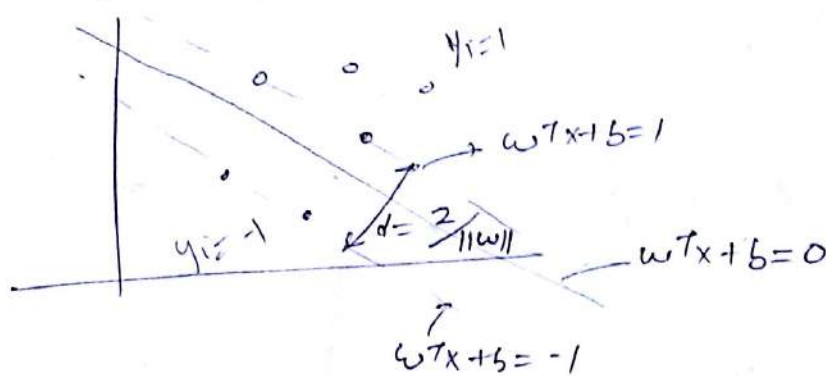


Q1

Let x_1, x_2, \dots, x_n be our data points
 & $y_i \in \{1, -1\}$ be the class label of x_i



$$y_i (w^T x_i + b) > 1 \quad \forall (x_i, y_i)$$

For $y_i \in \{1, -1\}$. Now decision boundary can be found as
 minimize $\frac{1}{2} \|w\|^2$
 subject to $y_i (w^T x_i + b) > 1$

So Lagrangian is

$$L = \frac{1}{2} w^T w + \sum_{i=1}^n \alpha_i (1 - y_i (w^T x_i + b))$$

$$\alpha_i \geq 0$$

Differentiating wrt. w & b

$$\frac{\partial L}{\partial w} = 0 \Rightarrow w - \sum_{i=1}^n \alpha_i y_i x_i = 0$$

$$\Rightarrow \bar{w} = \sum_{i=1}^n \alpha_i y_i x_i$$

$$\frac{\partial L}{\partial b} = 0 \Rightarrow \sum \alpha_i y_i = 0$$

So if we substitute $\bar{w} = \sum \alpha_i y_i x_i$ to L

we get

$$L = \frac{1}{2} \sum_{i=1}^n \alpha_i y_i x_i^T \sum_{j=1}^n \alpha_j y_j x_j + \sum \alpha_i (1 - y_i (\sum_{j=1}^n \alpha_j y_j x_j^T x_i + b))$$

$$= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j + \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \alpha_i y_i (\sum_{j=1}^n \alpha_j y_j x_j^T x_i + b)$$

$$= -\frac{1}{2} \sum_{i=1, j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j + \sum_{i=1}^n \alpha_i$$

$$\text{Since } \sum_{i=1}^n \alpha_i y_i = 0$$

So the dual problem is -

$$\begin{aligned} \max \quad w(\alpha) &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1, j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j \\ \text{Subject to } \alpha_i &\geq 0, \quad \sum_{i=1}^n \alpha_i y_i = 0 \end{aligned}$$

SVM solver for XOR data set,

Class 1: $x_1 = (0, 0)$, $x_4 = (1, 1)$ $\rightarrow y = -1$

Class 2: $x_2 = (1, 0)$, $x_3 = (0, 1)$ $\rightarrow y = 1$

Kernel function:—
polynomial of order 2
$$K(x, x') = (x^T x' + 1)^2$$

$$= (u_1 v_1 + u_2 v_2 + 1)^2$$

$$= 1 + 2u_2 v_2 + 2u_1 v_1 +$$

$\phi(x) = \phi([x_1, x_2])$
can be
$$= [1, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1x_2, x_1^2, x_2^2]$$

$\therefore \phi(x_1) = \phi([0, 0])$

$$= [1, 0, 0, 0, 0, 0]^T$$

$\phi(x_2) = \phi([1, 0]) = [1, \sqrt{2}, 0, 0, 1, 0]^T$

$\phi(x_3) = \phi([0, 1]) = [1, 0, \sqrt{2}, 0, 0, 1]^T$

$\phi(x_4) = \phi([1, 1]) = [1, \sqrt{2}, \sqrt{2}, \sqrt{2}, 1, 1]^T$

objective function is,

$$L_D(\alpha) = \alpha_1 + \alpha_2 + \alpha_3 + \alpha_4 - \frac{1}{2} \sum_{i=1}^4 \sum_{j=1}^4 \alpha_i \alpha_j y_i y_j K_{ij}$$

$$\text{s.t. } \sum \alpha_i y_i = 0 \Rightarrow \alpha_1 + \alpha_4 = \alpha_2 + \alpha_3$$

$$\alpha_i \geq 0$$

inner product is represented as

$$K = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 4 & 1 & 4 \\ 1 & 1 & 4 & 4 \\ 1 & 4 & 4 & 9 \end{bmatrix}$$

Now optimizing w.r.t. Lagrange multiplier leads to,

$$\frac{\partial L}{\partial \alpha_i} = 0 \Rightarrow$$

$$L(\alpha) = \sum \alpha_i - \frac{1}{2} (\alpha_1^2 + 4\alpha_2^2 + 4\alpha_3^2 + 9\alpha_4^2 - 2\alpha_1\alpha_2 - 2\alpha_1\alpha_3 + 2\alpha_1\alpha_4 + 2\alpha_2\alpha_3 - 8\alpha_2\alpha_4 - 8\alpha_3\alpha_4)$$

$$\frac{\partial L}{\partial \alpha_1} = 0 \Rightarrow \alpha_1 - \alpha_2 - \alpha_3 + \alpha_4 = 1$$

$$\frac{\partial L}{\partial \alpha_2} = 0 \Rightarrow 4\alpha_2 - \alpha_1 + \alpha_3 - 4\alpha_4 = 1$$

$$\frac{\partial L}{\partial \alpha_3} = 0 \Rightarrow 4\alpha_3 - \alpha_1 + \alpha_2 - 4\alpha_4 = 1$$

$$\frac{\partial L}{\partial \alpha_4} = 0 \Rightarrow 9\alpha_4 + \alpha_1 - 4\alpha_2 - 4\alpha_3 = 1$$

$$\frac{\partial L}{\partial \alpha_4} = 0 \Rightarrow$$

Solving these equations,

$$\alpha_1 = 13/3, \alpha_2 = 8/3, \alpha_3 = 8/3$$

$$\alpha_4 = 2.$$

Thus all the data points are support vectors in this case.

Q.E.D.

$$w = \sum_{i=1}^4 \alpha_i y_i \phi(x_i)$$

$$= \alpha_1 (-1) [1 \ 0 \ 0 \ 0 \ 0 \ 0]^T + \alpha_2 (1) [1 \ \sqrt{2} \ 0 \ 0 \ 1 \ 0]^T \\ + \alpha_3 (1) [1 \ 0 \ \sqrt{2} \ 0 \ 0 \ 1]^T + \alpha_4 (-1) [1 \ \sqrt{2} \ \sqrt{2} \ \sqrt{2} \ 1 \ 1]^T$$

$$= \begin{bmatrix} -1/3 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 8/3 \\ 8\sqrt{2}/3 \\ 0 \\ 0 \\ 8/3 \\ 0 \end{bmatrix} + \begin{bmatrix} 8/3 \\ 0 \\ 8\sqrt{2}/3 \\ 0 \\ 0 \\ 8/3 \end{bmatrix} + \begin{bmatrix} -2 \\ -2\sqrt{2} \\ -2\sqrt{2} \\ -2\sqrt{2} \\ -2 \\ -2 \end{bmatrix}$$

$$= \begin{bmatrix} -1 \\ 2\sqrt{2}/3 \\ 2\sqrt{2}/3 \\ -2\sqrt{2} \\ 2/3 \\ 2/3 \end{bmatrix}$$

So $g([x_1, x_2]) = w^T \phi(x)$

$$= w^T \begin{bmatrix} 1 & \sqrt{2} x_1 & \sqrt{2} x_2 & \sqrt{2} x_1 x_2 & x_1^2 & x_2^2 \end{bmatrix}^T$$

$$= \begin{bmatrix} -1 & 2\sqrt{2}/3 & 2\sqrt{2}/3 & -2\sqrt{2} & 2/3 & 2/3 \end{bmatrix} \begin{bmatrix} 1 \\ \sqrt{2} x_1 \\ \sqrt{2} x_2 \\ \sqrt{2} x_1 x_2 \\ x_1^2 \\ x_2^2 \end{bmatrix}$$

$$= -1 + 4/3 x_1 + 4/3 x_2 - 4 x_1 x_2 + 2/3 x_1^2 + 2/3 x_2^2$$

is the decision boundary,

(Q2) Bayesian Classifier: is a simple probabilistic classifier based on applying ~~bay~~ bayes' theorem. It is based on the idea that the role of class is to predict the values of features for the members of that class. If an agent know the class it can predict the values of other features. If it doesn't know Bayes rule can be applied to predict the class given some features. so if (X, Y) takes values in $\mathbb{R}^d \times \{1, 2, \dots, k\}$ Y is class label of X .
 $X|Y = r \sim P_r$ for $r = 1, \dots, k$
 \sim means "is distributed as" P_r is prob. distribution

A Classifier assigns to an observation $X = x$ a estimate of what unlabeled label $Y = r$ actually uses.

$$C^{\text{Bayes}}(x) = \underset{r \in \{1, 2, \dots, k\}}{\operatorname{argmax}} p(Y=r | X=x)$$

$$P(C | F_1, F_2, \dots, F_n) = \frac{P(C) P(F_1, \dots, F_n | C)}{P(F_1, \dots, F_n)}$$

$$\begin{aligned} &\propto P(C) P(F_1, \dots, F_n | C) \\ &= P(C) P(F_1 | C) P(F_2 | C, F_1) \\ &\quad P(F_3 | C, F_1, F_2) \dots \\ &\quad P(F_n | C, F_1, F_2, \dots, F_{n-1}) \end{aligned}$$

Q3

According to Fisher's Linear Discriminant methods-

If we have,

$$\bar{m}_1 = \frac{1}{N_1} \sum_{n \in C_1} \bar{x}_n \quad \text{and} \quad \bar{m}_2 = \frac{1}{N_2} \sum_{n \in C_2} \bar{x}_n$$

Simplest measure of separation of classes comes out to be when project onto w , is the separation of projected class means.

so $m_2 - m_1 = w^T (\bar{m}_2 - \bar{m}_1)$, $m_k = w^T \bar{m}_k$

and $w \propto \bar{m}_2 - \bar{m}_1$ (using Lagrange multiplier)

and within class variance is given by

$$s_k^2 = \sum_{n \in C_k} (y_n - m_k)^2$$

So Fisher criteria, $J(w) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2}$

$$J(w) = \frac{w^T S_B w}{w^T S_W w}$$

$$S_B = (m_2 - m_1)(m_2 - m_1)^T$$

$$S_W = \sum_{n \in C_1} (\bar{x}_n - \bar{m}_1)(\bar{x}_n - \bar{m}_1)^T + \sum_{n \in C_2} (\bar{x}_n - \bar{m}_2)(\bar{x}_n - \bar{m}_2)^T$$

$$\therefore \boxed{w \propto S_W^{-1} (\bar{m}_2 - \bar{m}_1)} \quad \text{--- (1)}$$

Now we show that if we use least square method using $t = x_1/x_2$ for class C_1
 $\hookrightarrow t = \frac{x_1}{x_2}$ for

So sum of squares error is given by

$$E = \frac{1}{2} \sum_{n=1}^N (\omega^T x_n + \omega_0 - t_n)^2$$

$$\frac{\partial E}{\partial \omega_0} = 0$$

\Rightarrow

$$\sum_{n=1}^N (\omega^T x_n + \omega_0 - t_n) = 0$$

$$\Rightarrow N \omega_0 = \cancel{\omega^T x_n} - \omega^T \cancel{x_n} \sum_{n=1}^N \cancel{x_n}$$

$$(\because \sum_{n=1}^N t_n = 0)$$

$$\Rightarrow \boxed{\omega_0 = -\omega^T \bar{m}}$$

$$= \frac{N_1}{N} \times N_1 - \frac{N_2}{N} \times N_2 = 0$$

$$\bar{m} = \frac{1}{N} \sum_{n=1}^N x_n = \frac{1}{N} (N_1 \bar{m}_1 + N_2 \bar{m}_2)$$

$$\text{Now } \frac{\partial E}{\partial \omega} = 0 \Rightarrow \sum_{n=1}^N (\omega^T x_n + \omega_0 - t_n) x_n = 0$$

Now by substituting $t_n = \frac{N_1}{N}$ for $n \in I_1$
 $t_n = \frac{N_2}{N}$ for $n \in I_2$

we get

$$\left(S \omega + \frac{N_1 N_2}{N} S \beta \right) \omega = N (m_1 - m_2)$$

Now $S \omega \propto m_1 - m_2$

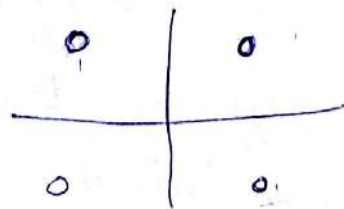
$\therefore \cancel{S \omega} \propto \cancel{S \omega}^{-1} (\bar{m}_2 - \bar{m}_1)$

$$S \omega = K (\bar{m}_1 - \bar{m}_2) \quad K \text{ is some constant}$$

$$\omega \propto S^{-1} (\bar{m}_1 - \bar{m}_2)$$

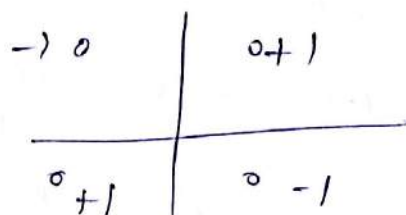
which is same as $\rightarrow (i)$

(05)



(a) Linear kernel

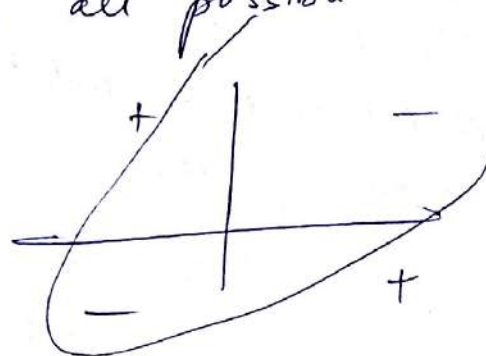
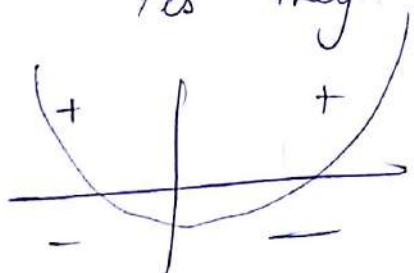
no it can't be shattered, since vc dim of the data given is 3. Clearly we can't find any linear boundary for say such 2 classes



Can't be shattered.

(b) Polynomial kernel of degree 2.

Yes they can shatter all possible labels



(c) Gaussian kernel

Now again a gaussian kernel can also shatter all possible labels.

