

# CS289–Spring 2017 — Homework 2 Solutions

Shuhui Huang, SID 3032129712

Collaborators:

## 1. Conditional Probability

(a) (i) Let  $X=1$ {an archer hits her target},  $Y=1$ {a gust of wind}.

$\therefore P(X=1|Y=1)=0.4$ ,  $P(X=1|Y=0)=0.7$ ,  $P(Y=1)=0.3$ .

$\therefore$  on a given shot there is a gust of wind and she hits her target:

$$P(X=1, Y=1)=P(X=1|Y=1) \cdot P(Y=1) = 0.3 \cdot 0.4 = 0.12$$

(ii) she hits the target with her first shot:

$$P(X=1)=P(X=1, Y=0)+P(X=1, Y=1)=P(X=1|Y=1) \cdot P(Y=1) + P(X=1|Y=0) \cdot P(Y=0) = 0.3 \cdot 0.4 + 0.7 \cdot 0.7 = 0.61$$

(iii) she hits the target exactly once in two shots:

$$P(\text{hits the target exactly once in two shots})=P(X_1=1, X_2=0) + P(X_1=0, X_2=1) = 2 \cdot 0.39 \cdot 0.61 = 0.4758$$

(iv) there was no gust of wind on an occasion when she missed:

$$P(Y=0|X=0)=\frac{P(X=0|Y=0) \cdot P(Y=0)}{P(X=0)} = \frac{0.3 \cdot 0.7}{0.39} = \frac{7}{13} = 0.53846$$

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

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

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

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

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

## 2.Positive Definiteness

(a) (1) first, prove (i)  $\implies$  (ii):

for a symmetric matrix A:

if  $A \succeq 0$ ,  $\forall x \in R^n, x^T A x \geq 0$

$\therefore$  for any x, we have  $x^T B^T A B x = (Bx)^T A (Bx)$

$\therefore$  according to definition, let  $x_{new} = (Bx)$ , we have  $x_{new} \in R^n$ , and  $x_{new}^T A x_{new} \geq 0$ , which means  $B^T A B \succeq 0$

Then, prove (ii)  $\implies$  (i):

if there exists invertible matrix  $B \in R^{n \times n}$ , such that  $B^T A B \succeq 0$

so for any  $x \in R^n, x^T B^T A B x \geq 0$

$\therefore$  for  $x_{new} = (Bx)$ , when B is invertible matrix,  $x \in R^n$ , we have  $x_{new} = (Bx) \in R^n$

$\therefore x_{new}^T A x_{new} \geq 0$ , for  $x_{new} \in R^n$

$\therefore A \succeq 0$

(2) Then, to prove (i)  $\implies$  (iii):

$\therefore x^T A x \geq 0$  for any x

$\therefore$  Consider v to be any eigenvector with  $Av = \lambda v$ . Then  $v^T A v = \lambda v^T v \geq 0$ .

$\therefore v^T v \geq 0$  (v is non-zero), we must have  $\lambda \geq 0$ .

(3) Then, to prove (iii)  $\implies$  (iv):

according to spectral theorem, we can obtain  $A = H^T D H$  with orthogonal H and diagonal matrix D with eigenvalues on the diagonal. (An orthogonal matrix U satisfies, by definition,  $H^T = H^{-1}$ )

Let  $E = D^{\frac{1}{2}} = E^T$  since all the eigenvalues of A are non-negative.

$\therefore$  we can get  $A = H^T E E H = (E H)^T E H$

Let  $U = E H$ , so there exists U, satisfy  $A = U^T U$

(4) Then, to prove (iv)  $\implies$  (i):

$\therefore$  there exists U such that  $A = U U^T$

$\therefore$  for any x, we always have  $x^T A x = x^T U U^T U x = (U x)^T U x \geq 0$  (since  $U x$  is non-zero)

$\therefore A \succeq 0$

(b) (i)  $\therefore x^T (A + \lambda I) x = x^T A x + x^T \lambda I x = x^T A x + \lambda \|x\|^2 > 0$

$\therefore$  we have  $A + \lambda I \succ 0$

(ii)  $\therefore x^T (A - \gamma I) x = x^T A x - x^T \gamma I x = x^T A x - \gamma \|x\|^2$

and there must exist  $\epsilon > 0$  such that  $x^T A x - \epsilon > 0$

so let  $\gamma = \frac{\epsilon}{\|x\|^2}$ , we have  $x^T (A - \gamma I) x = x^T A x - x^T \gamma I x = x^T A x - \gamma \|x\|^2 > 0$

$\therefore$  There exists a  $\gamma > 0$  such that  $A - \gamma I \succ 0$

(iii) Let  $e_1 = (1, 0, 0, \dots, 0)$  and so on, where  $e_i$  is a vector of all zeros, except for a 1 in the  $i$ th place.

Since A is positive definite, then  $x^T A x > 0$  for any non-zero vector  $x \in R^n$

so we always have  $e_i A e_i^T > 0$ , which means  $a_{ii} > 0$

So all the diagonal entries of A are positive.

(iv) Let  $x^T = (1, 1, 1, 1, 1 \cdots 1)_n$

We have  $x^T A x = \sum_{i=1}^n \sum_{j=1}^n A_{ij} > 0$ , since for any  $x \in R^n$ ,  $x^T A x > 0$

### 3. Derivatives and Norms

- (a) Let  $X = (x_1, x_2, \dots, x_n)^T$ , and  $a = (a_1, a_2, \dots, a_n)^T$

$$\text{Then } X^T a = \sum_{i=1}^n x_i a_i$$

$$\therefore \frac{\partial(X^T a)}{\partial X} = \left( \frac{\partial(X^T a)}{\partial x_1}, \frac{\partial(X^T a)}{\partial x_2}, \dots, \frac{\partial(X^T a)}{\partial x_n} \right)^T = (a_1, a_2, \dots, a_n)^T = a$$

- (b) Let  $A = [a_{ij}]_{n \times n}$ , so we have:

$$x^T A x = \sum_{i,j=1}^n a_{ij} x_i x_j$$

so for any k, we have :

$$\frac{\partial x^T A x}{\partial x_k} = \frac{\partial}{\partial x_k} (x_1 \sum_{j=1}^n a_{1j} x_j + x_2 \sum_{j=1}^n a_{2j} x_j + \dots + x_k \sum_{j=1}^n a_{kj} x_j + \dots + x_n \sum_{j=1}^n a_{nj} x_j)$$

$$= (x_1 a_{1k} + x_2 a_{2k} \dots + x_n a_{nk}) + (\sum_{j=1}^n a_{kj} x_j)$$

$$= (\sum_{i=1}^n a_{ik} x_i) + (\sum_{j=1}^n a_{kj} x_j)$$

$$\therefore \text{ we have: } \frac{\partial(x^T A x)}{\partial x} = (A + A^T)x$$

- (c) Let  $A = [a_{ij}]_{n \times n}$  and  $X = [x_{ij}]_{n \times n}$ , so we have:

$$\text{Trace}(XA) = \sum_{i=1}^n \sum_{j=1}^n x_{ij} a_{ji}$$

$\therefore$  we have:

$$\frac{\partial \text{Trace}(XA)}{\partial x_{ij}} = \frac{\partial}{\partial x_{ij}} (\sum_{i=1}^n \sum_{j=1}^n x_{ij} a_{ji}) = a_{ji}$$

$$\therefore \frac{\partial \text{Trace}(XA)}{\partial X} = [a_{ji}]_{n \times n} = A^T$$

- (d)  $f(x-y)$  should satisfy  $f(x-y) \geq f(x) - f(y)$

$$\text{Let } f(x) = (\sqrt{x_1} + \sqrt{x_2})^2, f(y) = (\sqrt{y_1} + \sqrt{y_2})^2, f(x-y) = (\sqrt{x_1 - y_1} + \sqrt{x_2 - y_2})^2$$

$$\therefore f(x) - f(y) = x_1 + x_2 - y_1 - y_2 + 2\sqrt{x_1 x_2} - 2\sqrt{y_1 y_2}$$

$$f(x-y) = x_1 + x_2 - y_1 - y_2 + 2\sqrt{(x_1 - y_1)(x_2 - y_2)}$$

$$\therefore (\sqrt{x_1 x_2} - \sqrt{y_1 y_2})^2 = x_1 x_2 + y_1 y_2 - 2\sqrt{(x_1 x_2 y_1 y_2)}, \text{ and } (\sqrt{(x_1 - y_1)(x_2 - y_2)})^2 = x_1 x_2 + y_1 y_2 - x_1 y_2 - x_2 y_1$$

$$\therefore x_1 y_2 + x_2 y_1 \geq 2\sqrt{x_1 x_2 y_1 y_2}$$

$$\therefore f(x-y) \leq f(x) - f(y) \text{ so } f(x) \text{ is not a norm. counterexample: } x_1 = 9, x_2 = 16, y_1 = 1, y_2 = 4$$

- (e)  $\therefore \|X\|_\infty = \max\{|x_1|, |x_2|, \dots, |x_n|\}$ , and  $\|X\|_2 = \sqrt{\sum_{i=1}^n x_i^2} \geq \sqrt{\max\{|x_1|, \dots, |x_n|\}^2}$

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

$$\text{and because } \|x\|_2 \leq \sqrt{n \cdot \max\{|x_1|, \dots, |x_n|\}^2} = \sqrt{n} \|x\|_\infty$$

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

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

$$= |x_1| + |x_2| + \dots + |x_n| = \|x\|_1 \therefore \|x\|_2 \leq \|x\|_1$$

$$\text{and } \therefore \|x\|_1 = (|x_1|, |x_2|, \dots, |x_n|) \cdot (1, 1, 1, \dots, 1)^T \leq \sqrt{\sum_{i=1}^n x_i^2} \cdot \sqrt{n} \text{ (according to Cauchy-Schwarz inequality)}$$

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

## 4. Eigenvalues

- (a)  $\because A$  is a symmetric matrix with  $A \succeq 0$ , and use the spectral theorem for symmetric matrix, we can get  $A = H^T D H$  with orthogonal  $H$  and diagonal matrix  $D$  with eigenvalues on the diagonal. (An orthogonal matrix  $U$  satisfies, by definition,  $H^T = H^{-1}$ )  
 $\therefore \|H\|_2 = 1$ , and the largest eigenvalue should be the largest element of  $D$ , so  $\lambda_{\max}(A) = \max_{\|x\|_2=1} x^T A x$
- (b) Similarly,  $\because A$  is a symmetric matrix with  $A \succeq 0$ , and use the spectral theorem for symmetric matrix, we can get  $A = H^T D H$  with orthogonal  $H$  and diagonal matrix  $D$  with eigenvalues on the diagonal. (An orthogonal matrix  $U$  satisfies, by definition,  $H^T = H^{-1}$ )  
 $\therefore \|H\|_2 = 1$ , and the smallest eigenvalue of  $A$  should be the smallest element of  $D$ , so  $\lambda_{\min}(A) = \min_{\|x\|_2=1} x^T A x$
- (c) Yes, they are convex. because  $\frac{\partial^2 x^T A x}{\partial x^2} \geq 0$ , so they are convex program.
- (d) if  $\lambda$  is an eigenvalue of  $A$ , so there is a vector  $v$  such that  $\lambda v = A v$   
 $\therefore A^2 v = \lambda A v = \lambda^2 v$ , which means  $\lambda^2$  is an eigenvalue of  $A^2$   
 since we have already proved in problem 2 that  $\lambda$  is non-negative,  
 so we can deduce that  $\lambda_{\max}(A^2) = \lambda_{\max}(A)^2$  and  $\lambda_{\min}(A^2) = \lambda_{\min}(A)^2$
- (e)  $\lambda_{\min}(A) = \min_{\|x\|_2=1} x^T A x \leq \|A x\|_2 = \sqrt{\sum \lambda^2 x^2} \leq \lambda_{\max}(A) = \max_{\|x\|_2=1} x^T A x$
- (f)  $\because \|x\|_2 = 1$   
 $\therefore \lambda_{\min}(A) = \lambda_{\min}(A) \|x\|_2 \leq \|A x\|_2 \leq \lambda_{\max}(A) \|x\|_2 = \lambda_{\max}(A)$

## 5.Gradient Descent

- (a) Let  $\frac{\partial(\frac{1}{2}x^T Ax - b^T x)}{\partial x} = Ax - b = 0$   
 $\therefore x^* = A^{-1}b$
- (b) w :arbitrary nonzero starting point (good choice is any  $A^{-1}b$ )  
 $R(w) > 0$   
 $V$  : set of indices i for which  $(A_i w - b_i) \cdot w < 0$   
w:  $w + \sum_{i \in V} (A_i w - b_i)$  return w
- (c)  $\therefore x^{(k)} = x^{(k-1)} - Ax^{k-1} + b = x^{(k-1)} - Ax^{k-1} + Ax^*$   
 $\therefore x^{(k)} - x^* = (I - A)(x^{(k-1)} - x^*)$
- (d)  $\|x^{(k)} - x^*\|_2 = \|(I - A)(x^{(k-1)} - x^*)\|_2$   
 $\therefore A$  is a symmetric matrix with  $0 < \lambda_{\min}(A)$  and  $\lambda_{\max}(A) < 1$   
 $\therefore \|x^{(k)} - x^*\|_2 = \|(I - A)(x^{(k-1)} - x^*)\|_2 \leq \rho \|x^{(k-1)} - x^*\|_2$
- (e) according to do d , $\|x^{(0)} - x^*\|_2 * \rho^k = \|x^{(k)} - x^*\|_2 \leq \epsilon$   
 $\therefore k \geq \log_{\rho} \frac{\epsilon}{\|x^{(0)} - x^*\|_2}$
- (f) the running time should be  $T(n) = n \cdot \log_{\rho} \frac{\epsilon}{\|x^{(0)} - x^*\|_2}$

**6.**

- (a) Define  $\lambda_{ij}l(f(x) = i, y = j)$  The risk of classifying a new data point as class  $i$  is:

$$R(\alpha_i|x) = \sum_j \lambda_{ij} P(\omega_j|x) = \lambda_s(1 - P(\omega_i|x))$$

and the risk of classifying the new data point as doubt is:

$$R(\alpha_{c+1}|x) = \lambda_r \sum_j P(\omega_j|x) = \lambda_r$$

For choosing doubt to be better than choosing any of the classes, the ratio of the risks must satisfy:

$$\frac{R(\alpha_{c+1}|x)}{R(\alpha_i|x)} = \frac{\lambda_r}{\lambda_s(1-P(\omega_i|x))} < 1$$

$$\therefore P(\omega_i|x) < 1 - \frac{\lambda_r}{\lambda_s}$$

$\therefore$  any particular  $i$  for which  $P(\omega_i|x) \geq 1 - \frac{\lambda_r}{\lambda_s}$  should not be assigned doubt

- (b) If  $\lambda_r = 0$ , then doubt will always be assigned, since for all  $i$ ,  $P(\omega_i|x) \geq 1 - \lambda_r/\lambda_s = 1$  is not satisfied unless  $P(\omega_i|x) = 1$ .

If  $\lambda_r > \lambda_s$ , then doubt will never be assigned, since for all  $i$ ,  $P(\omega_i|x) \geq 0 > 1 - \lambda_r/\lambda_s$ .

## 7. Gaussian Classification

- (a)  $\because P(\omega_1|x) = P(\omega_2|x) \rightarrow P(x|\omega_1)P(\omega_1) = P(x|\omega_2)P(\omega_2) \rightarrow P(x|\omega_1) = P(x|\omega_2) \rightarrow N(\mu_1, \sigma^2) = N(\mu_2, \sigma^2) \rightarrow (x - \mu_1)^2 = (x - \mu_2)^2 \rightarrow x = \frac{\mu_1 + \mu_2}{2}$   
 $\therefore$  The decision rule is to select  $\omega_1$  if  $x < \frac{\mu_1 + \mu_2}{2}$  and otherwise  $\omega_2$ .
- (b)  $\because P_e = \frac{1}{2} \int_{-\infty}^{\frac{\mu_1 + \mu_2}{2}} N(\mu_2, \sigma^2) + \frac{1}{2} \int_{\frac{\mu_1 + \mu_2}{2}}^{\infty} N(\mu_1, \sigma^2) du = 1 - \phi(\frac{\mu_2 - \mu_1}{2\sigma})$ , where  $\phi \sim N(0, 1)$   
 $\therefore P_e = 1 - \phi(\frac{\mu_2 - \mu_1}{2\sigma}) = \int_a^{\infty} \frac{1}{\sqrt{2\pi}} e^{-z^2/2} dz$ , where  $a = \frac{\mu_2 - \mu_1}{2\sigma}$



## 8. Maximum Likelihood Estimation

$$\begin{aligned}
 P(x_1, x_2, \dots, x_n) &= \frac{n!}{k_1!k_2!k_3!} p_1^{k_1} p_2^{k_2} p_3^{k_3} \\
 \therefore l(x_1, x_2, \dots, x_n) &= \log(P(x_1, x_2, \dots, x_n)) = \text{constant} + k_1 \log(p_1) + k_2 \log(p_2) + k_3 \log(p_3) \\
 &= \text{constant} + k_1 \log(p_1) + k_2 \log(p_2) + (n - k_1 - k_2) \log(1 - p_1 - p_2) \\
 \text{since } k_1 + k_2 + k_3 &= n, \quad p_1 + p_2 + p_3 = 1 \\
 \therefore \text{let } \frac{\partial l(x_1, x_2, \dots, x_n)}{\partial p_1} &= \frac{k_1}{p_1} + \frac{k_1 + k_2 - n}{1 - p_1 - p_2} = 0, \quad \frac{\partial l(x_1, x_2, \dots, x_n)}{\partial p_2} = \frac{k_2}{p_2} + \frac{k_1 + k_2 - n}{1 - p_1 - p_2} = 0 \\
 \therefore p_1 &= \frac{k_1(1 - p_2)}{n - k_2}, \quad p_2 = \frac{k_2(1 - p_1)}{n - k_1} \\
 \therefore \text{we get } \hat{P}_1 &= \frac{k_1}{n}, \hat{P}_2 = \frac{k_2}{n}, \hat{P}_3 = \frac{k_3}{n}, \text{ which is the MLE of } p_1, p_2, p_3
 \end{aligned}$$