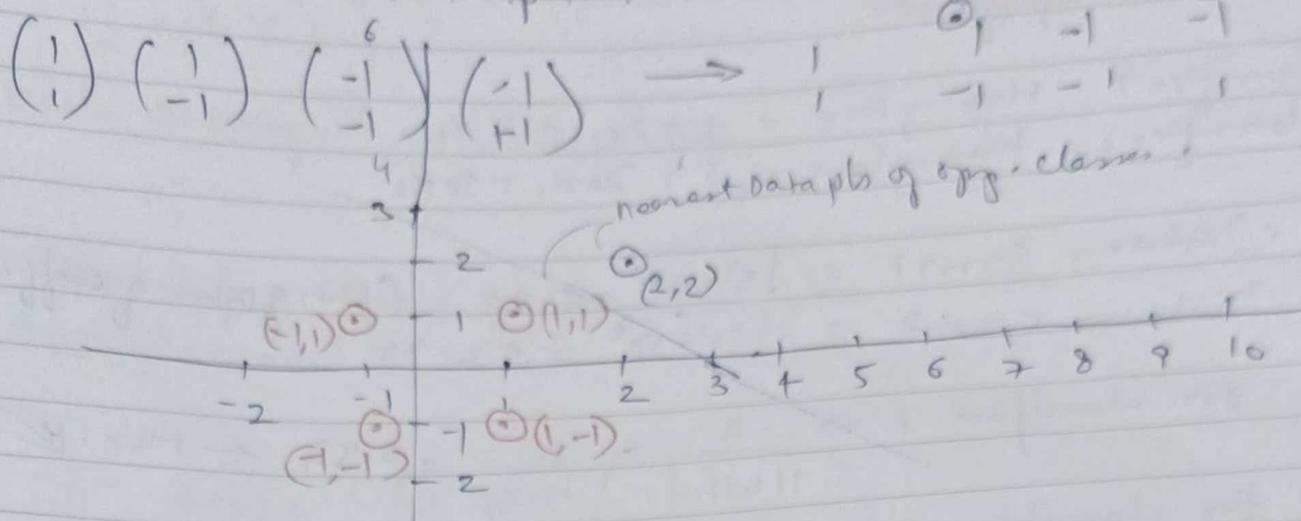
$$\therefore$$

$$\frac{6}{6}$$

$$\sqrt{(2)^2 + (2)^2} > 2 \quad \text{true}$$

$$\frac{6}{10}$$

For evenly binned data pts -



$$\therefore S_1 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, S_2 = \begin{pmatrix} 2 \\ 2 \end{pmatrix}$$

$$\text{Add Bias: } \tilde{S}_1 = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}, \tilde{S}_2 = \begin{pmatrix} 2 \\ 2 \\ 1 \end{pmatrix}$$

$$\alpha_1 \tilde{S}_1^\top \tilde{S}_1 + \alpha_2 \tilde{S}_2^\top \tilde{S}_1 = -1 \quad \dots \text{wst } S_1$$

$$\alpha_1 \tilde{S}_2^\top \tilde{S}_1 + \alpha_2 \tilde{S}_2^\top \tilde{S}_2 = +1 \quad \dots \text{wst } S_2$$

$$\therefore 3\alpha_1 + 5\alpha_2 = -1$$

$$5\alpha_1 + 9\alpha_2 = 1$$

$$\alpha_1 = -7, \alpha_2 = 4$$

$$\therefore \tilde{w} = \sum \alpha_i \tilde{s}_i =$$

$$\alpha_1 \tilde{S}_1 + \alpha_2 \tilde{S}_2$$

$$= -7 \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} + 4 \begin{pmatrix} 2 \\ 2 \\ 1 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$$

Q. Multivar Reg.

x_1	x_2	y_1	y_2
-------	-------	-------	-------

$$Y = [y_1, y_2], X = [x_1, x_2]$$

$$\Theta = [X^T X]^{-1} \cdot X^T Y$$

$$= \begin{pmatrix} -41 & 67 \\ 59 & 27 \\ 29 & -23 \end{pmatrix}$$

$$\therefore \hat{y}_1 = -41 + 59x_1 + 29x_2$$

$$\hat{y}_2 = 67 + 27x_1 - 23x_2$$

$$w = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$$

$$b = \text{bias} =$$

Margin \uparrow : better control over Generalization Error.
 \downarrow : overfitting hijayegi.

Selected Best Hyperplane:

Eg. Classifier 1 Hyperplane = $5 + 2x_1 + 5x_2$
Classifier 2 Hyperplane = $5 + 20x_1 + 50x_2$

\therefore Distance $\epsilon_{\text{margin}1} = \sqrt{5^2 + 2^2} = 5.39$ (sum of square of coeff.)

Distance $\epsilon_{\text{margin}2} = \sqrt{20^2 + 50^2} = 53.85$

for classifier 1 $\Rightarrow \frac{2}{\|w\|} = \frac{2}{5.39} = 0.37 \leftarrow \text{Max.} \checkmark$

for classifier 2 $\Rightarrow \frac{2}{\|w\|} = \frac{2}{53.85} = 0.037$

\therefore Hyperplane 1 is Best

age	income	student	credit Rating	bought Computer	Weight
1 Youth	high	no	fair	no	0.9403
2 Youth	lv	n	excellent	n	0.0597
3 Middle	h	n	f	yes	0.0597
4 Senior	medium	n	f	yes	0.0597
5 Senior	low	yes	f	yes	0.0597
6 Senior	l	y	e	n	0.0597
7 Middle	l	y	f	n	0.0597
8 Youth	m	n	f	y	0.0597
9 Youth	l	y	f	y	0.0597
10 Senior	m	y	f	y	0.0597
11 Youth	m	y	e	y	0.0597
12 Middle	m	n	e	y	0.0597
13 Middle	h	y	f	y	0.0597
14 Senior	m	n	e	n	0.0597

$$n(\text{yes}) = 9$$

$$n(\text{no}) = 5$$

age	2	3	5.	$-\frac{2}{5} \log \frac{2}{5} - \frac{3}{5} \log \frac{3}{5} = 0.971$
Youth	2	3	5.	$-\frac{3}{5} \log \frac{3}{5} - \frac{2}{5} \log \frac{2}{5} = 0.971$
Middleaged	4	0	4	
Senior	3	2	5.	
Total	y	n	weight	
Split-Info	$= -\frac{5}{14} \log \frac{5}{14} - \frac{4}{14} \log \frac{4}{14} - \frac{5}{14} \log \frac{5}{14} = 1.5774$			
$I_G =$	$0.971 \times \frac{5}{14}$	$+ 0.971 \times \frac{5}{14}$		$= 0.6934$
Gain	$0.9403 - 0.6934 = 0.2469$			

$$\text{Gain Ratio} =$$

$$\frac{0.2469}{1.5774} = 0.1559$$

income	2	2	4	1	$-\frac{4}{8} \log \frac{4}{8} - \frac{2}{8} \log \frac{2}{8} = 0.9183$
High	2	2	4	1	
Medium	4	2	6	0	
Low	3	1	4	1	$-\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4} = 0.8113$

$$\text{Split-Info} = -\frac{4}{14} \log \frac{4}{14} - \frac{4}{14} \log \frac{6}{14} - \frac{4}{14} \log \frac{5}{14} = 1.5566$$

$$I_G = 1 * \frac{4}{14} + 0.9183 * \frac{6}{14} + 0.8113 * \frac{4}{14} = 0.911$$

$$\text{Gain Ratio} = \frac{0.0293}{1.5566} = 0.0188$$

Student									
Yes	6	1	7	- $\frac{6}{14} \log \frac{6}{14}$	- $\frac{1}{14} \log \frac{1}{14}$	=	0.5917		
No	3	4	7	- $\frac{3}{14} \log \frac{3}{14}$	- $\frac{4}{14} \log \frac{4}{14}$	=	0.9852		

$$\text{Split Info} = -\left(\frac{6}{14} \log \frac{6}{14}\right) \times 2 = 1$$

$$IG = 0.5917 \times \frac{7}{14} + 0.9852 \times \frac{7}{14} = 0.7284$$

Gain Ratio = 0.1519

$$\text{Gain} = 0.9403 - 0.7284 = 0.1519$$

= 0.1519

Credit Rating									
Fair	6	2	8	- $\frac{6}{8} \log \frac{6}{8}$	- $\frac{2}{8} \log \frac{2}{8}$	=	0.8113		
Excellent	3	3	6	-	-				

$$IG = 0.8113 + \frac{8}{14} + 1 * \frac{6}{14} = 0.8922$$

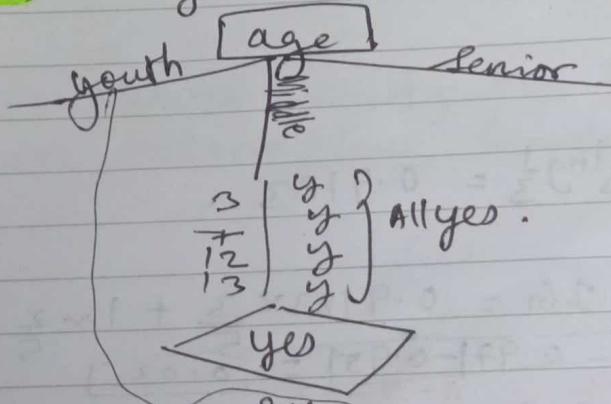
Gain Ratio = 0.0481

0.9852

= 0.0482

$$\text{Split Info} = -\frac{8}{14} \log \frac{8}{14} - \frac{6}{14} \log \frac{6}{14} = 0.9852$$

Max Gain ... age ... root ... max Gain Ratio ... Age



4	y
8	y
10	y
14	y

$$n(\text{yes}) = 3$$

$$n(\text{no}) = 2$$

$$\text{Gain} = -\frac{3}{5} \log \frac{3}{5} - \frac{2}{5} \log \frac{2}{5} = 0.971$$

$$\downarrow \text{Gain} = \frac{3}{5} \log \frac{3}{5} - \frac{2}{5} \log \frac{2}{5} = 0.971$$

income

$$IG = 1 * \frac{2}{5} = 0.4$$

$$\text{Gain} = 0.971 - 0.4 = 0.571$$

Gain Ratio = 0.571 / 0.5219 = 1.0971

$$\text{Split Info} = -\frac{2}{5} \log \frac{2}{5} - \frac{1}{5} \log \frac{1}{5} - \frac{2}{5} \log \frac{2}{5} = 1.5219$$

Student

y	2	0	2	0	$IG_1 = 0$
n	0	3	3	0	$\text{Gain} = 0.971 - 0 = 0.971$
$y \ n$	1	1	T		

split info = $-\frac{2}{5} \log \frac{2}{5} - \frac{3}{5} \log \frac{3}{5} = 0.9709$ Gain Ratio = $\frac{0.971}{0.9709} = 1$

Credit Rating

Fair	1	2	3	1	$IG = 1 \times \frac{2}{5} + \frac{3}{5} \times 0.9183 = 0.7443$
Excellent	1	1	2	1	
$y \ n$	1	1	T		

$\text{Gain} = 0.971 - 0.7443 = 0.2267$

split info = $-\frac{3}{5} \log \frac{3}{5} - \frac{2}{5} \log \frac{2}{5} = 0.9309$ Gain Ratio = $\frac{0.2267}{0.9709} = 0.2335$

Max Gain - Student

```

graph TD
    Root["student  
yes"] -- yes --> Leaf1
    Root -- no --> NodeNo["no"]
    NodeNo -- yes --> Leaf2
    NodeNo -- no --> Leaf3
  
```

Mex Gain Ratio - Student

for senior.

income

M	2	1	3	$-\frac{2}{3} \log \frac{2}{3} - \frac{1}{3} \log \frac{1}{3} = 0.9183$
L	1	1	2	1
H	0	0	0	0

$IG = 0.9183 \times \frac{3}{5} + 1 \times \frac{2}{5} = 0.951$

split Info = $-\frac{3}{5} \log \frac{3}{5} - \frac{2}{5} \log \frac{2}{5} = 0.9709$ Gain = $0.971 - 0.951 = 0.02$

Student

y	2	1	3	$-\frac{2}{3} \log \frac{2}{3} - \frac{1}{3} \log \frac{1}{3} = 0.9183$
n	1	1	2	1
$y \ n$	1	1	T	

$IG = 0.9183 \times \frac{3}{5} + 1 \times \frac{2}{5} = 0.92$

split ratio = $-\frac{3}{5} \log \frac{3}{5} - \frac{2}{5} \log \frac{2}{5} = 0.9709$

Gain Ratio = $\frac{0.02}{0.9709} = 0.0205$

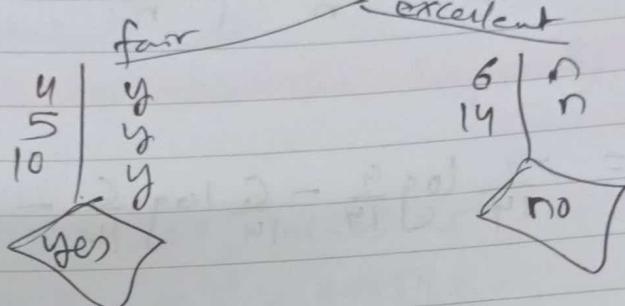
$$\text{Split Info} = -\frac{3}{5} \log \frac{3}{5} - \frac{2}{5} \log \frac{2}{5} = 0.9709$$

$$\text{Credit Rating} \quad \text{Gain ratio} = \frac{0.971}{0.9709} = 1$$

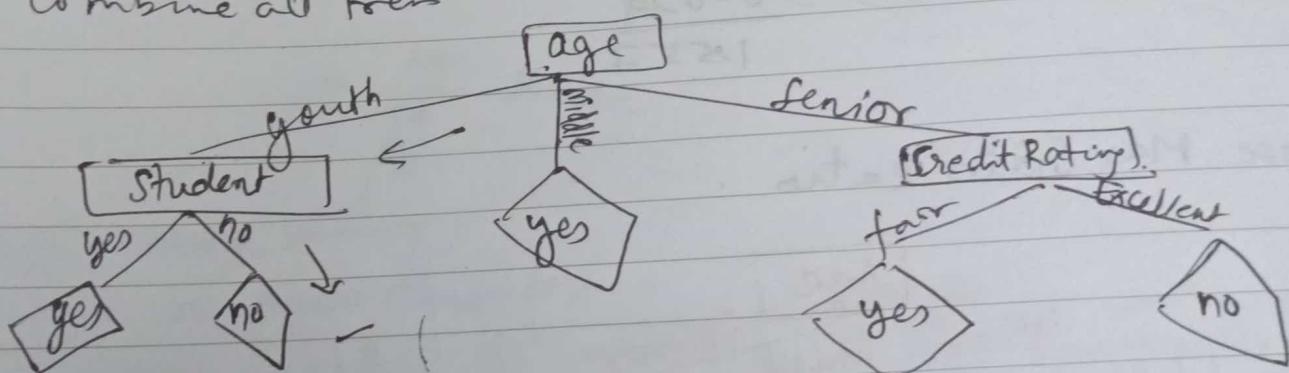
fair	3	0	3	0	$I_G = 0$
Excellent	0	2	2	0	$\text{Gain} = 0.971 - 0 = 0.971$

Max Gain ... credit rating

- Max Gain Ratio
Credit Rating



Combine all trees



Use this final tree to classify the record:

age = youth, income = low, Student = no, Credit Rating = Excellent

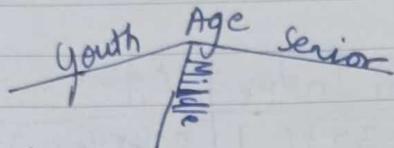
no

... same dataset.

$$n(\text{yes}) = 9, \quad n(\text{no}) = 5$$

$$\text{gini} = 1 - \left(\frac{9}{14}\right)^2 - \left(\frac{5}{14}\right)^2 = 0.459$$

$$\text{gini} = \frac{D_1}{\#n} \text{gini}(D_1) + \frac{D_2}{\#n} \text{gini}(D_2)$$



(Youth, middle, senior)

Consider all pairs (subsets)

$$\begin{aligned} \text{Youth, middle} & \left| \frac{9}{14} \left(1 - \left(\frac{6}{9}\right)^2 - \left(\frac{3}{9}\right)^2\right) + \frac{5}{14} \left(1 - \left(\frac{3}{5}\right)^2 - \left(\frac{2}{5}\right)^2\right) = 0.45 \right\} \text{Senior} \\ \text{Youth, senior} & \left| \frac{9}{14} \left(1 - \left(\frac{5}{9}\right)^2 - \left(\frac{4}{9}\right)^2\right) + \frac{5}{14} \left(1 - \left(\frac{4}{5}\right)^2 - \left(\frac{1}{5}\right)^2\right) = 0.357 \right\} \text{Youth, middle} \\ \text{middle, senior} & \left| \frac{9}{14} \left(1 - \left(\frac{2}{9}\right)^2 - \left(\frac{7}{9}\right)^2\right) + \frac{5}{14} \left(1 - \left(\frac{2}{7}\right)^2 - \left(\frac{3}{7}\right)^2\right) = 0.3936 \right. \end{aligned}$$

	income	#(ht+m) yes.	#(L) yes.
	high	medium	low
(h, m)	10	$\frac{10}{14} \left(1 - \left(\frac{2+2}{10}\right)^2 - \left(\frac{2+2}{10}\right)^2\right) + \frac{4}{14} \left(1 - \left(\frac{2}{4}\right)^2 - \left(\frac{1}{4}\right)^2\right) = 0.45$	
(m, l)	10	$\frac{10}{14} \left(1 - \left(\frac{3+4}{10}\right)^2 - \left(\frac{1+2}{10}\right)^2\right) + \frac{4}{14} \left(1 - \left(\frac{2}{5}\right)^2 - \left(\frac{2}{5}\right)^2\right) = 0.4428$	
(h, l)	8	$\frac{8}{14} \left(1 - \left(\frac{3+2}{8}\right)^2 - \left(\frac{1+2}{8}\right)^2\right) + \frac{6}{14} \left(1 - \left(\frac{4}{8}\right)^2 - \left(\frac{2}{8}\right)^2\right) = 0.4583$	
	+		

→ has difficulty when #classes ↑
 → equal-sized part favours + purity in both parts
 all yes

Income
 Yes
 no

$$(y, n) \Rightarrow y+n = 14$$

$$\frac{7}{14} \left(1 - \left(\frac{6}{7}\right)^2 - \left(\frac{1}{7}\right)^2 \right) + \frac{7}{14} \left(1 - \left(\frac{3}{7}\right)^2 - \left(\frac{4}{7}\right)^2 \right)$$

$$= 0.3673$$

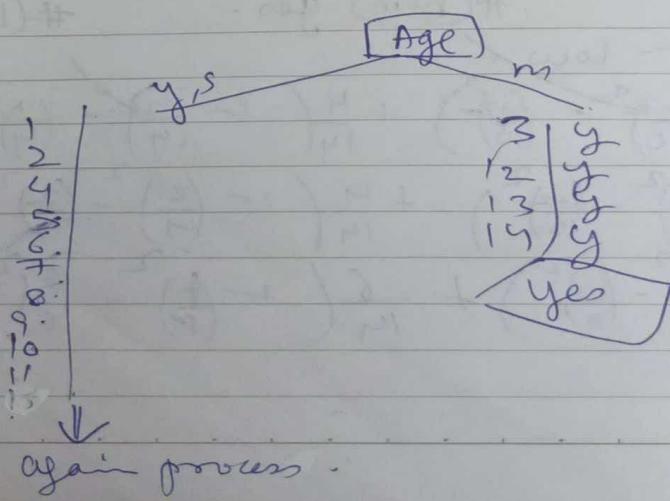
Credit Rating
 fair
 (f, e)
 excellent

$$f+e = 14$$

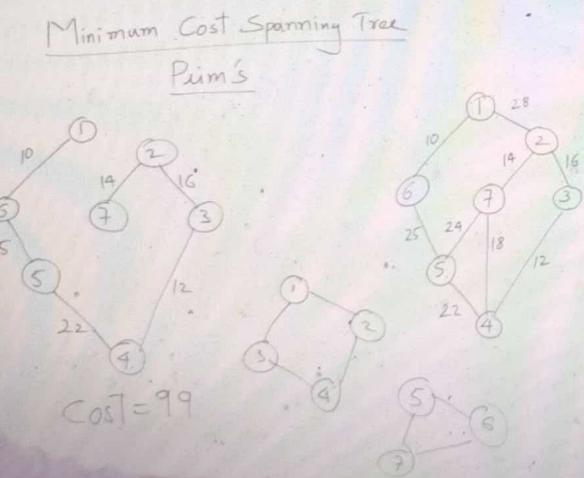
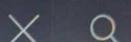
$$\frac{8}{14} \left(1 - \left(\frac{6}{8}\right)^2 - \left(\frac{2}{8}\right)^2 \right) + \frac{6}{14} \left(1 - \left(\frac{3}{6}\right)^2 - \left(\frac{3}{6}\right)^2 \right) = 0.4285$$

Attribute	split	Gini Index	$\Delta \text{Gini} =$
age	$(y, s) + m$	0.3571	$0.459 - 0.3571 = 0.1019$
income	$(m, b) + L$	0.4428	$0.459 - 0.4428 = 0.0162$
Student	Binary	0.3673	$0.459 - 0.3673 = 0.0917$
Credit Rating	Binary	0.4285	$0.459 - 0.4285 = 0.0305$

min. Gini Index
↑ Red: in Impurity



Kruskal's algorithm



*Untitled - Notepad

File Edit Format View Help
prim

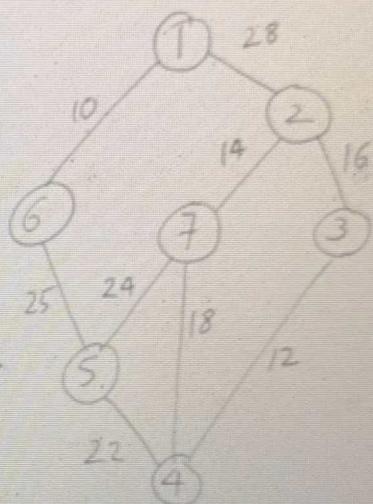
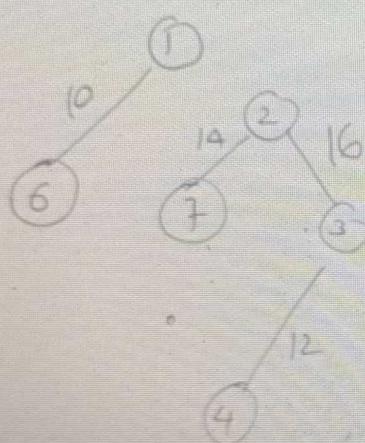
next min weight edge must be a neighbor so that the graph is connected at every step

kruskal

select most min only. its ok if its not adj.

Minimum Cost Spanning Tree

Kruskal's



12:38 / 20:11 • Kruskals... >



and Kruskals Algorithms - Greedy Method

Subscribe

39K

Share

ers

earch



Ex. 4.4.1 : Calculate the radial distance of point A having value of 6 with respect to point B and C having value of 8 and 12 respectively.

Soln. :

Assume $\gamma = 1$

Distance of point A with respect to B = $e^{-\gamma(x-y)^2} = e^{-1(6-8)^2} = e^{-4} = 0.018$

Distance of point A with respect to C = $e^{-\gamma(x-y)^2} = e^{-1(6-12)^2} = e^{-36} = 0$

So, point B has higher influence over point A than point C.

(Copyright No. L-98904/2021)

Practice Questions

Ex. 2.1.1: For the following data set, find the linear regression line. Predict the value of Y if X = 10.

X	Y
0	1
1	3
2	2
3	5
4	7
5	8
6	8
7	9
8	10
9	12

$$x_i - \bar{x} \quad y_i - \bar{y} \quad (x_i - \bar{x})(y_i - \bar{y}) \quad (x_i - \bar{x})^2$$

Soln. : Calculate the mean of X and Y and then calculate the other values as required.

X	Y	$X_i - \bar{X}$	$Y_i - \bar{Y}$	$(X_i - \bar{X})(Y_i - \bar{Y})$	$(X_i - \bar{X})^2$
0	1	-4.5	-5.5	24.75	20.25
1	3	-3.5	-3.5	12.25	12.25
2	2	-2.5	-4.5	11.25	6.25
3	5	-1.5	-1.5	2.25	2.25
4	7	-0.5	0.5	-0.25	0.25
5	8	0.5	1.5	0.75	0.25
6	8	1.5	1.5	2.25	2.25
7	9	2.5	2.5	6.25	6.25
8	10	3.5	3.5	12.25	12.25
9	12	4.5	5.5	24.75	20.25
\bar{X} (Mean of X) = 4.5		\bar{Y} (Mean of Y) = 6.5		Total = 96.5	Total = 82.5

Hence,

$$\beta_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$\beta_1 = \frac{96.5}{82.5} = 1.1696$$

$$\beta_0 = \bar{y} - \beta_1 \bar{x}$$

$$\beta_0 = 6.5 - 1.1696 \times 4.5 = 1.2368$$

$$Y = \beta_0 + \beta_1 X$$

$$Y = 1.2368 + 1.1696X$$

Hence, the regression line is

A sample plot of the regression line is as shown in Fig. P. 2.1.1.

Regression Line Plot

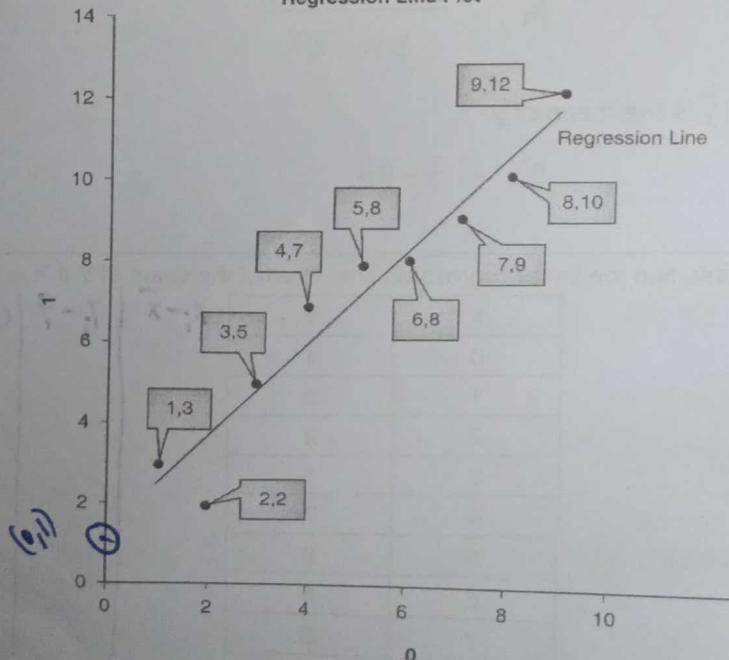


Fig. P. 2.1.1

Ex. 2.4.1 : Calculate the entropy of the following distribution.

Gender	Count
Male	9
Female	5

Soln. :

- $p(\text{Male}) = \frac{9}{14}$
- $p(\text{Female}) = \frac{5}{14}$

So, for calculating the entropy of the distribution you would use the following formula.

$$H_x = - \sum_{i=1}^n p(x_i) \log_2 p(x_i)$$

$$\text{Entropy} = - \left(\frac{5}{14} \log_2 \frac{5}{14} + \frac{9}{14} \log_2 \frac{9}{14} \right)$$

$$\text{Entropy} = - (- 0.53 - 0.40)$$

$$\text{Entropy} = 0.93$$

x. 2.4.3 : Calculate the entropy of the following distribution.

Fruit Colour	Taste	Count
Yellow	Sweet	10
Red	Sweet	5
Green	Sour	15
Orange	Sour	5

35

Soln. :

Note here that there are two input variables - colour and taste.

Assuming that the target of entropy calculation is on the basis of fruit colour. In this scenario, the taste column is then used as weight for entropy calculation.

$$p(\text{Yellow}) = \frac{10}{35}$$

$$p(\text{Red}) = \frac{5}{35}$$

$$p(\text{Green}) = \frac{15}{35}$$

$$p(\text{Orange}) = \frac{5}{35}$$

$$p(\text{Sweet}) = \frac{15}{35}$$

$$p(\text{Sour}) = \frac{20}{35}$$

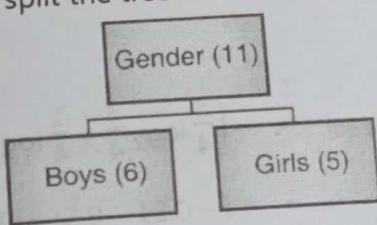
$$\begin{aligned} \text{Entropy} &= \text{Sweet} * E(\text{Yellow}) + \text{Sweet} * E(\text{Red}) + \\ &\quad \text{Sour} * E(\text{Green}) + \text{Sour} * E(\text{Orange}) \\ &= \left[\frac{15}{35} \times \frac{10}{35} \log \frac{10}{35} + \frac{15}{35} \times \frac{5}{35} \log \frac{5}{35} \right] + \dots \\ &= 0.92 \end{aligned}$$

Ex. 2.4.4 : Which attribute would you choose to split the following dataset into a decision tree?

Gender	Marks
Girl	65
Girl	46
Boy	56
Boy	43
Boy	53
Boy	49
Girl	42
Boy	84
Boy	44
Girl	42
Girl	40

Soln.: There are totally 11 datapoints. There are 5 girls and 6 boys. Also, there are 7 students with marks below 50 and 4 students with marks above 50. Now, you have to find out which of these two attributes (gender or marks) is better for splitting the dataset into a decision tree.

Let's take scenario 1 where you decide to split the tree based on gender.



Calculate entropy of the parent node (Gender).

$$p(\text{Boys}) = \frac{6}{11}$$

$$p(\text{Girls}) = \frac{5}{11}$$

$$\text{Entropy (Parent Node)} = -\left(\frac{6}{11} \log_2 \frac{6}{11} + \frac{5}{11} \log_2 \frac{5}{11}\right)$$

$$\text{Entropy (Parent Node)} = -(-0.48 - 0.52)$$

$$\text{Entropy (Parent Node)} = 1$$

After split,

$$\text{Entropy (Boys)} = -\left(\frac{6}{6} \log_2 \frac{6}{6}\right) = 0$$

$$\text{Entropy (Girls)} = -\left(\frac{5}{5} \log_2 \frac{5}{5}\right) = 0$$

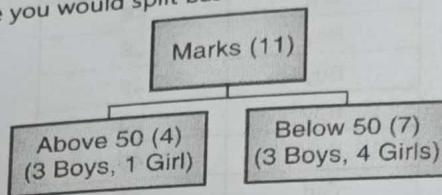
$$\text{Weighted entropy of child nodes} = \frac{6}{11} \times \text{Entropy (Boys)} + \frac{5}{11} \times \text{Entropy (Girls)} = 0$$

Hence,

$$\text{Information Gain, IG (Gender)} = \text{Parent Entropy} - \text{Child Entropy} = 1 - 0 = 1$$

Information Gain of 1 indicates that after split, pure sets of boys and girls are formed and splitting on gender makes sense.

Now, let's take scenario two where you would split based on marks scored (less than 50, 50 or above it).



After split,

Above 50

- $p(\text{Boys}) = \frac{3}{4}$
- $p(\text{Girls}) = \frac{1}{4}$

Below 50

- $p(\text{Boys}) = \frac{3}{7}$
- $p(\text{Girls}) = \frac{4}{7}$

$$\text{Entropy (Above 50)} = -\left(\frac{3}{4} \log_2 \frac{3}{4} + \frac{1}{4} \log_2 \frac{1}{4}\right)$$

$$= -(-0.31 - 0.5) = 0.81$$

$$\text{Entropy (Below 50)} = -\left(\frac{3}{7} \log_2 \frac{3}{7} + \frac{4}{7} \log_2 \frac{4}{7}\right)$$

$$= -(-0.52 - 0.46) = 0.98$$

Weighted entropy of child nodes

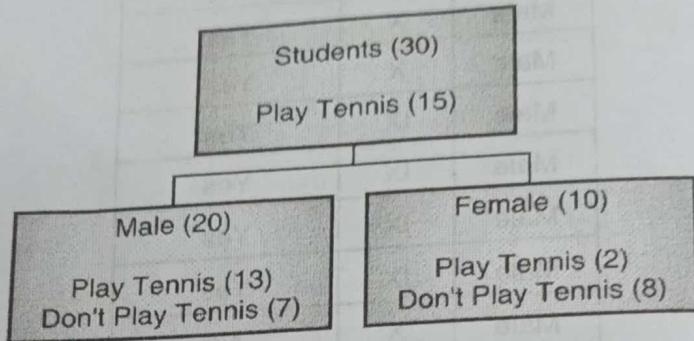
$$\begin{aligned} &= \frac{4}{11} \times \text{Entropy (Above 50)} + \frac{7}{11} \times \text{Entropy (Below 50)} \\ &= \frac{4}{11} \times 0.81 + \frac{7}{11} \times 0.98 \\ &= 0.29 + 0.62 = 0.91 \end{aligned}$$

Hence,

$$\begin{aligned} \text{Information Gain, IG (Marks)} &= \text{Parent Entropy} - \text{Child Entropy} \\ &= 1 - 0.91 = 0.09 \end{aligned}$$

As you see, the information gain on splitting the dataset by gender is higher than splitting by marks. Hence, suggested to split the given dataset on gender.

Scenario 1 : Split on Gender



$$\text{Gini for sub - nodes} = (P_i)^2$$

$\text{Gini}_{(\text{Male} | \text{Tennis})}$ = Square of probability of playing tennis + Square of probability of not playing tennis

$$\text{Gini}_{(\text{Male} | \text{Tennis})} = \left(\frac{13}{20}\right)^2 + \left(\frac{7}{20}\right)^2 = 0.42 + 0.12 = 0.54$$

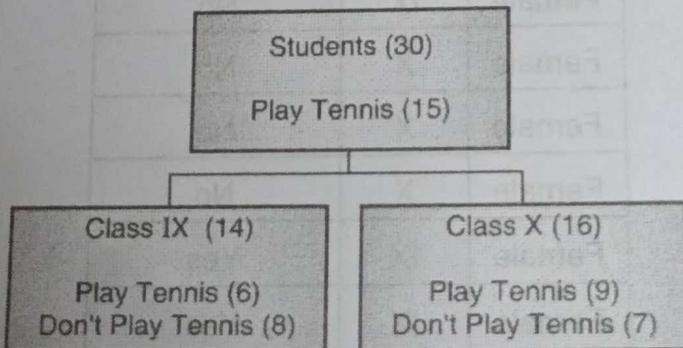
$\text{Gini}_{(\text{Female} | \text{Tennis})}$ = Square of probability of playing tennis + Square of probability of not playing tennis

$$\text{Gini}_{(\text{Female} | \text{Tennis})} = \left(\frac{2}{10}\right)^2 + \left(\frac{8}{10}\right)^2 = 0.04 + 0.64 = 0.68$$

$\text{Gini Index}_{(\text{Gender})}$ = Weight of sub-node (Male) \times $\text{Gini}_{(\text{Male} | \text{Tennis})}$ + Weight of sub-node (Female) \times $\text{Gini}_{(\text{Female} | \text{Tennis})}$

$$\text{Gini Index}_{(\text{Gender})} = \frac{20}{30} \times 0.54 + \frac{10}{30} \times 0.68 = 0.36 + 0.23 = 0.59$$

Scenario 2 - Split on Class



Date: / /

$$C = \begin{bmatrix} \text{cov}(x_1, y_1) & \text{cov}(x_1, y_2) \\ \text{cov}(y_1, x_1) & \text{cov}(y_1, y_2) \end{bmatrix}$$

$$C = \begin{bmatrix} 14 & -11 \\ -11 & 23 \end{bmatrix}$$

$$\lambda = \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix}$$

$$C - \lambda I = 0$$

$$\begin{bmatrix} 14 - \lambda & -11 \\ -11 & 23 - \lambda \end{bmatrix} = 0$$

$$(14 - \lambda)(23 - \lambda) - 121 = 0$$
$$(14)(23) - (14 + 23\lambda) + \lambda^2 - 121 = 0$$

$$\lambda^2 - 37\lambda + 201 = 0$$

$$\boxed{\lambda_1 = 30.3849}$$
$$\boxed{\lambda_2 = 6.615}$$

For $\lambda_1 = 30.3849$

$$(-\lambda_1 I)(x) = 0$$

$$\begin{bmatrix} 14 - \lambda_1 & -11 \\ -11 & 23 - \lambda_1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 0$$

$$(14 - \lambda_1)x_1 - 11x_2 = 0$$

$$-11x_1 + (23 - \lambda_1)x_2 = 0$$

$$\therefore \boxed{\frac{x_1}{11} = \frac{x_2}{14 - \lambda_1} = e}$$

$$e_1^T = \begin{bmatrix} 0.5574 & -0.8303 \end{bmatrix}$$

$$e_2^T = \begin{bmatrix} 0.8303 & 0.5574 \end{bmatrix}$$

$$P_{11} \quad P_{12} \quad P_{13} \quad P_{14}$$

PC1

$$P_{11} = e_1^T \begin{bmatrix} 4-8 \\ 11-8.5 \end{bmatrix} = -4.3052$$

$$P_{12} = e_1^T \begin{bmatrix} 8-8 \\ 4-8.5 \end{bmatrix} = 3.7363$$

$$P_{13} = 5.6930$$

$$P_{14} = -5.1240$$

$$\boxed{PC1 \approx \begin{bmatrix} -4.3052 & 3.7363 & 5.6930 & -5.1240 \\ P_{11} & P_{12} & P_{13} & P_{14} \end{bmatrix}}$$

2023

$$2-23-0-0506 \quad \left[\begin{array}{c} a_{21} \\ a_{22} \end{array} \right] = 0.$$

$$\left| \begin{array}{l} a_{21} = 0.79 \\ a_{22} = 0.61 \end{array} \right.$$

ments are given as
 $a_{12}x_2$

$$\left| z_2 = a_{21}x_1 + a_{22}x_2 \right.$$

PAGE
DATE

PC

$$C_1 \left(\begin{array}{l} 4-8 \\ 11-8.5 \end{array} \right)$$

$$8-8$$

$$4-8.5$$

$$13-8$$

$$5-8.5$$

$$7-8$$

$$14-8.5$$

Total

Predicted

		akiec	bcc	df	
		akiec	bcc	df	
Actual	akiec	1,1 ✓	1,2 ✓	1,3 ✓	
	bcc	2,1 ✓	2,2 ✓	2,3 ✓	
	df	3,1 ✓	3,2 ✓	3,3 ✓	
TP	1,1	2,2	3,3		
FP	(1,2)(1,3)	(2,1)(2,3)	(3,1)(3,2)		
TN	(2,2)(2,3) (3,2)(3,3)	(1,1)(1,3) (3,1)(3,3)	(1,1)(1,2) (2,1)(2,2)		
FN	(2,1)(3,1)	(1,2)(3,2)	(3,1)(3,3)		

+ Predicted -

Actual		+	TP	FN	
		-	FP	TN	
+	-				

TP	(1,1)
FP	(1,2)(1,3)
FN	(2,1)(3,1)
TN	(2,2)(2,3) (3,2)(3,3)