

HWS

Bhishan Poudel

## HW Assignment 5 (Due by 10:30am on Nov 2)

### 1 Theory (110 points)

1. [Properties of Linear Discriminants, 20 points]

We have proven in class that the distance between origin and the decision hyperplane  $h(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = 0$  is equal with  $-w_0 / \|\mathbf{w}\|$ . Prove that the margin between a point  $\mathbf{x}$  and the same decision hyperplane is equal with  $h(\mathbf{x}) / \|\mathbf{w}\|$ .

2. [Bonus, 20 points]

Prove the two properties above for the general  $n$ -dimensional case.

3. [Fisher Criterion and Least Squares, 30 points]

Show that the Fisher criterion can be written in the vectorized form shown below:

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2} = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}$$

4. [Fisher Criterion (\*), 20 points]

Reference the PRML Chapter 4 material available on Blackboard under Content.

Using the definitions of the between-class and within-class covariance matrices given by (4.27) and (4.28), respectively, together with (4.34) and (4.36) and the choice of target values described in Section 4.1.5, show that the expression (4.33) that minimizes the sum-of-squares error function can be written in the form (4.37).

5. [Perceptrons, 40 points]

Consider a training set that contains the following 8 examples:

$\mathbf{x}$	$x_1$	$x_2$	$x_3$	$t(x)$
$\mathbf{x}^{(1)}$	0	0	0	+1
$\mathbf{x}^{(2)}$	0	1	0	+1
$\mathbf{x}^{(3)}$	1.5	0	-1.5	+1
$\mathbf{x}^{(4)}$	1.5	1	-1.5	+1
$\mathbf{x}^{(5)}$	1.5	0	0	-1
$\mathbf{x}^{(6)}$	1.5	1	0	-1
$\mathbf{x}^{(7)}$	0	0	-1.5	-1
$\mathbf{x}^{(8)}$	0	1	-1.5	-1

(a) Prove that the perceptron algorithm does not converge on this dataset. Do not forget to include the bias.

(b) Consider a kernel perceptron that uses a polynomial kernel  $k(\mathbf{x}, \mathbf{y}) = (1 + \mathbf{x}^T \mathbf{y})^d$ . What is the smallest degree  $d$  for which the kernel perceptron would converge on this dataset?

6. [Perceptrons, 10 points]

A kernel perceptron for binary classification is run for a number of epochs  $E$  on a training dataset containing  $N$  examples, resulting in the dual parameters  $\alpha_1, \alpha_2, \dots, \alpha_N$ . What is the total number of mistakes that are made during training?

7. [Matrix Computations, 10 points]

Let  $U \in R_{k \times m}$  and  $X \in R_{n \times m}$ . Let  $u_i$  and  $x_i$  be the  $i$ -th columns of  $U$  and  $X$ , respectively, for  $1 \leq i \leq m$ . Prove that  $UX^T = \sum_{i=1}^m u_i x_i^T$ .

## 2 Submission

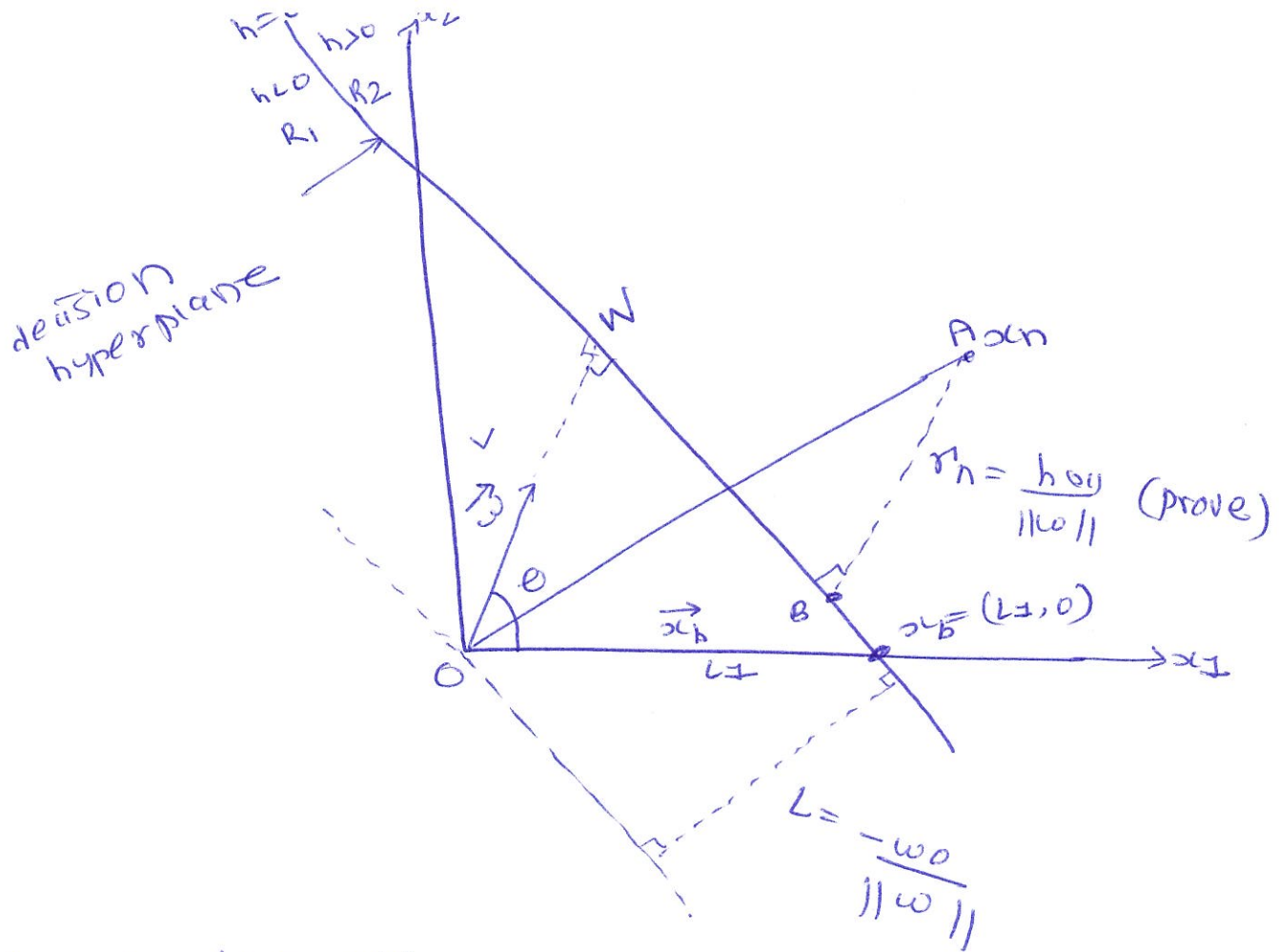
Turn in a hard copy of your homework report at the beginning of class on the due date. On this theory assignment, **clear and complete explanations and proofs of your results are as important as getting the right answer.**

HW 05  
Bhishan Poudel

QV1

prove:

$$\text{dist betn } \underline{\text{data point}} \text{ \& } \underline{\text{decision hyp.}} = \frac{h(x)}{\|w\|}$$



a) move  $L1 = -\frac{w_0}{\|w\|}$

Here,  $\vec{w} = [w_1, w_2]$

$\vec{x}_b = [x_1, x_2] = [L1, 0]$

$$\vec{w} \cdot \vec{x}_b + w_0 = 0 = w_1 x_1 + w_2 x_2 + w_0$$

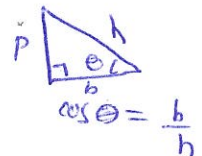
$$0 = w_1 L1 + w_2 \cdot 0 + w_0$$

$$\boxed{L1 = -\frac{w_0}{w_1}} \quad \text{--- ①}$$

Again,  $\boxed{\cos \theta = \frac{\text{base}}{\text{hypotenuse}} = \frac{L}{L1}}$

$$\cos \theta = \frac{L}{-w_0 / w_1}$$

$$\boxed{L = -\frac{w_0 \cos \theta}{w_1}} \quad \text{--- ②}$$



$$L = \frac{-w_0}{w_1} \cdot \frac{\vec{w} \cdot \vec{x}_b}{\|w\| \|x_b\|}$$

$$\therefore \vec{a} \cdot \vec{b} = \|a\| \|b\| \cos \theta$$

$$= \frac{-w_0}{w_1} \cdot \frac{(w_1, w_2) \cdot (1, 0)}{\|w\| \cdot 1}$$

$$= \frac{-w_0}{w_1} \frac{w_1}{\|w\|}$$

$$L = \frac{-w_0}{\|w\|} \quad \text{Ans}$$

Again,  
to find the co-ordinate of point B,

$$\vec{OB} = \vec{OA} - A\vec{B}$$

$$\vec{B} = \vec{x}_n - r_n \hat{w} \rightarrow \hat{w} \text{ is unit vector perpendicular to decision hyperplane}$$

$$B = x_n - r_n \frac{w}{\|w\|}$$

But, point B lies in decision boundary,  $\therefore$  so,

$$\vec{w} \cdot \vec{B} + w_0 = 0$$

$$w^T (x_n - r_n \frac{w}{\|w\|}) + w_0 = 0$$

$$\text{or, } w^T x_n - r_n \frac{w^T w}{\|w\|} + w_0 = 0$$

$$\text{or, } w^T x_n - r_n \|w\| + w_0 = 0$$

$$\therefore w^T w = \|w\|^2$$

$$\text{or } \omega^T x_n + w_0 = r_n \|\omega\|$$

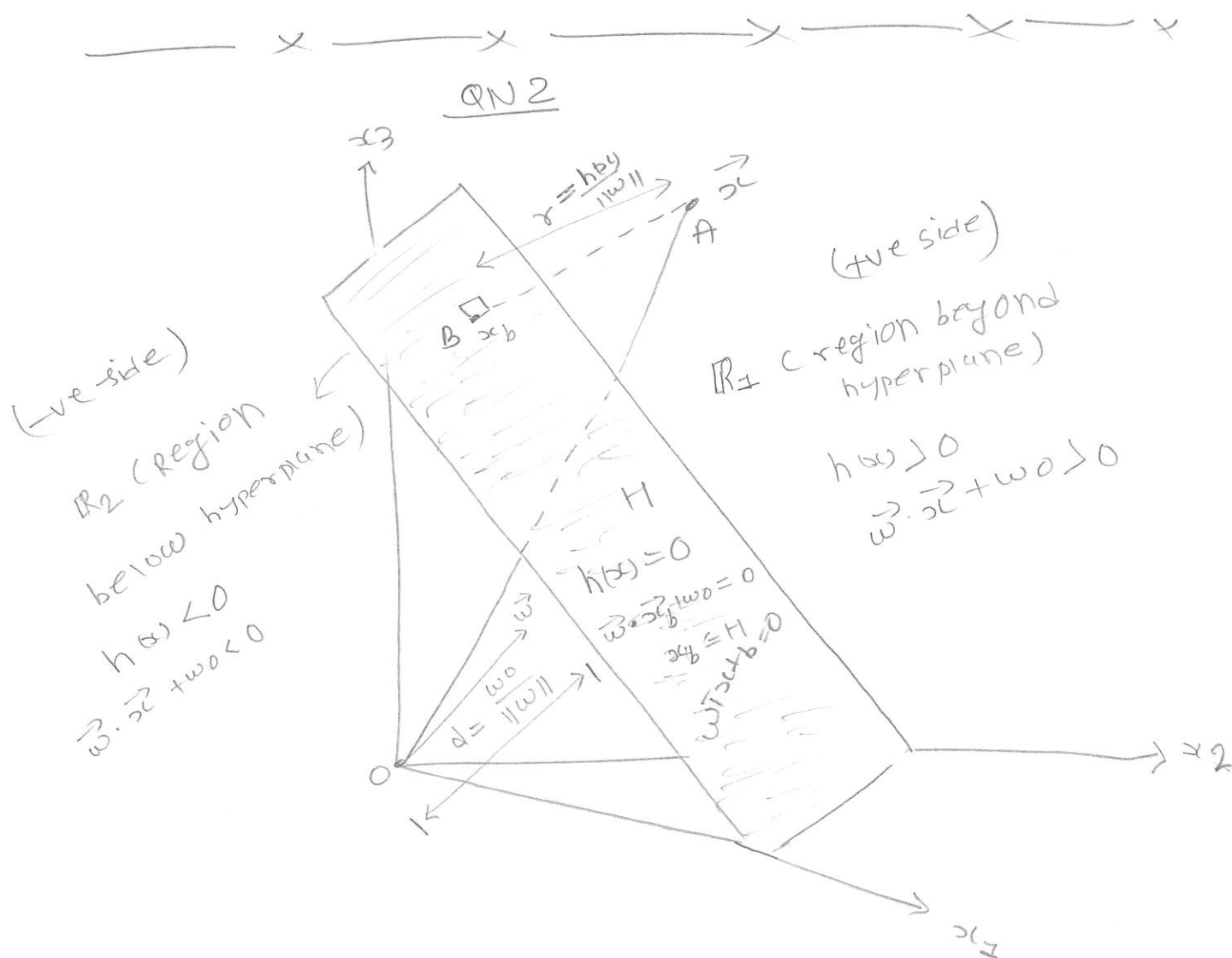
$$\therefore r_n = AB = \frac{\omega^T x_n + w_0}{\|\omega\|}$$

$$\text{where } h(x) = \omega^T x + w_0$$

$$r_n = \frac{h(x)}{\|\omega\|}$$

proved!

✓ 20





Q2

For n-dimensional case,  
prove:

$$\text{dist bet}^n \text{ origin and decision hyperplane} = \frac{-w_0}{\|w\|}$$

$$\& \text{ dist bet}^n \text{ datapoint and decision hyperplane} = \frac{h(w)}{\|w\|}$$

soln Let  $x$  be any point on the positive side of decision boundary (hyperplane), then,

$$x = x_b + r \frac{w}{\|w\|} \quad \text{where } \frac{w}{\|w\|} \text{ is unit vector perpendicular to decision hyperplane}$$

$$x_b = x - r \frac{w}{\|w\|} \quad \text{--- ①}$$

but,  $x_b$  lies in boundary, we need to find  $r$ ,

$$w \cdot x_b + w_0 = 0 \quad \text{--- ②}$$

$w$

$$\text{or, } w^T (x - r \frac{w}{\|w\|}) + w_0 = 0 = w^T x - r \frac{w^T w}{\|w\|} + w_0 = 0$$

$$\text{or, } w^T x + w_0 = r \|w\| = h(w)$$

$$\Rightarrow r = \frac{w^T x + w_0}{\|w\|}$$

distance from  
datapoint to  
hyperplane

$$r = \frac{h(w)}{\|w\|}$$

Ans

Also, we have to show the distance of hyperplane from the origin,

$$\vec{OB} = -\frac{w_0}{\|w\|} \hat{OB}$$

Here, the point  $x_b$  lies in decision hyperplane, so,

$$\vec{w} \cdot \vec{x}_b + w_0 = 0$$

$$w x_b + w_0 = 0$$

$$w x_b = -w_0$$

distance  
from

origin  
to

decision hyperplane

$$x_b = \frac{-w_0}{\|w\|}$$

Ans



Q13

Bishop  
ex 4.5

show that Fisher criterion

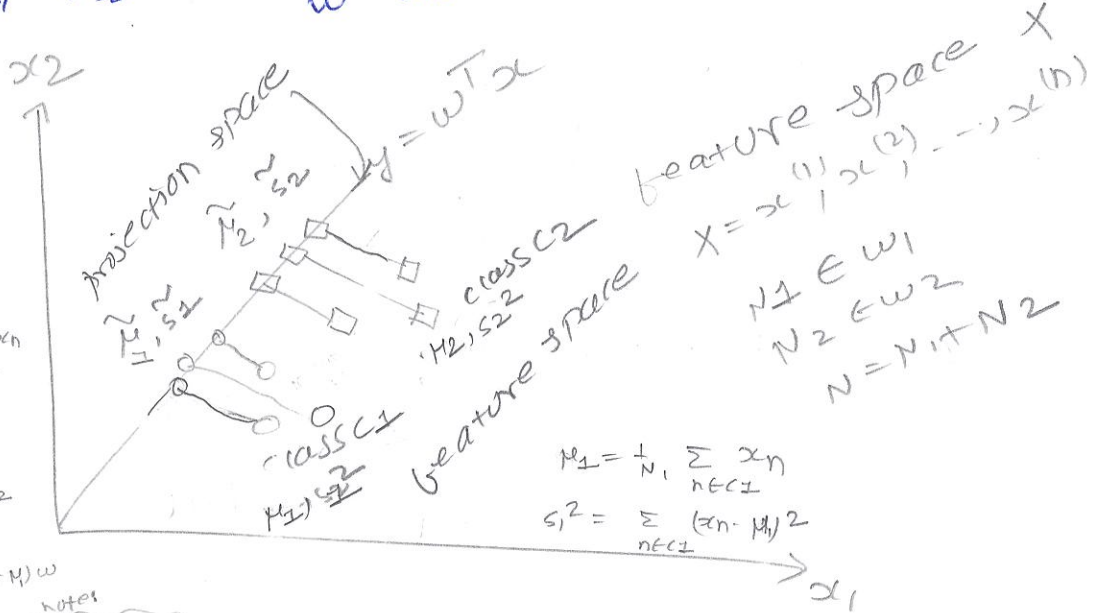
can be written as,

$$J(w) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2} = \frac{w^T S_B w}{w^T S_W w}$$

Solution

$$\begin{aligned} \tilde{\mu}_1 &= \frac{1}{N_1} \sum_{n \in c_1} w^T x_n \\ &= w^T \mu_1 \\ s_1^2 &= \sum_{n \in c_1} (x_n - \mu_1)^2 \\ &= \sum_{n \in c_1} (w^T x_n - w^T \mu_1)^2 \\ &= \sum_n w^T (x_n - \mu_1)(x_n - \mu_1)^T w \\ s_1^2 &= w^T S_1 w \end{aligned}$$

note:  
 $S_W = S_1 + S_2$



In feature space  $X$ , we have input data,

$$X = (x^{(1)}, x^{(2)}, \dots, x^{(N)})$$

feature space  $X$ , mean

$$\mu_i = \frac{1}{N_i} \sum_{x \in c_i} x$$

projection space  $Y$ , mean  $\tilde{\mu}_i = \frac{1}{N_i} \sum_{y \in c_i} y$

$$= \frac{1}{N_i} \sum_{x \in c_i} w^T x$$

$$\tilde{\mu}_i = w^T \mu_i$$

Also,

variance,  
(feature space)

$$\sigma_i^2 = \frac{1}{N_i} \sum_{x \in \mathcal{C}_i} (x - \mu_i)^2$$

variance

variance  
(projection  
space)

$$\tilde{\sigma}_i^2 = \frac{1}{N_i} \sum_{y \in \mathcal{C}_i} (y - \tilde{\mu}_i)^2$$

Also, scatter in feature space,

$$S_i^2 = \sum_{x \in \mathcal{C}_i} (x - \mu_i)^2 = \sum_{x \in \mathcal{C}_i} (x - \mu_i)(x - \mu_i)^T$$

scatter in projection space,

$$\tilde{S}_i^2 = \sum_{y \in \mathcal{C}_i} (y - \tilde{\mu}_i)^2$$

scatter

Also, within-class scatter  $S_W = S_1 + S_2$

Also,

$$\tilde{S}_W = \tilde{S}_1 + \tilde{S}_2$$

now, we define Fisher's criterion,

$$J(w) = \frac{|\tilde{\mu}_1 - \tilde{\mu}_2|^2}{\tilde{s}_1^2 + \tilde{s}_2^2}$$

Fisher's  
linear  
discriminant

now,

$$(\tilde{\mu}_1 - \tilde{\mu}_2)^2 = (w^T \mu_1 - w^T \mu_2)^2$$

$$= w^T (\mu_1 - \mu_2) (\mu_1 - \mu_2)^T w$$

$$(\tilde{\mu}_1 - \tilde{\mu}_2)^2 = w^T S_B w \quad \text{--- (A)}$$

where,  $S_B = (\mu_1 - \mu_2) (\mu_1 - \mu_2)^T$   
is between-class scatter.

also,

$$\tilde{s}_1^2 = \sum_{y \in c_1} (y - \tilde{\mu}_1)^2$$

$$= \sum_{x \in c_1} (w^T x - w^T \mu_1)^2$$

$$= \sum_{x \in c_1} w^T (x - \mu_1) (x - \mu_1)^T w$$

$$\tilde{s}_1^2 = w^T S_1 w$$

where, feature-space  
scatter

$$S_1 = \sum_{x \in c_1} (x - \mu_1) (x - \mu_1)^T$$

$$\Rightarrow \tilde{s}_1^2 = w^T S_1 w$$

$$\tilde{s}_2^2 = w^T S_2 w$$

$$\tilde{s}_1^2 + \tilde{s}_2^2 = w^T (S_1 + S_2) w = w^T S_W w$$

where  $S_W = S_1 + S_2$  is  
within-class scatter.

$$\Rightarrow \boxed{\tilde{x}_1^2 + \tilde{x}_2^2 = \omega^T S \omega} \quad \text{--- (B)}$$

using equations (A) & (B), we have,

$$J(\omega) = \frac{|\tilde{x}_1 - \tilde{x}_2|^2}{\tilde{x}_1^2 + \tilde{x}_2^2}$$

$$\boxed{J(\omega) = \frac{\omega^T S_B \omega}{\omega^T S \omega}}$$

Q.E.D.!

30

(Q4)

In Bishop Book, (p. 190 & 208)

Bishop  
ex 4.6  
p.

use equations: 4.27, 4.28, 4.34, 4.36, 4.33

to prove eq 4.37,

derive:  $\left( S_W + \frac{N_1 N_2}{N} S_B \right) w = N (m_1 - m_2)$  — (4.37)

Here, given, SSE is given by

$$E(w) = \frac{1}{2} \sum_{n=1}^N (w^T x_n + w_0 - t_n)^2 \quad \text{--- (4.31) p. 190/208}$$

$$\frac{\partial E}{\partial w} = 0 = \sum_{n=1}^N (w^T x_n + w_0 - t_n) x_n \quad \text{--- (4.33)}$$

(4.34)

$$w_0 = -w^T m$$

$$t_n = \begin{cases} \frac{N}{N_1} & \text{if } x_n \in C_1 \\ -\frac{N}{N_2} & \text{if } x_n \in C_2 \end{cases}$$

so that  $\sum_{n=1}^N t_n = N_1 \frac{N}{N_1} - N_2 \frac{N}{N_2}$

mean of the total dataset (population mean)

(4.36)  $m = \frac{1}{N} \sum_{n=1}^N x_n = \frac{1}{N} N_1 m_1 + \frac{1}{N} N_2 m_2$

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

where, (sample mean of class)

(4.24) p. 208

$$m_1 = \frac{1}{N_1} \sum_{n \in C_1} x_n \quad (\text{mean of } x_n \text{ belonging to class } C_1)$$
$$m_2 = \frac{1}{N_2} \sum_{n \in C_2} x_n$$

$$\Rightarrow \sum_{n \in C_1} x_n = N_1 m_1$$

$$\Rightarrow \sum_{n \in C_2} x_n = N_2 m_2$$



we define within-class scatter as,

eq 4.28  
p. 207

$$S_W = \sum_{n \in C_1} (x_n - m_1)(x_n - m_1)^T + \sum_{n \in C_2} (x_n - m_2)(x_n - m_2)^T$$

and, between-class scatter as,

eq 4.29  
p. 207

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

we have to prove,

TO PROVE

$$(S_W + \frac{N_1 N_2}{N} S_B) w = N(m_1 - m_2) \quad (\text{eq 4.37})$$

$$\Rightarrow \frac{N(m_1 - m_2)}{w} = S_W + \frac{N_1 N_2}{N} S_B \quad \text{--- (1)}$$

we start from gradient zero equation (4.33),

$$0 = \sum_{n=1}^N (w^T x_n - \underset{\rightarrow}{w_0} - t_n) x_n$$

$$= \sum_{n=1}^N (w^T x_n - w^T m - t_n) x_n$$

$$= \sum_{n=1}^N w^T (x_n x_n^T - m x_n) - \sum_{n=1}^N t_n x_n$$

$$= \sum_{n \in C_1} (x_n x_n^T - x_n m^T) w - \sum_{n \in C_1} t_n x_n$$

$$+ \sum_{n \in C_2} (x_n x_n^T - x_n m^T) w - \sum_{n \in C_2} t_n x_n$$

$$m_2 = \frac{1}{N_2} \sum_{n \in C_2} x_n \Rightarrow \sum_{n \in C_2} x_n = N_2 m_2$$

use  $t_n = \begin{cases} \frac{N_1}{N} & \text{for } C_1 \\ -\frac{N_2}{N} & \text{for } C_2 \end{cases}$

$$0 = \left( \sum_{n \in C_1} x_n x_n^T - N_1 m_1 m_1^T \right) \omega - \frac{N}{N_1} \mu_1 m_1$$

$$\left( \sum_{n \in C_2} x_n x_n^T - N_2 m_2 m_2^T \right) \omega + \frac{N}{N_2} \mu_2 m_2$$

$$0 = \left( \sum_{n \in C_1} x_n x_n^T + \sum_{n \in C_2} x_n x_n^T - (N_1 m_1 + N_2 m_2) m^T \right) \omega - N(m_1 - m_2)$$

$$\frac{N(m_1 - m_2)}{\omega} = \underbrace{\sum_{n \in C_1} x_n x_n^T + \sum_{n \in C_2} x_n x_n^T - (N_1 m_1 + N_2 m_2) m^T}_{\text{write in terms of } S_W} \quad \text{②}$$

$m = \frac{N_1 m_1 + N_2 m_2}{N}$

now, we expand the first term in  $S_W$ ,

$$\begin{aligned} \sum_{n \in C_1} (x_n - m_1)(x_n - m_1)^T &= \sum_{n \in C_1} (x_n x_n^T - x_n m_1^T - m_1 x_n^T + m_1 m_1^T) \\ &= \sum_{n \in C_1} x_n x_n^T - m_1^T \sum_{n \in C_1} x_n - m_1 \sum_{n \in C_1} x_n^T + \sum_{n \in C_1} m_1 m_1^T \\ &= \sum_{n \in C_1} x_n x_n^T - N_1 m_1 m_1^T - \cancel{N_1 m_1 m_1^T} + \cancel{N_1 m_1 m_1^T} \end{aligned}$$

$$\boxed{\sum_{n \in C_1} (x_n - m_1)(x_n - m_1)^T = \sum_{n \in C_1} x_n x_n^T - N_1 m_1 m_1^T}$$

$$\text{So, } S_W = \sum_{n \in C_1} (x_n - m_1)(x_n - m_1)^T + \sum_{n \in C_2} (x_n - m_2)(x_n - m_2)^T$$

$$\boxed{S_W = \sum_{n \in C_1} x_n x_n^T + \sum_{n \in C_2} x_n x_n^T - N_1 m_1 m_1^T - N_2 m_2 m_2^T}$$

$$\Rightarrow \sum_{n \in C_1} x_n x_n^T + \sum_{n \in C_2} x_n x_n^T = S_W + N_1 m_1 m_1^T + N_2 m_2 m_2^T$$

Then eqn (2) becomes,

$$\frac{N(m_1 - m_2)}{\omega} = s\omega + N_1 \underbrace{m_1}_{\uparrow} \underbrace{m_1^T}_{\uparrow} + N_2 \underbrace{m_2}_{\uparrow} \underbrace{m_2^T}_{\uparrow} - (N_1 m_1 + N_2 m_2) \cdot \left( \frac{N_1 \overset{\downarrow}{m_1} + N_2 \overset{\downarrow}{m_2}}{N} \right)^T$$

$$= s\omega + m_1 m_1^T \left( N_1 - \frac{N_1^2}{N} \right) + m_2 m_2^T \left( N_2 - \frac{N_2^2}{N} \right)$$

$$- \frac{N_1 N_2}{N} (m_1 m_2^T + m_2 m_1^T)$$

$$N = N_1 + N_2$$

$$= s\omega + \frac{(N_1 + N_2)N_1 - N_1^2}{N} m_1 m_1^T + \frac{(N_1 + N_2)N_2 - N_2^2}{N} m_2 m_2^T$$

$$- \frac{N_1 N_2}{N} (m_1 m_2^T + m_2 m_1^T)$$

$$= s\omega + \frac{N_1 N_2}{N} (m_1 m_1^T + m_2 m_2^T)$$

$$- \frac{N_1 N_2}{N} (m_1 m_2^T + m_2 m_1^T)$$

$$= s\omega + \frac{N_1 N_2}{N} (m_1 m_1^T + m_2 m_2^T - m_1 m_2^T - m_2 m_1^T)$$

$$\frac{N(m_1 - m_2)}{\omega} = s\omega + \frac{N_1 N_2}{N} s_B$$

$$\begin{aligned} s_B &= (m_2 - m_1)(m_2 - m_1)^T \\ &= m_2 m_2^T - m_2 m_1^T - m_1 m_2^T + m_1 m_1^T \\ &= m_1 m_1^T + m_2 m_2^T - m_1 m_2^T - m_2 m_1^T \end{aligned}$$

$$\therefore \left( s\omega + \frac{N_1 N_2}{N} s_B \right) \omega = N(m_1 - m_2)$$

Q.E.D.

20

✓

QNS

Given training set,

$x$	$x_1$	$x_2$	$x_3$	$t(x)$
$x^{(1)}$	0	0	0	+1
$x^{(2)}$	0	1	0	+1
$x^{(3)}$	1.5	0	-1.5	+1
$x^{(4)}$	1.5	1	-1.5	+1
$x^{(5)}$	1.5	0	0	-1
$x^{(6)}$	1.5	1	0	-1
$x^{(7)}$	0	0	-1.5	-1
$x^{(8)}$	0	1	-1.5	-1

a) prove that the perceptron algorithm does not converge on this dataset. (Note: Include bias term)

b) consider a kernel perceptron that uses a polynomial kernel  $K(x, y) = (1 + x^T y)^d$ .

what is the smallest degree  $d$  for which the kernel perceptron would converge on this dataset?

So for a given example, example belongs to class  $\pm 1$  if  
 $w_0 + \vec{w} \cdot \vec{x}_n \geq 0$

and, example belongs to class  $-1$  if  
 $w_0 + \vec{w} \cdot \vec{x}_n < 0$

now,  $x^{(1)} = \begin{matrix} w_1 & w_2 & w_3 \\ [0, 0, 0] \end{matrix} \quad t^{(1)} = \pm 1$

$$w_0 + w_1 \cdot 0 + w_2 \cdot 0 + w_3 \cdot 0 \geq 0$$

$$w_0 \geq 0$$

similarly,

+ve	①	$w_0 + 0 + 0 + 0 \geq 0$	①	
	②	$w_0 + 0 + w_2 + 0 \geq 0$	②	
	③	$w_0 + 1.5w_1 + 0 - 1.5w_3 \geq 0$	③	④ - ③ $\Rightarrow w_2 \geq 0$
	④	$w_0 + 1.5w_1 + w_2 - 1.5w_3 \geq 0$	④	
-ve	⑤	$w_0 + 1.5w_1 + 0 + 0 < 0$	⑤	⑥ - ⑤ $\Rightarrow w_2 < 0$
	⑥	$w_0 + 1.5w_1 + w_2 + 0 < 0$	⑥	
	⑦	$w_0 + 0 + 0 - 1.5w_3 < 0$	⑦	
	⑧	$w_0 + 0 + w_2 - 1.5w_3 < 0$	⑧	

here, ④ - ③ gives  $w_2 \geq 0$

but ⑥ - ⑤ gives  $w_2 < 0$

this is contradictory  
 $\Rightarrow$   
 perceptron fail.



(5b) Here, the given dataset is not linearly separable in  $x$ -space.

we shall propose a perceptron that uses polynomial kernel,

the polynomial kernel is,

$$K(x, x') = (a^T x + b)^d$$

for  $a = b = 1$

$$K(x, x') = (1 + x^T x')^d$$

$a$  = scaling factor

$b$  = bias term

$d$  = degree of polynomial

$x$  = augmented dataset  
of shape  $N \times (m+1)$

$x'$  = one example of dataset  
with  $m+1$  features

when  $d=1$  : the kernel is linear classifier, but  
our dataset is NON linear, so it cannot  
classify correctly. 12

when  $d=2$  : when degree is 2, quadratic kernel  
can classify linearly non-separable  
dataset. CAN

we find a vector to  
separate the data set.

when  $d \geq 2$  : higher degree polynomial kernels can  
fit the linearly non-separable dataset  
good, but, they may overfit the  
test dataset. may overfit

$\therefore$  smallest degree  $d=2$  Ans

qiv 6

For a kernel perceptron,  
from lecture note, (Lecture 05)  
the algorithm is

1.  $f(x) = w^T x = \sum_n \alpha_n t_n \phi_n^T x = \sum_n \alpha_n t_n \underbrace{K(x, x_n)}_{\text{kernel}}$

2. initialize  $\alpha$  :  $\alpha_n = 0$

for  $e$  in  $E$  epochs

3. for  $n = 1 \dots N$

$$h_n = \text{sgn}(f(x_n))$$

if  $h_n \neq t_n$  then

$$\alpha_n = \alpha_n + 1$$

Now, from step 3, if we run perceptron  
 $E$  number of epochs,

then total number of mistakes made

so far is,

$$M = \sum_{n=1}^N \alpha_n$$

Ans

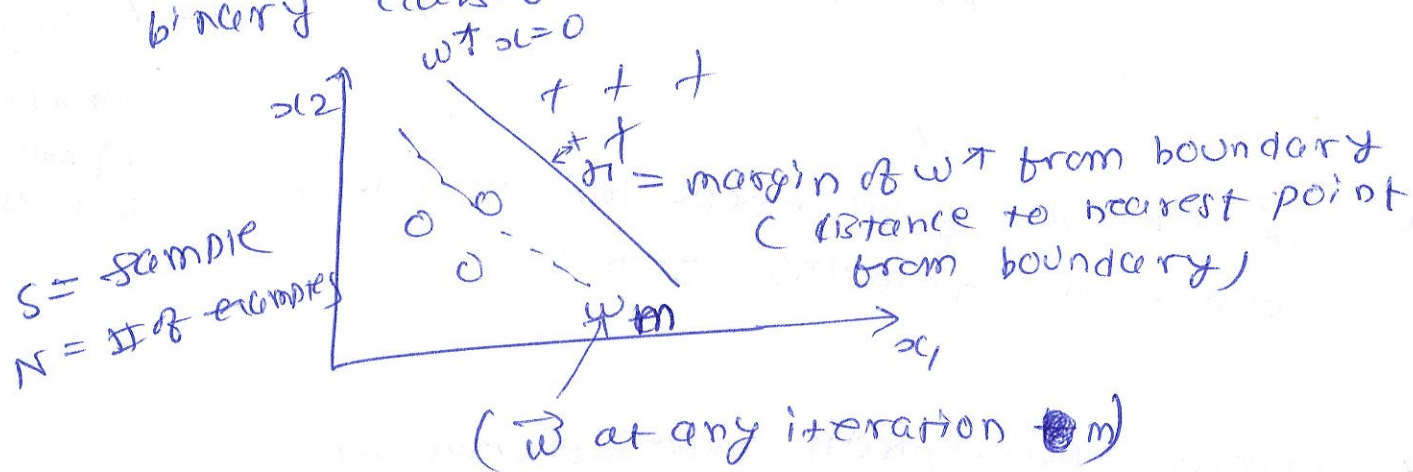
/D

✓

QNB Method 2 (Extra practice)  
A Kernel perceptron for binary

maximum error bound  
 Classification is run for a number of epochs  $E$  on a training dataset containing  $N$  examples, resulting in the dual parameters  $\alpha_1, \alpha_2, \dots, \alpha_N$ .  
 what is the total number of mistakes that are made during training?

Soln Consider a kernel perceptron for binary classification.



Let there is a unit vector  $\vec{w}^T$  that can separate  $N$  training examples into two classes with margin  $\gamma$ .

Let  $R = \text{norm of maximum value of } \vec{x}$  in dataset

$$\|\vec{R}\| = \max_{x \in S} \|\vec{x}\|$$

# perceptron algorithm:

1. initialize  $w$  to zero.  $w_1 = [0 \ 0 \ 0 \ 0]$

2. for  $e$  in epochs

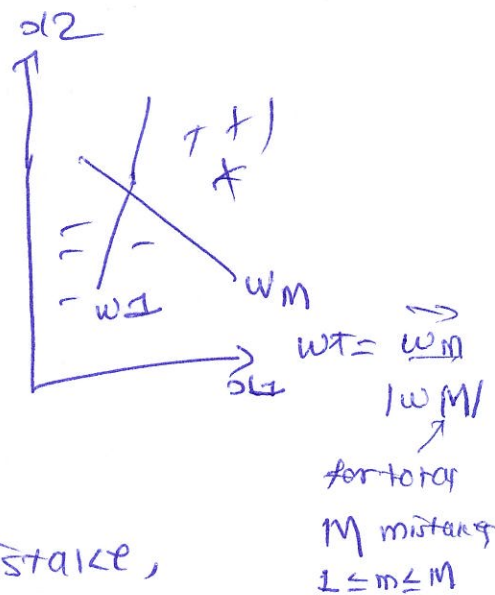
for  $n = 1$  to  $N$

if  $t(x_n) \cdot w \cdot x_n < 0$

$w = w + t x_n$

if converged  $w$ :

break



\*  $\vec{w}$  changes if we encounter mistake, otherwise remains same.

\* Suppose perceptron makes a mistake at iteration  $m$ , then,

$$\vec{w}_m \cdot w^* = (\vec{w}_{m-1} + t \vec{x}) \cdot w^*$$

$$= \vec{w}_{m-1} \cdot w^* + t(\vec{x} \cdot w^*)$$

$$w_m \cdot w^* \geq \vec{w}_{m-1} \cdot w^* + \gamma$$

from definition of  $\gamma$ .  
each data point  $x$   
is  $\geq \gamma$  from the  
boundary  $w^*$



$$\therefore w_m \cdot w^T \geq w_{m-1} \cdot w^T + \gamma \quad \text{--- (A)}$$

Also,

$$\begin{aligned} \|w_m\|^2 &= \|w_{m-1} + t\alpha\|^2 \\ &= \|w_{m-1}\|^2 + 2t(w_{m-1} \cdot \alpha) + \underbrace{\|t\alpha\|^2}_{\text{⑥ maximum of } \|w\| \text{ is } R} \end{aligned}$$

⑥ Here we have mistake at iteration  $m$ , so the previous iteration has hypothesis  $t(w_{m-1} \cdot \alpha) < 0$

$$\|w_m\|^2 \leq \|w_{m-1}\|^2 + R^2 \quad \text{--- (B)} \quad \left( \begin{array}{l} \text{from} \\ \text{⑥ \& ⑤} \end{array} \right)$$

Here, ~~both~~ inequalities ⑥ & ③ hold true for all the iterations  $m$ ,

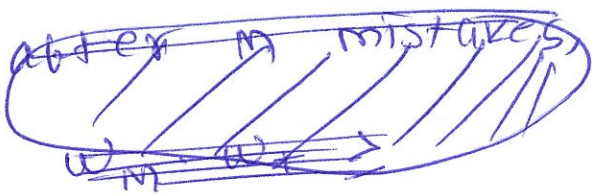
$$w_m \cdot w^T \geq w_{m-1} \cdot w^T + \gamma$$

$$\|w_m\|^2 \leq \|w_{m-1}\|^2 + R^2$$

$$\|w_m\|^2 - \|w_{m-1}\|^2 \leq R^2$$

( If there is a mistake then square of weight vectors is less than  $R^2$  ),  
 ( for each mistake  $w^2 - w_{m-1}^2$  cannot be larger than  $R^2$  )





For all iterations,

$$w_m \cdot w_T \geq w_{m-1} \cdot w_T + \gamma$$

after  $M$  mistakes,

$$\boxed{\vec{w}_M \cdot \vec{w}_T \geq \gamma M} \quad \text{--- (C)}$$

also, for all iterations,

$$\|w_m\|^2 \leq \|w_{m-1}\|^2 + R^2$$

after  $M$  mistakes

$$\|w_M\|^2 \leq R^2 M$$

$$\boxed{\|w_M\| \leq R \sqrt{M}} \quad \text{--- (D)}$$

so, from (C) & (D),

$$\gamma M \leq \vec{w}_M \cdot \vec{w}_T$$

$$\leq \|\vec{w}_M\| \|\vec{w}_T\|$$

$$\leq \|w_M\|$$

$$\gamma M \leq R \sqrt{M}$$

$$\gamma \sqrt{M} \leq R$$

$$\gamma^2 M \leq R^2$$

$$\therefore \vec{a} \cdot \vec{b} = ab \cos \theta$$

~~cos theta is less than or equal to 1~~

cos  $\theta$  is less than or equal to 1

$\downarrow$   $w_T$  is unit vector of  $\vec{w}_M$

$$\|w_T\| = 1$$

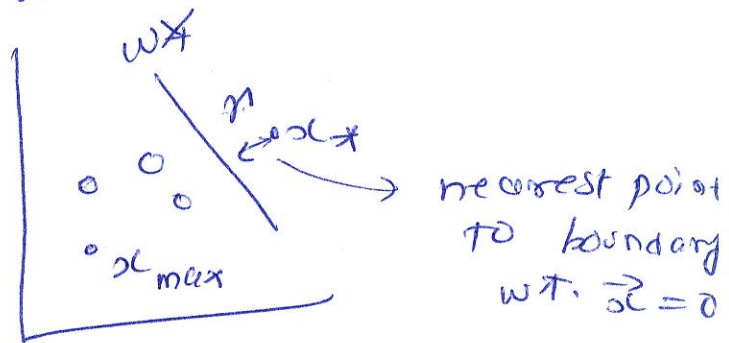
$$\boxed{M \leq \left(\frac{R}{\gamma}\right)^2}$$

Ans

Hence,

Let the data (training) sample  $S$  is,

$x_1$	$t_1$
$x_2$	$t_2$
$x_3$	$t_3$
$\vdots$	$\vdots$
$x_N$	$t_N$



$$\|x_{\max}\| = R$$

example,

$x_1:$	0 0 1 2 3	$t_1 = +$
$x_2:$	5 0 3 2 6	$t_2 = -$
$x_N:$	10 20 30 40	$t_N = +$

then, if we use Perceptron to train our model, the total number of mistakes ( $M$ ) made is,

$$M \leq \left( \frac{R}{\gamma} \right)^2$$

QNT

Let  $U \in \mathbb{R}^{k \times m}$

$X \in \mathbb{R}^{n \times m}$

$u_i, x_i = i\text{th column of } U \text{ and } X \text{ resp.}$   
 $1 \leq i \leq m.$

prove  $UX^T = \sum_{i=1}^m u_i x_i^T$

Solution:

$$U = U_{k,m} = \begin{pmatrix} u_{11} & u_{12} & \dots & u_{1m} \\ u_{21} & u_{22} & \dots & u_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ u_{k1} & u_{k2} & \dots & u_{km} \end{pmatrix}_{k,m} = \begin{pmatrix} u_1 & u_2 & \dots & u_m \end{pmatrix}_{1,m}$$

$\uparrow$   
 $k \text{ rows}$   
 $m \text{ cols}$

$$X = X_{n,m} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1m} \\ x_{21} & x_{22} & \dots & x_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nm} \end{pmatrix}_{n,m} = \begin{pmatrix} x_1 & x_2 & \dots & x_m \end{pmatrix}_{1,m}$$

$\uparrow$   
 $x_i \text{ column vector}$

$n \text{ rows}$   
 $m \text{ cols}$

$$V = X^T = \begin{pmatrix} x_{11} & x_{21} & \dots & x_{n1} \\ x_{12} & x_{22} & \dots & x_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ x_{1m} & x_{2m} & \dots & x_{nm} \end{pmatrix} = \begin{pmatrix} x_{11}^T \\ x_{12}^T \\ \vdots \\ x_{1m}^T \end{pmatrix}_{m,1}$$

From definition of matrix multiplication,

$$\text{if } U = U_{km}$$

$$\text{and } V = V_{mn}$$

then,

$$(UV)_{pq} = \sum_{i=1}^m U_{pi} V_{iq}$$

$$1 \leq p \leq k$$

$$1 \leq q \leq n$$

$$\forall 1 \leq p \leq k \text{ and } 1 \leq q \leq n$$

$$\text{where } X = \begin{pmatrix} x_{11} & x_{12} & x_{1m} \\ x_{21} & x_{22} & x_{2m} \\ x_{n1} & x_{n2} & x_{nm} \end{pmatrix}_{n,m}$$

n rows  
m cols

$$x_i = \begin{pmatrix} x_{i1} \\ x_{i2} \\ x_{im} \end{pmatrix}_{1,2}$$

$$= \begin{pmatrix} U_{11} & U_{12} & U_{1m} \\ U_{21} & U_{22} & U_{2m} \\ \vdots & \vdots & \vdots \\ U_{k1} & U_{k2} & U_{km} \end{pmatrix}_{k,m} \quad \downarrow \quad \begin{pmatrix} x_{11} & x_{21} & \dots & x_{n1} \\ x_{12} & x_{22} & & x_{n2} \\ & & & x_{nm} \end{pmatrix}_{m,n}$$

$$= \begin{pmatrix} U_{11}x_{11} + U_{12}x_{12} + U_{1m}x_{1m} & U_{11}x_{21} + U_{12}x_{22} + U_{1m}x_{2m} & U_{11}x_{n1} + U_{12}x_{n2} + U_{1m}x_{nm} \\ U_{21}x_{11} + U_{22}x_{12} + U_{2m}x_{1m} & U_{21}x_{21} + U_{22}x_{22} + U_{2m}x_{2m} & U_{21}x_{n1} + U_{22}x_{n2} + U_{2m}x_{nm} \\ \vdots & \vdots & \vdots \\ U_{k1}x_{11} + U_{k2}x_{12} + U_{km}x_{1m} & U_{k1}x_{21} + U_{k2}x_{22} + U_{km}x_{2m} & U_{k1}x_{n1} + U_{k2}x_{n2} + U_{km}x_{nm} \end{pmatrix}_{k,n}$$

$$= \begin{pmatrix} U_{11}x_{11} & U_{11}x_{21} & U_{11}x_{n1} \\ U_{21}x_{11} & U_{21}x_{21} & U_{21}x_{n1} \\ \vdots & \vdots & \vdots \\ U_{k1}x_{11} & U_{k1}x_{21} & U_{k1}x_{n1} \end{pmatrix}_{k,n} + \begin{pmatrix} U_{12}x_{12} & U_{12}x_{22} & U_{12}x_{n2} \\ U_{22}x_{12} & U_{22}x_{22} & U_{22}x_{n2} \\ \vdots & \vdots & \vdots \\ U_{k2}x_{12} & U_{k2}x_{22} & U_{k2}x_{n2} \end{pmatrix}_{k,n} + \begin{pmatrix} U_{1m}x_{1m} & U_{1m}x_{2m} & U_{1m}x_{nm} \\ U_{2m}x_{1m} & U_{2m}x_{2m} & U_{2m}x_{nm} \\ \vdots & \vdots & \vdots \\ U_{km}x_{1m} & U_{km}x_{2m} & U_{km}x_{nm} \end{pmatrix}_{k,n}$$

$$= U_1 X_1^T + U_2 X_2^T + U_m X_m^T$$

10 ✓

since,

$$U_1 X_1^T = \begin{pmatrix} u_{11} \\ u_{21} \\ \vdots \\ u_{k1} \end{pmatrix} \begin{pmatrix} x_{11} \\ x_{21} \\ \vdots \\ x_{n1} \end{pmatrix}^T$$

$$= \begin{matrix} \rightarrow \\ \begin{pmatrix} u_{11} \\ u_{21} \\ \vdots \\ u_{k1} \end{pmatrix} \end{matrix} \begin{matrix} \downarrow \\ (x_{11} \ x_{21} \ x_{n1}) \end{matrix}_{1,n} \begin{matrix} \\ k,1 \end{matrix}$$

$$U_1 X_1^T = \begin{pmatrix} u_{11}x_{11} & u_{11}x_{21} & u_{11}x_{n1} \\ u_{21}x_{11} & u_{21}x_{21} & u_{21}x_{n1} \\ \vdots & \vdots & \vdots \\ u_{k1}x_{11} & u_{k1}x_{21} & u_{k1}x_{n1} \end{pmatrix}_{k,n} \quad \text{--- (a)}$$

similarly,

$$U_2 X_2^T = \begin{pmatrix} u_{12}x_{12} & u_{12}x_{22} & u_{12}x_{n2} \\ u_{22}x_{12} & u_{22}x_{22} & u_{22}x_{n2} \\ \vdots & \vdots & \vdots \\ u_{k2}x_{12} & u_{k2}x_{22} & u_{k2}x_{n2} \end{pmatrix}_{k,n} \quad \text{--- (b)}$$

and,

$$U_m X_m^T = \begin{pmatrix} u_{1m}x_{1m} & u_{1m}x_{2m} & u_{1m}x_{nm} \\ u_{2m}x_{1m} & u_{2m}x_{2m} & u_{2m}x_{nm} \\ \vdots & \vdots & \vdots \\ u_{km}x_{1m} & u_{km}x_{2m} & u_{km}x_{nm} \end{pmatrix} \quad \text{--- (c)}$$

$$\therefore U X^T = U_1 X_1^T + U_2 X_2^T + \dots + U_m X_m^T$$

$$\boxed{U X^T = \sum_{i=1}^m U_i X_i^T} \quad \text{Q.E.D.}$$