王井飞 2215335

Hw4

1.(a) $\{y_i(x_i \cdot w + \psi) \geq 0, \forall i=1\cdots n\}$    as ~~the~~ a set of $n$ constraints
$f(x_i)$
     to ensure that all $n$ points are correctly classified.

(b) We have known $H = \{x \in R^d : x \cdot w + \psi = 0\}$ and $x_i \in R^d$

$\therefore$ we assume another two points : $\begin{cases} V_1 \text{ is the point which is the projection of } x_i \text{ on } H \\ V_2 \text{ is any point on the } H \ (V_2 \neq V_1) \end{cases}$

And then the angle between $\overrightarrow{x_i V_1}$ and $\overrightarrow{x_i V_2}$ is $\theta$

$\Rightarrow d = \|\overrightarrow{x_i V_1}\| = \|\overrightarrow{x_i V_2}\| \cdot \cos\theta$   (1)   and   $\cos\theta = \dfrac{\overrightarrow{x_i V_1} \cdot \overrightarrow{x_i V_2}}{\|\overrightarrow{x_i V_1}\| \|\overrightarrow{x_i V_2}\|}$   (2)

$\overset{(1)(2)}{\therefore\Rightarrow} d = \dfrac{\overrightarrow{x_i V_1} \cdot \overrightarrow{x_i V_2}}{\|\overrightarrow{x_i V_1}\|}$    (5)

And then because $V_1$ and $V_2$ both are on the $H$

$\therefore \Rightarrow \begin{cases} V_1 w + \psi = 0 \\ V_2 w + \psi = 0 \end{cases}$    $\Rightarrow (V_1 - V_2) w = 0$     (3)
                        $\overrightarrow{V_1 V_2} \cdot w = 0$

And then because $\overrightarrow{x_i V_1} \perp \overrightarrow{V_1 V_2}$   $\therefore \overrightarrow{x_i V_1} \cdot \overrightarrow{V_1 V_2} = 0$   (4)

$\therefore \overset{(3)(4)}{\Rightarrow} \overrightarrow{x_i V_1} = k \cdot w$  (6)  because $\overrightarrow{V_1 V_2} \neq \vec{0}$

$\therefore \overset{(5)(6)}{\Rightarrow} d = \dfrac{|k| \cdot |w \cdot \overrightarrow{x_i V_2}|}{|k| \|w\|} = \dfrac{|w \cdot (x_i - V_2)|}{\|w\|} \overset{w V_2 + \psi = 0}{=} \dfrac{|w x_i + \psi|}{\|w\|}$

$\therefore$ Finally, we get it $d = \dfrac{|w x_i + \psi|}{\|w\|}$

(c) Because all the points have been classified correctly, then we can calculate all the distances between points and decision boundary ($H$).

$\Rightarrow \dfrac{|w x_i + \psi|}{\|w\|} \geq r_w , \forall x_i \in \{x_1 \cdots x_n\}$. where $r_w$ is the min distance
                                         between points and $H$ given $w$.

(d) Because support vector is the training points closest to the decision boundary then we can use these points to calculate the distance.

That is $\dfrac{|w x_j + \psi|}{\|w\|}$ , in which $x_j$ is $\in$ closet set (support vector)

(e) $\max\limits_{w \in R^d} r_w = \max\limits_{w \neq 0} \min\limits_{x_i \in \{x_1 \cdots x_n\}} \dfrac{|w x_i + \psi|}{\|w\|}$ , s.t. $\forall x_i \in \{x_1 \cdots x_n\}$ $y_i(x_i w + \psi) \geq 0$.

**2** ① Assume all the observations in $X$ have been normalized, which is let $X_{new\,i} = X_i - \mu$ for every $X_i$ in $X$. And $\mu = \frac{1}{n}\sum_{i}^{n} X_i$.
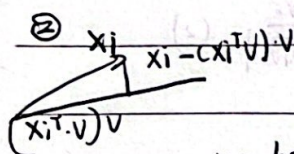
Then we can have
$$\max_{V} \frac{1}{n}\sum_{i=1}^{n}(X_i^T V - \mu^T V)^2 \quad s.t.\ \|V\|=1$$
$$= \max_{V} \frac{1}{n} V^T\left(\sum_{i=1}^{n}(X_i-\mu)(X_i-\mu)^T\right)\cdot V$$
$$= \max_{V} \frac{1}{n} V^T X X^T V \quad \longleftarrow$$

$$\left(\begin{array}{l} \because n \text{ is scalar (constant number)} \\ \therefore \Rightarrow \max_{V} V^T X X^T V \quad s.t.\ V^T V = 1 \end{array}\right) \text{ max projected variance}$$

②



$X_j$, $X_i - (X_i^T V)V$, $(X_i^T V)V$

we have this projection relation. So from the Pythagoras theorem, we have $\|X_i\|^2 = \|(X_i^T V)V\|^2 + \|X_i - (X_i^T V)V\|^2$.

Therefore, we can sum it up and have $\boxed{\frac{1}{n}\sum_{i}^{n}\|X_i\|^2} = \frac{1}{n}\sum_{i=1}^{n}\|(X_i^T V)V\|^2 + \frac{1}{n}\sum_{i}^{n}\|X_i - (X_i^T V)V\|^2$
and divide by $n$

↓ obviously this is a constant     ↓ max projected variance     ↓ our goal: min projected error

$\therefore$ Finally we get it because $\underline{\text{max projected var}} + \underline{\text{our goal}} = \text{constant}$

$\underset{\text{(}||\text{)}}{|||}$   $\frac{1}{n}\sum_{i}^{n}\|X_i - (X_i^T V)V\|^2$

$\therefore$   $\min \frac{1}{n}\sum_{i}^{n}\|X_i - (X_i^T V)V\|^2$  //

3. (a) compute $\mu = \frac{1}{n}\sum_{i}^{n} x_i = \frac{1}{6}\begin{bmatrix} 0+0+1+1+2+2 \\ 0+1+0+2+1+2 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$  $\dot{X} = X - \mu = \begin{bmatrix} -1 & -1 \\ -1 & 0 \\ 0 & -1 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix}$
$R^{6\times2}$

(b) $\dot{X}^T\dot{X} = \begin{bmatrix} -1 & -1 & 0 & 0 & 1 & 1 \\ -1 & 0 & -1 & 1 & 0 & 1 \end{bmatrix}\begin{bmatrix} -1 & -1 \\ -1 & 0 \\ 0 & -1 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} = \begin{bmatrix} 4 & 2 \\ 2 & 4 \end{bmatrix}$

$|\dot{X}^T\dot{X} - \lambda I| = \begin{vmatrix} 4-\lambda & 2 \\ 2 & 4-\lambda \end{vmatrix} = (4-\lambda)^2 - 4 = 16 + \lambda^2 - 8\lambda - 4 = \lambda^2 - 8\lambda + 12$
$= (\lambda-2)(\lambda-6)$

① when $\lambda_1 = 2$  $\begin{bmatrix} 2 & 2 \\ 2 & 2 \end{bmatrix} = \dot{X}^T\dot{X} - \lambda I$

$\Rightarrow \begin{bmatrix} 1 & 1 \\ 0 & 0 \end{bmatrix}v_1 = 0 \Rightarrow v_1 = \begin{bmatrix} +\frac{\sqrt{2}}{2} \\ -\frac{\sqrt{2}}{2} \end{bmatrix}$ or $\begin{bmatrix} -\frac{\sqrt{2}}{2} \\ +\frac{\sqrt{2}}{2} \end{bmatrix}$ $\quad \because ||v_1|| = ||v_2|| = 1$

② when $\lambda_2 = 6$  $\begin{bmatrix} -2 & 2 \\ 2 & -2 \end{bmatrix} = \dot{X}^T\dot{X} - \lambda I$

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

(c) substitude $v_1$ and $v_2$

$\Rightarrow \begin{cases} \frac{1}{n}v_{11}^T \cdot \dot{X}^T\dot{X} \cdot v_{11} = \frac{2}{3} & \frac{1}{n}v_{12}^T \cdot \dot{X}^T\dot{X} \cdot v_{12} = \frac{2}{3} \\ \frac{1}{n}v_{21}^T \cdot \dot{X}^T\dot{X} \cdot v_{21} = 2 & \frac{1}{n}v_{22}^T\dot{X}^T\dot{X} \cdot v_{22} = 2 \end{cases}$
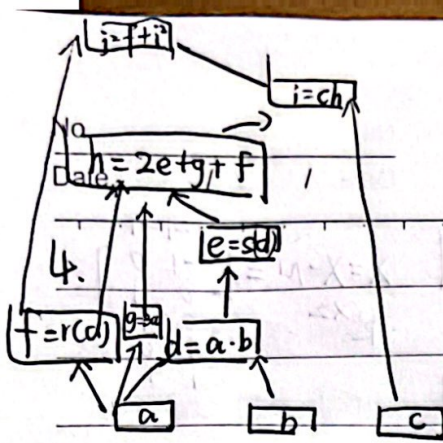
$\therefore$ ① we may choose $\lambda_2 = 6$ as the main principal component

$\rightarrow$ ② we can use the value of eignvalue $\lambda$ to choose the best. The ~~more~~ bigger $\lambda$, the better principal component.

③ Because it will max the projected variance and then it can make the information loss minimize during the dimension reduce process.

(or you can calculate one by one)

Wengu

4.

$$j=i+i$$
$$i=ch$$
$$h=2e+g+f$$
$$e=s(d)$$
$$f=r(d)$$
$$d=a\cdot b$$
$$g=a$$
$$a \qquad b \qquad c$$

① $\dfrac{\partial j}{\partial e} = \dfrac{\partial j}{\partial h}\dfrac{\partial h}{\partial e} = 2ic\cdot 2 = 4ic$

② $\dfrac{\partial j}{\partial f} = \dfrac{\partial j}{\partial i}\dfrac{\partial i}{\partial h}\dfrac{\partial h}{\partial f} + \dfrac{\partial j}{\partial f} = 2ic\cdot 1+1 = 2ic+1$

③ $\dfrac{\partial j}{\partial g} = \dfrac{\partial j}{\partial i}\dfrac{\partial i}{\partial h}\dfrac{\partial h}{\partial g} = 2ic\cdot$

④ $\dfrac{\partial j}{\partial d} = \dfrac{\partial j}{\partial i}\dfrac{\partial i}{\partial h}\dfrac{\partial h}{\partial e}\dfrac{\partial e}{\partial d} = 4ic\cdot\left(\dfrac{1}{1+e^d}\right)' = 4ic\cdot\dfrac{+e^{-d}}{(1+e^d)^2} = +4ic\cdot s(d)(1-s(d))$

⑤ $\dfrac{\partial j}{\partial c} = \dfrac{\partial j}{\partial i}\dfrac{\partial i}{\partial c} = 2ih$

⑥ $\dfrac{\partial j}{\partial b} = \dfrac{\partial j}{\partial i}\dfrac{\partial i}{\partial h}\dfrac{\partial h}{\partial e}\dfrac{\partial e}{\partial d}\dfrac{\partial d}{\partial b} = +4ic\cdot s(d)(1-s(d))\cdot a$
$= +4ica\, s(d)(1-s(d))$

⑦ $\dfrac{\partial j}{\partial a} = \dfrac{\partial j}{\partial f}\dfrac{\partial f}{\partial a} + \dfrac{\partial j}{\partial g}\cdot\dfrac{\partial g}{\partial a} + \dfrac{\partial j}{\partial d}\dfrac{\partial d}{\partial a}$

$= (2ic+1)\cdot \text{Ⅱ}(a>0) + 2ic\cdot 3 + (+4ic\cdot s(d)(1-s(d)))\cdot b$
$= (2ic+1)\text{Ⅱ}(a>0) + 6ic + 4icb\, s(d)(1-s(d))$