

# Question 1: Conditional Probability

## • Part (a)

$$(i) P(\text{windy}, \text{hits}) = P(\text{hit} | \text{windy}) P(\text{windy}) = (0.4)(0.3) = \underline{0.12}$$

$$(ii) P(\text{hit}) = P(\text{hit}, \text{windy}) + P(\text{hit}, \neg \text{windy}) = P(\text{hit} | \text{windy}) P(\text{windy}) + P(\text{hit} | \neg \text{windy}) P(\neg \text{windy}) = 0.4(0.3) + 0.7(0.7) = \underline{0.61}$$

$$(iii) \frac{P(\text{hit then not hit}) + P(\text{not hit then hit})}{P(\text{hit, hit}) + P(\text{hit, not hit}) + P(\text{not hit, hit}) + P(\text{not hit, not hit})} = \frac{0.61(0.39) + 0.39(0.61)}{0.61^2 + 2(0.61)(0.39) + 0.39^2} = \underline{0.4758}$$

$$(iv) P(\neg \text{windy} | \text{miss}) = \frac{P(\neg \text{windy}, \text{miss})}{P(\text{miss})} = \frac{P(\text{miss} | \neg \text{windy}) P(\neg \text{windy})}{1 - P(\text{hit})} = \frac{(1 - P(\text{hit} | \neg \text{windy}))(1 - P(\text{windy}))}{1 - P(\text{hit})} = \frac{0.3(0.7)}{0.39} = \underline{0.538}$$

## • Part (b)

$$P(A|B, C) = \frac{P(A, B, C)}{P(B, C)} = \frac{P(C|A, B) P(A, B)}{P(C|B) P(B)} = \frac{P(C|A, B)}{P(C|B)} \frac{P(A, B)}{P(B)}$$

$$P(A|B) = \frac{P(A, B)}{P(B)}$$

Since we are given that  $P(A|B, C) > P(A|B)$  we can write:

$$\frac{P(C|A, B)}{P(C|B)} \frac{P(A, B)}{P(B)} > \frac{P(A, B)}{P(B)} \longrightarrow \frac{P(C|A, B)}{P(C|B)} > 1 \longrightarrow P(C|A, B) > P(C|B)$$

we can rewrite this as:

(\*)

$$1 - P(C^c|A, B) > 1 - P(C^c|B) \longrightarrow \underline{P(C^c|A, B) < P(C^c|B)}$$

$$P(A|B, C^c) = \frac{P(C^c|A, B) P(A, B)}{P(C^c|B) P(B)} = \frac{P(C^c|A, B)}{P(C^c|B)} \frac{P(A, B)}{P(B)} = \frac{P(C^c|A, B)}{P(C^c|B)} P(A|B)$$

Since we know (\*), and probabilities are positive,  $\frac{P(C^c|A, B)}{P(C^c|B)} < 1$  and thus,  $P(A|B, C^c) < P(A|B)$  given that  $P(A|B, C) > P(A|B)$ .

## Question 2: Positive Definiteness

Part (a)

(i)  $\Rightarrow$  (ii)

$$A \geq 0$$

$$x^T A x \geq 0$$

$(By)^T A (By) \geq 0$  Say that  $x = By$  where  $By \in \mathbb{R}^n$  and  $B$  is invertible in  $\mathbb{R}^{n \times n}$  since  $x \in \mathbb{R}^n$

$$y^T (B^T A B) y \geq 0 \quad (*)$$

since  $B$  is invertible,  $y \in \mathbb{R}^n$  and since  $(*)$  is in the form of the definition of a positive semidefinite matrix, we can see that  $B^T A B \geq 0$  for some invertible matrix  $B \in \mathbb{R}^{n \times n}$

(ii)  $\Rightarrow$  (i)

$$B^T A B \geq 0$$

$$x^T B^T A B x \geq 0$$

$$(Bx)^T A (Bx) \geq 0$$

$$y^T A y \geq 0, \text{ where } y = Bx \quad (*)$$

since  $\forall x \in \mathbb{R}^n$   $x^T B^T A B x \geq 0$ , we can see that applying the invertible  $\mathbb{R}^{n \times n}$  matrix  $B$  on  $x$  maps  $y \in \mathbb{R}^n$ . Therefore,  $A$  in  $(*)$  is in the form of the definition of a positive semidefinite matrix and  $A \geq 0$ .

(i)  $\Rightarrow$  (iii)

(i)  $A \geq 0$

$$x^T A x \geq 0$$

$\lambda^T (\lambda v) \geq 0$  using the definition of eigenvalue  $\lambda$  and eigenvector  $v \in \mathbb{R}^n$ :  $Av = \lambda v$ , for some  $\lambda$  and its  $v$

$$\lambda v^T v \geq 0$$

$\lambda \|v\|_2^2 \geq 0$  for this to hold true,  $\lambda \geq 0$ . Therefore, all eigenvalues of  $A$  must be non-negative

(iii)  $\Rightarrow$  (iv)

$A = U \Lambda U^T$  by the Spectral Theorem for Symmetric Matrices

$A = U \Lambda^{1/2} \Lambda^{1/2} U^T$  since  $\Lambda$  is a diagonal matrix and all  $\lambda \geq 0$

$$A = U \Lambda^{1/2} (\Lambda^{1/2})^T U^T$$

$$A = (U \Lambda^{1/2}) (U \Lambda^{1/2})^T$$

$$A = M M^T, \quad M = U \Lambda^{1/2}$$

Therefore there exists some matrix  $M \in \mathbb{R}^{n \times n}$  s.t.  $A = M M^T$

(iv)  $\Rightarrow$  (i)

$$A = U U^T$$

$$x^T A x = x^T U U^T x$$

$$x^T A x = (U^T x)^T (U^T x)$$

$x^T A x = \|U^T x\|_2^2 \geq 0$  Therefore if there exists some matrix  $U \in \mathbb{R}^{n \times n}$  s.t.  $A = U U^T$ , then  $A \geq 0$

$$x^T A x \geq 0$$

$$A \geq 0$$

Therefore, we conclude that (i), (ii), (iii), and (iv) are equivalent.

Part (b)

$$(i) \underline{x}^T \underline{x} = \underline{x}^T \underline{x} = \|\underline{x}\|_2^2 > 0 \quad \text{Therefore } \lambda I > 0 \text{ since } \lambda > 0$$

$$\underline{x}^T (A + \lambda I) \underline{x}$$

$$\underline{x}^T A \underline{x} + \underline{x}^T \lambda I \underline{x}$$

$$\underline{x}^T A \underline{x} + \lambda \|\underline{x}\|_2^2$$

$$\text{Since } A > 0, \underline{x}^T A \underline{x} > 0 \text{ and since } \lambda > 0, \lambda \|\underline{x}\|_2^2 > 0$$

$$\text{Therefore, } \underline{x}^T A \underline{x} + \lambda \|\underline{x}\|_2^2 > 0 \text{ and } A + \lambda I > 0$$

$$(ii) A - \gamma I > 0$$

$$\underline{x}^T (A - \gamma I) \underline{x} > 0$$

$$\underline{x}^T A \underline{x} - \underline{x}^T \gamma I \underline{x} > 0$$

$$\underline{x}^T A \underline{x} - \gamma \|\underline{x}\|_2^2 > 0$$

$$\gamma < \frac{\underline{x}^T A \underline{x}}{\|\underline{x}\|_2^2}$$

$$\gamma > 0 \text{ since } \underline{x}^T A \underline{x} > 0 \text{ and } \|\underline{x}\|_2^2 > 0$$

$$\text{Therefore there exists some } \gamma \text{ s.t. } A - \gamma I > 0$$

(iii) Say that not all the diagonal entries of  $A$  are positive s.t. the  $i$ -th diagonal element  $d_i < 0$ .  
 construct  $\underline{x}$  as a vector with 1 as its  $i$ -th element and 0 for the rest.  $\underline{x}$  must still be in  $\mathbb{R}^n - \{0\}$ .  
 Then we can see that  $\underline{x}^T A \underline{x} = d_i$ , which is a contradiction. Thus, all the diagonal entries of  $A$  must be positive.

$$(iv) \text{ Say that } \underline{x} \in \mathbb{R}^n \text{ and } \underline{x} = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}$$

$$\underline{x}^T A \underline{x} = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}^T [a_{11} \dots a_{nn}] \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} = \sum_{i=1}^n \sum_{j=1}^n a_{ij} > 0 \text{ by definition of } A > 0$$

$$\text{Therefore, } \sum_{i=1}^n \sum_{j=1}^n A_{ij} > 0 \text{ if } A > 0 \in \mathbb{R}^{n \times n}$$

### Question 3: Derivatives and Norms

Part (a)

$$\nabla_x (a^T x) = \nabla_x (\sum a_i x_i) = \underline{a}$$

Part (b)

$$\begin{aligned} \nabla_x (x^T A x) &= \nabla_x ([x^T a_1 \dots x^T a_n] x) = \nabla_x (\sum_{i=1}^n x_i \sum_{j=1}^n x_j A_{ij}) = \nabla_x (\sum_{i=1}^n \sum_{j=1}^n x_i x_j A_{ij}) \\ &= \begin{bmatrix} 2x_1 A_{11} + \sum_{j=2}^n x_j (A_{1j} + A_{j1}) \\ \vdots \\ 2x_n A_{nn} + \sum_{j=1}^{n-1} x_j (A_{nj} + A_{jn}) \end{bmatrix} = \begin{bmatrix} \sum_{j=1}^n x_j (A_{1j} + A_{j1}) \\ \vdots \\ \sum_{j=1}^n x_j (A_{nj} + A_{jn}) \end{bmatrix} = (A + A^T) x \end{aligned}$$

if  $A$  is symmetric then  $A = A^T$ . Thus if  $A$  is symmetric  $\nabla_x (x^T A x) = 2Ax$

Part (c)

$$A^T x = [a_1^T x, \dots, a_n^T x]$$

$$\text{trace}(A^T x) = \sum_{i=1}^n (A^T x)_{ii} = \sum_{i=1}^n \sum_{j=1}^n A_{ij}^T x_{ij} = \sum_{i=1}^n \sum_{j=1}^n A_{ij} x_{ji} = t \rightarrow \nabla_x(t) = \begin{bmatrix} \partial t / \partial x_{11} & \dots & \partial t / \partial x_{1n} \\ \vdots & & \vdots \\ \partial t / \partial x_{n1} & & \partial t / \partial x_{nn} \end{bmatrix}$$

$$\nabla_x (\text{trace}(A^T x)) = A$$

Part (d)

$$\text{Say } x = \begin{bmatrix} a^2 \\ b^2 \end{bmatrix} \text{ and } y = \begin{bmatrix} b^2 \\ a^2 \end{bmatrix}, a, b \in \mathbb{R}$$

$$f(x+y) = (\sqrt{a^2+b^2} + \sqrt{b^2+a^2})^2 = (2\sqrt{a^2+b^2})^2 = 4(a^2+b^2)$$

$$f(x) + f(y) = (\sqrt{a^2} + \sqrt{b^2})^2 + (\sqrt{b^2} + \sqrt{a^2})^2 = 2(a+b)^2$$

clearly  $f(x+y) > f(x) + f(y)$ . Counterexample so  $f(x)$  is not a norm for vectors  $x \in \mathbb{R}^2$ .

Part (e)

$$\|x\|_\infty = \max_i |x_i|$$

$$\|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$$

It's easy to see that  $\|x\|_\infty \leq \|x\|_2$  since we are summing only positive squares for  $\|x\|_2$ , of which  $\max_i |x_i|$  is included. The max value of  $\|x\|_2$  would be a uniform vector  $\underline{u} = \begin{bmatrix} u \\ \vdots \\ u \end{bmatrix}$ . In this case,

$$\|\underline{u}\|_2 = \sqrt{\sum_{i=1}^n u^2} = \sqrt{n u^2} = \sqrt{n} u$$

$$\text{since } u \text{ is the max element in } \underline{u}, \|\underline{u}\|_2 = \sqrt{n} \max_i |x_i| = \sqrt{n} \|x\|_\infty$$

$$\therefore \|x\|_\infty \leq \|x\|_2 \leq \sqrt{n} \|x\|_\infty$$

Part (f)

By Cauchy-Schwarz:  $|\langle x, y \rangle| \leq \|x\|_2 \|y\|_2$

$$|\langle x, \underline{1} \rangle| \leq \|x\|_2 \|\underline{1}\|_2$$

$$\sum_{i=1}^n |x_i| \leq \sqrt{\sum_{i=1}^n x_i^2} \sqrt{\sum_{i=1}^n 1}$$

$$\|x\|_1 \leq \sqrt{n} \|x\|_2$$

$$\|x\|_1^2 = \left( \sum_{i=1}^n |x_i| \right)^2 = (|x_1| + |x_2| + \dots + |x_n|)^2$$

$$\|x\|_2^2 = \sum_{i=1}^n x_i^2 = (x_1^2 + x_2^2 + \dots + x_n^2)$$

clearly  $\|x\|_1^2 \geq \|x\|_2^2$  and since  $\|x\|_1 \geq 0$ , then

$$\|x\|_1 \geq \|x\|_2$$

$$\therefore \|x\|_2 \leq \|x\|_1 \leq \sqrt{n} \|x\|_2$$

## Question 4: Eigenvalues

Part (a)

$\underline{x} = \alpha_1 \underline{v}_1 + \dots + \alpha_n \underline{v}_n$  since eigenvalues span the vector space

$$\underline{x}^T A \underline{x} = (\alpha_1 \underline{v}_1^T + \dots + \alpha_n \underline{v}_n^T) A (\alpha_1 \underline{v}_1 + \dots + \alpha_n \underline{v}_n)$$

$$= \alpha_1^2 \underline{v}_1^T A \underline{v}_1 + \dots + \alpha_n^2 \underline{v}_n^T A \underline{v}_n$$

$$= \alpha_1^2 \underline{v}_1^T (\lambda_1 \underline{v}_1) + \dots + \alpha_n^2 \underline{v}_n^T (\lambda_n \underline{v}_n) \quad \text{since } A \underline{v} = \lambda \underline{v}$$

$$= \alpha_1^2 \lambda_1 \|\underline{v}_1\|_2 + \dots + \alpha_n^2 \lambda_n \|\underline{v}_n\|_2$$

$$\max_{\|\underline{x}\|_2=1} \underline{x}^T A \underline{x} = \alpha_1^2 \lambda_1 + \dots + \alpha_n^2 \lambda_n$$

Since we are constrained by  $\|\underline{x}\|_2 = 1$ ,  $\sum_{i=1}^n |\alpha_i|^2 = 1$

clearly, to maximize  $\underline{x}^T A \underline{x}$  we should set  $\alpha_i = 1$  where  $\lambda_{\max}(A) = \lambda_i$  so that  $\underline{x} = \underline{v}_i$

$$\therefore \max_{\|\underline{x}\|_2=1} \underline{x}^T A \underline{x} = \lambda_{\max}(A)$$

Part (b)

We use a similar logic as in 4(a)

$$\min_{\|\underline{x}\|_2=1} \underline{x}^T A \underline{x} = \alpha_1^2 \lambda_1 + \dots + \alpha_n^2 \lambda_n > 0 \quad \text{since } \lambda \text{ and } \alpha_i^2 \text{ are non-negative (by 2(c))}$$

since we are constrained by  $\|\underline{x}\|_2 = 1$ ,  $\sum_{i=1}^n |\alpha_i|^2 = 1$

clearly, to minimize  $\underline{x}^T A \underline{x}$  we should set  $\alpha_i = 1$  where  $\lambda_{\min}(A) = \lambda_i$  so that  $\underline{x} = \underline{v}_i$

$$\therefore \min_{\|\underline{x}\|_2=1} \underline{x}^T A \underline{x} = \lambda_{\min}(A)$$

Part (c)

Both optimization problems in (a) and (b) aren't convex programs. Because their constraint  $\|\underline{x}\|_2 = 1$  isn't convex.

That is,  $\|\underline{x}\|_2 = 1$  creates a circle, which fails the 2 points test.

Part (d)

$$A = V \Lambda V^T$$

$$A^2 = (V \Lambda V^T)(V \Lambda V^T) = V \Lambda^2 V^T$$

Since  $\Lambda$  is a diagonal matrix of eigenvalues,  $\Lambda^2$  is the diagonal matrix of eigenvalues squared. Thus if  $\lambda$  is an eigenvalue of  $A$ , then  $\lambda^2$  is an eigenvalue of  $A^2$  as the diagonal of  $\Lambda^2$  are the eigenvalues of  $A^2$ .

$$\therefore \lambda_{\max}(A^2) = \lambda_{\max}(A)^2 \text{ and } \lambda_{\min}(A^2) = \lambda_{\min}(A)^2$$

Part (e)

$$\|A \underline{x}\|_2 = \sqrt{(A \underline{x})^T (A \underline{x})} = \sqrt{\underline{x}^T A^T A \underline{x}} = \sqrt{\underline{x}^T A^2 \underline{x}} = \sqrt{\underline{x}^T A \sqrt{A} \underline{x}}, \text{ which is in the form } \underline{y}^T A \underline{y} \text{ where } \underline{y} = \sqrt{A} \underline{x}$$

by 4(a) we can see that  $\max_{\|\underline{y}\|_2=1} \underline{y}^T A \underline{y} = \lambda_{\max}(A)$  thus  $\|A \underline{x}\|_2 \leq \lambda_{\max}(A)$

by 4(b) we can see that  $\min_{\|\underline{y}\|_2=1} \underline{y}^T A \underline{y} = \lambda_{\min}(A)$  thus  $\|A \underline{x}\|_2 \geq \lambda_{\min}(A)$

$$\therefore \lambda_{\min}(A) \leq \|A \underline{x}\|_2 \leq \lambda_{\max}(A)$$

Part (f)

$$\lambda_{\min}(A) \leq \|A \underline{y}\|_2 \leq \lambda_{\max}(A) \quad \text{from 4(e) where } \underline{y} = \frac{1}{\|\underline{x}\|_2} \underline{x}, \|\underline{x}\|_2 \neq 0 \text{ and } \|\underline{y}\|_2 = 1$$

$$\lambda_{\min}(A) \leq \frac{1}{\|\underline{x}\|_2} \|A \underline{x}\|_2 \leq \lambda_{\max}(A)$$

$$\lambda_{\min}(A) \|\underline{x}\|_2 \leq \|A \underline{x}\|_2 \leq \lambda_{\max}(A) \|\underline{x}\|_2$$

# Question 5: Gradient Descent

## Part (a)

the first-order optimality conditions  $\rightarrow$  setting the gradient of the objective function to 0 then solving for  $\underline{x}$

$$\nabla_{\underline{x}} \left( \frac{1}{2} \underline{x}^T A \underline{x} - \underline{b}^T \underline{x} \right) = \nabla_{\underline{x}}^2 \left( \frac{1}{2} \underline{x}^T A \underline{x} - \underline{b}^T \underline{x} \right) = A \geq 0 \text{ since its eigenvectors are non-negative } (0 < \lambda_{\min}(A) \text{ and } \lambda_{\max}(A) < 1)$$

$$\frac{1}{2} (2A \underline{x}^*) - \underline{b} = 0 \quad \text{by 2(a)}$$

$$\underline{x}^* = A^{-1} \underline{b}$$

Therefore the optimization problem is convex.

## Part (b)

$$\underline{x}_{i+1} \leftarrow \underline{x}_i - \eta \nabla_{\underline{x}} F(\underline{x}_i)$$

$$\underline{x}_{i+1} \leftarrow \underline{x}_i - \eta \nabla_{\underline{x}} \left( \frac{1}{2} \underline{x}^T A \underline{x} - \underline{b}^T \underline{x} \right)$$

$$\underline{x}_{i+1} \leftarrow \underline{x}_i - \eta (A \underline{x}_i - \underline{b})$$

$$\underline{x}_{i+1} \leftarrow (I - \eta A) \underline{x}_i + \eta \underline{b}$$

## Part (c)

$$\underline{x}^{(k)} - \underline{x}^* = (I - A) \underline{x}^{(k-1)} - \underline{b} - \underline{x}^*$$

Say  $\underline{x}^* = \underline{x}^{(\infty)}$  since we have a quadratic program as our loss function

$$\underline{x}^{(k)} - \underline{x}^{(\infty)} = (I - A) \underline{x}^{(k-1)} - \underline{b} - ((I - A) \underline{x}^{(\infty)} - \underline{b})$$

$$\underline{x}^{(k)} - \underline{x}^{(\infty)} = (I - A) \underline{x}^{(k-1)} + (I - A) \underline{x}^{(\infty)}$$

$$\underline{x}^{(k)} - \underline{x}^{(\infty)} = (I - A) (\underline{x}^{(k-1)} - \underline{x}^{(\infty)})$$

$$\underline{x}^{(k)} - \underline{x}^* = (I - A) (\underline{x}^{(k-1)} - \underline{x}^*)$$

## Part 5 (d)

$$A \underline{v} = \lambda \underline{v}$$

$$\underline{v} - A \underline{v} = \underline{v} - \lambda \underline{v}$$

$$(I - A) \underline{v} = (1 - \lambda) \underline{v} \text{ since we are given that } 0 < \lambda_{\min}(A) \text{ and } \lambda_{\max}(A) < 1, 1 - \lambda > 0$$

Since  $I$  is the identity matrix with a diagonal of ones  $\begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{bmatrix}$  and  $1 - \lambda_i > 0$ , then all the eigenvalues of  $I - A$  are positive and it's easy to see that  $I - A$  is also symmetric ( $I$  identity matrix consists of 1's on diagonal so it doesn't affect the symmetry along the diagonal).

$$\underline{x}^{(k)} - \underline{x}^* = (I - A) (\underline{x}^{(k-1)} - \underline{x}^*) \text{ from 5(c)}$$

$$\underline{x}^{(k)} - \underline{x}^* = B \underline{y}, \quad B = I - A, \quad \underline{y} = \underline{x}^{(k-1)} - \underline{x}^*$$

$$\|\underline{x}^{(k)} - \underline{x}^*\|_2 = \|B \underline{y}\|_2 \leq \lambda_{\max}(B) \|\underline{y}\|_2 \text{ by 4(f)}$$

$$\|\underline{x}^{(k)} - \underline{x}^*\|_2 \leq \lambda_{\max}(I - A) \|\underline{x}^{(k-1)} - \underline{x}^*\|_2$$

Since all  $\lambda_i$  of  $I - A$  are in the range  $(0, 1)$ , then

$$\|\underline{x}^{(k)} - \underline{x}^*\|_2 \leq \rho \|\underline{x}^{(k-1)} - \underline{x}^*\|_2, \quad 0 < \rho < 1$$

Part 5(e)

$$\|x^{(k)} - x^*\|_2 \leq \rho \|x^{(k-1)} - x^*\|_2 \quad \text{by 5(d)}$$

$$\|x^{(k)} - x^*\|_2 \leq \rho \|x^{(k-1)} - x^*\|_2 \leq \rho^2 \|x^{(k-2)} - x^*\|_2, \quad \text{using the recurrence relation in 5(d)}$$

$$\|x^{(k)} - x^*\|_2 \leq \rho^k \|x^{(0)} - x^*\|_2$$

Since we want to find the  $k$  where we are  $\varepsilon > 0$  close to  $x^*$

$$\|x^{(k)} - x^*\|_2 \leq \rho^k \|x^{(0)} - x^*\|_2 \leq \varepsilon$$

$$k \log \rho \leq \log \varepsilon - \log \|x^{(0)} - x^*\|_2$$

$$k \geq \frac{\log \varepsilon - \log \|x^{(0)} - x^*\|_2}{\log \rho}$$

since  $0 < \rho < 1$  so  $\log \rho < 0$

Part 5(f)

By 5(e) we know that it takes  $k \geq \frac{\log \varepsilon - \log \|x^{(0)} - x^*\|_2}{\log \rho}$  steps to get  $\varepsilon > 0$  close to  $x^*$

H takes  $(2n-1)n$  for matrix-vector multiplication. Thus the running time is  $k(2n-1)n = \frac{\log \varepsilon - \log \|x^{(0)} - x^*\|_2}{\log \rho} (2n-1)n$

$$\text{or } \frac{\log \varepsilon - \log \|x^{(0)} - x^*\|_2}{\log \rho} O(n^2)$$

# Question 6(a): Classification

Part (a)

$$f(x) = \begin{cases} i & \text{if } P(Y=i|x) \geq P(Y=j|x) \forall j \text{ and } P(Y=i|x) \geq 1 - \lambda_r/\lambda_s \\ c+1 & \end{cases}$$

case 1: if  $f(x) = c+1$

$$\begin{aligned} R(f(x)=i|x) &= \sum_{j=1}^c L(f(x)=i, y=j) P(Y=j|x) \\ &= \lambda_r \sum_{j=1}^c P(Y=j|x) \\ &= \lambda_r \end{aligned}$$

Now say we have some policy  $g: \mathbb{R}^d \rightarrow \{1, \dots, c+1\}$ . We will consider when  $g(x) \in \{1, \dots, c\}$

$$\begin{aligned} R(g(x)=i|x) &= \sum_{j=1}^c L(g(x)=i, y=j) P(Y=j|x) \\ &= \left[ \sum_{j=1, j \neq i}^c L(g(x)=i, y=j) P(Y=j|x) \right] + L(g(x)=i, y=i) P(Y=i|x) \\ &= \sum_{j=1, j \neq i}^c \lambda_s P(Y=j|x) + 0 P(Y=i|x) \\ &= \lambda_s (1 - P(Y=i|x)) \end{aligned}$$

We can rewrite this as  $P(Y=i|x) = 1 - \frac{1}{\lambda_s} R(g(x)=i|x)$

Since our policy  $f$  chose  $c+1$ , we know that  $P(Y=i|x) < 1 - \lambda_r/\lambda_s$  because there must be a greatest subset s.t.  $\exists i \in \{1, \dots, c\} P(Y=i|x) \geq P(Y=j|x) \forall j$ .

$$1 - \frac{1}{\lambda_s} R(g(x)=i|x) < 1 - \lambda_r/\lambda_s$$

$$R(g(x)=i|x) < \lambda_r$$

Thus, for some policy  $g$  to classify  $i \in \{1, \dots, c\}$ ,  $R(g(x)=i|x) < \lambda_r$  and therefore, to classify with the doubt class  $c+1$ ,  $R(g(x)=i|x) \geq \lambda_r$ .

$\therefore f$  obtains minimal risk for  $c+1$



case 2: if  $f(x) \in \{1, \dots, c\}$

$$\begin{aligned} R(f(x)=i|x) &= \sum_{j=1}^c L(f(x)=i, y=j) P(Y=j|x) \\ &= \left[ \sum_{j=1, j \neq i}^c L(f(x)=i, y=j) P(Y=j|x) \right] + L(f(x)=i, y=i) P(Y=i|x) \\ &= \sum_{j=1, j \neq i}^c \lambda_s P(Y=j|x) + 0 \\ &= \lambda_s (1 - P(Y=i|x)) \end{aligned}$$

We can rewrite this as  $P(Y=i|x) = 1 - \frac{1}{\lambda_s} R(f(x)=i|x)$

since our policy  $f$  chose  $i$ , we know that  $P(Y=i|x) \geq 1 - \lambda_r/\lambda_s$

$$1 - \frac{1}{\lambda_s} R(f(x)=i|x) \geq 1 - \lambda_r/\lambda_s$$

$$R(f(x)=i|x) < \lambda_r$$

Now say we have some policy  $g: \mathbb{R}^d \rightarrow \{1, \dots, c\}$  and check when  $g(x) = c+1$

$$\begin{aligned} R(g(x)=i|x) &= \sum_{j=1}^c L(g(x)=i, y=j) P(Y=j|x) \\ &= \lambda_r \sum_{j=1}^c P(Y=j|x) \\ &= \lambda_r \end{aligned}$$

Thus we can see that to classify as  $c+1$ ,  $R(g(x)=i|x) \geq \lambda_r$ . Thus, our policy  $f$  is minimal in classifying  $i \in \{1, \dots, c\}$

We conclude then that  $f$  obtains minimum risk

Part (b)

If  $\lambda_r = 0$ , then it is beneficial to choose the doubt class  $c+1$  since there would be no loss for choosing doubt, just like for choosing the correct class. If  $\lambda_r > \lambda_s$  then we would never classify doubt because we would have a higher loss for choosing doubt than choosing an incorrect class. That is, we would always classify and never choose doubt. This is consistent with intuition since we have shown in 6(a) that to choose doubt  $R(g(x)=i|x) \geq \lambda_r$

## Question 7: Gaussian Classification

Part (a)

$P(w_1|x) = P(w_2|x)$  definition of Bayes' optimal decision boundary

$$\frac{P(x|w_1)P(w_1)}{P(x)} = \frac{P(x|w_2)P(w_2)}{P(x)}$$

$$\frac{\frac{1}{2}P(x|w_1)}{P(x)} = \frac{\frac{1}{2}P(x|w_2)}{P(x)}$$

$$\frac{1}{\sqrt{2\sigma^2}\pi} e^{-\frac{(x-\mu_1)^2}{2\sigma^2}} = \frac{1}{\sqrt{2\sigma^2}\pi} e^{-\frac{(x-\mu_2)^2}{2\sigma^2}}$$

$$(x-\mu_1)^2 = (x-\mu_2)^2$$

$$x^2 - 2\mu_1 x + \mu_1^2 = x^2 - 2\mu_2 x + \mu_2^2$$

$$2x(\mu_2 - \mu_1) = \mu_2^2 - \mu_1^2$$

$$x_b = \frac{\mu_2^2 - \mu_1^2}{2(\mu_2 - \mu_1)} = \frac{\mu_2 + \mu_1}{2} \quad \text{decision boundary}$$

decision rule: choose class  $w_2$  if  $x > x_b$  else  $w_1$ , since  $\mu_2 > \mu_1$ .

Part (b)

$$P_e = P(w_1|w_2)P(w_2) + P(w_2|w_1)P(w_1)$$

$$P(w_1|w_2) = P\left(x < \frac{\mu_1 + \mu_2}{2} \mid x \sim N(\mu_2, \sigma^2)\right) = P\left(\frac{x - \mu_2}{\sigma} < \frac{\mu_1 - \mu_2}{2\sigma}\right)$$

$$P(w_2|w_1) = P\left(x \geq \frac{\mu_1 + \mu_2}{2} \mid x \sim N(\mu_1, \sigma^2)\right) = 1 - P\left(\frac{x - \mu_1}{\sigma} < \frac{\mu_2 - \mu_1}{2\sigma}\right) = 1 - \Phi\left(\frac{\mu_2 - \mu_1}{2\sigma}\right) = P(w_1|w_2)$$

$$P_e = \frac{1}{2}P(w_1|w_2) + \frac{1}{2}P(w_2|w_1) = P(w_1|w_2) = 1 - \Phi\left(\frac{\mu_2 - \mu_1}{2\sigma}\right)$$

$$= 1 - \frac{1}{\sqrt{2\pi}} \int_{-\infty}^a e^{-z^2/2} dz = \frac{1}{\sqrt{2\pi}} \int_a^{\infty} e^{-z^2/2} dz, \quad a = \frac{\mu_2 - \mu_1}{2\sigma}$$

### Question 8: Maximum Likelihood Estimation

Multinomial Distribution:  $P(k_1, k_2, k_3) = \frac{3!}{\prod_{i=1}^3 (k_i!)} \prod_{i=1}^3 p_i^{k_i}$

$$\ell(p_1, p_2, p_3) = \ln(P(k_1, k_2, k_3)) = \ln(3!) - \sum_{i=1}^3 \ln(k_i!) + \sum_{i=1}^3 k_i \ln p_i$$

$$\mathcal{L}(p_1, p_2, p_3, \lambda) = \ell(p_1, p_2, p_3) + \lambda(1 - \sum_{i=1}^3 p_i)$$

$$\begin{cases} \frac{\partial \mathcal{L}}{\partial p_1} = 0 \rightarrow \frac{k_1}{p_1} - \lambda = 0 \rightarrow p_1 = \frac{k_1}{\lambda} \\ \frac{\partial \mathcal{L}}{\partial p_2} = 0 \rightarrow \frac{k_2}{p_2} - \lambda = 0 \rightarrow p_2 = \frac{k_2}{\lambda} \\ \frac{\partial \mathcal{L}}{\partial p_3} = 0 \rightarrow \frac{k_3}{p_3} - \lambda = 0 \rightarrow p_3 = \frac{k_3}{\lambda} \\ \frac{\partial \mathcal{L}}{\partial \lambda} = 0 \rightarrow 1 - (p_1 + p_2 + p_3) = 0 \rightarrow p_1 + p_2 + p_3 = 1 \end{cases}$$

$$\frac{k_1}{\lambda} + \frac{k_2}{\lambda} + \frac{k_3}{\lambda} = 1$$

$\lambda = k_1 + k_2 + k_3 = n$  since  $k_i$  are counts and we are summing counts from all classes  $\{1, 2, 3\}$

$$\therefore p_1 = \frac{k_1}{n}, \quad p_2 = \frac{k_2}{n}, \quad p_3 = \frac{k_3}{n}$$

$$H = \nabla_p^2 = \begin{bmatrix} -k_1/p_1^2 & 0 & 0 \\ 0 & -k_2/p_2^2 & 0 \\ 0 & 0 & -k_3/p_3^2 \end{bmatrix}$$

for  $\underline{x} \in \mathbb{R}^3$

$$\underline{x}^T H \underline{x} = \left[ -\frac{k_1}{p_1^2} x_1, -\frac{k_2}{p_2^2} x_2, -\frac{k_3}{p_3^2} x_3 \right] \underline{x}$$

$$\underline{x}^T H \underline{x} = \left[ -\frac{k_1}{p_1^2} x_1^2, -\frac{k_2}{p_2^2} x_2^2, -\frac{k_3}{p_3^2} x_3^2 \right]$$

$\forall \underline{x} \in \mathbb{R}^3 - \{0\}$ ,  $\underline{x}^T H \underline{x} < 0$  since  $k_i > 0$ , probabilities are positive, and we square  $x_i$ . Therefore, the Hessian is negative definite and  $\mathcal{L}$  is concave.