

$$M_1 = 10.664 \quad M_2 = 0.485$$

$$M_3 = 0.2, 0.92.$$

Three lines (hyperplane)

$$\text{eq ① } 2x_1 + 3x_2 - 5 = 0 \quad \beta_0 = -5 \quad \beta_1 = 2 \quad \beta_2 = 3$$

$$\text{for } x_1 = 3, x_2 = 4 \quad d_1 = -5 + 2(3) + 3(4) - 5 = 13$$

Dataset

	x_1	x_2	y
1	3	4	+1
2	2	3	+1
3	1	-1	-1
4	-2	+1	-1

$$d_2 = -5 + 2(2) + 3(3) - 5 = 13$$

$$= -5 + 6 + 12 = \frac{13}{3.6} = \cancel{6.38} \quad 3.6$$

$$\text{for } x_1 = 1, x_2 = -1$$

$$d_3 = -5 + 2(1) + 3(-1) - 5 = 13$$

$$= -5 + 2 - 5 = \frac{8}{3.6} = \cancel{2.22}$$

$$\text{for } x_1 = -2, x_2 = +1$$

$$d_4 = -5 + 2(-2) + 3(1) - 5 = \frac{10}{3.6} = 2.7$$

$$\text{for } x_1 = -2, x_2 = -1$$

$$d_5 = -5 + 2(-2) + 3(-1) - 5 = \frac{-16}{3.6} = \cancel{-4.44}$$

$$\Rightarrow \underline{\underline{1.6}}$$

$$28 \mu A + H = 1.8 \text{ m} \rightarrow H$$

$$28 \mu A + H = 1.8 \text{ m}$$

$$M_1 = 1.66$$

(most probable result)

$$\text{eq. } ②. -x_1 + 4x_2 + 7 = 0$$

$$\text{for } x_1 = 3, x_2 = 4$$
$$\frac{-3 + 4(4) + 7}{\sqrt{-1 + (4)^2}} = \frac{-3 + 16 + 7}{\sqrt{-1 + 16}} = \frac{20}{\sqrt{15}} = \frac{20}{3.8} = 5.2$$

$$\text{for } x_1 = 2, x_2 = 3$$
$$\frac{-2 + 4(3) + 7}{\sqrt{-1 + (3)^2}} = \frac{-2 + 12 + 7}{\sqrt{-1 + 9}} = \frac{17}{\sqrt{8}} = \frac{17}{3.8} = 4.47$$

$$\text{for } x_1 = 1, x_2 = -1$$
$$\frac{-1 + 4(-1) + 7}{\sqrt{-1 + (-1)^2}} = \frac{-1 - 4 + 7}{\sqrt{-1 + 1}} = \frac{2}{\sqrt{0}} = \frac{2}{0} = \infty$$
$$\frac{-4 + -4 + 7}{\sqrt{-1 + 0^2}} = \frac{-8 + 7}{\sqrt{-1 + 0}} = \frac{-1}{\sqrt{-1}} = \frac{-1}{i} = 1$$

$$\text{for } x_1 = -2, x_2 = 1$$
$$\frac{-(-2) + 4(1) + 7}{\sqrt{-1 + (1)^2}} = \frac{2 + 4 + 7}{\sqrt{-1 + 1}} = \frac{13}{\sqrt{0}} = \frac{13}{0} = \infty$$
$$\frac{2 + 4 + 7}{\sqrt{-1 + 0^2}} = \frac{13}{\sqrt{-1 + 0}} = \frac{13}{i} = -13$$

Evaluation metrics

Actual values
not predicted ones
red circle → positive class
white circle → negative class
threshold = 0.5

Confusion Matrix

	Original +ve	orig. -ve
+ve	9 (TP)	1 (FP)
-ve	1 (FN)	8 (TN)
T = True		
P = +ve		
N = orig. -ve		
F = False.	-ve	

$$\begin{aligned} \text{Accuracy} &= (TP + TN) / \text{Total} \\ &= (9+8) / 20 \\ &= \underline{\underline{0.85}}. \end{aligned}$$

$$\text{F}_1 \text{ score} = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

$$\text{precision} = \frac{TP}{TP + FP}$$

$$\text{recall} = \frac{TP}{TP + FN}$$

precision = tp / (tp + fp) = 4 / (4 + 2) = 2/3 ≈ 0.67

—/—

Observation	Actual score	Predicted score	Predicted label
1	1	0.85	1
2	0	0.60	1 * FP
3	1	0.70	1
4	1	0.40	0 * FN
5	0	0.55	1 * PP
6	1	0.50	1
7	0	0.65	1 * FP
8	0	0.35	0
9	1	0.60	1
10	0	0.20	0

positive + predicted

Compute the following : 1) Accuracy

2) Precision

threshold = 0.5 3) Recall or sensitivity

4) Specificity

Accuracy = $(TP + TN) / \text{Total}$ 5) False Positivity

= 6/10 = 0.6

Confusion Matrix

6) F1-score

	neg-svc	sig-rev
pred-neg	✓ 4	✓ 3
pred-pos	1	✓ 2

$$\begin{aligned} \text{Accuracy} &= (4+2)/10 \\ &= 6/10 \\ &= 0.6 \end{aligned}$$

11

2) Precision = $\frac{TP}{TP + FP} = \frac{4}{4+3} = \frac{4}{7} = 0.57$

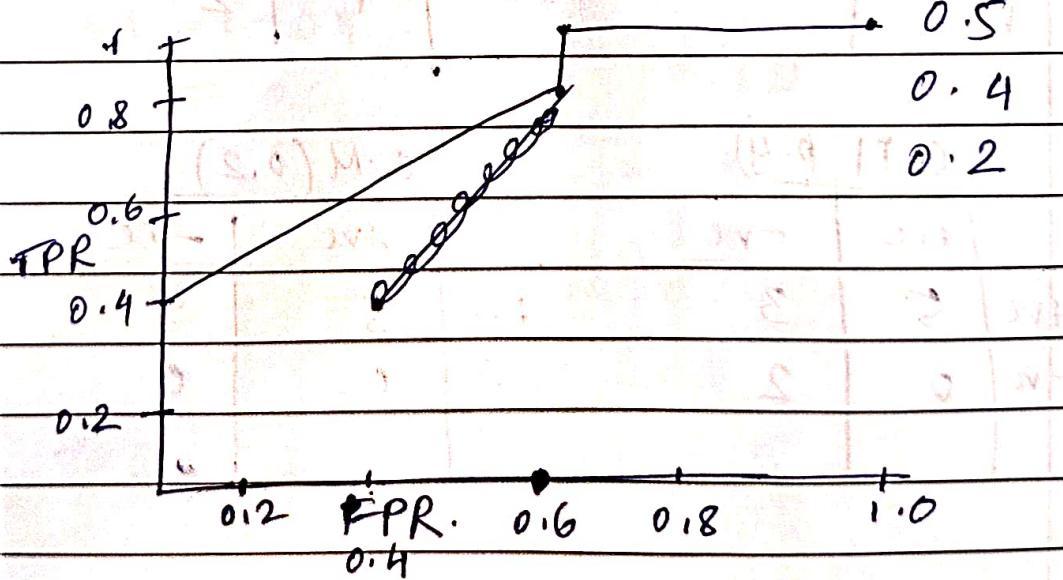
3) Recall = $\frac{TP}{TP + FN} = \frac{4}{4+1} = \frac{4}{5} = 0.8$

4) Specificity = Negative Recall = $TN = \frac{2}{TN + FP} = \frac{2}{2+3} = \frac{2}{5} = 0.4$

5) False Positivity =

6) F₁ score = $2 \times \frac{P \times R}{P+R} = 2 \times \frac{0.57 \times 0.66}{0.57 + 0.66} = 2 \times \frac{0.648}{1.23} = 2 \times 0.52 = 0.648$

7) Roc curve \rightarrow Threshold values $\rightarrow 0.7, 0.5, 0.4, 0.2$



Actual score Predicted score 0.7 0.5 0.4 0.2

1	0.85	1	1	1	1
0	0.60	0	1	1	1
1	0.70	1	1	1	1
0	0.40	0	0	1	1
0	0.55	0	1	1	1
1	0.50	0	1	1	1
0	0.65	0	1	1	1
0	0.35	0	0	0	1
1	0.60	0	1	1	1
0	0.20	0	0	0	1

CONFUSION MATRIX (C.M.)

C.M. (0.7) C.M. (0.5)

pred		org		C.M. (0.7)		C.M. (0.5)	
		+ve	-ve	+ve	-ve	+ve	-ve
+ve	2	0	-ve	4	3	-ve	-ve
-ve	3	5	-ve	1	2	-ve	-ve

C.M. (0.4)

C.M. (0.2)

		C.M. (0.4)		C.M. (0.2)	
		+ve	-ve	+ve	-ve
+ve	5	3	+ve	5	5
-ve	0	2	-ve	0	0

0.2

$$FPR = \frac{0}{0+5} = 1 - \frac{0}{5} = 1 - 0 = 1$$

1/1

$$FPR = \frac{5}{5+0} = 5/5 = 1$$

Threshold	TPR	FPR
0.7	0.4	0
0.5	0.8	0.6
0.4	1	0.6
0.2	(1-0.2) * (1-0.4) = 0.6 * 0.6 = 0.36	0.6

0.7

$$\frac{TPR}{FPR} = 1 - \text{specificity}$$

$$1 - \frac{TN}{TN+FP} \Rightarrow 1 - \frac{5}{5+0} = 1 - \frac{5}{5} = 1 - 1 = 0$$

TPR/FPR = Sensitivity; Recall

$$= \frac{TP}{TP+FN} = \frac{2}{2+3} = \frac{2}{5} = 0.4$$

= 0.285

0.5

$$FPR = 1 - 0.9$$

$$TPR = \frac{TP}{TP+FN}$$

0.4

$$= 0.6$$

$$= \frac{4}{4+1} = \frac{4}{5} = 0.8$$

0.4

$$FPR = 1 - \frac{TN}{TN+FP} = 1 - \frac{2}{2+3} = 1 - \frac{2}{5} = \underline{\underline{0.6}}$$

$$TPR = \frac{TP}{TP+FN} = \frac{5}{5+0} = \frac{5}{5} = 1$$

$$AUC = \sum_{j=1}^4 (FPR_i - FRR_{i-1}) * (TPR_i + TPR_{i-1})$$

$\frac{1}{2}$ | P_C | $\frac{1}{2}$

$\frac{0.7}{2} \Rightarrow 0$ | 0.8 | 0.5

$$\underline{0.5} \Rightarrow \frac{(0.6 - 0)}{2} \times (0.8 + 0.4)$$

$$= \frac{0.6}{2} \times \frac{1.2}{2} = \frac{0.72}{2} = 0.36$$

$$\underline{0.4} \Rightarrow (0.6 - 0.6) \times \frac{0.4}{2} = 0.$$

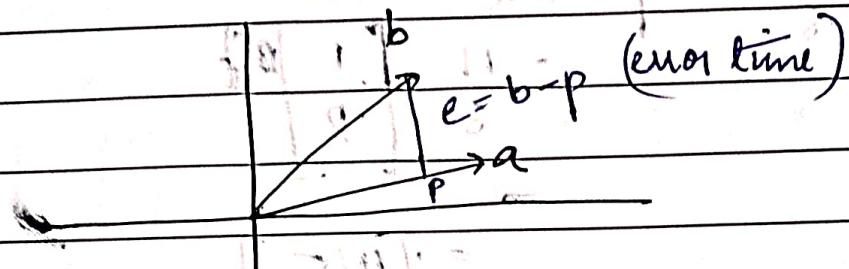
$$\underline{0.2} \Rightarrow (1 - 0.6) \times \frac{(1+1)}{2} = 0.4$$

$$= \frac{0.4}{2} \times 2 = \frac{0.8}{2} = 0.4$$

$$0.36 + 0.4 = 0.76$$

a is perpendicular to e

$$a \perp^+ e \quad a^T e = 0 \quad (e = b - p)$$



$$p = x a$$

$$a^T (b - p) = 0$$

$$a^T (b - x a) = 0$$

$$a^T b = x a^T a$$

if a is unit vector

$$a^T a = 1$$

$$x = \frac{a^T b}{a^T a}$$

scalar

$$p = x a \Rightarrow a^T b = x a^T a$$

if you want to

project b onto a \rightarrow do a projection matrix of a and multiply it with b .

$$P = A (A^T A)^{-1} A^T b$$

$$b = \begin{pmatrix} 3 \\ 4 \end{pmatrix} \quad a = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$$

A diagram illustrating the decomposition of a vector b into a component parallel to a line (vector p) and a component perpendicular to the line (vector e). The vector b is shown originating from the same point as vector a . A dashed line represents the direction of a . Vector p is the projection of b onto this line, and vector $e = b - p$ is the part of b that is perpendicular to the line.

find p

$$x = \frac{a^T b}{a^T a} \Rightarrow \begin{bmatrix} 1 & 2 \end{bmatrix} \begin{bmatrix} 3 \\ 4 \end{bmatrix} = \begin{bmatrix} 3+8=11 \\ 1+4=5 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix} = \begin{bmatrix} 3+6=9 \\ 1+4=5 \end{bmatrix}$$

x_0

such that $x_0 \in \mathbb{R}^n$

$$(q-d) \cdot x_0 = q \cdot x_0 - d \cdot x_0$$

$$x = \frac{11}{8}$$

$$p = Xa$$

(want max)

$$x = \frac{11}{5}$$

$$\begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

$$= \begin{bmatrix} 11/5 \\ 22/5 \end{bmatrix}$$

$$0 = (q-d) \cdot x_0$$

$$X \cdot x = \begin{bmatrix} 22/5 \end{bmatrix}$$

$$0 = (q-d) \cdot x_0$$

$$\vec{a}^T \vec{x} = \vec{a}^T \vec{a} + \vec{a}^T \vec{b} =$$

$$\vec{a}^T \vec{a} = \lambda$$

$$\begin{bmatrix} 1 \\ 2 \end{bmatrix} \begin{bmatrix} 1 & 2 \end{bmatrix} = \begin{bmatrix} 3 \\ 4 \end{bmatrix}$$

$$\begin{bmatrix} 3 \\ 4 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 \\ 2 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix} = \begin{bmatrix} 5 \\ 10 \end{bmatrix}$$

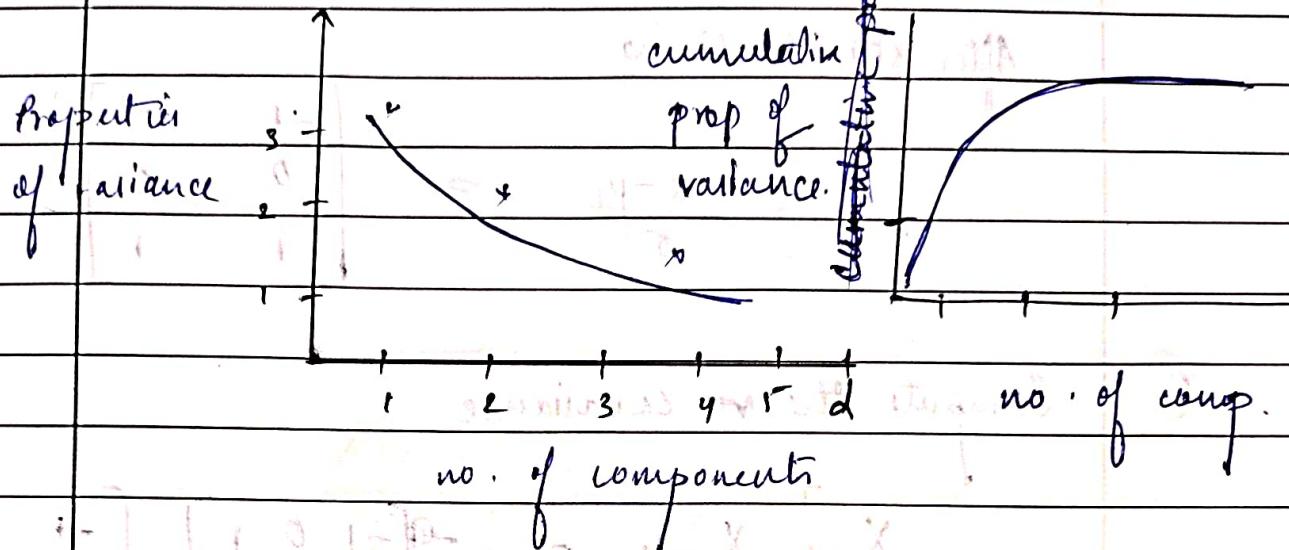
Max 2nd col

$$\Rightarrow \begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix} = \begin{bmatrix} 3 \\ 4 \end{bmatrix}$$

$$\begin{bmatrix} 3 \\ 4 \end{bmatrix}$$

$$\Rightarrow \underline{\underline{3}} \quad \underline{\underline{3}} \quad \underline{\underline{3+8}} \quad \underline{\underline{6+16}}$$

$$\begin{bmatrix} 1 \\ 2 \end{bmatrix} = \begin{bmatrix} 11 \\ 22 \end{bmatrix} = \begin{bmatrix} 11/5 \\ 22/5 \end{bmatrix}$$



$$X = \begin{bmatrix} 2 & 100 \\ -4 & 200 \\ 6 & 300 \end{bmatrix} \quad \text{apply PCA} \rightarrow 2 = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

① Standardize each column $\Rightarrow x - \bar{x} / \sigma$

$$2 + 4 + 6 = 12/3 = 4. \quad \mu_2 = 200$$

$$\mu_1 = 4 \quad \sigma = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

~~$$\sigma = \sqrt{\frac{2+4+6}{2}} = \sqrt{12} = 2\sqrt{3}$$~~

$$\mu_1 = 4$$

$$\mu_2 = 200$$

$$\sigma_1 = 2$$

$$\sigma_2 = 100$$

$$\sigma_2 = 100$$

$$= 4$$

$$0$$

$$4$$

$$\frac{8}{8}$$

$$(100)^2 = 10000$$

$$(100)^2 = 10000$$

$$\frac{10000}{20000}$$

$$= \sqrt{100}$$

After standardising

$$\frac{x - \mu_1}{\sigma} = \begin{bmatrix} -1 & -1 \\ 0 & 0 \\ 1 & 1 \end{bmatrix}$$

②

Compute the covariance.

$$X_{\text{std}} X_{\text{std}}^T = \begin{bmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} -1 & -1 \\ 0 & 0 \\ 1 & 1 \end{bmatrix}$$

$$= \begin{bmatrix} 1+0+1 & 1+0+1 \\ 1+0+1 & 1+0+1 \end{bmatrix} = \begin{bmatrix} 2 & 2 \\ 2 & 2 \end{bmatrix} / n-1$$

$$\Rightarrow \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

Ax ③

Compute eigen values & eigen vectors.

needed

④

Decide the % of variance and choose k principle components.

⑤

Transform X into new space of principle components.

⑥

Use these 2 to train a model.

$$\lambda_1 = 0, \lambda_2 = 2$$

$$\det(A - \lambda I) = 0$$

$$A = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \quad I = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$A = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \quad A - \lambda I = \begin{bmatrix} 1-\lambda & 0 \\ 0 & 1-\lambda \end{bmatrix}$$

$$A - \lambda I = \begin{bmatrix} 1-\lambda & 1 \\ 1 & 1-\lambda \end{bmatrix}$$

$$(1-\lambda)(1-\lambda) - 1 = 0$$

$$\Rightarrow 1 - 2\lambda + \lambda^2 - 1 = 0$$

$$\Rightarrow -2\lambda + \lambda^2 = 0$$

$$\lambda(-2 + \lambda) = 0$$

$$\lambda = 0 \quad \lambda = 2$$

$$A - 2I = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} - 2 \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} - \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$

$$= \begin{bmatrix} -1 & 1 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$-x + y = 0$$

$$\begin{matrix} x \\ y \end{matrix}$$

$\lambda_1 = 2, \lambda_2 = 1$

$\sqrt{\text{length}^2}$

$$0 = (1 + \frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}})$$

$$\begin{bmatrix} -1 & 1 \\ 1 & 0 \end{bmatrix} \underline{v_1} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\begin{bmatrix} -1 & 1 \\ 1 & 0 \end{bmatrix} v_2 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

(5)

$$z_{\text{PCA}} = \begin{bmatrix} 8 & 6 \\ 4 & 3 \end{bmatrix} \times v_2$$

$$= \begin{bmatrix} -1 & -1 \\ 0 & 0 \\ 1 & 1 \end{bmatrix} \times \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$= \begin{bmatrix} -1 & -1 \\ 0 & 0 \\ 1 & 1 \end{bmatrix} \times \begin{bmatrix} i/\sqrt{2} \\ 1/\sqrt{2} \end{bmatrix}$$

$$= \begin{bmatrix} \frac{-1+i}{\sqrt{2}} \\ 0 \\ \frac{1+i}{\sqrt{2}} \end{bmatrix} = \begin{bmatrix} -2 \\ 0 \\ 2 \end{bmatrix}$$

Eigen vector capturing the largest eigen value will give the highest variance.

a) for each of K-clusters, compute the cluster centroid.

b) Assign each observations to the cluster whose centroid is closest (Euclidean distance).

Sample no.	x_1	x_2	$R=2$	Sample no.	2 dimensions
1	1	1		1)	$C_1 = \{2, 3, 4, 6\}$
2	1	2			$C_2 = \{1, 5\}$
3	2	2			Centroid d
4	8	8	$C_4 = \text{avg} \rightarrow 3 + 5$		$\begin{bmatrix} 5 \\ 5 \end{bmatrix}$
5	8	9			
6	9	8	take the sample no. \rightarrow add them	$C_2 = \begin{bmatrix} 4.5 \\ 5 \end{bmatrix}$	
			and take that divide it by the no. of sample.		

$$d = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \leftarrow 1^{\text{st}} \text{ sample. } C_1 - d = \begin{bmatrix} 4 \\ 4 \end{bmatrix} = 8.5 \sqrt{2}$$

$$C_2 - d = \begin{bmatrix} 3.5 \\ 4 \end{bmatrix} = 5.3$$

After 1st iteration

$$\text{or } C_2 = \{1\}$$

$$d = \begin{bmatrix} 1 \\ 2 \end{bmatrix} = C_1 - d = \begin{bmatrix} 4 \\ 3 \end{bmatrix} = 5 \quad C_2 - d = \begin{bmatrix} 3.5 \\ 3 \end{bmatrix} = 4.6$$

standard set $C_2 = \{1, 2\}$ with 21. p. values and 100

1. bivariate

standard set $d = \begin{bmatrix} 2 \\ 2 \end{bmatrix}$ of $\Rightarrow C_2 - d = \begin{bmatrix} 3 \\ 3 \end{bmatrix} \Rightarrow 4.2$

(2. bivariate) $\Rightarrow \begin{bmatrix} 3 \\ 3 \end{bmatrix}$ subtract

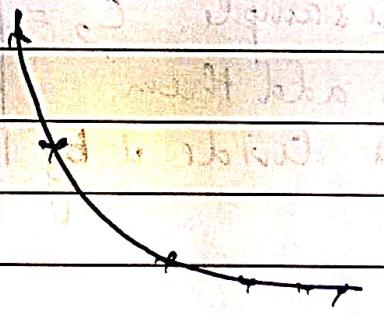
$C_2 - d = \begin{bmatrix} 2.5 \\ 2 \end{bmatrix} = 3.2$

$C_2 = \{1, 2, 3\}$

bivariate

for practical \rightarrow EDA insights to be written.

Score



-2.00

[1] [2] [3] - 70% 100 \Rightarrow [1] - 6

[1] K. [1] [1]

[1] [2] [3] - 70% 100 \Rightarrow [1] - 6

[1] [2] [3] - 70% 100 \Rightarrow [1] - 6

1. bivariate

Score

1. bivariate

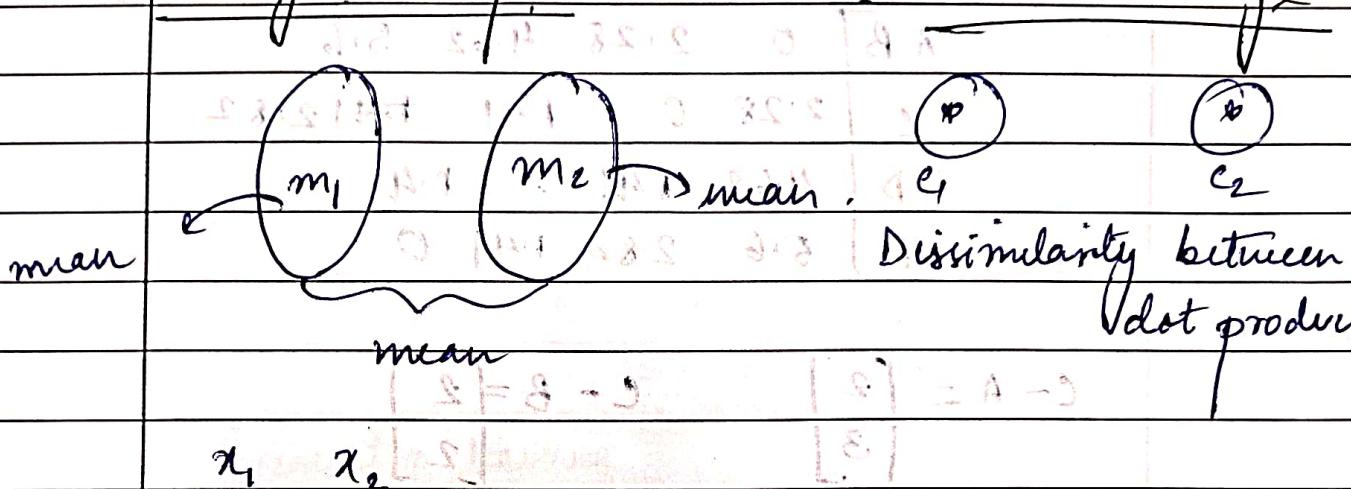
use the min. number to fuse

for complete linkage max of two}

min, max should be chosen on the basis of measure
and data: (SA) (in which way)

Average linkage

Centroid linkage



A

1

1.1+1.1 = Single linkage

B

2

1.8+1.8 = Euclidean distance as

C

3

4.3+3.3 = the measure.

D

4

Step 1

E

5

4.5

A

0

1.1+1.1 = 2.2

B

1

0.8+0.8 = 1.6

C

3.6

2.8+0.8 = 3.6

D

5

4.2+4.2 = 8.4

E

6.4

5.6+2.8 = 8.4

$$\Rightarrow \frac{120}{2 \times 6} = \frac{120}{12} = 10$$

$$(x_2 - x_1)^2 = 1 - 2$$

$$A = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad B = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$

$$C = \begin{bmatrix} 4 \\ 3 \end{bmatrix} = 3 = 9 + 4$$

$$D = \begin{bmatrix} 5 \\ 4 \end{bmatrix} \Rightarrow \begin{bmatrix} 3 \\ 2 \end{bmatrix} = 2$$

start of iteration number 21

(first) new partition diagram

Merge A & B

New clusters are (AB) C D E

updated clustering		AB	C	D	E	rank
		0	2.28	4.62	5.6	
(*)	C	2.28	0	1.41	2.82	
*	D	4.62	1.41	0	1.41	
wanted clustering	E	5.6	2.82	1.41	0	

Euclidean dist.

$$C - A = \begin{bmatrix} 2 \\ 3 \end{bmatrix} \quad C - B = \begin{bmatrix} 2 \\ 2 \end{bmatrix}$$

$$= \sqrt{4+9} = \sqrt{13}$$

$$\text{dist AB} = \sqrt{13}, \quad \text{dist BC} = \sqrt{16}$$

$$\text{dist AC} = \sqrt{13 + 16} = \sqrt{29}$$

~~$$D - A = \begin{bmatrix} 3 \\ 4 \end{bmatrix}$$~~

$$D - B = \begin{bmatrix} 3 \\ 3 \end{bmatrix}$$

$$= \sqrt{9+16} = \sqrt{25} = 5$$

$$= \sqrt{18} = \sqrt{18}$$

$$= 5 + 4.24 = 9.24$$

$$= \sqrt{9+9} = \sqrt{18}$$

$$= 4.24 + 4.24 = 8.48$$

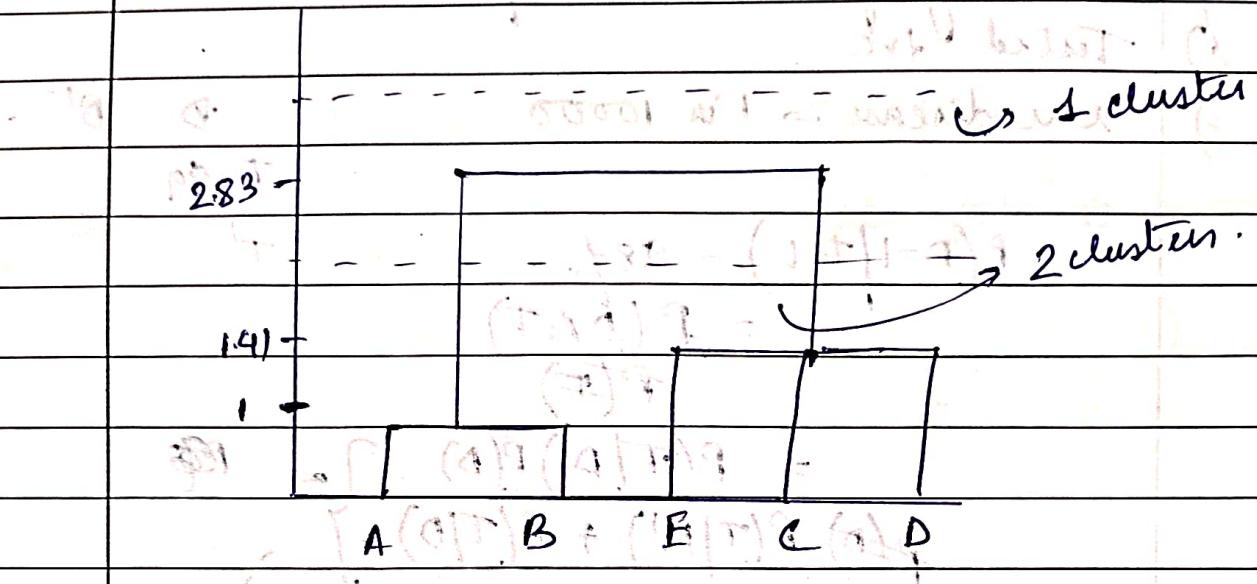
$$= 4.24 + 4.24 = 8.48$$

$$= 4.24 + 4.24 = 8.48$$

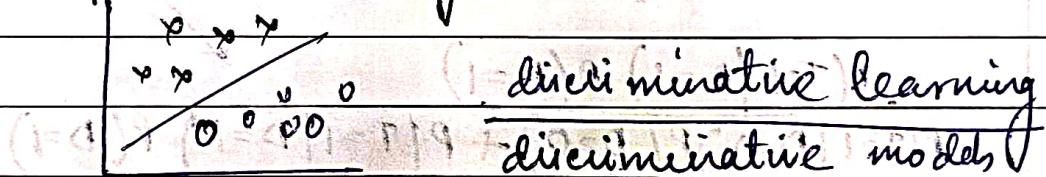
11

Decision boundary - separates a person with

high blood pressure.



Generative learning



capable of generating new data points x_1, x_2, x_3 (not so much into prediction) \rightarrow studying objects of different classes!

Accuracy of the Test = 99%.

Tested true.

rare disease \rightarrow 1 in 10000

D D'

T 99

$$P(D=1 | T=1) = 98\%$$

$$= \frac{P(D \cap T)}{P(T)}$$

$$= P(T|D)P(D)$$

$$\cancel{P(D)} P(T|D') + P(T|D)$$

$$= (0.99)(0.0001)$$

D = 0

= 0.0001

$$P(T=1 | D=1) P(D=1)$$

$$P(T=1 | D=0) P(D=0) + P(T=1 | D=1) P(D=1)$$

$$\Rightarrow 0.0001 (0.99) (0.0001)$$

$$\cancel{0.9999 + 0.0001}$$

$$(0.00009)(0.9999) + (0.99)(0.0001)$$

$$= 0.0098 / 0.09$$

Naive Bayes Classifier

- Spam detection deployed this model in production vector representations have to be there for model

If test is input then $\{1.2, 8.4, \dots\}$

"To by this gradient"

Semantic distance.

between words.

at	0	0	0	1	1	0	0	0	8.2
able	2000	0	0	1	1	0	0	0	10x1
buy		1	1	1	1	0	0	0	clan
:	1	0	0	1	1	0	0	1	spam
more		0	0	1	1	1	0	0	spam
:		0	0	1	1	1	0	0	spam
play	1	0	0	0	0	0	1	0	no spam
5000	500000	1	0	0	0	0	0	0	no spam
spaces				1	1	1	1	1	no spam
lowered				0	1	0	0	0	spam
?				1	0	0	0	0	spam
				0	1	0	0	0	spam
				1	0	0	0	0	spam
				1	1	0	0	0	No spam

~~Preprocessing~~ ~~$2^3 = 8$ combination~~

x_1	x_2	x_3	prob $y=1$ (spam)	prob $y=0$ (no spam)
0	0	0		
0	0	1		
0	1	0	Learning from the data (prior)	
0	1	1	Like how many times 0 0 0 has been a spam in the previous table (probability)	
1	0	1		
1	1	0		

this solution
is not
practical

theta as parameters	2/7	5/7	0
	2/7	5/7	0
	1/7	6/7	0
	2/7	5/7	1/4
	0	1	1/4

→ theta and 1-theta

a new sample has 1 1 0 (it is a spam or not
it is a \leftarrow spam)

$$\text{spam} \log \frac{2/7}{5/7} > \underline{\underline{1/4}}$$

Naive Bayes assumption

x_i 's are independent given y .

↳ Conditional Independent
assumption

- ~~Naive Bayes classifier~~
- 1) Free win now. spam
 - 2) Win a prize. spam
 - 3) Hello How are you? not-spam
 - 4) Let's win it. not-spam
 - 5) free lunch today. not-spam

	word "free"	word "win"	Label
1)	Yes	Yes	spam
2)	No	Yes	spam
3)	No	No	not-spam
4)	No	Yes	not-spam
5)	Yes	No	not-spam

Test message: Free win

Step 1: Calculate prior.

probability of spam | test message

or probability of not-spam | test message.

$$P(y=1|x) = \frac{P(x|y=1) P(y=1)}{P(x|y=1) P(y=1) + P(x|y=0) P(y=0)}$$

$$\text{spam} = 2 \quad P(\text{spam} = 1) = \frac{2}{5} = 0.4$$

$$\text{not spam} = 3 \quad P(\text{spam} = 0) = \frac{3}{5} = 0.6$$

2. Compute likelihood:

$$\text{free} = \text{yes}, \text{mim} = \text{yes}$$

for "spam" class free appears $\frac{1}{2}$ of 2 spam class

$$P(\text{free} = \text{yes} | \text{spam} = 1) = \frac{1}{2}$$

$$P(\text{mim} = \text{yes} | \text{spam} = 1) = \frac{2}{2} = 1$$

For "not spam" class

$$P(\text{free} = \text{yes} | \text{not spam}) = \frac{1}{3}$$

$$P(\text{mim} = \text{yes} | \text{not spam}) = \frac{1}{3}$$

Inference phase

free = yes, mim = yes

$$P(\text{spam} | \mathbf{x}) = P(\mathbf{x} | \text{spam}) \cdot P(\text{spam}).$$

$$P(\mathbf{x} | \text{spam}) = P(x_1, x_2 | \text{spam})$$

$$= P(x_1 | \text{spam}) \cdot P(x_2 | \text{spam}).$$

$$P(\text{spam} | x) = p(x_1 | \text{spam}) p(x_2 | \text{spam}) \cdot p(\text{spam})$$

$$= p(\text{Free} = \text{yes} | \text{spam}) p(\text{win} = \text{yes} | \text{spam}) \\ p(\text{spam})$$

$$= \left(\frac{1}{2}\right)(1)\left(\frac{2}{3}\right) = \frac{2}{10} = \frac{1}{5} = 0.2$$

$$P(\text{not spam} | x) = p(x_1 | \text{not spam}) p(x_2 | \text{not spam}) \\ p(\text{not spam})$$

$$= P(\text{Free} = \text{no} | \text{not spam})$$

$$= P(x_1, x_2 | \text{Not spam})$$

$$= P(x_1 | \text{not spam}), P(x_2 | \text{not spam})$$

$$= P(x_1 | \text{not spam}), P(x_2 | \text{not spam})$$

$$= 1/3 \times 1/3$$

$$P(\text{not spam} | x) = \frac{1}{3} \times \frac{1}{3} = \frac{1}{9}$$

$$= \frac{1}{18} = 0.05555555555555555$$

$$\nabla_{\theta} J(\theta) = 0$$

$$X^T X \theta - X^T y = 0$$

$$X^T y = X^T X \theta$$

$$\theta = (X^T X)^{-1} X^T y$$

Steps first go row wise

1) Compute the predictions (first go row wise)

$$h_{\theta} x^{(1)} = \theta_0 x_0^{(1)} + \theta_1 x_1^{(1)} + \theta_2 x_2^{(1)}$$

$$= 0 \times 1 + 0 \times 1 + 0 \times 2$$

$$= 0 + 0 + 0$$

$$= 0$$

$$h_{\theta} x^{(2)} = 0 \times 1 + 0 \times 2 + 0 \times 1$$

$$= 0 + 0 + 0$$

$$= 0$$

$$h_{\theta} x^{(3)} = 0 \times 1 + 0 \times 3 + 0 \times 1$$

$$= 0$$

2) Compute gradients.

$$e^{(1)} = h_{\theta}(x^{(1)}) - y^{(1)} = 0 - 3 = -3$$

$$e^{(2)} = h_{\theta}(x^{(2)}) - y^{(2)} = 0 - 4 = -4$$

$$e^{(3)} = h_{\theta}(x^{(3)}) - y^{(3)} = 0 - 5 = -5$$

go column wise

3) Calculate gradients.

$$\frac{dJ}{d\theta_0} = \sum_{i=1}^n (h_{\theta} x^{(i)} - y^{(i)}) x_0^{(i)}$$

$$= (-3 \times 1) + (-4 \times 1) + (-5 \times 1)$$

$$= -3 + (-4) + (-5)$$

$$= -3 - 9 = -12$$

$$\frac{\partial J}{\partial \theta_1} = \sum_{i=1}^n (h_\theta x^{(i)} - y^{(i)}) x_1^{(i)}$$

$$= (-3)1 + (-4)2 + (-5)3$$

$$= -3 + (-8) + (-15)$$

$$= -3 + (-23)$$

$$= -26$$

$$\frac{\partial J}{\partial \theta_2} = \sum_{i=1}^n (h_\theta x^{(i)} - y^{(i)}) x_2^{(i)}$$

$$= (-3)2 + (-4)1 + (-5)3$$

$$= -6 + -4 + -15$$

$$= -25$$

3) Update Parameters

$$\theta_0 = \theta_0 - \alpha \frac{\partial J}{\partial \theta_0} \Rightarrow 0 - (0.1)(-12)$$

$$\theta_1 = 0 - \alpha \frac{\partial J}{\partial \theta_1} \Rightarrow 0 - (0.1)(-26)$$

$$\theta_2 = 0 - \alpha \frac{\partial J}{\partial \theta_2} = 0 - (0.1)(-25)$$

To find J

$$J_\theta = \frac{1}{2} \sum_{i=1}^n (h_\theta x^{(i)} - y^{(i)})^2$$

$$= \frac{1}{2} [(-3)^2 + (-4)^2 + (-5)^2]$$

$$= \frac{1}{2} [9 + 16 + 25]$$

$$= 0.5 \cdot \frac{50}{2} = 25$$

1st iteration

$$\theta = \begin{bmatrix} 1.2 \\ 2.6 \\ 2.5 \end{bmatrix}$$

Step 1) Compute the predictions

$$\begin{aligned} h_{\theta}(x^{(1)}) &= \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 \\ &= 1.2 + 1.2 \times 1 + 2.6 \times 2 \\ &= 1.2 + 1.2 + 2.6 \times 1 + 2.5 \times 2 \\ &= 8.8 \end{aligned}$$

$$= 1.2 + 2.6 + 5.0$$

$$= 8.8$$

$$h_{\theta}(x^{(2)}) = 1.2 + (1.2 \times 1) + (2.6 \times 2) + (2.5 \times 1)$$

$$= 1.2 + 5.2 + 2.5$$

$$= 8.9$$

$$h_{\theta}(x^{(3)}) = (1.2 \times 1) + (2.6 \times 3) + (2.5 \times 2)$$

$$= 1.2 + 7.8 + 7.5$$

$$= 16.5$$

Step 2) Compute gradients

$$e^{(1)} = h_{\theta}(x^{(1)}) - y^{(1)} = 8.8 - 3 = 5.8$$

$$e^{(2)} = 8.9 - 4 = 4.9$$

$$e^{(3)} = 16.5 - 5 = 11.5$$

Step 3) Calculate gradients

$$\frac{dJ}{d\theta_0} = \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$= 5.8 \times 1 + 4.9 \times 1 + 11.5 \times 1$$

$$= 22.2$$

$$\frac{dJ}{d\theta_1} = \sum_{i=1}^n (h_\theta x^{(i)} - y^{(i)}) x_1^{(i)}$$

$$= (5.8 \times 1) + (4.9 \times 2) + (11.5 \times 3)$$

$$= 5.8 + 9.8 + 34.5$$

$$= \underline{\underline{50.1}}$$

$$\frac{dJ}{d\theta_2} = \sum_{i=1}^n (h_\theta x^{(i)} - y^{(i)}) x_2^{(i)}$$

$$= (5.8 \times 2) + (4.9 \times 1) + (11.5 \times 3)$$

$$= 11.6 + 4.9 + 34.5$$

$$= \underline{\underline{51}}$$

3) Update Parameters

$$\theta_0 = \theta_0 - \alpha \frac{dJ}{d\theta_0} \Rightarrow \theta_0 - (0.1)(22.2)$$

$$\theta_1 = \theta_1 - \alpha \frac{dJ}{d\theta_1} \Rightarrow 2.6 - (0.1)(50.1)$$

$$\theta_2 = \theta_2 - \alpha \frac{dJ}{d\theta_2} \Rightarrow -2.41$$

$$\theta_2 = \theta_2 - \alpha \frac{dJ}{d\theta_2} \Rightarrow 2.5 - (0.1)(51)$$

$$\Rightarrow -2.6$$

To find J

$$J_\theta = \frac{1}{2} \sum_{i=1}^n (h_\theta x^{(i)} - y^{(i)})^2$$

$$= \frac{1}{2} [(5.8)^2 + (4.9)^2 + (11.5)^2]$$

$$= \frac{1}{2} \times 22.2 \Rightarrow \underline{\underline{11.1}}$$

K-fold cross validation

$x_1 \mid x_2 \mid y$

1	-1	1
2	1.2	0
3	2	1
4	-2	1
5	0.1	0

3-fold CV

Acc

std

$$F_1 \{0_1, 0_2\} = \{1.8, 2.8\}$$

$$F_2 \{0_1, 0_2\} = \{2.1, 3.1\}$$

$$F_3 \{0_1, 0_2\} = \{19, 4\}$$

put dates into 3 splits

fold 1 100%

	x_1	x_2	y
1	-1	1	1
2	1.2	0	0

$$\text{Train set} = \{1, 2, 3, 4\}$$

$$\text{Test set} = \{5\}$$

$$\text{Acc. Test set} = \{5\}$$

fold 2

50%

	x_1	x_2	y
1	2	1	1
2	-1	1	1

$$\text{Train set} = \{3, 4, 5\}$$

$$\text{Test set} = \{1, 2\}$$

Acc.

fold 3

50%

	x_1	x_2	y
4	0.1	0	0

$$\text{Train set} = \{4\}$$

$$\text{Test set} = \{1, 2, 3\}$$

$$\text{Train set} = \{1, 2, 5\}$$

$$\text{Test set} = \{3, 4\}$$

$$\begin{aligned}
 h_0 x &= \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 \\
 &= 0 \cdot 1 + 0 \cdot 2 + 0 \cdot (-1) \quad (\theta_0 = 0) \\
 &= 0 + 0 + 0 \\
 &= 0
 \end{aligned}$$

$$\begin{aligned}
 h_{00} x &= \theta_1 x_1 + \theta_2 x_2 \\
 &= -1.8 \cdot 2 + 2.8 \cdot (-1) \\
 &= -3.6 + (-2.8) \\
 &= -3.6 - 2.8 \Rightarrow \frac{6.4}{0.6} \Rightarrow \frac{1}{1 + e^{6.4}}
 \end{aligned}$$

$$\begin{aligned}
 h_{01} x &= \theta_1 x_1 + \theta_2 x_2 \\
 &= -1.8 \cdot (-5) + (2.8 \cdot 1.2)
 \end{aligned}$$

$$\begin{aligned}
 h_{02} x &= \theta_1 x_1 + \theta_2 x_2 \\
 &= -1.8 \cdot 1 + (2.8 \cdot 2) = \frac{10}{1 + 4.2} \quad (e^{-12.36} = 4.2)
 \end{aligned}$$

$$\begin{aligned}
 &= -1.8 + (5.6) = \frac{1}{5.2} = 0.052
 \end{aligned}$$

$$\begin{aligned}
 &\Rightarrow 8.8 = \frac{1}{1 + 44.701} = \frac{1}{45.701} = 0.021
 \end{aligned}$$

$$h_0 x = \theta_1 x_1 + \theta_2 x_2$$

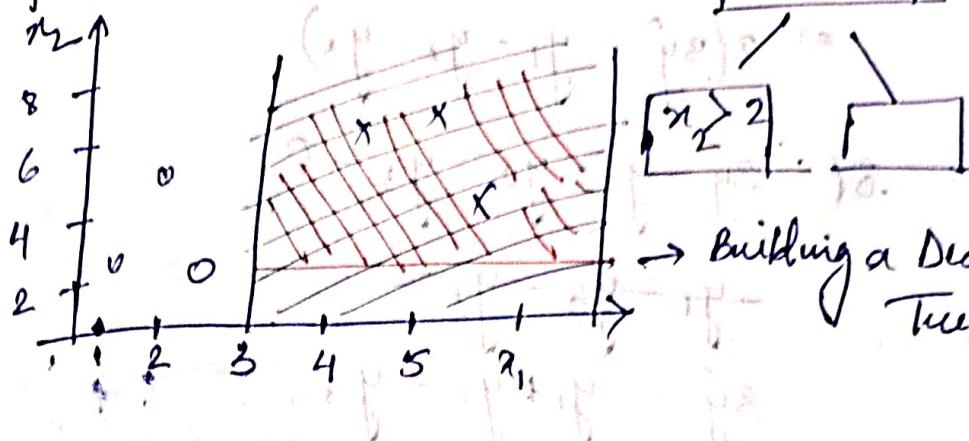
$$= (-1.8 \cdot -3) + (2.8 \cdot 2) \Rightarrow 5.4 + 5.6$$

$$\Rightarrow 11 = \frac{1}{1 + 1.22} \Rightarrow \frac{1}{2.22} = 0.022$$

$$= 2 \cdot e^{-(-0.2)} = 1.22$$

$$\Rightarrow \frac{1}{1 + 1.22} \Rightarrow \frac{1}{2.22} = 0.022$$

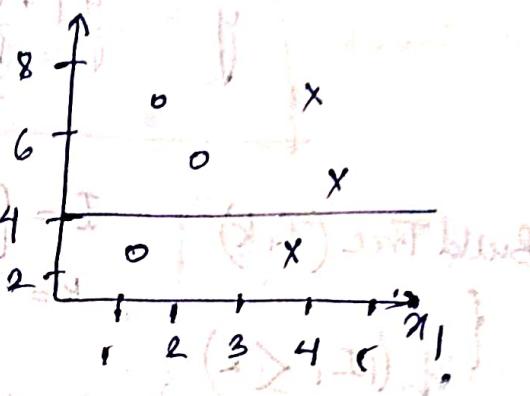
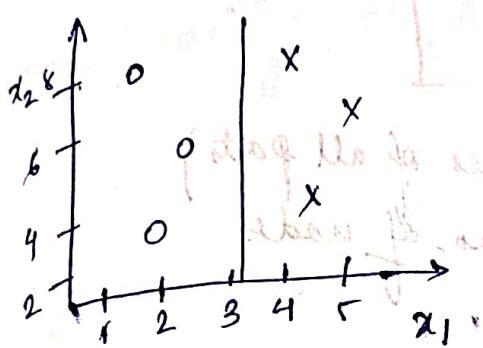
Partition Space



$x_1 \geq 3$ | 6 points

$x_2 \geq 2$

→ Building a Decision Tree.



$x_1 \geq 3$

$x_2 \geq 2$

This is better because it's splitting into homogeneous set of cases.

$$y = \begin{bmatrix} -2 \\ 3 \\ 4 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

$$f = \sum_{i=1}^n (\hat{y} - y^{(i)})^2$$

$$\hat{y} \cdot \alpha f = (\hat{y} - y_1)^2 + (\hat{y} - y_2)^2 + (\hat{y} - y_3)^2$$

$$\min_{\hat{y}} \frac{df}{d\hat{y}} = 0$$

$$\frac{df}{d\hat{y}} = \frac{d(\hat{y} - y_1)^2}{d\hat{y}} + \frac{d(\hat{y} - y_2)^2}{d\hat{y}} + \frac{d(\hat{y} - y_3)^2}{d\hat{y}}$$

$$= 2(\hat{y} - y_1) + 2(\hat{y} - y_2) + 2(\hat{y} - y_3)$$

$$= 2(\hat{y} - 2y_1 - 2y_2 + 2\hat{y} - 2y_3 + 2\hat{y} - 2y_3)$$

$$\Rightarrow 6\hat{y} - 2y_1 - 2y_2 - 2y_3$$

$$\Rightarrow 2(3\hat{y} - y_1 - y_2 - y_3)$$

$$B6: 3\hat{y} - y_1 - y_2 - y_3 = 0$$

$$\cancel{-y_1 \quad y_2 \quad y_3}$$

$$3\hat{y} = y_1 + y_2 + y_3$$

$$\hat{y} = \frac{y_1 + y_2 + y_3}{3}$$

Build Tree (I, k) $I = \{\text{indices of all parts}\}$

$k = \min \text{ no. of nodes}$

If $|I| \leq k$

Set $\hat{y} = \text{Average } (y^{(i)})$

return leaf (label = \hat{y})

$$2.6 \times 10^{-3}$$

else
for each split-dim $j \rightarrow \text{feature}$

for each split value $S \rightarrow \text{value}$

// inner loop

$$I_{j,s}^+ = \{i \in I \mid x_j^{(i)} \geq S\}$$

$$I_{j,s}^- = \{i \in I \mid x_j^{(i)} < S\}$$

$\hat{y} = \text{average of } (y^{(i)})$

$$\begin{cases} \hat{y}^+ & i \in I_{j,s}^+ \\ \hat{y}^- & i \in I_{j,s}^- \end{cases}$$

Based on the
chosen feature
and value
we will do
partition.

mean on
both sides

(the chance of overfitting is less.)

$$C_{\alpha}(T) = \sum_{i=1}^n L(T(x^{(i)}, y^{(i)})) \rightarrow \text{Build tree procedure}$$

Pruning of tree
→ removing leaves.

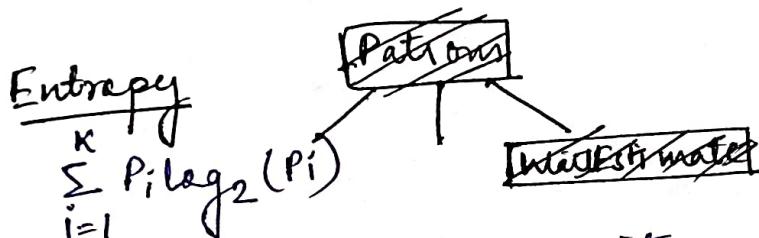
$\boxed{x|T|}$ + leaves need not match with the training data

Regularization cont. → hyperparameter

Classification Tree

CART

Leaves are always be the target values
 \downarrow
 Labels



To measure homogeneity we can measure the entropy and ensure if the tree is homogeneous or not.

$$p=6, n=6 \quad H\left(\frac{p}{p+n}, \frac{n}{p+n}\right)$$

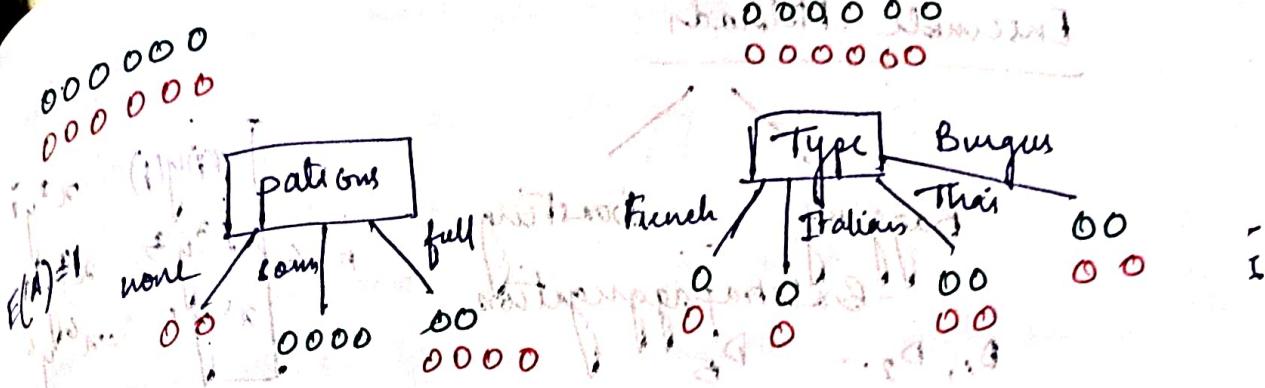
$$= -\frac{p}{p+n} \log_2 \frac{p}{p+n} - \frac{n}{p+n} \log_2 \frac{n}{p+n}$$

$$= -\sum_{i=1}^k p_i \log_2 p_i$$

$$H\left(\frac{6}{12}, \frac{6}{12}\right) = -\frac{6}{12} \log_2 \frac{6}{12} - \frac{6}{12} \log_2 \frac{6}{12}$$

$$EH(\text{pattern}) = \sum_{i=1}^3 -\frac{p_{ini}}{p+n} H\left(\frac{p_i}{p+n_i}, \frac{n_i}{p+n_i}\right)$$





		Type				Burgers
		French		Italian		Thai
		00	00	00	00	00
00	00	00	00	00	00	00
00	00	00	00	00	00	00

$$E(H) = \sum_{i=1}^K -\left(\frac{P_i + n_i}{P+n}\right) H\left(\frac{P_i}{P+n}, \frac{n_i}{P+n}\right)$$

$$\text{Information gain} = H\left(\frac{P}{P+n}, \frac{n}{P+n}\right) - E(H)$$

To maximize $I(A)$ lowers $E(H)$

$$\textcircled{4} \quad P_i, n_i \rightarrow \frac{2}{12} H(0,1) + \frac{4}{12} H(1,0) + \frac{6}{12} H\left(\frac{2}{6}, \frac{4}{6}\right)$$

$$\left[\frac{2}{12} H(0,1) + \frac{4}{12} H(1,0) + \frac{6}{12} H\left(\frac{2}{6}, \frac{4}{6}\right) \right]$$

$$\frac{2}{12}(-0 \log_2 0 - 1 \log_2 1) + \frac{4}{12}(1 \log_2 1 - 0 \log_2 0)$$

$$\frac{2}{12} \left(\frac{2}{6} \log_2 \frac{2}{6} - \frac{4}{6} \log_2 \frac{4}{6} \right)$$

$$= 0.455 \rightarrow \text{Initial entropy} = 1$$

$$I_G(\text{Patrons}) = H(\text{Patrons}) - E(H(\text{Patrons}))$$

$$= 1 - 0.455$$

$$I_G(\text{Patron}) = 0.541$$

Is this justified? → Yes

Is this justified? → Yes

Kernel

$$K(x, y) = \langle \phi(x), \phi(y) \rangle$$

transformed into higher dimensional X
and y respectively.

~~dot prod~~ $\langle \cdot \rangle$ symbol for multiplication.

$$x = (x_1, x_2, x_3)$$

$$y = (y_1, y_2, y_3)$$

is a vector
not a feature

$$\phi(x) = \begin{bmatrix} x_1 & x_1 \\ x_1 & x_2 \\ x_1 & x_3 \\ x_2 & x_1 \\ x_2 & x_2 \\ x_2 & x_3 \\ x_3 & x_1 \\ x_3 & x_2 \\ x_3 & x_3 \end{bmatrix}$$

$$x = [1, 2, 3]$$

$$y = [4, 5, 6]$$

$$\langle \phi x, \phi y \rangle$$

$$\phi x = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 2 \\ 4 \\ 3 \end{bmatrix} \quad \phi y = \begin{bmatrix} 16 \\ 20 \\ 24 \\ 20 \\ 25 \\ 30 \end{bmatrix}$$

$$\phi x \cdot \phi y = 16 + 40 + 72 + 40 + 100 + 180 + 72 + 180 + 324 \\ = 1024$$

$$K(x, y) = \langle x, y \rangle^2 \quad (\text{squaring on original data})$$

$$= \cancel{(1024)^2} = \cancel{1048576} \quad 32 \times 32 = \cancel{1024}$$

gradient descent on higher dimensional space
is computationally expensive hence Kernel
helps to compute on a lower dimensional space.