

## Exercise 2

### Question 1

Suppose  $A, B \in \mathbb{R}^{n \times n}$ ,  $i, j \in \{1, \dots, n\}$

(a)(1)  $(AB)_{ij} = \sum_{k=1}^n A_{ik} B_{kj} = \sum_{k=1}^n B_{jk}^T A_{ki}^T = (B^T A^T)_{ji}$ . Hence  $(AB)^T = B^T A^T$ .

(2)  $(A+B)_{ij}^T = (A+B)_{ji} = A_{ji} + B_{ji} = A_{ij}^T + B_{ij}^T$ . Hence  $(A+B)^T = A^T + B^T$ .

(3) Choose  $A = \begin{bmatrix} 1 & 2 \\ 0 & 3 \end{bmatrix}$ ,  $B = \begin{bmatrix} 2 & 3 \\ 1 & 2 \end{bmatrix}$ , then  $AB = \begin{bmatrix} 4 & 7 \\ 3 & 6 \end{bmatrix}$  with  $10 = \text{tr}(AB) \neq \text{tr}(A) \text{tr}(B) = 4 \cdot 4 = 16$ .

(4)  $\text{tr}(AB) = \sum_{i=1}^n \sum_{k=1}^n A_{ik} B_{ki} = \sum_{k=1}^n \sum_{i=1}^n B_{ki} A_{ik} = \text{tr}(BA)$

(5)  $\text{tr}(A+B) = \sum_{i=1}^n (A+B)_{ii} = \sum_{i=1}^n (A_{ii} + B_{ii}) = \sum_{i=1}^n A_{ii} + \sum_{i=1}^n B_{ii} = \text{tr} A + \text{tr} B$

(6)  $(AB)B^{-1}A^{-1} = I$  and  $B^{-1}A^{-1}(AB) = I$ . Hence  $(AB)^{-1} = B^{-1}A^{-1}$

(7) Choose  $A = \begin{bmatrix} 1 & 2 \\ 0 & 3 \end{bmatrix}$ ,  $B = \begin{bmatrix} 2 & 3 \\ 1 & 2 \end{bmatrix}$ , then  $A^{-1} = \frac{1}{3} \begin{bmatrix} 3 & -2 \\ 0 & 1 \end{bmatrix}$ ,  $B^{-1} = \begin{bmatrix} 2 & -3 \\ -1 & 2 \end{bmatrix}$ , while  $A+B = \begin{bmatrix} 3 & 5 \\ 1 & 5 \end{bmatrix}$ ,  $(A+B)^{-1} = \frac{1}{10} \begin{bmatrix} 5 & -5 \\ -1 & 3 \end{bmatrix} \neq A^{-1} + B^{-1}$

(b) Let  $u, v \in \ker A$ ,  $c \in \mathbb{R}^n$ , then  $A(u+v) = Au + Av = 0 + 0 = 0$ ,  $A(cv) = cAv = 0$ , and  $0 \in \ker A$ . Hence,  $\ker A$  is a subspace of  $\mathbb{R}^n$ . Let  $u, v \in \text{Im } A$  then  $\exists u', v' \in \mathbb{R}^n$  s.t.  $Au' = u$ ,  $Av' = v$ .  $A(u' + v') = Au' + Av' = u + v \in \text{Im } A$  (closure under addition) and  $A(cv') = cAv' = cv \in \text{Im } A$  (closure under scalar multiplication). Since  $A0 = 0 \in \text{Im } A$  we get that  $\text{Im } A$  is a subspace of  $\mathbb{R}^n$ .

(c) Suppose  $A \in \mathbb{R}^{n \times p}$  with  $T \in L(\mathbb{R}^p, \mathbb{R}^n)$  its linear map. Let  $v \in \mathbb{R}^n$  and  $u \in \text{null } T$ , then  $\langle T^*v, u \rangle = \langle v, Tu \rangle = \langle v, 0 \rangle = 0$ .

This shows that  $\text{Im}(T^*) \subseteq (\text{null } T)^\perp$ . But  $\dim(\text{null } T)^\perp = \text{rank } T = \dim \text{Im}(T^*)$ , hence  $\text{Im}(T^*) = (\text{null } T)^\perp$ .

### Question 2

(a) Suppose  $P_u \in \mathbb{R}^{n \times n}$  is an orthogonal projection matrix, projecting onto the subspace  $U \subseteq \mathbb{R}^n$ . Let  $u_1, \dots, u_m$  an orthonormal basis of  $U$  and extend it into an orthonormal basis  $u_1, \dots, u_m, \dots, u_n$  of  $\mathbb{R}^n$ .  $\forall v \in \mathbb{R}^n \exists ! u \in U = \text{span}(\{u_1, \dots, u_m\})$

$u' \in U^\perp = \text{span}(\{u_{m+1}, \dots, u_n\})$  s.t.  $v = u + u'$ .  $P_u(v) = 1u + 0u'$ , implying that 1 and 0 are  $P_u$ 's only eigenvalues, with  $\dim E(1, P_u) = m = \text{rank } P_u$ , and  $\dim E(0, P_u) = n - m = \dim \text{null } P_u$ .

(b) Suppose  $X \in \mathbb{R}^{n \times p}$  is a full rank matrix, with  $T \in L(\mathbb{R}^p, \mathbb{R}^n)$  its linear map. Let  $v \in \mathbb{R}^p$ , then  $T^*T v = 0 \rightarrow$

$\langle T^*T v, v \rangle = 0 \rightarrow \langle T v, T v \rangle = 0 \rightarrow \|T v\| = 0 \rightarrow v = 0$  and  $T^*T$  is injective. Since it is an operator on  $\mathbb{R}^p$  it is also invertible. In addition,  $(T^*T)^* = T^*T$ , meaning  $T^*T$  is self-adjoint. Let  $P_X = T(T^*T)^{-1}T^*$ , then

$P_X^2 = T(T^*T)^{-1}T^*T(T^*T)^{-1}T^* = T(T^*T)^{-1}T^* = P_X$  and  $P_X$  is idempotent, meaning it is a projection.

$P_X^* = (T(T^*T)^{-1}T^*)^* = T(T^*T)^{-1}T^* = P_X$ , and  $P_X$  is self-adjoint. Let  $v \in \mathbb{R}^n$ , then  $\langle P_X v, (I - P_X)v \rangle = \langle P_X v, v \rangle - \langle P_X v, P_X v \rangle = \langle P_X v, v \rangle - \langle P_X^* P_X v, v \rangle = \langle P_X v, v \rangle - \langle P_X^2 v, v \rangle = 0$ , and  $P_X$  is an orthogonal projection.

$\forall v \in \mathbb{R}^n P_X = T[(T^*T)^{-1}T^*v]$ , implying that  $P_X v \in \text{Im } T$ .

(c)  $P_X$  is self-adjoint, hence according the real spectral theorem, it is diagonalizable w.r.t. some orthonormal basis

$u_1, \dots, u_n$ . Let  $U = [u_1, \dots, u_n]$  then  $P_X = U \Lambda U^{-1}$  where  $\Lambda$  is an  $n \times n$  diagonal matrix with  $p$  ones and  $n-p$  zeros in its diagonal, corresponding to  $P_X$ 's eigenvalues.  $\text{tr}(P_X) = \text{tr}(U \Lambda U^{-1}) = \text{tr}(U^{-1} U \Lambda) = \text{tr}(\Lambda) = p$ .

### Question 3 part 1

(a) Suppose  $A \in \mathbb{R}^{n \times p}$  and  $T \in L(\mathbb{R}^p, \mathbb{R}^n)$  its linear map. Let  $v \in \mathbb{R}^n$  s.t.  $T^* T v = \lambda v$  for some  $\lambda \in \mathbb{R}$ , then

$T T^* T v = \lambda T v$ , meaning that  $T v$  is an eigenvector of  $T T^*$ , with  $\lambda$  being its eigenvalue. In fact,

$T: E(\lambda, T^* T) \rightarrow E(\lambda, T T^*)$ . For  $\lambda \neq 0$  we get  $T v \neq 0$ , meaning that  $T v|_{E(\lambda, T^* T)}$  is injective and

$\dim E(\lambda, T^* T) \leq \dim E(\lambda, T T^*)$ . Similarly, for  $\lambda \neq 0$ ,  $T^*: E(\lambda, T T^*) \rightarrow E(\lambda, T^* T)$  is injective and

$\dim E(\lambda, T T^*) \leq \dim E(\lambda, T^* T)$ . This implies that  $E(\lambda, T^* T) = E(\lambda, T T^*)$ . Let  $\lambda \in \mathbb{R}$  s.t.  $T T^* v = \lambda v$  for some  $v \in \mathbb{R}^n$ ,

then  $\lambda \langle v, v \rangle = \langle T T^* v, v \rangle = \langle T^* v, T^* v \rangle = \|T^* v\|^2 \geq 0$ , implying that  $\lambda \geq 0$ . To conclude,  $T^* T$  and  $T T^*$  share the same non-negative eigenvalues, with the positive eigenvalues having the same geometric multiplicities.

(b) Let  $r = \text{rank } A = \text{rank } A^T$  and assume that  $A$  can be decomposed s.t.  $A = U S V^T$  with  $U \in \mathbb{R}^{n \times n}$  and  $V \in \mathbb{R}^{p \times p}$  being orthogonal and  $S \in \mathbb{R}^{n \times p}$ . Then  $A A^T = U S V^T (U S V^T)^T = U S V^T V S^T U^T = U S S^T U^T$ , and we've got a familiar spectral decomposition

of a symmetric matrix  $A A^T$  with its eigenvalues in the diagonal of  $S S^T \in \mathbb{R}^{n \times n}$ . Sort these eigenvalues in a descending order:

$\lambda_1, \dots, \lambda_r, 0, \dots, 0$  and set  $U$ 's columns accordingly. Similarly,  $A^T A = (U S V^T)^T U S V^T = V S^T S V^T$  which is the spectral

decomposition of the symmetric matrix  $A^T A$  with its eigenvalues on the diagonal of  $S^T S \in \mathbb{R}^{p \times p}$ . According to (a),  $A A^T$  and  $A^T A$

share the same eigenvalues with positive ones having the same geometric multiplicity. Therefore  $\lambda_1, \dots, \lambda_r, 0, \dots, 0$  are also the diagonal

entries of  $S^T S$ , maybe with different number of trailing zeros, and we can sort them in a descending order and arrange  $V$ 's

columns accordingly. We can now define  $D \in \mathbb{R}^{r \times r}$  to be a diagonal matrix with  $\sqrt{\lambda_1}, \dots, \sqrt{\lambda_r}$  in its diagonal, and  $S \in \mathbb{R}^{n \times p}$  as

$\begin{bmatrix} D & 0 \\ 0 & 0 \end{bmatrix}$ . In order to find the singular values  $s_1, \dots, s_r$  we take the square-roots of  $A^T A$  eigenvalues.  $V$ 's columns  $v_1, \dots, v_p$  are

constructed by turning the eigenvectors into an orthonormal list, following the Gram-Schmidt procedure. Let  $U = [u_1, \dots, u_n]$  and

$V = [v_1, \dots, v_p]$ , then  $u_1, \dots, u_n$  are the orthonormal eigenvectors of  $T T^*$  and  $v_1, \dots, v_p$  are the orthonormal eigenvectors of  $T^* T$ ,

both are relative to the same sorted eigenvalue list  $\lambda_1, \dots, \lambda_r, 0, \dots, 0$  (up to the number of zeros). As shown in (a),  $\forall i=1, \dots, r$

$T|_{E(\lambda_i, T^* T)}$  is an isomorphism between  $E(\lambda_i, T^* T)$  and  $E(\lambda_i, T T^*)$ . But  $\|T v_i\|^2 = \langle T v_i, T v_i \rangle = \langle v_i, T^* T v_i \rangle = \lambda_i \|v_i\|^2 = \lambda_i$

while  $\|u_i\|^2 = 1$ . Hence  $u_i = \frac{T v_i}{\sqrt{\lambda_i}} = \frac{T v_i}{s_i}$ .

### part 2

(a) The first principal component, defined by the unit vector  $w \in \mathbb{R}^p$  that maximizes the variance of the data points  $x_1, \dots, x_n \in \mathbb{R}^p$  is

formally defined as  $PC_1^{var} = \arg \max_{w \in \mathbb{R}^p, \|w\|=1} \|X w\|^2$  using the standard inner product space over  $\mathbb{R}$  with  $X = [x_1, \dots, x_n]^T$ . Let  $T \in L(\mathbb{R}^p, \mathbb{R}^n)$

the corresponding linear map of  $X$ , then  $T^* T \in L(\mathbb{R}^p)$  is self adjoint thus according to the spectral theorem it is diagonalizable

w.r.t. some orthonormal basis. Sort  $T^* T$ 's eigenvalues and their corresponding orthonormal eigenvectors s.t.  $\lambda_1 > \dots > \lambda_p \geq 0$ .

$w = \sum_{i=1}^p \langle w, u_i \rangle u_i$  and  $T^* T_w = \sum_{i=1}^p \lambda_i \langle w, u_i \rangle u_i$ . Hence,  $\langle T_w, T_w \rangle = \langle w, T^* T_w \rangle = \sum_{i=1}^p \lambda_i \langle w, u_i \rangle^2$ . In addition,  $\|w\|^2 = 1 \rightarrow \left\langle \sum_{i=1}^p \langle w, u_i \rangle u_i, \sum_{i=1}^p \langle w, u_i \rangle u_i \right\rangle = \sum_{i=1}^p \langle w, u_i \rangle^2 = 1$ . Therefore  $\sum_{i=1}^p \lambda_i \langle w, u_i \rangle^2 \leq \sum_{i=1}^p \lambda_1 \langle w, u_i \rangle^2 = \lambda_1$  and this happens when  $\langle w, u_1 \rangle = 1 \rightarrow w = u_1$ . Furthermore,  $\sum_{i=1}^p \lambda_i \langle w, u_i \rangle^2 = \sum_{i=1}^p \lambda_1 \langle w, u_i \rangle^2 \rightarrow \sum_{i=1}^p (\lambda_i - \lambda_1) \langle w, u_i \rangle^2 = 0$  which is true only if  $\langle w, u_i \rangle = 0 \quad \forall i=2, \dots, p$ .

$$(b) \sum_{i=1}^n \text{dist}(x_i, w)^2 = \sum_{i=1}^n \langle x_i - \langle w, x_i \rangle w, x_i - \langle w, x_i \rangle w \rangle = \sum_{i=1}^n \langle x_i, x_i \rangle - 2 \langle w, x_i \rangle^2 + \langle w, x_i \rangle^2 = \sum_{i=1}^n \|x_i\|^2 - \langle w, x_i \rangle^2 = \|X\|_F^2 - \sum_{i=1}^n |w^T x_i|^2.$$

Therefore  $w = u_1$  that maximizes  $\sum_{i=1}^n |w^T x_i|^2$  also minimizes  $\sum_{i=1}^n \text{dist}(x_i, w)^2$  and  $PC_1^{\text{var}} = PC_1^{\text{LS}}$ .

#### Question 4

(a) Suppose  $Y$  is a random variable with finite expectation and variance. Let  $a \in \mathbb{R}$  and define  $M(a) = E[(Y-a)^2]$ . Then

$$M(a) = E[Y^2] - 2aE[Y] + a^2 \quad \frac{\partial}{\partial a} M(a) = -2E[Y] + 2a = 0 \rightarrow a = E[Y]. \text{ Alternatively, we can treat } Y \text{ and } a1$$

as vectors in a Hilbert space  $L^2$  with some  $(\Omega, \mathcal{F}, P)$  with the inner product defined as  $\langle A, B \rangle = E[AB]$ . Thus, we're

seeking to minimize  $Y - a1$  by finding the projection of  $Y$  onto  $\text{span}(1)$ . Such a projection  $\hat{a}1$  satisfies  $\langle Y - \hat{a}1, 1 \rangle = 0$

$$\rightarrow \langle Y, 1 \rangle = \hat{a} \langle 1, 1 \rangle \rightarrow \hat{a} = E[Y].$$

(b) Suppose  $X, Y$  are random variables with finite expectation and variance. Define  $MSE = E[(Y - f(x))^2]$  for some function  $f: \mathbb{R} \rightarrow \mathbb{R}$ .

According to the law of total expectation  $MSE = E[E[(Y - f(x))^2 | X]]$ . Based on (a), for all  $x$  in the support of  $X$

$E[(Y - f(x))^2 | X = x]$  is minimized when  $\hat{f}(x) = E(Y | X = x)$ , hence, defining  $\hat{f}(x) = E(Y | X = x)$  would minimize

the MSE.