# Introduction to Machine Learning (67577)

## Exercise 1
## Estimation Theory & Mathematical Background

### 2 Theoretical Part

#### 2.1 Mathematical Background

#### 2.1.1 Linear Algebra

1. Prove that orthogonal matrices are isometric transformations. That is, let $T : V \mapsto W$ be some linear transformation and $A$ the corresponding matrix. Show that if $A$ is an orthogonal matrix then $\forall x \in V \ ||Ax|| = ||x||$.

נוכיח כי $\forall x \in V$, $\forall$ כי (N) מטריצה $A$ היא איזומטרית נגדיר

$$||Ax|| = \sqrt{<Ax, Ax>} = \sqrt{x^T A^T A x} = \sqrt{x^T x}$$

$$= \sqrt{<x,x>} = ||x||$$

$$A^T A = I \downarrow$$

כי $A$ אורתוגונלית

2. Calculate the SVD of the following matrix $A$. That is, find the matrices $U, \Sigma, V^\top$ where $U, V$ are orthogonal matrices and $\Sigma$ diagonal.

$$A = \begin{bmatrix} 1 & 1 & 0 \\ 1 & -1 & 2 \end{bmatrix}$$

Recall, that to find the SVD of $A$ we can calculate $A^\top A$ to deduce $V, \Sigma$ and then calculate $AA^\top$ to deduce $U$. Equivalently, once we deduced $V, \Sigma$ we can fine $U$ using the equality $AV = U\Sigma$.

$$A^T A = \begin{bmatrix} 1 & 1 \\ 1 & -1 \\ 0 & 2 \end{bmatrix} \begin{bmatrix} 1 & 1 & 0 \\ 1 & -1 & 2 \end{bmatrix} = \begin{bmatrix} 2 & 0 & 2 \\ 0 & 2 & -2 \\ 2 & -2 & 4 \end{bmatrix}$$

נחשב את ה $A$:

$$\det(A^T A - I\lambda) = \begin{vmatrix} 2-\lambda & 0 & 2 \\ 0 & 2-\lambda & -2 \\ 2 & -2 & 4-\lambda \end{vmatrix}$$

$$= (2-\lambda)\begin{vmatrix} 2-\lambda & -2 \\ -2 & 4-\lambda \end{vmatrix} + 2\begin{vmatrix} 0 & 2-\lambda \\ 2 & -2 \end{vmatrix}$$

$$= (2-\lambda)\left(4 - 6\lambda + \lambda^2\right) - 4\left(2 - \lambda\right)$$

$$= (2-\lambda)\left(\lambda^2 - 6\lambda\right) = -\lambda(\lambda-2)(\lambda-6)$$

$$\boxed{\lambda = 0, 2, 6} \qquad : \overset{=}{rr} \quad \text{إذ}$$

$\underline{\lambda = 0}$:

$$\begin{bmatrix} 2 & 0 & 2 \\ 0 & 2 & -2 \\ 2 & -2 & 4 \end{bmatrix}\begin{bmatrix} x \\ y \\ z \end{bmatrix} = \bar{0} \longrightarrow \begin{array}{l} x = -z \\ y = z \\ x - y + 2z = 0 \end{array}$$

$$: \overset{=}{r}_1 \qquad \overset{=}{\text{is}}$$

$$\frac{1}{\sqrt{3}}\begin{bmatrix} 1 \\ -1 \\ -1 \end{bmatrix}$$

$\underline{\lambda = 2}$
$$\begin{bmatrix} 0 & 0 & 2 \\ 0 & 0 & -2 \\ 2 & -2 & 2 \end{bmatrix}\begin{bmatrix} x \\ y \\ z \end{bmatrix} = \bar{0} \longrightarrow \begin{array}{l} z = 0 \\ x = y \end{array}$$

$$\frac{1}{\sqrt{2}}\begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}$$

<div dir="rtl">סעיף ī</div>

$$\lambda = 6: \quad \begin{bmatrix} -4 & 0 & 2 \\ 0 & -4 & -2 \\ 2 & -2 & -2 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \bar{0} \rightarrow \begin{array}{l} z = 2x \\ z = -2y \end{array}$$

<div dir="rtl">סעיף īī</div>

$$\frac{1}{\sqrt{6}}\begin{bmatrix} 1 \\ -1 \\ 2 \end{bmatrix}$$

<div dir="rtl">סעיף</div>

$$V = \begin{bmatrix} \frac{1}{\sqrt{6}} & \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{3}} \\ \frac{-1}{\sqrt{6}} & \frac{1}{\sqrt{2}} & \frac{-1}{\sqrt{3}} \\ \frac{2}{\sqrt{6}} & 0 & \frac{-1}{\sqrt{3}} \end{bmatrix} \qquad D = \begin{bmatrix} 6 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

<div dir="rtl">סעיף</div>

$$\sigma_1 = \sqrt{6}, \quad \sigma_2 = \sqrt{2}, \quad \sigma_3 = 0$$

<div dir="rtl">נחשב את V:</div>

$$AA^T = \begin{bmatrix} 1 & 1 & 0 \\ 1 & -1 & 2 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & -1 \\ 0 & 2 \end{bmatrix}$$

$$= \begin{bmatrix} 2 & 0 \\ 0 & 6 \end{bmatrix} \qquad \lambda_{1,2} = 2, 6$$

$$\lambda = 2: \quad \begin{bmatrix} 0 & 0 \\ 0 & 4 \end{bmatrix} \begin{bmatrix} x \\ -y \end{bmatrix} = \bar 0 \rightarrow \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$\lambda = 6: \qquad\qquad\qquad\qquad \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$AA^T = U \begin{bmatrix} 6 & 0 \\ 0 & 2 \end{bmatrix} U^T$$

$$U = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

$$A = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} \sqrt6 & 0 & 0 \\ 0 & \sqrt2 & 0 \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt6} & \frac{-1}{\sqrt6} & \frac{2}{\sqrt6} \\ \frac{1}{\sqrt2} & \frac{1}{\sqrt2} & 0 \\ \frac{1}{\sqrt3} & \frac{-1}{\sqrt3} & \frac{-1}{\sqrt3} \end{bmatrix}$$

3. In this question we prove the Power-Iteration algorithm for finding the SVD of a matrix. Let $A \in \mathbb{R}^{m \times n}$ and define $C_0 = A^\top A$. Denote $\lambda_1 \geq \ldots \geq \lambda_n$ the eigenvalues of $C_0$, with the corresponding normalized eigenvectors $v_1, \ldots, v_n$.

Let us assume the $\lambda_1 > \lambda_2$. Define $b_k \in \mathbb{R}$ as follows:

$$b_0 = \sum_{i=1}^{n} a_i v_i, \quad b_{k+1} = \frac{C_0 b_k}{\|C_0 b_k\|}$$

where $a_1 \neq 0$. Show that: $\lim_{k \to \infty} b_k = \pm v_1$.

$$b_{k+1} = \frac{C_0 b_k}{\|C_0 b_k\|} = C_0 \cdot \frac{C_0 b_{k-1}}{\|C_0 b_{k-1}\|} \cdot \frac{1}{\left\| \frac{C_0 \cdot C_0 b_{k-1}}{\|C_0 b_{k-1}\|} \right\|}$$

$$= C_0^2 \frac{b_{k-1}}{\|C_0 b_{k-1}\|} \cdot \frac{\|C_0 b_{k-1}\|}{\|C_0^2 b_{k-1}\|} = \frac{C_0^2 b_{k-1}}{\|C_0^2 b_{k-1}\|}$$

$$b_k = \frac{C_0^k b_0}{\|C_0^k b_0\|} = \frac{C_0^k \sum_i a_i v_i}{\|C_0^k \sum_i a_i v_i\|} \qquad \text{الى}$$

$$= \frac{(V \Sigma^\top \Sigma V^\top)^k \sum_i a_i v_i}{\|(V \Sigma^\top \Sigma V^\top)^k \sum_i a_i v_i\|} = \frac{V D^k V^\top a_i v_i}{\|V D^k V^\top a_i v_i\|}$$

$$= \frac{\sum_i a_i \lambda_i^k V e_i}{\|\sum_i a_i \lambda_i^k V e_i\|} = \frac{\sum_i a_i \lambda_i^k v_i}{\|\sum_i a_i \lambda_i^k v_i\|}$$

$$= \frac{\lambda_1^k\left(\sum a_i\left(\frac{\lambda_i}{\lambda_1}\right)^k v_i\right)}{\lambda_1^k\left\|\sum a_i\left(\frac{\lambda_i}{\lambda_1}\right)^k v_i\right\|} = \frac{a_1 v_1}{\sqrt{\|a_1 v_1\|}} = v_1$$

כאשר $k \to \infty$

$\longrightarrow v_i$

אנחנו נקבל .

### 2.1.2 Multivariate Calculus

4. Let $x \in \mathbb{R}^n$ be a fixed vector and $U \in \mathbb{R}^{n \times n}$ a fixed orthogonal matrix. Calculate the Jacobian of the function $f : \mathbb{R}^n \to \mathbb{R}^n$:

$$f(\sigma) = U \cdot \operatorname{diag}(\sigma) U^\top x$$

Where $\operatorname{diag}(\sigma)$ is an $n \times n$ matrix where

$$\operatorname{diag}(\sigma)_{ij} = \begin{cases} \sigma_i & i = j \\ 0 & i \neq j \end{cases}$$

$$f(\sigma) = U \operatorname{diag}(\sigma) U^\top x = U \sum \sigma_i u_i^\top x$$

נגדיר
$$\frac{\partial f}{\partial \sigma_i} = U$$
ונקבל

$$\boxed{\operatorname{Jac}(f) = U \operatorname{diag}(U^\top x)}$$

5. Use the chain rule to calculate the gradient of $h(\sigma) = \frac{1}{2}\|f(\sigma) - y\|^2$

$$h(\sigma) = \frac{1}{2}\|f(\sigma) - y\|^2 =$$

$$\nabla h(\sigma) = \left(f(\sigma) - y\right)^T \nabla f(\sigma)$$

לפי כלל השרשרת:

6. Calculate the Jacobian of the softmax function $S : \mathbb{R}^d \to [0,1]^k$

$$S(\mathbf{x})_j = \frac{e^{x_j}}{\sum_{l=1}^{k} e^{x_l}}$$

$$\frac{\partial S(x)_j}{\partial x_i} = - \frac{e^{x_j + x_i}}{\left(\sum_l e^{x_l}\right)^2} \quad \Downarrow \quad i \neq j$$

$$\frac{\partial S(x)_j}{\partial x_i} = \frac{e^{x_i}}{\sum_l e^{x_l}} - \frac{e^{x_j + x_i}}{\left(\sum_l e^{x_l}\right)^2} \quad \Downarrow \quad i = j$$

$$J(S(x))_{ji} = \begin{cases} & i \neq j \\ & i = j \end{cases}$$

לפי

7. Let $f : \mathbb{R}^d \to \mathbb{R}$ be defined as $f(x,y) = x^3 - 5xy - y^5$. Calculate the Hessian of $f$.

$$H(f) = \begin{array}{cc} \dfrac{\partial f}{\partial x^2} & \dfrac{\partial f}{\partial x \partial y} \\[4mm] \dfrac{\partial f}{\partial y \partial x} & \dfrac{\partial f}{\partial y^2} \end{array}$$

$\dfrac{\partial f}{\partial x} = 3x^2 - 5y$

$\dfrac{\partial f}{\partial x \partial y} = \dfrac{\partial f}{\partial y \partial x} = -5$

$\dfrac{\partial f}{\partial y} = -5x - 5y^4$

$\dfrac{\partial f}{\partial x^2} = 6x$

$\dfrac{\partial f}{\partial y^2} = -20y^3$

$$H(f) = \begin{array}{cc} 6x & -5 \\[3mm] -5 & -20y^3 \end{array}$$

ل.م

## 2.2 Estimation Theory

8. Let $x_1, x_2, \ldots \overset{iid}{\sim} \mathcal{P}$ be a sample of infinity size drawn from some probability distribution function $\mathcal{P}$ with finite expectation and variance. Show that the sample mean estimator $\hat{\mu}_n = \frac{1}{n}\sum x_i$ calculated over the first $n$ samples is a consistent estimator. Hint: for any given fixed value of $n \in \mathbb{N}$ bound from above the probability of deviating more than $\varepsilon$.

$$\hat{\mu}_n = \frac{1}{n}\sum x_i$$

$$p\left(\hat{\mu}_n - \mu \geq \varepsilon\right) \leq \frac{1}{\sqrt{4n}}$$

חישבנו בדרכים שונות את 3

לכן אנו צריכים כי בסך הכל n עם גבול 4 אכן מתכנס

לכל ערך ב ε הוא קטן לכ ל 0

9. Let $\mathbf{x}_1, \ldots, \mathbf{x}_m \overset{iid}{\sim} \mathcal{N}(\mu, \Sigma)$ be $m$ observations sampled i.i.d from a multivariate Gaussian with expectation of $\mu \in \mathbb{R}^d$ and a covariance matrix $\Sigma \in \mathbb{R}^{d \times d}$. Derive the log-likelihood function of $\mathcal{N}(\mu, \Sigma)$. Hint: follow the approach used to derive the likelihood function for the univariate case.

נסמן $\bar{x}_i$ בתור שמות נקודה כ:

$$f_{(\mu, \Sigma)}(x_i) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left(-\frac{1}{2}(\bar{x}_i - \mu)^T \Sigma^{-1}(\bar{x}_i - \mu)\right)$$

נכן

$$L(x, \mu, \Sigma) = \prod_{i=1}^{m} f_{(\mu, \Sigma)}(\bar{x}_i)$$

המחשב log likelihood:

$$L(x, M, \Sigma) = \log \prod_{i=1}^{m} f_{(M, \Sigma)}(\bar{x}_i)$$

$$= \sum_{i=1}^{m} -\frac{d}{2} \log(2\pi) - \frac{1}{2} \log(|\Sigma|) - \frac{1}{2}(x_i - M)^T \Sigma^{-1}(x_i - M)$$

$$= -\frac{md}{2} \log(2\pi) - \frac{m}{2} \log(|\Sigma|) - \frac{1}{2} \sum_{i=1}^{m} (x_i - M)^T \Sigma^{-1}(x_i - M)$$