

$$1a) i. (.4)(.3) = .12$$

$$ii. P(T) = (.4)(.3) + (.7)(.7) = .61$$

$$iii. P(A_1)P(A_2^c) + P(A_1^c)P(A_2) = (.61)(.39) + (.39)(.39) = .4758$$

$$iv. P(W^c|T^c) = \frac{(.3)(.7)}{.39} = .538.$$

$$1b) P(A|B,C) > P(A|B)$$

$$P(A,B,C) / P(B,C) > P(A|B)$$

$$P(C|A,B)P(A,B) / P(B,C) > P(A,B) / P(B)$$

$$P(C|A,B) > P(B,C) / P(B)$$

$$1 - P(C^c|A,B) > P(B,C) / P(B)$$

$$-P(C^c|A,B) > P(B,C) / P(B) - 1$$

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

$$\frac{P(C^c,A,B)}{P(A,B)} < 1 - \frac{P(B,C)}{P(B)}$$

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

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

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

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

$$P(A|B,C^c) < P(A|B)$$

2a. $\boxed{ii \Rightarrow i}$ We have $B^T A B \geq 0$, and can transform it s.t:
 $x^T B^T A B x \geq 0$ } By definition of positive semidefinite.

$$\text{let } y = Bx, \text{ and } y^T = (Bx)^T = x^T B^T$$

From this we can reformulate our equation as:

$$y^T y \geq 0$$

which follows the form of (i). Therefore we have shown
that $ii \Rightarrow i$

2a) i. \Leftrightarrow ii We can take the transpose of (ii) such that:

$$(B^T A B)^T = (A B)^T (B)^T = B^T A B \rightarrow (B^T A B)^T = B^T A B.$$

Because (i) tells us that A is symmetric and positive semi-definite, we can say that because $B^T A B$ is symmetric for some invertible $B \in \mathbb{R}^{n \times n}$, it must also be positive semi-definite.

iii \Leftrightarrow iv. $A = P D P^T = P \sqrt{D} \sqrt{D} P^T$ (valid b/c of iii's implications)

$$\text{let } U = P \sqrt{D} \text{ and } U^T = (P \sqrt{D})^T = \sqrt{D}^T P^T = \sqrt{D} P^T \text{ (valid b/c } D \text{ is diagonal)}$$

$$\therefore A = U U^T$$

i \Leftrightarrow iii The definition of an eigenvalue and eigenvector is as follows:

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

The definition of SPD is: $\vec{x}^T A \vec{x} \geq 0 \quad \forall \vec{x} \in \mathbb{R}^n$. If \vec{x} is then an eigenvector of A , it follows that:

$$\vec{x}^T \vec{x} \lambda \geq 0$$

Because $\vec{x}^T \vec{x}$ is positive, in order for $\vec{x}^T A \vec{x}$ to be ≥ 0 , we know that λ must be ≥ 0 . (We know $\vec{x}^T \vec{x}$ is positive because it is the sum of squares)

iv. \Leftrightarrow i Because U is symmetric $\in \mathbb{R}^{n \times n}$, it holds that:

$$A = U U^T \Rightarrow \vec{x}^T A \vec{x} = \vec{x}^T U U^T \vec{x} \Rightarrow (U^T \vec{x})^T (U^T \vec{x}) = \|U^T \vec{x}\|_2^2 \geq 0$$

Therefore $U U^T$ is positive semi-definite.

2b) (i) We know that $\vec{x}^T A \vec{x} \geq 0$ by the definition of symm. positive definite, $\vec{x}^T \vec{x} > 0$ because it is composed of the sum of squared, and $\lambda > 0$ as given in the problem. From this we can also ascertain that $A > 0$.

$$\vec{x}^T A \vec{x} + \vec{x}^T \lambda \vec{x} > 0$$

$$\vec{x}^T (A + \lambda I) \vec{x} > 0$$

$\vec{x}^T (A + \lambda I) \vec{x} > 0$ } holds b/c we know λ is positive

$$\therefore A + \lambda I > 0$$

$$(ii) x^T A x - x^T \gamma x > 0$$

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

$$A - \gamma I > 0$$

$$A > \gamma I \text{ or } A > \gamma I.$$

Therefore for $A - \gamma I > 0$, $A > \gamma I$ where $\gamma \neq 0$.

(iii) Say there is a vector e_i ; that is all zeros except for a 1 in the i^{th} position. By the definition of A being SPD, it holds that $e_i^T A e_i > 0$ for all i . If the i^{th} entry of A wasn't positive, then $e_i^T A e_i = 0 + \dots + a_{ii} + \dots + 0$, which means that $e_i^T A e_i$ would be ≤ 0 , a direct violation of the definition of SPD.

Therefore in order for $A > 0 \in \mathbb{R}^{n \times n}$, $a_{ii} > 0$ for $i = 1, \dots, n$.

(iv) Say our vector x is $[1]$. By the definition of SPD, we know that $x^T A x > 0$. With our definition of x , we then have:

$$[1] A [1] > 0$$

Because x is composed of 1's, the sum of A_{ij} for all $i, j = 1, \dots, n$ must be > 0 . We see this in the following example:

$$[1 \ 1] \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} [1] = [a_{11} + a_{12} \ a_{21} + a_{22}] [1] > 0$$

$$\begin{bmatrix} a_{11} + a_{12} \\ a_{21} + a_{22} \end{bmatrix} > 0, \text{ where } a_{11}, a_{22} \text{ are component of } A \in \mathbb{R}^{2 \times 2}$$

In the 2×2 matrix where $i=2$ and $j=2$ we see that $\sum_{i=1}^2 \sum_{j=1}^2 A_{ij} > 0$,

therefore this expands to the greater case where $i=n$ and $j=n$, such that $\sum_{i=1}^n \sum_{j=1}^n A_{ij} > 0$.

3.

a) $f(x) = \bar{a}^T x$

$$\nabla_x f(x) = a$$

b) $f(x) = x^T A x = \sum_{i=1}^n \sum_{j=1}^n x_i a_{ij} x_j$
 $= \sum_{i=1}^n a_{ii} x_i^2 + \sum_{i \neq j} \sum_{j=1}^n a_{ij} x_i x_j$
 $\frac{\partial f(x)}{x_k} = 2 \sum_{i=1}^n a_{ki} x_i \quad \text{where } x_k \text{ for } k \in \{1, 2, \dots, n\}$

Therefore $\nabla f(x) = 2Ax$

c) $\text{tr}(A^T x) = \sum_{j=1}^n \sum_{i=1}^n a_{ij} x_{ij}$

$$\frac{\partial (\sum_{j=1}^n \sum_{i=1}^n a_{ij} x_{ij})}{\partial x_{ij}} = a_{ij} \Rightarrow \nabla_x \text{tr}(A^T x) = A$$

d) $(\sqrt{|x_1 - y_1|} + \sqrt{|x_2 - y_2|})^2 + (\sqrt{|x_1|} + \sqrt{|y_1|})^2 \leq (\sqrt{|y_1|} + \sqrt{|y_2|})^2$

Counterexample: $x = \begin{bmatrix} 1 \\ 3 \end{bmatrix} \quad y = \begin{bmatrix} 4 \\ 7 \end{bmatrix}$

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

$$\begin{aligned}\|x\|_{\infty} &= \max_{1 \leq i \leq n} |x_i| \leq \sqrt{\max_{1 \leq i \leq n} (x_i)^2 + \sum_{j=1, j \neq i}^n x_j^2} \\ &= \sqrt{\sum_{i=1}^n x_i^2} \\ &= \|x\|_2\end{aligned}$$

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

$$\begin{aligned}\|x\|_2 &= \sqrt{\sum_{j=1}^n x_j^2} \leq \sqrt{n \cdot \max_{1 \leq i \leq n} (x_i^2)} \quad \left. \begin{array}{l} \text{valid because } x_i \text{ is} \\ \text{greatest from } x_i \rightarrow n, \\ \text{therefore } n \cdot x_i \geq x_1 + \dots + x_j \end{array} \right\} \\ &= \sqrt{n} \|x\|_{\infty}\end{aligned}$$

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

Therefore we have shown that $\|x\|_{\infty} \leq \|x\|_2 \leq \sqrt{n} \|x\|_{\infty}$

$$\begin{aligned}f) \|x\|_2 &= \sqrt{x_1^2 + x_2^2 + \dots + x_n^2} \leq \sqrt{x_1^2 + \dots + x_n^2 +} \\ &\quad = \sqrt{(|x_1| + |x_2| + \dots + |x_n|)^2} \\ &= |x_1| + |x_2| + \dots + |x_n| \\ &= \|x\|_1\end{aligned}$$

$$\therefore \|x\|_2 \leq \|x\|_1$$

$$\|x\|_1 = |x_1| + \dots + |x_n| \quad \left. \begin{array}{l} \text{let this be vector } x \\ \text{=} x^T * 1 \end{array} \right\}$$

$$\|x\|_1 \leq \|x\|_2 \|1\|_2 \quad \left. \begin{array}{l} \text{By Cauchy-Schwarz, where } \|x\|_2 = \|x\|_1 \\ \text{and } \|1\|_2 = \sqrt{n} \end{array} \right.$$

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

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

Therefore we have shown that $\|x\|_2 \leq \|x\|_1 \leq \sqrt{n} \|x\|_2$

4

a) Because A is symmetric, we can use the Spectral Thm. such that:

$$A = UDU^T \Rightarrow A = U\lambda U^T \quad \text{where } U \text{ and } U^T \text{ are orthogonal}$$

We can apply this to our original question:

$$\begin{aligned} \lambda_{\max}(A) &= \max_{\|x\|_2=1} x^T U \lambda U^T x \quad \text{let } Y = x^T U, Y^T = U^T x \\ &= \max_{\|Y\|_2=1} Y \lambda Y^T \\ &= \lambda_{\max} Y_1^2 + \lambda Y_2^2 + \dots + \lambda_{\min} Y_n^2 \leq \lambda_{\max} (Y_1^2 + \dots + Y_n^2) \\ &= \max_{\|Y\|_2=1} Y \lambda_{\max} Y^T \end{aligned}$$

From here we know that $YY^T = I$ because $\|x\|_2 = 1$, and if we reduce A to the diagonal case we get the following:

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

b) Using spectral theorem in our original equation, we get:

$$\begin{aligned} \lambda_{\min}(A) &= \min_{\|x\|_2=1} x^T U \lambda U^T x \quad \text{let } Y = x^T U, Y^T = U^T x \\ &= \min_{\|Y\|_2=1} Y \lambda Y^T \\ &= \lambda_{\min} Y_1^2 + \dots + \lambda_{\max} Y_n^2 \geq \lambda_{\min} (Y_1^2 + \dots + Y_n^2) \\ &= \min_{\|Y\|_2=1} Y \lambda_{\min} Y^T \end{aligned}$$

From here we know that $YY^T = I$, because $\|x\|_2 = 1$, and if we reduce A to the diagonal case we get:

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

c) Neither are convex programs. Since A is positive semi-definite, the objective fn is convex, proving that the maximization problem is not convex. The minimization is also not convex because there exists no linear combination of vectors in the feasible set that satisfies $\|x\|_2 = 1$.

$$d) Ax = \lambda x \Rightarrow A^2 x = A\lambda x \Rightarrow A^2 x = \lambda^2 x$$

From this we see that λ is an eigenvalue of A , and λ^2 is an eigenvalue of A^2 . From this we can use similar logic in part a to deduce that:

$$\lambda_{\max}(A^2) = \max_{\|x\|_2=1} x^T A^2 x = (\max_{\|x\|_2=1} x^T Ax)^2 = \lambda_{\max}(A)^2 \quad] \lambda \text{ are } > 0$$

We can make a similar argument to prove that:

$$\lambda_{\min}(A^2) = \lambda_{\min}(A)^2$$

$$e) \lambda_{\min}(A) \leq \|Ax\|_2 \leq \lambda_{\max}(A)$$

$$\sqrt{(\lambda_{\min}(A))^T Ax} \leq \sqrt{(Ax)^T Ax} \leq \sqrt{(\lambda_{\max}(A))^T Ax}$$

$$\sqrt{\lambda_{\min}(A)^T Ax} \leq \|Ax\|_2 \leq \sqrt{\lambda_{\max}(A)^T Ax}$$

$$\sqrt{\lambda_{\min}(A)^2} \leq \|Ax\|_2 \leq \sqrt{\lambda_{\max}(A)^2}$$

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

$$f) \lambda_{\min}(A) \|x\|_2 \leq \|x\|_2 \sqrt{A^2(x/\|x\|_2, \cdot x/\|x\|_2)} \leq \lambda_{\max}(A) \|x\|_2$$

$$\lambda_{\min}(A) \leq \sqrt{A^2(x/\|x\|_2, \cdot x/\|x\|_2)} \leq \lambda_{\max}(A) \quad] \|x\|_2 = 1$$

$$\lambda_{\min}(A) \leq \sqrt{(Ax)^T Ax} \leq \lambda_{\max}(A)$$

Because $\|x\|_2$ is the unit vector, we simplify our original equation to the one that we proved in 4e.

5

a)

Because A is symmetric and positive semi-definite, we can also infer that A is convex. This tells us that A must have a global minimum. We can find this minimum by taking the derivative of the optimization problem and setting it equal to 0.

$$f(x) = x^T A x - b^T x$$

$$\nabla f(x) = 2Ax - b \quad [\text{we sub this back into the original problem}]$$

$$\frac{1}{2}(2Ax) - b = 0 \Rightarrow Ax - b = 0 \Rightarrow Ax = b \Rightarrow x^* = b/A$$

Therefore we know that $x^* = b/A$.

$$\begin{aligned} b) \quad x &= x - \left(\frac{1}{2} \nabla x (x^T A x - b^T x) \right) \\ &= x - \left(\frac{1}{2} (2Ax) - b \right) \\ &= x - Ax + b. \end{aligned}$$

$$\begin{aligned} c) \quad x^{(k)} &= x^{(k-1)} - Ax^{(k-1)} + b \\ x^{(k)} - x^* &= x^{(k-1)} - Ax^{(k-1)} + b - b/A \quad [\text{sub out } b \text{ w/ } Ax^*] \\ &= x^{(k-1)} - Ax^{(k-1)} + Ax^* - b/A \\ &= (I - A)(x^{(k-1)} - x^*) \quad : \end{aligned}$$

$$d) \quad \|x^{(k)} - x^*\|_2 \leq p \|x^{(k-1)} - x^*\|_2$$

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

$$\|(I - A)\|_2 \|x^{(k-1)} - x^*\|_2 \leq p \|x^{(k-1)} - x^*\|_2 \quad [\text{holds by matrix-norm rules}]$$

$$\|I - A\|_2 \leq p$$

Because the $\lambda_{\min}(A) > 0$ and $\lambda_{\max}(A) < 1$, we know that $\|I - A\|_2$ will also be between 0 and 1.

From this we can deduce that $\exists p$ where $0 < p < 1$ s.t.

$$\|x^{(k)} - x^*\| \leq p \|x^{(k-1)} - x^*\|_2$$

$$e) \|x^{(k)} - x^*\| \leq p \|x^{(k-1)} - x^*\|_2 \leq p^k \|x^{(0)} - x^*\|_2$$

$$\epsilon = p^k \|x^{(0)} - x^*\|_2$$

$$k = \log_p \frac{\epsilon}{\|x^{(0)} - x^*\|_2}$$

$$f) n^2 \log_p \frac{\epsilon}{\|x^{(0)} - x^*\|_2}$$

This is because the matrix-vector multiplication of Ax takes up n^2 time.

6.

a) Suppose $\exists i, j$ s.t. $P(Y=i|x) \geq P(Y=j|x)$. From this we can tell that $P(Y=j|x) \leq 1/c$, given $j \neq i$ and $i \neq 1+c$. We can then apply this to the given risk fn s.t.

$$R \leq (c-1)(1/c)\lambda_s \Rightarrow R \leq (1-1/c)\lambda_s$$

On the other hand, $\forall j$ where $P(Y=i|x) < P(Y=j|x)$, $P(Y=j|x) > 1/c$. We can then apply this to the given risk fn s.t.:

$$R > (c-1)(1/c)\lambda_s \Rightarrow R > (1-1/c)\lambda_s$$

From this we can see that the policy outlined, where we select class i if $P(Y=i|x) \geq P(Y=j|x)$, is indeed the minimum, as R increases otherwise. For the second part of the policy we can apply the following transformation:

$$P(Y=i|x) \geq 1 - \lambda_r / \lambda_s$$

$$1 - P(Y=i|x) \geq -\lambda_r / \lambda_s$$

$$(1 - P(Y=i|x))\lambda_s \geq -\lambda_r$$

$$(1 - P(Y=i|x))\lambda_s \leq \lambda_r$$

This transformation shows that we face an equal or smaller amount of risk by guessing on a classification rather than choosing doubt.

b) IF $\lambda_r = 0$, you are better off choosing doubt rather than misclassifying, as reflected by the loss function. IF $\lambda_r > \lambda_s$, you are better off misclassifying a data point over choosing doubt - this is, once again, reflected by the risk and loss functions. We can see this reflected in the policy equations as well. For $\lambda_r = 0$, the policy becomes:

$$P(Y=i|x) \geq 1 - \delta/\lambda_s \Rightarrow P(Y=i|x) \geq 1$$

and $\lambda_r > \lambda_s$ results in:

$P(Y=i|x) \geq b$, where b is some negative number, because λ_r/λ_s will always be > 1 , so $1 - \lambda_r/\lambda_s$ is necessarily negative.

7.

$$a) P(w_1 | x) = P(w_2 | x)$$

$$P(x|w_1) P(w_1) = P(x|w_2) P(w_2)$$

$$N(\mu_1, \sigma^2) = N(\mu_2, \sigma^2) \rightarrow (x - \mu_1)^2 = (x - \mu_2)^2$$

$$x = \frac{\mu_1 + \mu_2}{2}$$

Therefore, if $x < \frac{\mu_1 + \mu_2}{2}$ select w_1 ,

else select w_2 .

$$b) P_e = \frac{1}{2} P\left(\frac{x - \mu_2}{\sigma^2} \leq \frac{\mu_1 - \mu_2}{\sigma^2} \mid x = w_2\right) + \frac{1}{2} P\left(\frac{x - \mu_1}{\sigma^2} \geq \frac{\mu_2 - \mu_1}{\sigma^2} \mid x = w_1\right)$$

$$= P\left(\frac{x - \mu_1}{\sigma^2} \geq \frac{\mu_2 - \mu_1}{\sigma^2} \mid x = w_1\right)$$

From here we can use the PDF to show the association

such that:

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

8. Using the method of maximum likelihood, we can estimate p_1 , p_2 , and p_3 as follows:

$$L(p) = \prod p_x^{k_x}$$

We use this equation because x is multinomially distributed. Because we only have n observations of x , our $L(p)$ has the constraint $\sum_x p_x = 1$, where $x \in \{1, 2, 3\}$. To maximize L , we take the gradient of L and our constraint, C and show that they are colinear such that:

$$\frac{\partial}{\partial p_x} L(p) = \lambda \frac{\partial}{\partial p_x} C(p) \quad \exists \lambda \forall x$$

After taking the derivative we get:

$$k_x / p_x L(p) = \lambda$$

From this we see that p_x and k_x are proportional.

The original problem states that $p_1 + p_2 + p_3 = 1$ (also stated in $C(p)$). With this, we know that

$$\hat{p}_x = k_x / (k_1 + k_2 + k_3) \text{. For each } p \text{ we then have:}$$

$$\hat{p}_1 = k_1 / (k_1 + k_2 + k_3)$$

$$\hat{p}_2 = k_2 / (k_1 + k_2 + k_3)$$

$$\hat{p}_3 = k_3 / (k_1 + k_2 + k_3)$$

Collaborators: Rushni Patel, Sydney McMurdoch, Lynn Kong
Problem 8 Citation: Math Stack Exchange # 421105.

1a) i. $(.4)(.3) = .12$

ii. $P(T) = (.4)(.3) + (.7)(.7) = .61$

iii. $P(A_1)P(A_2^c) + P(A_1^c)P(A_2) = (.61)(.39) + (.39)(.39) = .4758$

iv. $P(W^c|T^c) = \frac{(.3)(.7)}{.39} = .538$.

1b) $P(A|B, C) > P(A|B)$

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

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

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

$$1 - P(C^c|A, B) > \frac{P(B, C)}{P(B)} / 1$$

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

$$\frac{P(C^c, A, B)}{P(A, B)} < 1 - \frac{P(B, C)}{P(B)}$$

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

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

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

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

$$P(A|B, C^c) < P(A|B).$$