שלו ראקובלה

# Introduction to Machine Learning (67577)

# Exercise 1
# Estimation Theory & Mathematical Background

Second Semester, 2023

## Contents

# 1   Submission Instructions

Please make sure to follow the general submission instructions available on the course website. In addition, for the following assignment, submit a single `ex1_ID.tar` file containing:
- An `Answers.pdf` file with the answers for all theoretical and practical questions (include plotted graphs *in* the PDF file).
- The following python files (without any directories): `gaussian_estimators.py`, `fit_gaussian_estimators.py`

The `ex1_ID.tar` file must be submitted in the designated Moodle activity prior to the date specified *in the activity*.
- Late submissions will not be accepted and result in a zero mark.
- Plots included as separate files will be considered as not provided.

# 2   Theoretical Part

## 2.1   Mathematical Background

### 2.1.1   Linear Algebra

Based on Recitation 1

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 $\forall x \in V \ ||Ax|| = ||x||$.
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 $AV = U\Sigma$.
3. Show that the outer product of two vectors $\mathbf{v} \in \mathbb{R}^n, \mathbf{u} \in \mathbb{R}^m$, which is denoted by $\mathbf{v} \otimes \mathbf{u}$ or $\mathbf{v} \cdot \mathbf{u}^\top$ is a matrix $A \in \mathbb{R}^{n \times m}$ with $rank(A) = 1$. That is, show that all rows (or columns) in $A$ are linearly dependent.
4. Show that for any otrthonormal basis $(\mathbf{u}_1, ..., \mathbf{u}_n)$ and any aribtrary vector $\mathbf{x} \in \mathbb{R}_n$ such that $\mathbf{x} = \sum_{i=1}^n a_i \cdot \mathbf{u}_i$, it holds that $a_i = \langle \mathbf{x}, \mathbf{u}_i \rangle$ for any $i \in [1, n]$. That is, show that the i'th coefficient of representing $\mathbf{x}$ in the basis $(\mathbf{u}_1, ..., \mathbf{u}_n)$, is the inner product between $\mathbf{x}$ and $\mathbf{u}_i$.

### 2.1.2   Multivariate Calculus

Based on Recitation 2

5. 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 $f : \mathbb{R}^n \to \mathbb{R}^n$:

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

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

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

6. Use the chain rule to calculate the gradient of $h(\sigma) = \frac{1}{2}||f(\sigma) - y||^2$
7. Calculate the Jacobian of the softmax function $S : \mathbb{R}^d \to [0,1]^k$

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

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

### 2.1.3 convexity

Based on Recitation 2

9. Prove that the intersection $C := \bigcap_{i \in I} C_i$ for $\{C_i : i \in I\}$ a collection of convex sets is convex.
10. Prove that the vector sum $C_1 + C_2 := \{c_1 + c_2 : c_1 \in C_1, c_2 \in C_2\}$ of two convex sets is convex.
11. Prove that the set $\lambda C := \{\lambda c : c \in C\}$ is convex, for any convex set $C$, and every scalar $\lambda$.

## 2.2 Estimation Theory

Based on Lecture 1

12. 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 $\varepsilon$.

13. 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 $\Sigma \in \mathbb{R}^{d \times d}$. Provide an expression for the log-likelihood function of $\mathcal{N}(\mu, \Sigma)$. Develop the expression as much as you can. Hint: follow the approach used to derive the likelihood function for the univariate case.

## 3  Practical Part

Before starting the practical part please make sure to have cloned/downloaded the IML.HUJI GitHub repository and set up a working virtual environment. Write the necessary code in the files specified in the questions.

## 3.1  Univariate Gaussian Estimation

Based on lecture 1

Implement the `UnivariateGaussian` class in the `learners.gaussian_estimators.py` file. Follow details specified in class and function documentation.

1. Using `numpy.random.normal` draw 1000 samples $x_1, \ldots, x_{1000} \overset{iid}{\sim} \mathcal{N}(10, 1)$ and fit a univariate Gaussian. Print the estimated expectation and variance. Output format should be `(expectation, variance)`.

2. Over previously drawn samples, fit a series of models of increasing samples size: 10, 20,...,100, 110,...1000. Plot the absolute distance between the estimated- and true value of the expectation, as a function of the sample size. Provide meaningful axis names and title.

3. Compute the PDF of the previously drawn samples using the model fitted in question 1. Plot the empirical PDF function under the fitted model. That is, create a scatter plot with the ordered sample values along the x-axis and their PDFs (using the `UnivariateGaussian.pdf` function) along the y-axis. Provide meaningful axis names and title. What are you expecting to see in the plot?

## 3.2  Multivariate Gaussian Estimation

Based on Lecture 1

Implement the `Multivariate` class in the `learners.gaussian_estimators.py` file. Follow details specified in class and function documentation.

NOTICE: When implementing the `log_likelihood` function you are required to use the expression developed in the q13 above. That is, the expression for $\ell(\mu, \Sigma | \mathbf{x}_1, \dots, \mathbf{x}_m)$.

4. Using `numpy.random.multivariate_normal` draw 1000 samples $\mathbf{x}_1, \dots, \mathbf{x}_{1000} \overset{iid}{\sim} \mathcal{N}(\mu, \Sigma)$

$$\mu = \begin{bmatrix} 0 \\ 0 \\ 4 \\ 0 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} 1 & 0.2 & 0 & 0.5 \\ 0.2 & 2 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0.5 & 0 & 0 & 1 \end{bmatrix}$$

Fit a multivariate Gaussian and print the estimated expectation and covariance matrix. Print each in a separate line.

5. Using the samples drawn in the question above calculate the log-likelihood for models with expectation $\mu = [f_1, 0, f_3, 0]^\top$ and the true covariance matrix defined above, where $f_1, f_3$ get values returned from `np.linspace(-10, 10, 200)`. Plot a heatmap of $f1$ values as rows, $f_3$ values as columns and the color being the calculated log likelihood. Provide meaningful axis names and title. What are you able to learn from the plot?

6. Of all values tested in question 5, which model (pair of values for feature 1 and 3) achieved the maximum log-likelihood value? Round to 3 decimal places

## 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 $\forall x \in V \; ||Ax|| = ||x||$.

$$||Ax||^2 = (Ax)^t \, Ax = x^t A^t \, Ax = x^t x = ||x||^2$$

transpose חוק

$A$ אורתוגונלי
ולכן $A^t A = I_n$

(א

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 $AV = U\Sigma$.

(2    בכדי למצוא את ה-SVD של $A$, נוכל למצוא את ה-EVD של $A^t A$:

$$A^t_A = \begin{pmatrix} 1 & 1 \\ 1 & -1 \\ 0 & 2 \end{pmatrix} \begin{pmatrix} 1 & 1 & 0 \\ 1 & -1 & 2 \end{pmatrix} = \begin{pmatrix} 2 & 0 & 2 \\ 0 & 2 & -2 \\ 2 & -2 & 4 \end{pmatrix}$$

נמצא את הערכים העצמיים $3 \times 3$ של $A^T A$ המטריצה המתאימה ל-$A^t A$:

$$\det(\lambda I_n - A^t A) = \det \begin{pmatrix} \lambda-2 & 0 & -2 \\ 0 & \lambda-2 & 2 \\ -2 & 2 & \lambda-4 \end{pmatrix} = (\lambda-2) \begin{pmatrix} \lambda-2 & 2 \\ 2 & \lambda-4 \end{pmatrix} -2 \begin{pmatrix} 0 & -2 \\ \lambda-2 & 2 \end{pmatrix} = (\lambda-2)\left[(\lambda-2)(\lambda-4)-4\right]$$

$$\underbrace{-2\left(0-(-2)(\lambda-2)\right)}_{\underbrace{-2 \cdot 2(\lambda-2)}_{-4(\lambda-2)}} = (\lambda-2)\underbrace{\left(\lambda^2-6\lambda\right)}_{\text{הכפלה וצמצום}} = \underbrace{(\lambda-2)\lambda(\lambda-6)}_{\text{הערכים העצמיים של } A^t A} \Rightarrow \text{ לפי כלל אפסים: } 2,0,6$$

נמצא את הוקטורים העצמיים של $A^t A$ בכדי למצוא את $V$: נציב לכל ערך עצמי במטריצה $\lambda I_n - A^t A$ ונפתור:

$$\begin{pmatrix} 0 & 0 & -2 \\ 0 & 0 & 2 \\ -2 & 2 & -2 \end{pmatrix} \longrightarrow \begin{pmatrix} 1 & -1 & 1 \\ 0 & 0 & 2 \\ 0 & 0 & -2 \end{pmatrix} \rightarrow \begin{pmatrix} 1 & -1 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & -1 \end{pmatrix} \longrightarrow \begin{pmatrix} 1 & -1 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{pmatrix}$$

למשל 2=λ:

$x_1 = t \;,\; x_3 = 0 \;,\; x_2 = t$ לכן $V_2 = \left\{ \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix} t \;\middle|\; t \in \mathbb{R} \right\}$ , נפוט $\begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix}$ לכן לא תלוי 2.

$$\begin{pmatrix} -2 & 0 & -2 \\ 0 & -2 & 2 \\ -2 & 2 & -4 \end{pmatrix} \longrightarrow \begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & -1 \\ 1 & -1 & 2 \end{pmatrix} \longrightarrow \begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & -1 \\ 0 & -1 & 1 \end{pmatrix} \longrightarrow \begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & -1 \\ 0 & 0 & 0 \end{pmatrix}$$

$x_3 = t$, $x_1 = -t$, $x_2 = t$ נקבל או $\begin{pmatrix} -1 \\ 1 \\ 1 \end{pmatrix}$ ע"י $k$ ל $0$.

$$\begin{pmatrix} 4 & 0 & -2 \\ 0 & 4 & 2 \\ -2 & 2 & 2 \end{pmatrix} \longrightarrow \begin{pmatrix} 2 & 0 & -1 \\ 0 & 2 & 1 \\ -1 & 1 & 1 \end{pmatrix} \longrightarrow \begin{pmatrix} 0 & 2 & 1 \\ 0 & 0 & 0 \\ -1 & 1 & 1 \end{pmatrix} \longrightarrow \begin{pmatrix} 1 & -1 & -1 \\ 0 & 1 & \frac{1}{2} \\ 0 & 0 & 0 \end{pmatrix}$$

ל $\lambda = 4$:

$x_3 = t$, $x_2 = -\frac{1}{2}t$, $x_1 = \frac{1}{2}t$ נקבל או $\begin{pmatrix} 1 \\ -1 \\ 2 \end{pmatrix}$ ע"י $k$ ל $6$.

סדר נפנה ע"י $\begin{pmatrix} 1 \\ -1 \\ 2 \end{pmatrix}, \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix}, \begin{pmatrix} -1 \\ 1 \\ 1 \end{pmatrix}$ של הם $0, 2, 6$ בהתאמה. נראה מ $A^t A$ הם וקטורים עצמיים ע"י של $\lambda$

שונים ולכן בל"ת הקטורים ולכן נאורמנל כדי שיהיו אורתונ' ב $R^3$.

$$\left\| \begin{pmatrix} 1 \\ 1 \\ 2 \end{pmatrix} \right\| = \sqrt{1+1+4} = \sqrt{6}, \quad \left\| \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix} \right\| = \sqrt{2}, \quad \left\| \begin{pmatrix} -1 \\ 1 \\ 1 \end{pmatrix} \right\| = \sqrt{3} \Leftarrow (v_1, v_2, v_3) = \left( \begin{pmatrix} 1 \\ -1 \\ 2 \end{pmatrix} \frac{1}{\sqrt{6}}, \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix} \frac{1}{\sqrt{2}}, \begin{pmatrix} -1 \\ 1 \\ 1 \end{pmatrix} \frac{1}{\sqrt{3}} \right)$$

בסיס אורתונ' ל $R^3$.

או נקבל: $A^t A = \begin{pmatrix} | & | & | \\ v_1 & v_2 & v_3 \\ | & | & | \end{pmatrix} \begin{pmatrix} 6 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} | & | & | \\ v_1 & v_2 & v_3 \\ | & | & | \end{pmatrix}^t$

$= V^T$

לפי נוסחת ה- $V^T = \begin{pmatrix} 1/\sqrt{6} & -1/\sqrt{6} & 2/\sqrt{6} \\ 1/\sqrt{2} & 1/\sqrt{2} & 0 \\ -1/\sqrt{3} & 1/\sqrt{3} & 1/\sqrt{3} \end{pmatrix}$

הערכים העצמיים הם הריבועים של ערכי או

$A^t A$ חיוביים $0$, וזהו $\sqrt{6}, \sqrt{2}$. ולכן $\Sigma = \begin{pmatrix} \sqrt{6} & 0 & 0 \\ 0 & \sqrt{2} & 0 \end{pmatrix}$ ($\Sigma$ הם הערכים הסינגולריים של $A$).

נמצא את $AA^t$ כדי למצוא את הוקטורים הסינגולריים השמאליים-

$AA^t = \begin{pmatrix} 1 & 1 & 0 \\ 1 & -1 & 2 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 1 & -1 \\ 0 & 2 \end{pmatrix} = \begin{pmatrix} 2 & 0 \\ 0 & 6 \end{pmatrix}$

נקבל ש- $(e_2, e_1)$ בסיס אורתונ' של ע"י (כבר ידוע ולכן הערכים העצמיים) נראה ש- $AA^t e_1 = 2e_1$, $AA^t e_2 = 6e_2$,

$e_2 \perp e_1$, $\|e_1\| = \|e_2\| = 1$. ולכן $V = (e_2, e_1) = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$.

או: $A = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} \sqrt{6} & 0 & 0 \\ 0 & \sqrt{2} & 0 \end{pmatrix} \begin{pmatrix} 1/\sqrt{6} & -1/\sqrt{6} & 2/\sqrt{6} \\ 1/\sqrt{2} & 1/\sqrt{2} & 0 \\ -1/\sqrt{3} & 1/\sqrt{3} & 1/\sqrt{3} \end{pmatrix}$.

$\qquad\qquad\qquad\qquad V \qquad\qquad \Sigma \qquad\qquad\qquad V^T$

3. Show that the outer product of two vectors $\mathbf{v} \in \mathbb{R}^n, \mathbf{u} \in \mathbb{R}^m$, which is denoted by $\mathbf{v} \otimes \mathbf{u}$ or $\mathbf{v} \cdot \mathbf{u}^\top$ is a matrix $A \in \mathbb{R}^{n \times m}$ with $rank(A) = 1$. That is, show that all rows (or columns) in $A$ are linearly dependent.

**Definition 1.4** For two vectors $\mathbf{v} \in \mathbb{R}^n, \mathbf{u} \in \mathbb{R}^m$, the *outer product* of $\mathbf{v}$ and $\mathbf{u}$, which is denoted by $\mathbf{v} \otimes \mathbf{u}$ or $\mathbf{vu}^\top$ is an $n \times m$ matrix with entries:

מוכר $n \times 1$    $1 \times m \Rightarrow n \times m$
בהכפלה
הרגילה

$$[\mathbf{v} \otimes \mathbf{u}]_{ij} = v_i \cdot u_j, \quad \mathbf{v} \otimes \mathbf{u} = \begin{bmatrix} v_1 u_1 & v_1 u_2 & \cdots & v_1 u_m \\ \vdots & \vdots & \ddots & \vdots \\ v_n u_1 & v_n u_2 & \cdots & v_n u_m \end{bmatrix} = A$$

נעיין רק אם נסתכל על המטריצה $A = \begin{pmatrix} v_1 \cdot u^t \\ v_2 \cdot u^t \\ \vdots \\ v_n \cdot u^t \end{pmatrix}$, נוכל לראות שכל שורה ניתן לכתוב כמכפלה של אותה שורה (השורה) בסקלר, ולכן השורות תלויות ליניארית.

לכן, $u \neq 0$ כי, $0 \neq N$ שורה (כלשהי) $v_i$   $v = \begin{pmatrix} v_1 \\ \vdots \\ v_n \end{pmatrix}$ קיים בהכרח קורדינטה $\neq 0$ אם $u, v \neq 0$ אכן, לכן.    לדוגמה: $v_2 u^t = v_2 u^t \cdot \frac{v_2}{v_n}$

$v_i u^t \neq 0$ ולכן     $rank(A) = 1$.

4. Show that for any otrthonormal basis $(\mathbf{u}_1, ..., \mathbf{u}_n)$ and any aribtrary vector $\mathbf{x} \in \mathbb{R}_n$ such that $\mathbf{x} = \sum_{i=1}^n a_i \cdot \mathbf{u}_i$, it holds that $a_i = \langle \mathbf{x}, \mathbf{u}_i \rangle$ for any $i \in [1, n]$. That is, show that the i'th coefficient of representing $\mathbf{x}$ in the basis $(\mathbf{u}_1, ..., \mathbf{u}_n)$, is the inner product between $\mathbf{x}$ and $\mathbf{u}_i$.

יהי $1 \leq i \leq n$. אכן -

$$\langle x, u_i \rangle = \langle \sum_{j=1}^n a_j u_j, u_i \rangle = \sum_{j=1}^n a_j \langle u_j, u_i \rangle = a_i$$

לינאריות
של המכפלה הפנימית
דלתא קרונקר

מכיוון שכל $u$-ים הם אורתונורמליים המכפלה הפנימית $\langle u_i, u_j \rangle = 0$ כאשר $i \neq j$ ו $1$ אחרת, כי הם אורתונורמליים

**Definition 2.4** Let $\mathbf{f} : \mathbb{R}^d \to \mathbb{R}^m$ where $\mathbf{f}(\mathbf{x}) = (f_1(\mathbf{x}), ..., f_m(\mathbf{x}))^\top$. The *Jacobian* of $f$ is the $m \times d$ matrix of all partial derivatives:

נזכור: המטריצה המסמלת את ההשפעה של השינוי בכל $x_i$ על $f$ וכן את הקשר של שינוי. ית.

$$J_\mathbf{x}(\mathbf{f}) := \begin{bmatrix} \frac{\partial f_1(\mathbf{x})}{\partial x_1} & \cdots & \frac{\partial f_1(\mathbf{x})}{\partial x_d} \\ \vdots & & \vdots \\ \frac{\partial f_m(\mathbf{x})}{\partial x_1} & \cdots & \frac{\partial f_m(\mathbf{x})}{\partial x_d} \end{bmatrix}$$

(3

Based on Recitation 2

5. 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 $f : \mathbb{R}^n \to \mathbb{R}^n$:

$$f(\sigma) = \underbrace{U}_{n \times n} \cdot \underbrace{\operatorname{diag}(\sigma)}_{n \times n} \underbrace{U^{\top}}_{n \times n} \underbrace{x}_{n \times 1}$$

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

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

$$U = \begin{pmatrix} | & & | \\ u_1 & \dots & u_n \\ | & & | \end{pmatrix} \qquad \text{נגדיר} \quad (5)$$

$$f(\sigma) = U \operatorname{diag}(\sigma) U^{\top} x = U \begin{pmatrix} \sigma_1 & & 0 \\ & \sigma_2 & \\ 0 & & \ddots & \\ & & & \sigma_n \end{pmatrix} \begin{pmatrix} u_1^t x \\ \vdots \\ u_n^t x \end{pmatrix} = U \begin{pmatrix} \sigma_1 u_1^t x \\ \vdots \\ \sigma_n u_n^t x \end{pmatrix} = \overbrace{\sum_{i=1}^{n}}^{\text{סקלר}} \underbrace{(\sigma_i u_i^t x)}_{1 \times n} \underbrace{u_i}_{n \times 1} =$$

$$= \sum_{i=1}^{n} u_i (u_i^t \sigma_i x) = \sum_{i=1}^{n} \sigma_i u_i u_i^t x$$

$$\left[ J_\sigma(F) \right]_{i,j} = \frac{\partial f_i(\sigma)}{\partial \sigma_i} = \left[ u_i u_i^t x \right]_j \qquad \Longrightarrow \quad J_\sigma(F) = \begin{pmatrix} | & & | \\ u_1 & \dots & u_n \\ | & & | \end{pmatrix} \begin{pmatrix} u_1^{\top} x & & 0 \\ & \ddots & \\ 0 & & u_n^t x \end{pmatrix} = U \operatorname{diag}(U^{\top} x)$$

$u_i u_i^t x$ הקואורדינטה ה-$j$ של

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

(6)

$h(\sigma) = \frac{1}{2}\|f(\sigma)-y\|^2 = \frac{1}{2}\langle f(\sigma)-y, f(\sigma)-y\rangle = \frac{1}{2}\left(\|f(\sigma)\|^2 + \|y\|^2 - 2\langle f(\sigma), y\rangle\right) = \frac{1}{2}\|f(\sigma)\|^2 + \frac{1}{2}\|y\|^2 - \langle f(\sigma), y\rangle =$

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

לפי הגדרה

$\nabla h(\sigma) = \frac{\partial h(\sigma)}{\partial \sigma} = \frac{\partial h(\sigma)}{\partial f(\sigma)} \cdot \frac{\partial f(\sigma)}{\partial \sigma} = \frac{\partial\left(\frac{1}{2}\|f(\sigma)\|^2 - f(\sigma)^t y + \frac{1}{2}\|y\|^2\right)}{\partial f(\sigma)} \cdot \frac{\partial f(\sigma)}{\partial \sigma} = \left(f(\sigma)-y\right)^t \cdot \frac{\partial f(\sigma)}{\partial \sigma} =$

$= \left(f(\sigma)-y\right)^t \underset{h\times n}{J_\sigma(f)} \underset{6 \text{ שאלה לפי}}{=} \left(f(\sigma)-y\right)^t \underset{h\times h}{U \, diag}\left(U^T x\right)$

נזכיר כי נגזרות ב-$\mathbb{R}^h$ מקיימות $x \cdot y$ לכן $\frac{\partial x^t y}{\partial x_i} = \sum_i x_i \partial_i = y_i \Rightarrow \frac{\partial x^t y}{\partial x} = y$, $\frac{\partial x^t x}{\partial x_i} = \frac{\sum_i x_i^2}{\partial x_i} = 2x_i \Rightarrow \frac{\partial \|x\|^2}{\partial x} = 2x$ ✳

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

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

כלל המנה: $\left(\dfrac{f}{g}\right)' = \dfrac{f'g - fg'}{g^2}$

$$[J_K(s)]_{i,j} = \frac{\partial S(x)_i}{\partial x_j} = \frac{\partial}{\partial x_j}\frac{e^{x_i}}{\sum_{\ell=1}^{k} e^{x_\ell}} = \begin{cases} \dfrac{e^{x_i}\sum_{\ell=1}^{k} e^{x_\ell} - e^{x_i}\cdot e^{x_j}}{\left(\sum_{\ell=1}^{k} e^{x_\ell}\right)^2} \underset{\text{נוציא}}{\overset{\text{גורם}}{=}} \dfrac{e^{x_i}}{\sum_{\ell=1}^{k} e^{x_\ell}}\cdot\dfrac{\left(\sum_{\ell=1}^{k} e^{x_\ell} - e^{x_j}\right)}{\sum_{\ell=1}^{k} e^{x_\ell}} & i=j \end{cases}$$

$$S(x)_i\cdot\left(\frac{\sum_{\ell=1}^{k} e^{x_\ell}}{\sum_{\ell=1}^{k} e^{x_\ell}} - \frac{e^{x_i}}{\sum_{\ell=1}^{k} e^{x_\ell}}\right) = S(x)_i\cdot(1 - S(x)_i)$$

$$\frac{0\cdot\sum_{\ell=1}^{k} e^{x_\ell} - e^{x_i}\cdot e^{x_j}}{\left(\sum_{\ell=1}^{k} e^{x_\ell}\right)^2} = -\frac{e^{x_i}}{\sum_{\ell=1}^{k} e^{x_\ell}}\cdot\frac{e^{x_j}}{\sum_{\ell=1}^{k} e^{x_\ell}} = -S(x)_i\, S_j(x) = S(x)_i\,(-S(x)_j) \quad i\ne j$$

כלומר $[J_x(f)]_{i,j} = S(x)_i\cdot(\delta_{ij} - S(x)_i)$ כאשר $\delta_{ij} = \begin{cases} 1 & i=j \\ 0 & i\ne j \end{cases}$.

ובכתיבה מטריצונית: $J_x(f) = \text{diag}(S) - SS^T$ נשים לב שכן $[SS^T]_{ij} = S(x)_i\cdot S(x)_j$, וכן $[\text{diag}(S)]_{ij} = \begin{cases} S(x)_i & i=j \\ 0 & i\ne j \end{cases}$

וכן (להשלים) $S(x)_i - S(x)_i\,S(x)_j = S(x)_i\,(1-S(x)_j)$, $S(x)_i\,S(x)_j = -S(x)_i\,S(x)_j$, $0 - S(x)_i\,S(x)_j$

---

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

(8)

$$\frac{\partial f}{\partial x} = 3x^2 - 5y \qquad\qquad \frac{\partial f}{\partial y} = -5x - 5y^4 \qquad\qquad \frac{\partial f}{\partial x\partial y} = -5$$

$$\qquad\qquad\qquad \frac{\partial f}{\partial^2 y} = -20y^3 \qquad\qquad \frac{\partial f}{\partial y\partial x} = -5$$

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

אחר נחברם:

$$M(f)_{(x,y)} = \begin{pmatrix} \dfrac{\partial^2 f}{\partial^2 x} & \dfrac{\partial^2 f}{\partial x\partial y} \\ \dfrac{\partial^2 f}{\partial y\partial x} & \dfrac{\partial^2 f}{\partial^2 y} \end{pmatrix} = \begin{pmatrix} 6x & -5 \\ -5 & -20y^3 \end{pmatrix}$$

## 2.1.3 convexity

Based on Recitation 2

9. Prove that the intersection $C := \bigcap_{i \in I} C_i$ for $\{C_i : i \in I\}$ a collection of convex sets is convex.
10. Prove that the vector sum $C_1 + C_2 := \{c_1 + c_2 : c_1 \in C_1, c_2 \in C_2\}$ of two convex sets is convex.
11. Prove that the set $\lambda C := \{\lambda c : c \in C\}$ is convex, for any convex set $C$, and every scalar $\lambda$.

(9) יהיו $u,v \in \bigcap_{i \in I} C_i$. נראה שמתקיים שמתקיים $\alpha u + (1-\alpha)v \in C$ לכל $\alpha \in [0,1]$.

לכל $i \in I$ מתקיים $\bigcap_{i \in I} C_i \subseteq C_i$ מאחר ומתקיים $\alpha u + (1-\alpha)v \in C_i$ שכן $u,v \in C_i$ ומכאן שהקבוצה קמורה. ולכן מאחר וכל $C_i$ קמורה. לכן מתקיים $\alpha u + (1-\alpha)v \in C$

(10) יהי $u,v \in C_1 + C_2$ כלומר קיימים $c_1, c_3 \in C_1$, $c_2, c_4 \in C_2$ כך ש- $u = c_1 + c_2$, $v = c_3 + c_4$ מאחר ו- $C_1, C_2$ קבוצות קמורות מתקיים לכל $v \in C$

$$\alpha u + (1-\alpha)v = \alpha(C_1 + C_2) + (1-\alpha)(C_3 + C_4) = (\alpha C_1 + (1-\alpha)C_3) + (\alpha C_2 + (1-\alpha)C_4)$$

נשים לב ש- $\alpha C_1 + (1-\alpha)C_3 \in C_1$, $\alpha C_2 + (1-\alpha)C_4 \in C_2$ כי הקבוצות קמורות. לכן מתקבלת צורה קמורה.

(11) יהי $u,v \in \lambda C$. כלומר קיים $c_1, c_2 \in C$ כך ש- $u = \lambda c_1$, $v = \lambda c_2$ כך-

$$\alpha u + (1-\alpha)v = \lambda \alpha C_1 + (1-\alpha)\lambda C_2 = \lambda(\alpha C_1 + (1-\alpha)C_2) \Rightarrow \alpha u + (1-\alpha)v \in \lambda C$$
$$c \supseteq c \in C \text{ קמור קבוצה} \qquad \text{מהגדרה}$$

## 2.2 Estimation Theory

Based on Lecture 1 מ

12. 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 $\varepsilon$.

13. 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 $\Sigma \in \mathbb{R}^{d \times d}$. Provide an expression for the log-likelihood function of $\mathcal{N}(\mu, \Sigma)$. Develop the expression as much as you can. Hint: follow the approach used to derive the likelihood function for the univariate case.

(12)

נסמן $\mu$ בתוחלת של $\mathcal{P}$. קלא. לפי הגדרת $\varepsilon>0$:

$$P(|\hat{\mu}_n - \mu| > \varepsilon) \leq P(|\mu_n - E[\mu_n]| \geq \varepsilon) \leq \frac{\text{Var}[\mu_n]}{\varepsilon^2} = \frac{\text{Var}(x)}{n\varepsilon^2} \xrightarrow[n\to\infty]{} 0 \Rightarrow$$

מתכנס בהסתברות $\mu_n$ אשר רצינו להוכיח.

$E(\frac{1}{n}\sum_{i=1}^{n} x_i) = \frac{1}{n} \cdot n E(x_i) = E(x)$
$x \sim \mathcal{P}$

משוויון צ'בישב
בגלל כי קבוע
ומתקיים ש־$\varepsilon>0$-ו

$\text{Var}[\frac{1}{n}\sum_{i=1}^{n} x_i] = \frac{1}{n^2}\sum_{i=1}^{n} \text{Var}(x_i) = \frac{\text{Var}(x)}{n}$
כי $x_1, \ldots, x_n$    $x \sim \mathcal{P}$

$$f(X) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left\{-\frac{1}{2}(X - \mu)^\top \Sigma^{-1} (X - \mu)\right\}$$

$X = \begin{pmatrix} | & & | \\ x_1 & \cdots & x_m \\ | & & | \end{pmatrix} \in M_{d \times m}$

נסמן (13)

$$L(\mu, \Sigma \mid x_1, \ldots, x_m) = f_\mu(x_1, \ldots, x_m) = \prod_{i=1}^{m} f_\mu(x_i) = \prod_{i=1}^{m} \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left(-\frac{1}{2}\overset{t}{(x_i - \mu)}\overset{-1}{\Sigma}(x_i - \mu)\right)$$

$1 \times d \quad d \times d \quad d \times 1$

$$= \left((2\pi)^d |\Sigma|\right)^{-\frac{m}{2}} \cdot \exp\left(\sum_{i=1}^{m}\left(-\frac{1}{2}(x_i - \mu)^t \Sigma^{-1}(x_i - \mu)\right)\right) = \left((2\pi)^d |\Sigma|\right)^{-\frac{m}{2}} \cdot \exp\left(-\frac{1}{2}\sum_{i=1}^{m}(x_i - \mu)^t \Sigma^{-1}(x_i - \mu)\right)$$

$e^{x+y} = e^x \cdot e^y$

נעבור לאת $\log L(\mu, \Sigma \mid x_1, \ldots, x_m)$:

$$\log\left(\left((2\pi)^d |\Sigma|\right)^{-\frac{m}{2}} \cdot \exp\left(-\frac{1}{2}\sum_{i=1}^{m}(x_i - \mu)^t \Sigma^{-1}(x_i - \mu)\right)\right) = -\frac{m}{2}\log\left((2\pi)^d |\Sigma|\right) - \frac{1}{2}\sum_{i=1}^{m}(x_i - \mu)^t \Sigma^{-1}(x_i - \mu)$$

$\log x \cdot y = \log x + \log y$
$\log a^r = r \log a$
$\log(\exp(x)) = x$

$$= -\frac{m}{2}\left(\log((2\pi)^d) + \log(|\Sigma|)\right) - \frac{1}{2}\sum_{i=1}^{m}(x_i - \mu)^t \Sigma^{-1}(x_i - \mu) =$$

$\log x \cdot y = \log x + \log y$
$\log a^r = r \log a$

$$= -\frac{m}{2} \cdot d\log(2\pi) - \frac{m}{2}\log(|\Sigma|) - \frac{1}{2}\sum_{i=1}^{m}(x_i - \mu)^t \Sigma^{-1}(x_i - \mu) = -\frac{1}{2}\left(m(d\log(2\pi) + \log(|\Sigma|)) + \sum_{i=1}^{m}(x_i - \mu)^t \Sigma^{-1}(x_i - \mu)\right)$$
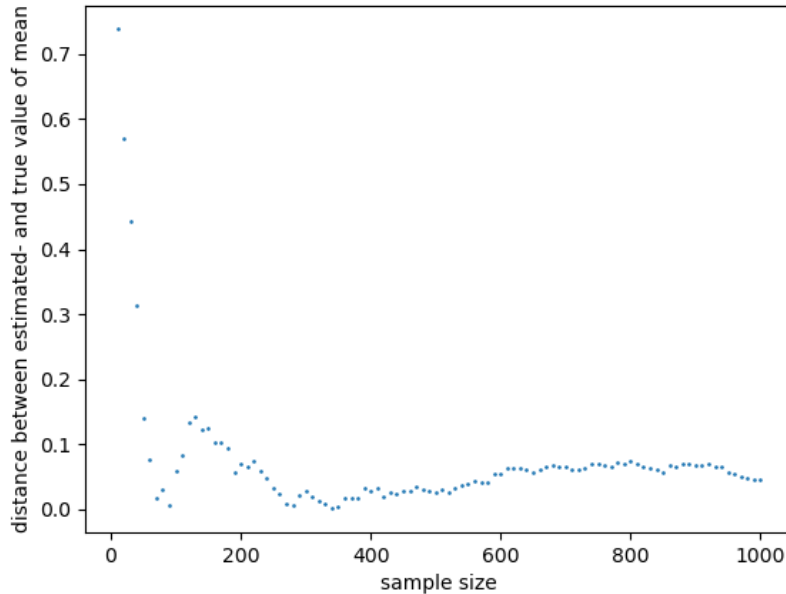
$r\log a = \log a^r$

(1

(2

Question 2- deviation of estimated expectation from true expectation
as function of sample sizefor samples from N(10,1) distribution



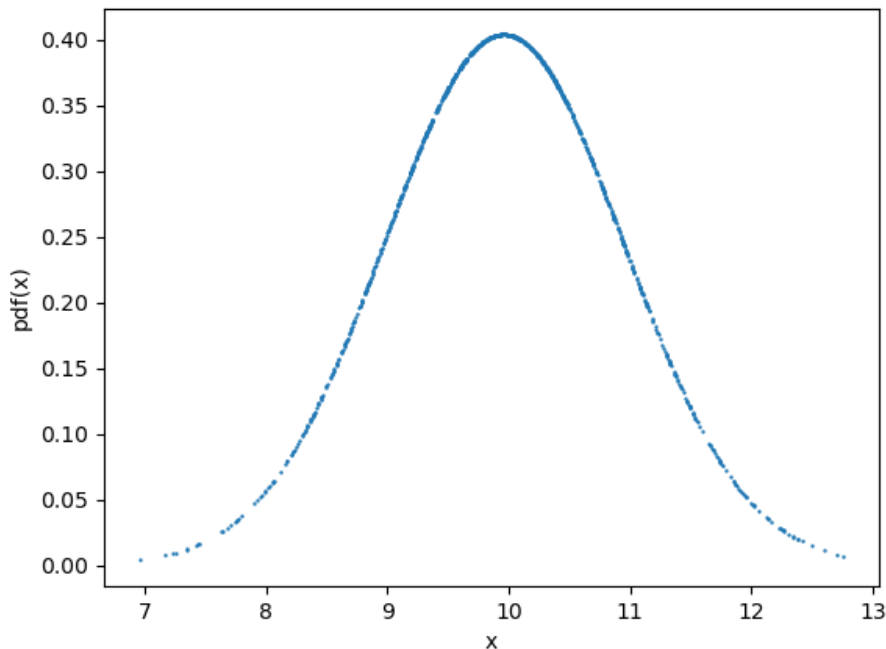what are you expecting to see in the plot?

(3

הגרלנו את וקטור ה-X מתוך התפלגות (10,1)N ולכן נצפה לראות צורת פעמון באופן
דומה לפונקציות צפיפות נורמליות, עם מרכוז סביב 10 (התוחלת) וצפיפות של
נקודות בעיקר סביב התוחלת. זאת מכיוון שכזכור מקורס הסתברות שלקחנו בקירוב
68% מהדגימות בהתפלגות נורמלית נמצאות במרחק סטיית תקן אחת לכל היותר
מהתוחלת, ובערך 95% במרחק שתי סטיות תקן מהתוחלת לכל היותר. מכיוון שהשונות
קטנה (1) נצפה לראות נקודות בעיקר בתחום [9,11].

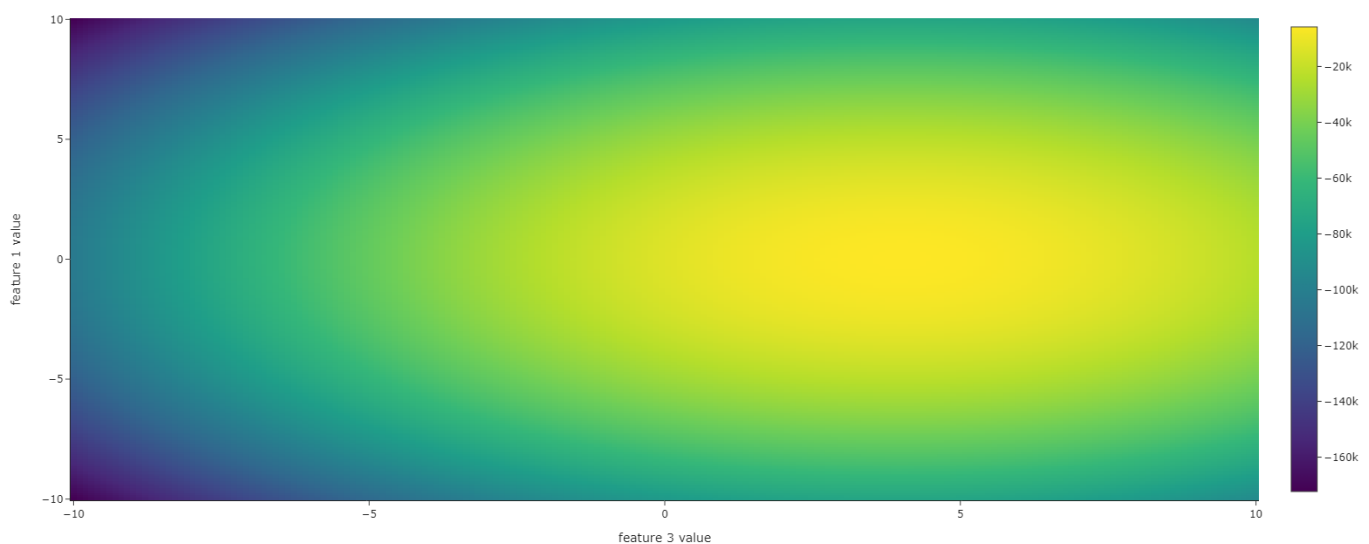Question 3- PDF values of samples from N(10,1) distribution

mean =

COV =

```
[-0.023 -0.043  3.993 -0.02 ]
[[ 0.917  0.166 -0.03   0.463]
 [ 0.166  1.974 -0.006  0.046]
 [-0.03  -0.006  0.98  -0.02 ]
 [ 0.463  0.046 -0.02   0.973]]
```

(5

Question 5- Log likelihood of Multivariate Gaussian as function of mean's features 1,3



What are you able to learn from the plot?

נבחין שהאזור הצהוב, בו ה-log- likelihoodגבוה יותר, הוא האזור שבו ה-features קרובים יותר
לערך האמת שלהם בתוחלת (f1= 0, f3= 4).
כלומרניתן ללמוד מכך ששערוך לפי עיקרון הלמידה של מקסום פונקציית ה-log-likelihood היא
בחירה טובה(נזכור שהגענו לחפש מקסימום לפונקציית ה- log-likelihood מתוך מטרה לחפש את ה
MLE-).

(6

```
features 1,3 that maximize the log-likelihood are:
f1= -0.05 ,f3= 3.97
```