

247 HW1.

Tiwen Zhang. UID: 805899489

## 1. Linear Algebra.

(a). Q is a real orthogonal matrix.

(i). Since Q is orthogonal, then  $Q^T Q = Q Q^T = I$ .

For  $Q^T$ :  
|  $(Q^T)^T \cdot Q^T = Q \cdot Q^T = I$ , then  $Q^T$  is orthogonal.  
|  $Q^T \cdot (Q^T)^T = Q^T Q = I$

For  $Q^{-1}$ :  
|  $(Q^{-1})^T \cdot Q^{-1} = (Q^T)^{-1} \cdot Q^{-1} = (Q^T Q)^{-1} = I^{-1} = I$ , then  $Q^{-1}$  is orthogonal.  
|  $Q^T (Q^{-1})^T = Q^T (Q^T)^{-1} = (Q \cdot Q^T)^{-1} = I^{-1} = I$

(ii). For eigenvalue  $\lambda$ . eigenvector  $x$  for A: (Any random  $\lambda$  and  $x$ ).

$$Q \cdot x = \lambda x$$

$$\Rightarrow \text{Premultiply } Q^T: Q^T Q x = \lambda \cdot Q^T x$$

$$\Rightarrow I \cdot x = \lambda \cdot Q^T x$$

$L_2$  norm of both sides:  $\|x\|_2 = \|\lambda \cdot Q^T x\|_2 = \|\lambda\| \cdot \|Q^T x\|_2$  (Since  $\lambda$  is a scalar).

$$\Rightarrow \|\lambda\|^2 \cdot \|Q^T x\|_2^2 \text{ and } \|Q^T x\|_2^2 = (Q^T x)^T Q^T x = x^T Q Q^T x = x^T x = \|x\|_2^2.$$

$$\Rightarrow \|\lambda\|^2 = 1 \Rightarrow \|\lambda\| = 1 \Rightarrow \text{all eigenvalues of } A \text{ have norm of 1.}$$

(iii). Since  
|  $\det(Q^T Q) = \det(I) = 1$   
|  $\det(Q^T) \cdot \det(Q) = \det(Q)^2 = \det(Q^T Q)$

$$\Rightarrow \det(Q)^2 = 1 \Rightarrow \det(Q) = \pm 1$$

(iv). For  $\forall b \in \mathbb{R}^n$ ,  $\|A \cdot b\|_2^2 = (A \cdot b)^T \cdot A \cdot b = b^T \cdot A^T A \cdot b = b^T b = \|b\|_2^2$ .

$\Rightarrow$  After transformation of A, the transformed vector still has the same length as the original vector.

(b). (i). For  $A \in \mathbb{R}^{m \times n}$ , we can have SVD decomposition as:

$$A = U\Sigma V^T, \text{ where } U \in \mathbb{R}^{m \times m}, \Sigma \in \mathbb{R}^{m \times n}, V \in \mathbb{R}^{n \times n}$$

Left singular vectors of  $A$  = Columns of  $U$  = Eigenvectors (Orthonormal) of  $A A^T$   
 Right singular vectors of  $A$  = Columns of  $V$  = Eigenvectors (Orthonormal) of  $A^T A$ .

(ii).

$$\text{Let } A = U\Sigma V^T, \Sigma = \begin{bmatrix} S & 0 \\ 0 & 0 \end{bmatrix}, S = \text{diag}(a_1, \dots, a_r).$$

$$A^T = V\Sigma U^T.$$

$$\Rightarrow A A^T = U \Sigma V^T V \Sigma U^T = U \Sigma \Sigma U^T = \underbrace{U \Delta_1 U^T}_{\substack{\text{Orthogonal} \\ (\text{since } V \text{ is orthogonal})}} \underbrace{\Delta_1}_{\substack{\text{Eigenvalue} \\ \text{Decomposition}}} = \begin{bmatrix} S & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} S & 0 \\ 0 & 0 \end{bmatrix}^T = \begin{bmatrix} S^2 & 0 \\ 0 & 0 \end{bmatrix} = \text{diag}(a_1^2, \dots, a_r^2, 0)$$

$$A^T A = V \Sigma^T U^T U \Sigma V^T = V \Sigma^T \Sigma V^T = \underbrace{V \Delta_2 V^T}_{\substack{\text{Orthogonal} \\ (\text{since } U \text{ is orthogonal})}} \underbrace{\Delta_2}_{\substack{\text{Eigenvalue} \\ \text{Decomposition}}} = \Delta_2 = \text{diag}(a_1^2, \dots, a_r^2, 0, \dots, 0).$$

$$\Rightarrow \text{Singular value of } A = \sqrt{\text{Eigenvalue of } A A^T \text{ or } A^T A}$$

(c). (i). False.

Because eigenvalues can be equal to each other.

i.e. For  $I \in \mathbb{R}^{n \times n}$ ,  $I = U \cdot \Delta \cdot U^T$ , where  $\Delta = I$ .  $U$  is orthogonal.

Here  $I$  has  $n$  same eigenvalues: 1.

(ii). False.

$$\text{i.e. } A \cdot v_1 = \lambda_1 \cdot v_1. \text{ Let } v_3 = av_1 + bv_2.$$

$$A \cdot v_2 = \lambda_2 \cdot v_2. \quad A \cdot v_3 = A \cdot av_1 + A \cdot bv_2$$

$$= a\lambda_1 \cdot v_1 + b\lambda_2 \cdot v_2$$

If  $A v_3 = \lambda_3 v_3$ , then  $a(\lambda_1 - \lambda_3) v_1 + b(\lambda_2 - \lambda_3) v_2$ , In general, this doesn't always hold.

(iii). True.

$$\text{For } Av_1 = \lambda_1 v_1 \Rightarrow v_1^T A v_1 = \lambda_1 \cdot v_1^T v_1.$$

$$\text{Since } b^T A b \geq 0 \Rightarrow \lambda_1 \cdot v_1^T v_1 = \lambda_1 \|v_1\|_2^2 \geq 0, \text{ where } \|v_1\|_2^2 > 0$$
$$\Rightarrow \lambda_1 \geq 0.$$

(iv). False.

"r" represents the number of linear independent eigenvectors,  
so it couldn't be larger than the number of distinct eigenvalues.

(v). True.

$$\text{For } Av_1 = \lambda v_1, \text{ let } v_3 = av_1 + bv_2.$$

$$Av_2 = \lambda v_2$$

$$\begin{aligned} Av_3 &= A(a \cdot v_1 + b \cdot v_2) = a \cdot \lambda v_1 + b \cdot \lambda v_2 \\ &= \lambda(av_1 + bv_2) \\ &= \lambda v_3. \text{ Still an eigenvector.} \end{aligned}$$

## 2. Probability.

(a). Suppose we define following events:

H5: Get coins of H50

H6: Get coins of H60

D: Heads up

T: Tails up.

$$(i). P(H5|T) = \frac{P(T \cap H5)}{P(T)} = \frac{P(H5) \cdot P(T|H5)}{\sum P(T|Hi) \cdot P(Hi)} = \frac{\frac{1}{2} \times \frac{1}{2}}{\frac{1}{2} \times \frac{1}{2} + \frac{1}{2} \times \frac{2}{5}} = \frac{0.25}{0.25 + 0.2} = \frac{5}{9}$$

$$(ii) P(H5|THHH) = \frac{P(H5 \cap THHH)}{P(THHH)} \quad P(H5 \cap THHH) = P(H5) \cdot P(THHH|H5) = \frac{1}{32} = 0.03125$$
$$\Rightarrow P(H5|THHH) = \frac{0.03125}{0.07445} \approx 0.4197$$

$$\begin{aligned} P(THHH) &= P(THHH|H5) \cdot P(H5) + P(THHH|H6) \cdot P(H6) \quad \Rightarrow P(H5|THHH) = \frac{0.03125}{0.07445} \\ &= \frac{1}{2} \times (\frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2}) + \frac{1}{2} \times (\frac{2}{5} \times \frac{3}{5} \times \frac{3}{5} \times \frac{3}{5}) \\ &= \frac{1}{32} + \frac{27}{625} = 0.03125 + 0.432 = 0.07445 \end{aligned}$$

$$(iii). P(H50|H9T1) = \frac{P(H50 \cap H9T1)}{P(H9T1)} = \frac{P(H50) \cdot P(H9T1|H50)}{P(H9T1)}$$

$$P(H55|H9T1) = \frac{P(H55) \cdot P(H9T1|H55)}{P(H9T1)}$$

$$P(H60|H9T1) = \frac{P(H60) \cdot P(H9T1|H60)}{P(H9T1)}$$

$$\begin{aligned} P(H9T1) &= \sum P(H_i) \cdot P(H9T1|H_i) \\ &= P(H50) \cdot P(H9T1|H50) + P(H55) \cdot P(H9T1|H55) + P(H60) \cdot P(H9T1|H60) \\ &= \frac{1}{3} \times \left(\frac{1}{2}\right)^9 \times \left(\frac{1}{2}\right) + \frac{1}{3} \times (0.55)^9 \times 0.45 + \frac{1}{3} \times (0.6)^9 \times 0.4 \\ &= a + b + c. \text{ where } a = \frac{1}{3} \times \left(\frac{1}{2}\right)^9, b = \frac{1}{3} \times 0.55^9 \times 0.45, c = \frac{1}{3} \times 0.6^9 \times 0.4 \end{aligned}$$

$$\Rightarrow P(H50|H9T1) = \frac{a}{a+b+c}$$

$$P(H55|H9T1) = \frac{b}{a+b+c}$$

$$P(H60|H9T1) = \frac{c}{a+b+c}$$

(b). Suppose we define the following events:

S: student from Science.  $P(S) = 15\%$ .

H: student from Health care.  $P(H) = 21\%$ .

L: student from liberal arts.  $P(L) = 24\%$ .

E: student from Engineering.  $P(E) = 40\%$ .

A: student likes the lecture.

$$P(S|A) = \frac{P(S \cap A)}{P(A)}$$

$$\begin{aligned} P(A) &= \sum_{X \in \{S, H, L, E\}} P(A|X) \cdot P(X) = 15\% \times 90\% + 21\% \times 18\% + 24\% \times 0 + 40\% \times 10\% \\ &= 0.135 + 0.0378 + 0.04 \\ &= 0.2128 \end{aligned}$$

$$P(S \cap A) = P(S) \cdot P(A|S) = 15\% \times 90\% = 0.135.$$

$$\Rightarrow P(S|A) = \frac{0.135}{0.2128} \approx 0.634$$

(c). Create a table like this:

P	Preg	Nonpreg	$P(\text{preg}) = 1\%$
Positive	99%	1%	
Negative	1%	99%	

$$P(\text{Preg} | \text{pos}) = \frac{P(\text{preg} \cap \text{pos})}{P(\text{pos})} = \frac{P(\text{preg}) \times P(\text{pos} | \text{preg})}{P(\text{pos} | \text{preg}) \times P(\text{preg}) + P(\text{pos} | \text{unpreg}) \times P(\text{unpreg})}$$

$$= \frac{1\% \times 99\%}{1\% \times 99\% + 99\% \times 1\%} = \frac{1}{11}$$

It means the testing result couldn't define that people really get pregnant or not. because in all women, probability of pregnant is just 1%.

To get more precise answer, people need to do a second test.

$$(d). E(Ax+b) = E(Ax) + E(b)$$

$$= A \cdot E(x) + E(b) = A \cdot E(x) + b$$

If vector  $x$  has mean  $m$ ,  $m \in \mathbb{R}^n$ , then  $E(Ax+b) = A \cdot m + b \in \mathbb{R}^n$ .

$$(e). \text{cov}(x) = E((x - E(x))(x - E(x))^T), \text{ Let } y = Ax + b$$

$$\begin{aligned} \text{cov}(y) &= E((Ax+b - E(Ax+b))(y^T - E(y))^T) \\ &= E((Ax - A \cdot E(x))(x^T A^T + b^T - (A \cdot E(x) + b)^T)) \\ &= E((Ax - A \cdot E(x))(x^T A^T + b^T - E(x^T \cdot A^T - b^T))) \\ &= E(A(x - E(x)) \cdot (x^T - E(x^T) \cdot A^T)) \\ &= A \cdot E((x - E(x))(x - E(x))^T) \cdot A^T \\ &= A \cdot \text{cov}(x) \cdot A^T \end{aligned}$$

### 3. Multivariate derivatives.

(a)  $x \in \mathbb{R}^n$ ,  $y \in \mathbb{R}^m$ ,  $A \in \mathbb{R}^{n \times m}$ . Find  $\nabla_x x^T A y$ .

$$f(x) = x^T A y \in \mathbb{R}^{1 \times 1} \text{ is a scalar, } \nabla_x f(x) \in \mathbb{R}^n.$$

Since  $\frac{\partial x^T a}{\partial x} = \frac{\partial a^T x}{\partial x} = a$ , (referring to "matrixcook" and lecture notes).  
 then  $\nabla_x f(x) = \frac{\partial x^T (Ay)}{\partial x} = \frac{\partial (Ay)^T x}{\partial x} = A^T y \in \mathbb{R}^n$ .

(b). Find  $\nabla_y x^T A y$ .

$$\text{Since } \frac{\partial \theta^T x}{\partial x} = \theta, \text{ then } \nabla_y x^T A y = \frac{\partial (A^T x)^T y}{\partial y} = A^T x \in \mathbb{R}^{m \times 1}.$$

(c). Find  $\nabla_a x^T A y$ .

$$\text{Since } \frac{\partial a^T x b}{\partial x} = a b^T \text{ (referring to matrixcook book).}$$

$$\text{Then } \nabla_a x^T A y = x y^T \in \mathbb{R}^{n \times m}.$$

(d).  $x \in \mathbb{R}^n$ ,  $A \in \mathbb{R}^{n \times n}$ ,  $f = x^T A x + b^T x$ , find  $\nabla_x f(x)$ .

$f$  is a scalar,  $x \in \mathbb{R}^n$ .  $\nabla_x f(x) \in \mathbb{R}^n$ .

$$\begin{aligned}\nabla_x f(x) &= \nabla_x (x^T A x + b^T x) \\ &= \nabla_x (x^T A x) + \nabla_x (b^T x) \\ &= (A + A^T)x + b \in \mathbb{R}^n.\end{aligned}$$

(e).  $A, B \in \mathbb{R}^{n \times n}$ ,  $f = \text{tr}(AB)$ , find  $\nabla_A f$ .

$$\text{Since } \frac{\partial \text{Tr}(XA)}{\partial x} = A^T \text{ (referring to matrixcook).}$$

$f$  is scalar,  $A \in \mathbb{R}^{n \times n}$ , so  $\nabla_A f \in \mathbb{R}^{n \times n}$ ,

$$\nabla_A f = \nabla_A \text{Tr}(AB) = B^T.$$

(f).  $f = \text{tr}(BA + A^T B + A^2 B)$ , find  $\nabla_A f$ .

Referring to  $\frac{\partial \text{tr}(AXB)}{\partial x} = ATB^T$ ,  $\frac{\partial \text{tr}(x^TA)}{\partial x} = A$ ,  $\frac{\partial \text{tr}(x^2A)}{\partial x} = (xA + Ax)^T$ .

$$\nabla_A f = \nabla_A (\text{tr}(BAI + A^T B + A^2 B))$$

$$= \nabla_A (\text{tr}(B \cdot A \cdot I)) + \nabla_A (\text{tr}(A^T B)) + \nabla_A (\text{tr}(A^2 B))$$

$$= B^T \cdot I^T + B + (AB + BA)^T$$

$$= B^T + B + (AB + BA)^T$$

(g).  $f = \|A + \lambda B\|_F^2$ . find  $\nabla_A f$ .

Referring to  $\frac{\partial}{\partial x} \|x\|_F^2 = \frac{\partial}{\partial x} \text{Tr}(xx^H) = 2x$ .

then  $\frac{\partial f}{\partial A} = \frac{\partial}{\partial A} \text{Tr}(A + \lambda B)(A + \lambda B)^T$  (since A, B are both real matrix).

$$= \frac{\partial}{\partial A} \text{Tr}(AA^T + \lambda AB^T + \lambda BA^T + \lambda^2 BB^T)$$

$$= \frac{\partial}{\partial A} \text{Tr}(AA^T) + \frac{\partial}{\partial A} \text{Tr}(\lambda AB^T) + \frac{\partial}{\partial A} \text{Tr}(\lambda BA^T) + \frac{\partial}{\partial A} \text{Tr}(\lambda^2 BB^T)$$

$$= 2A + \lambda B + \lambda B + 0$$

$$= 2(A + \lambda B) \in \mathbb{R}^{n \times n}$$

#### 4. Derive least-square with matrix derivatives.

$$\min \frac{1}{2} \sum_{i=1}^n \|y^{(i)} - w \cdot x^{(i)}\|^2 \quad (*).$$

Let  $y^{(i)}, x^{(i)} \in \mathbb{R}^m$  for  $i \in [n]$ , then  $w \in \mathbb{R}^{m \times m}$

$$\|y^{(i)} - w \cdot x^{(i)}\|^2 = (y^{(i)} - w \cdot x^{(i)})^T (y^{(i)} - w \cdot x^{(i)})$$

$$\text{Make } \tilde{Y} = [y_1, y_2 \dots y_n], \quad \tilde{X} = [w x_1, w x_2 \dots w x_n].$$

$$\text{Then } \tilde{Y} - \tilde{X} = [y_1 - w x_1, y_2 - w x_2, \dots y_n - w x_n] \in \mathbb{R}^{n \times n}$$

$$(\tilde{Y} - \tilde{X})^T = \begin{bmatrix} (y_1 - w x_1)^T \\ \vdots \\ (y_n - w x_n)^T \end{bmatrix} \in \mathbb{R}^{n \times n}$$

$$\Rightarrow A = (\tilde{Y} - \tilde{X})^T \cdot (\tilde{Y} - \tilde{X}) = \begin{bmatrix} (y_1 - w x_1)^T (y_1 - w x_1) & \cdots & (y_1 - w x_1)^T (y_n - w x_n) \\ (y_2 - w x_2)^T (y_1 - w x_1) & \ddots & \vdots \\ \vdots & \cdots & (y_n - w x_n)^T (y_n - w x_n) \end{bmatrix} \in \mathbb{R}^{n \times n}$$

$$\text{Then } \min \frac{1}{2} \sum_i \|y^{(i)} - w x^{(i)}\|^2 = \min \frac{1}{2} \text{tr}(A).$$

$$\text{where } A = (Y - w X)^T (Y - w X) = Y^T Y - Y^T w X - X^T w^T Y + X^T w^T w X.$$

$$\begin{aligned} \frac{\partial \text{tr}(A)}{\partial w} &= \frac{\partial}{\partial w} (\text{tr}(Y^T Y) - \text{tr}(Y^T w X) - \text{tr}(X^T w^T Y) + \text{tr}(X^T w^T w X)) \\ &= 0 - Y X^T - Y X^T + \frac{\partial}{\partial w} \text{tr}(w X X^T w^T) \quad (\text{tr}(AB) = \text{tr}(B^T A)) \\ &= w (X X^T + X^T X) - 2 Y X^T \end{aligned}$$

$$\Rightarrow w = Y X^T \cdot (X X^T)^{-1}$$

#### 5. Regularized least squares.

$$L(\theta) = \underbrace{\frac{1}{2} \sum_i (y^{(i)} - \theta^T x^{(i)})^2}_{L_1(\theta)} + \underbrace{\frac{\lambda}{2} \|\theta\|_2^2}_{L_2(\theta)}, \text{ need to minimize } L(\theta).$$

$$\frac{\partial L(\theta)}{\partial \theta} = \frac{\partial L_1(\theta)}{\partial \theta} + \frac{\partial L_2(\theta)}{\partial \theta}, \quad \frac{1}{2} \|\theta\|_2^2 = \frac{\lambda}{2} \theta^T \theta. \Rightarrow \frac{\partial}{\partial \theta} \left( \frac{\lambda}{2} \theta^T \theta \right) = \frac{\lambda}{2} \frac{\partial}{\partial \theta} (\theta^T \theta).$$

$$\text{let } S = \theta^T \theta = \sum_{i=1}^n \theta_i^2, \quad \frac{\partial S}{\partial \theta_1} = 2 \theta_1, \quad \dots \quad \frac{\partial S}{\partial \theta_n} = 2 \theta_n \Rightarrow \frac{\partial S}{\partial \theta} = \begin{bmatrix} \frac{\partial S}{\partial \theta_1} & \dots & \frac{\partial S}{\partial \theta_n} \end{bmatrix} = 2 \theta.$$

$$\Rightarrow \frac{\partial L(\theta)}{\partial \theta} = X^T X \theta - X^T Y + \frac{\lambda}{2} \cdot 2 \theta = (X^T X + \lambda I) \theta - X^T Y$$

$$\text{Let } \frac{\partial L(\theta)}{\partial \theta} = 0, \Rightarrow \theta = (X^T X + \lambda I)^{-1} \cdot X^T Y.$$