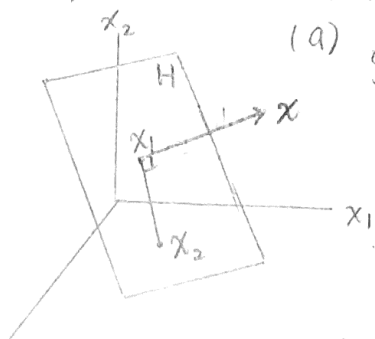


1.



$$(a) \quad g(x) = w_0 + \underline{w}^T \underline{x}$$

Assume x_1 and x_2 are the points on H

$$g(x_1) = w_0 + \underline{w}^T \underline{x}_1 = 0$$

$$g(x_2) = w_0 + \underline{w}^T \underline{x}_2 = 0$$

$$g(x_1) - g(x_2) = \underline{w}^T (\underline{x}_1 - \underline{x}_2) = 0 \Rightarrow \underline{w}^T \cdot (\underline{x}_1 - \underline{x}_2) = 0$$

$\Rightarrow \underline{w}$ is normal to H .

$(\underline{x}_1 - \underline{x}_2)$ is a vector on H .

(b) Assume there is a point x on $g(x) > 0$ side, and $x = x_1 + a\underline{w}$ ($a > 0$)
the vector from x_1 to x is $(x - x_1) = a\underline{w}$

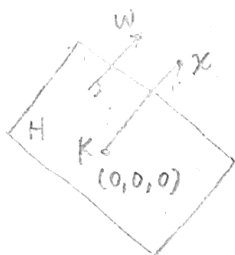
$$(w, (x - x_1)) = \underline{w}^T \cdot (x - x_1) = \underline{w}^T \cdot a\underline{w} = a \|\underline{w}\| \|\underline{w}\| \cos 0 = a \underline{w}^2 > 0$$

\therefore the inner product of vector w and $(x - x_1)$ is bigger than 0

\Rightarrow They point on the same side of H

$\Rightarrow \underline{w}$ points to the positive side of H

(c)



\therefore in "augmented feature space" $\Rightarrow H$ must pass through the origin $K(0,0,0)$

$$r = \frac{g(x)}{\|\underline{w}\|} = \frac{\underline{w}^T \underline{x}}{\|\underline{w}\|}, \text{ Assume there's a point } x \text{ on the positive side of } H.$$

$$\underline{x} = a\underline{w} \text{ take into } K, \quad r = \frac{\underline{w}^T \underline{x}}{\|\underline{w}\|} = \frac{\underline{w}^T (a\underline{w})}{\|\underline{w}\|} = \frac{a \|\underline{w}\|^2 \cos 0}{\|\underline{w}\|} = a \|\underline{w}\| > 0$$

$$\underline{x} = a\underline{w}$$

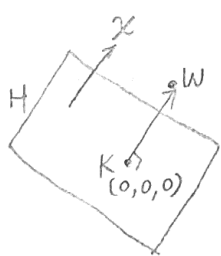
$\therefore x$ is on the positive side of $H(x)$

when $a > 0$, \underline{w} and \underline{x} point to $r > 0$ side

$a < 0$, $a\underline{w}$ points to $r < 0$ side

$\} \underline{w}$ points to the $r > 0$ side of H .

(d)



\therefore "augmented weight space" $\Rightarrow H$ must pass through the origin $K(0,0,0)$

$$r = \frac{g(w)}{\|\underline{x}\|} = \frac{\underline{w}^T \underline{x}}{\|\underline{x}\|} \text{ Assume there's a point } w \text{ on the positive side of } H$$

$$\underline{w} = a\underline{x} \text{ take into } K, \quad r = \frac{\underline{w}^T \underline{x}}{\|\underline{x}\|} = \frac{\underline{w}^T (a\underline{x})}{\|\underline{x}\|} = \frac{a \|\underline{x}\|^2 \cos 0}{\|\underline{x}\|} = a \|\underline{x}\| > 0$$

$$\underline{w} = a\underline{x}$$

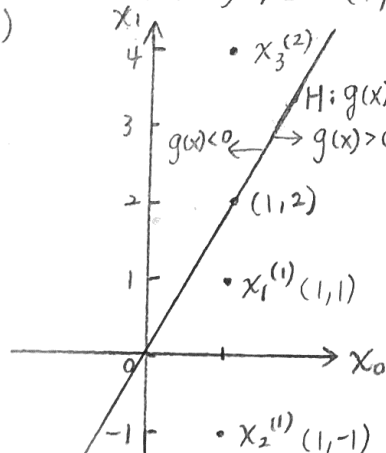
$\therefore w$ is on the positive side of H

when $a > 0$, \underline{w} and \underline{x} point to the positive side of H ($r > 0$)

$a < 0$, \underline{w} points to the negative side of H . ($r < 0$)

$$2. \quad x_1^{(1)} = (1, 1), \quad x_2^{(1)} = (1, -1), \quad x_3^{(2)} = (1, 4)$$

(a)



$$\text{Mean of class 1: } (1, \frac{1+(-1)}{2}) = (1, 0)$$

$$\text{Mean of class 2: } (1, 4)$$

The decision boundary need to pass through $(0,0)$ and the center of mean 1 and mean 2.

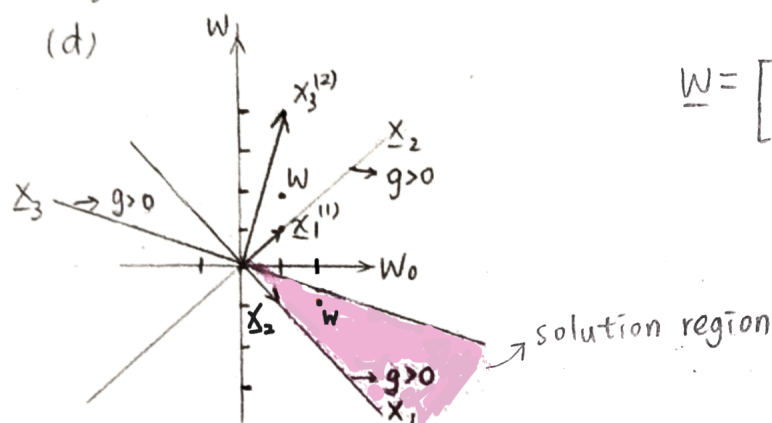
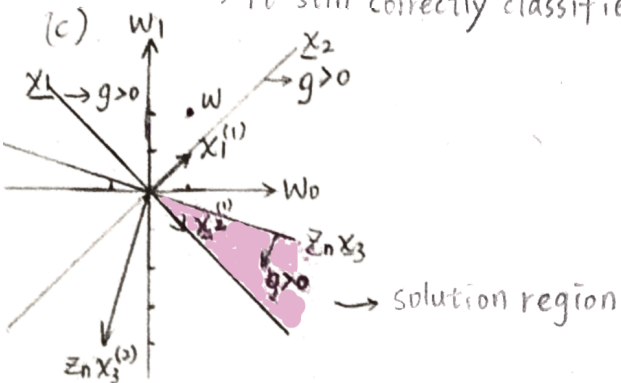
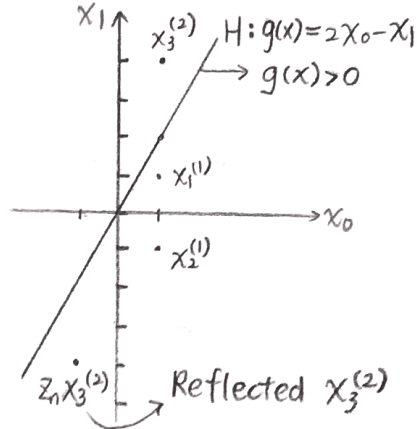
The side of class 1 is positive.

(b) Original: class 1 $\Rightarrow g(x) > 0$
 class 2 $\Rightarrow g(x) < 0$

Reflected: $z_n^{(k)} \begin{cases} 1, & k=1 \\ -1, & k=2 \end{cases}$

$$\underline{w}^T z_n^{(k)} \underline{x}^{(k)} > 0$$

\Rightarrow it still correctly classifies the points.



$\underline{w} = \begin{bmatrix} 2 \\ -1 \end{bmatrix} \Rightarrow \underline{w}$ is in the solution region.

3. (a) $f(t) = \begin{bmatrix} x_1(t) \\ x_2(t) \\ \vdots \\ x_D(t) \end{bmatrix} \Rightarrow f(p) = \begin{bmatrix} x_1(p) \\ x_2(p) \\ \vdots \\ x_D(p) \end{bmatrix} \quad p(x) = \begin{bmatrix} p_1(x) \\ p_2(x) \\ \vdots \\ p_D(x) \end{bmatrix}$

$$\nabla_x f[p(x)] = \begin{bmatrix} \nabla_x p(x) \cdot \frac{d}{dp} x_1(p) \\ \nabla_x p(x) \cdot \frac{d}{dp} x_2(p) \\ \vdots \\ \nabla_x p(x) \cdot \frac{d}{dp} x_D(p) \end{bmatrix} = \begin{bmatrix} \frac{d}{dp} x_1(p) \\ \frac{d}{dp} x_2(p) \\ \vdots \\ \frac{d}{dp} x_D(p) \end{bmatrix} \quad \nabla_x p(x) = \frac{d}{dp} \begin{bmatrix} x_1(p) \\ x_2(p) \\ \vdots \\ x_D(p) \end{bmatrix} \cdot \nabla_x p(x) = \left(\frac{d}{dp} f(p) \right) \nabla_x p(x) \quad \text{Q.E.D.}$$

(b) $\frac{\partial}{\partial x} [x^T M x] = (M + M^T) x$

$\nabla_x (x^T x) = \nabla_x (x^T I x) = (I + I^T) x = 2I x = 2x$

(c) $\underline{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix} \quad \nabla_x (x^T x) = \nabla_x \left([x_1 \ x_2 \ x_3 \ \dots \ x_D] \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix} \right) = \begin{bmatrix} \frac{\partial (x_1^2 + x_2^2 + \dots + x_D^2)}{\partial x_1} \\ \frac{\partial (x_1^2 + x_2^2 + \dots + x_D^2)}{\partial x_2} \\ \vdots \\ \frac{\partial (x_1^2 + x_2^2 + \dots + x_D^2)}{\partial x_D} \end{bmatrix} = \begin{bmatrix} 2x_1 \\ 2x_2 \\ \vdots \\ 2x_D \end{bmatrix} = 2x$

(d) $f(t) = t^3$, $p = x^T x \Rightarrow \nabla_x p(x) = 2x$

$$\nabla_x [(x^T x)^3] = \nabla_x f[p(x)] = \frac{df(p)}{dp} \cdot \nabla_x p(x) = 3p^2 \cdot 2x = 3(x^T x)^2 \cdot 2x = 3 \|x\|_2^4 \cdot 2x = 6 \|x\|_2^4 \cdot x$$

4. (a) $w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_D \end{bmatrix}$ $\|w\|_2 = \sqrt{w_1^2 + w_2^2 + w_3^2 + \dots + w_D^2}$
 let $p = w^T w = w^2$, $f(t) = t^{\frac{1}{2}}$

$$\nabla_w \|w\|_2 = \nabla_w f[p(w)] = \left(\frac{df(p)}{dp} \right) \nabla_w p(w) = \frac{1}{2} p^{-\frac{1}{2}} \cdot 2w = \frac{1}{2} \cdot (\|w\|_2^2)^{-\frac{1}{2}} \cdot 2w = \frac{w}{\|w\|_2}$$

(b) $p(w) = (Mw - b)^T (Mw - b) = w^T M^T M w - w^T M^T b - b^T M w + b^T b$

$$\nabla_w p(w) = 2M^T M w - M^T b - M^T b = 2M^T M w - 2M^T b$$

$$\begin{aligned} \nabla_w \|Mw - b\| &= \frac{1}{2} [(Mw - b)^T (Mw - b)]^{-\frac{1}{2}} \cdot (2M^T M w - 2M^T b) \\ &= \frac{M^T M w - M^T b}{\|Mw - b\|_2} \end{aligned}$$

5. a) Show that total linear separability implies linear separability.

For points which are totally linear separable, for each class n there exists a hyperplane

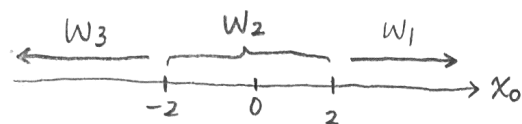
$H: g(x) = 0$ that separates the samples belonging to class n from the rest.

$$g(x) > 0 \text{ for } x \in W_n$$

$$g(x) < 0 \text{ for } x \notin W_n$$

b) Show that linear separability doesn't necessarily imply total linear separable.

Take a counterexample to illustrate:



$$\begin{cases} W_1: x_0 > 2 \\ W_2: -2 < x_0 < 2 \\ W_3: x_0 < -2 \end{cases}$$

W_1, W_2 and W_3 are linear separable, but they are

not totally linear separable because there's no discriminant function $g(x)$ that can separate W_2 samples from the rest.