



67577 נמי פאראכט נמי:

1 נמי נמי:

2 נמי נמי:

211975453 נמי:

Introduction to Machine Learning (67577)

Exercise 1 Estimation Theory & Mathematical Background

Second Semester, 2023

Contents

1	Submission Instructions	2
2	Theoretical Part	2
2.1	Mathematical Background	2
2.1.1	Linear Algebra	2
2.1.2	Multivariate Calculus	2
2.1.3	convexity	3
2.2	Estimation Theory	3
3	Practical Part	3
3.1	Univariate Gaussian Estimation	3
3.2	Multivariate Gaussian Estimation	4

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 $\forall x \in V \quad \|Ax\| = \|x\|$.
2. Calculate the SVD of the following matrix A . That is, find the matrices U, Σ, V^\top where U, V are orthogonal matrices and Σ 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, Σ and then calculate AA^\top to deduce U . Equivalently, once we deduced V, Σ we can find 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 $\text{rank}(A) = 1$. That is, show that all rows (or columns) in A are linearly dependent.
4. Show that for any orthonormal basis $(\mathbf{u}_1, \dots, \mathbf{u}_n)$ and any arbitrary 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 $(\mathbf{u}_1, \dots, \mathbf{u}_n)$, 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 $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$:

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

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 $S : \mathbb{R}^d \rightarrow [0, 1]^k$

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

8. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ 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 λ .

2.2 Estimation Theory

Based on Lecture 1

12. Let $x_1, x_2, \dots \stackrel{iid}{\sim} \mathcal{P}$ 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* (find the definition in the course book, page 14, Definition 1.1.10 under "Consistency"). Hint: for any given fixed value of $n \in \mathbb{N}$ bound from above the probability of deviating more than ε .
13. Let $\mathbf{x}_1, \dots, \mathbf{x}_m \stackrel{iid}{\sim} \mathcal{N}(\mu, \Sigma)$ 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 $x_1, \dots, x_{1000} \stackrel{iid}{\sim} \mathcal{N}(10, 1)$ 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 $\mathbf{x}_1, \dots, \mathbf{x}_{1000} \stackrel{iid}{\sim} \mathcal{N}(\mu, \Sigma)$

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

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 f_1 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

הוכחה 1

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

A - ה $T : V \rightarrow W$ הוא איזומטריה, כלומר $\|Ax\| = \|x\|$.

$\forall x \in V \quad \|Ax\| = \|x\|$ משום ש A אורתוגונלית.

הוכחה:

$$\|Ax\|^2 = \langle Ax | Ax \rangle = (Ax)^T \cdot (Ax) = \underbrace{x^T}_{\text{טבורי}} \underbrace{A^T A}_{\text{טבורי}} \underbrace{(Ax)}_{\text{טבורי}} = x^T \cdot I \cdot x = x^T \cdot x = \langle x | x \rangle = \|x\|^2$$

□

$\forall x \in V \quad \|Ax\| = \|x\|$, $\|Ax\|^2 = \|x\|^2$.

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

.2

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

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

: A le SVD \rightarrow jk den

: $A^\top A$ le EVD \rightarrow jk den \Rightarrow svd

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

: $A^\top A$ le svd jk k3v

$$\det(A^\top A - \lambda I) = \det \begin{pmatrix} 2-\lambda & 0 & 2 \\ 0 & 2-\lambda & -2 \\ 2 & -2 & 4-\lambda \end{pmatrix} = -\lambda^3 + 8\lambda^2 - 12\lambda = 0 \Leftrightarrow \lambda_1 = 0, \lambda_2 = 2, \lambda_3 = 6$$

: svd jk k3v

$$(A^\top A - \lambda I) \cdot V = 0 \quad : \text{svd} \quad V = \begin{pmatrix} V_1 \\ V_2 \\ V_3 \end{pmatrix} \quad \text{jk k3v}$$

$$(A^\top A - \lambda I) = (A^\top A - 0 \cdot I) = A^\top A = \begin{bmatrix} 2 & 0 & 2 \\ 0 & 2 & -2 \\ 2 & -2 & 4 \end{bmatrix} \quad : \lambda = 0 \quad \text{svd}$$

$$\begin{bmatrix} 2 & 0 & 2 \\ 0 & 2 & -2 \\ 2 & -2 & 4 \end{bmatrix} \begin{pmatrix} V_1 \\ V_2 \\ V_3 \end{pmatrix} = \begin{pmatrix} 2V_1 + 2V_3 \\ 2V_2 - 2V_3 \\ 2V_1 - 2V_2 + 4V_3 \end{pmatrix} = 0 \Leftrightarrow \begin{cases} V_1 = -V_3 \\ V_2 = V_3 \end{cases} = V_1 = -V_2$$

. $\lambda = 0 \quad -\delta \quad \text{svd} \quad \text{svd} \quad \text{svd} \quad \begin{pmatrix} 1 \\ -1 \\ 1 \end{pmatrix} \quad \text{svd} \quad \text{svd}$

$$(A^\top A - \lambda I) = (A^\top A - 2 \cdot I) = \begin{bmatrix} 2-2 & 0 & 2 \\ 0 & 2-2 & -2 \\ 2 & -2 & 4-2 \end{bmatrix} = \begin{bmatrix} 0 & 0 & 2 \\ 0 & 0 & -2 \\ 2 & -2 & 2 \end{bmatrix} \quad : \lambda = 2 \quad \text{svd}$$

$$\begin{bmatrix} 0 & 0 & 2 \\ 0 & 0 & -2 \\ 2 & -2 & 2 \end{bmatrix} \begin{pmatrix} V_1 \\ V_2 \\ V_3 \end{pmatrix} = \begin{pmatrix} 2V_3 \\ -2V_3 \\ 2V_1 - 2V_2 + 2V_3 \end{pmatrix} = 0 \Leftrightarrow \begin{cases} V_3 = 0 \\ 2V_1 - 2V_2 + 2V_3 = 0 \end{cases} \Leftrightarrow \begin{cases} V_3 = 0 \\ V_1 = V_2 \end{cases}$$

2. סעיפים 2 ו-3

$$\lambda = 2 - \sqrt{6} \quad \text{לפניהם}$$

$$(A^T A - \lambda I) = (A^T A - 6 \cdot I) = \begin{bmatrix} 2-6 & 0 & 2 \\ 0 & 2-6 & -2 \\ 2 & -2 & 4-6 \end{bmatrix} = \begin{bmatrix} -4 & 0 & 2 \\ 0 & -4 & -2 \\ 2 & -2 & -2 \end{bmatrix} \quad : \lambda = 6 \quad \text{ריבועי}$$

$$\begin{bmatrix} -4 & 0 & 2 \\ 0 & -4 & -2 \\ 2 & -2 & -2 \end{bmatrix} \begin{pmatrix} V_1 \\ V_2 \\ V_3 \end{pmatrix} = \begin{pmatrix} -4V_1 + 2V_3 \\ -4V_2 - 2V_3 \\ 2V_1 - 2V_2 - 2V_3 \end{pmatrix} = 0 \iff \begin{cases} V_1 = -V_2 \\ 2V_1 - 2V_2 - 2V_3 = 0 \end{cases}$$

$$\lambda = 6 \quad -\sqrt{6} \quad \text{לפניהם}$$

$$V_3 = \begin{pmatrix} 1/\sqrt{3} \\ -1/\sqrt{3} \\ 1/\sqrt{3} \end{pmatrix}, \quad V_2 = \begin{pmatrix} 1/\sqrt{2} \\ 1/\sqrt{2} \\ 0 \end{pmatrix}, \quad V_1 = \begin{pmatrix} 1/\sqrt{6} \\ -1/\sqrt{6} \\ 2/\sqrt{6} \end{pmatrix} \quad \text{מכיון ש } V_i \text{ יתקיים}$$

$$A^T A = V D V^T \quad \text{בנוסף} \quad D = \begin{bmatrix} 6 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 0 \end{bmatrix}, \quad V = \begin{bmatrix} 1/\sqrt{6} & 1/\sqrt{2} & 1/\sqrt{3} \\ -1/\sqrt{6} & 1/\sqrt{2} & -1/\sqrt{3} \\ 2/\sqrt{6} & 0 & 1/\sqrt{3} \end{bmatrix} \quad : \text{הוכן}$$

$$u_i = \frac{Av_i}{\sqrt{\lambda_i}} \quad \text{כאות כוונתית ב-SVD}$$

$$u_1 = \frac{Av_1}{\sqrt{\lambda_1}} = \frac{1}{\sqrt{6}} \cdot \begin{pmatrix} 0 \\ 6/\sqrt{6} \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \quad : \text{רמז}$$

$$u_2 = \frac{Av_2}{\sqrt{\lambda_2}} = \frac{1}{\sqrt{2}} \cdot \begin{pmatrix} 2/\sqrt{2} \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$$

(u_3 פיך פשוט ו-ORTHOGONAL)

$$A = U \Sigma V^T \quad \text{בנוסף} \quad \Sigma = \begin{bmatrix} \sqrt{6} & 0 & 0 \\ 0 & \sqrt{2} & 0 \\ 0 & 0 & 0 \end{bmatrix} \quad -1 \quad U = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}, \quad \text{הוכן}$$

ב證據 A ב-SVD ->, מוכח

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 $\text{rank}(A) = 1$. That is, show that all rows (or columns) in A are linearly dependent. .3

- ב' $\exists A \in \mathbb{R}^{n \times m}$ גורם $\mathbf{v} \in \mathbb{R}^n$, $\mathbf{u} \in \mathbb{R}^m$ כך ש- \mathbf{v} ו- \mathbf{u} מוגדרים על ידי $\mathbf{v} \otimes \mathbf{u}$ או $\mathbf{v} \cdot \mathbf{u}^\top$ ן.3

$$\text{rank}(A) = 1$$

לפיכך כ' מכך נובע

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

: ב' מילא י' כ' כך כי $i \neq j$ $\exists k \in [n]$ כך ש- \mathbf{v}_k ו- \mathbf{u}_j מוגדרים על ידי $\mathbf{v} \cdot \mathbf{u}^\top$ ן.3

$$(v_j)_k = v_k u_j, \quad (v_i)_k = v_k \cdot u_i : k \in [n] \quad \text{סב' סעיפים ו' ו' נובע}$$

$$\text{, ב' מילא } \alpha_i - 1 \quad \alpha_j \text{ כ' מילא } \alpha_j = v \cdot u_j, \quad \alpha_i = v \cdot u_i \quad \text{ונכון}$$

$$\text{. סב' } \text{rank}(A) = 1 \quad / \sigma /$$

■

Show that for any orthonormal basis $(\mathbf{u}_1, \dots, \mathbf{u}_n)$ and any arbitrary 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 $(\mathbf{u}_1, \dots, \mathbf{u}_n)$, is the inner product between \mathbf{x} and \mathbf{u}_i .

4

$$\text{Show that } \mathbf{x} = \sum_{i=1}^n a_i \mathbf{u}_i \text{ for any } \mathbf{x} \in \mathbb{R}^n \text{ and } (\mathbf{u}_1, \dots, \mathbf{u}_n) \text{ orthonormal basis.}$$

$$\forall i \in [n] \quad a_i = \langle \mathbf{x}, \mathbf{u}_i \rangle$$

$$\forall i \in [n] \quad \text{Show that } \mathbf{x} = \sum_{i=1}^n a_i \mathbf{u}_i \text{ for any } \mathbf{x} \in \mathbb{R}^n \text{ and orthonormal basis.}$$

$$\begin{aligned} \langle \mathbf{x}, \mathbf{u}_i \rangle &= \left\langle \sum_{j=1}^n a_j \mathbf{u}_j, \mathbf{u}_i \right\rangle = a_1 \langle \mathbf{u}_1, \mathbf{u}_i \rangle + a_2 \langle \mathbf{u}_2, \mathbf{u}_i \rangle + \dots + a_n \langle \mathbf{u}_n, \mathbf{u}_i \rangle = a_i \langle \mathbf{u}_i, \mathbf{u}_i \rangle = a_i \cdot 1 = a_i \end{aligned}$$

$$\therefore a_i = \langle \mathbf{x}, \mathbf{u}_i \rangle$$

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 $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$:

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

.5

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

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

$$(\mathcal{J}_x(f))_{i,j} = \frac{\partial f_i(x)}{\partial x_j} : 2 \quad \text{לפונקציית גזע}$$

$$(U \cdot D \cdot U^\top)x = \sum_{k=1}^n \lambda_k \langle x | u_k \rangle u_k \quad \text{המוגדר ב-1}$$

$$D = \text{diag}(\sigma) \quad \text{-1. גורם גזע } \lambda_1, \dots, \lambda_n \quad \text{בנוסף, } U \cdot D \cdot U^\top \text{ מוגדר ב-1}$$

$$f(\sigma) = (U \cdot \text{diag}(\sigma) \cdot U^\top)x = \sum_{k=1}^n \lambda_k \langle x | u_k \rangle u_k \quad \text{המוגדר}$$

$\text{diag}(\sigma)$ הוא גורם גזע נורמלי. $\lambda_1, \dots, \lambda_n$ נורמיים. $D = \text{diag}(\sigma)$ הוא גורם גזע נורמי.

$$f(\sigma) = (U \cdot \text{diag}(\sigma) \cdot U^\top)x = \sum_{k=1}^n \sigma_k \langle x | u_k \rangle u_k$$

$$(\mathcal{J}_\sigma(f))_{i,j} = \frac{\partial f_i(x)}{\partial \sigma_j} = \frac{\partial}{\partial \sigma_j} \left(\sum_{k=1}^n \sigma_k \langle x | u_k \rangle u_k \right)_i = \underbrace{[\langle x | u_j \rangle u_j]_i}_{\substack{k \neq j \\ \text{אך}}} = \underbrace{[(x^\top u_j) \cdot u_j]_i}_{\substack{\text{ו-} \\ \text{אך}}}$$

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

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

$$f: \mathbb{R}^n \rightarrow \mathbb{R}^n \quad \text{def} \quad f(\sigma) = (\nabla \cdot \text{diag}(\sigma) \nabla^t) x \quad \text{so } \sigma \text{ is a } 1 \times k$$

$$\cdot S(x) = \frac{1}{2} \|x\|^2, \quad T(\sigma) = f(\sigma) - y \quad \text{def}$$

$$h(\sigma) = (S \circ T)(\sigma) = S(f(\sigma) - y) = \frac{1}{2} \|f(\sigma) - y\|^2 \quad \text{def}$$

$$\text{def} \quad J_x(S \circ T)_{i,j} = \frac{\partial}{\partial y_k(x)} S_i(T(x)) \cdot \frac{\partial}{\partial x_j} T(x) \quad \text{so we can use the chain rule to get}$$

$$\cdot J_x(S \circ T) = J_{T(x)}(S) \cdot J_x(T)$$

we can do this:

$$S'(x) = \left(\frac{1}{2} \|x\|^2 \right)' = \frac{1}{2} \cdot 2x^t = x^t \quad \text{so}, \quad 2x^t \text{ is the } \|x\|^2 \text{ function}$$

$$\cdot J_{T(x)}(S) = (T(x))^t \quad \text{def}$$

$$\therefore J_{\sigma}(f) = \nabla \cdot \text{diag}(\nabla^t x) \quad \text{so we get to}$$

$$J_x(T) = T'(\sigma) = (f(\sigma) - y)' = \underbrace{f'(\sigma)}_{\text{def}} = \nabla \cdot \text{diag}(\nabla^t x)$$

$$J_{\sigma}(h) = J_{\sigma}(S \circ T) = J_{T(\sigma)}(S) \cdot J_{\sigma}(T) = (T(x))^t \cdot \nabla \cdot \text{diag}(\nabla^t x) \quad \text{def}$$

$$\therefore \nabla h(\sigma) = (J_{\sigma}(h))^t \quad \text{so we can now calculate } \nabla h(\sigma)$$

$$\nabla h(\sigma) = (J_{\sigma}(h))^t = ((T(x))^t \cdot \nabla \cdot \text{diag}(\nabla^t x))^t = (\nabla \cdot \text{diag}(\nabla^t x))^t \cdot T(x) = (J_{\sigma}(f))^t \cdot f(\sigma) - y$$

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

-7

$$S : \mathbb{R}^d \rightarrow [0, 1]^d \quad \text{softmax} \leftarrow \quad S(\mathbf{x})_j = \frac{e^{x_j}}{\sum_{l=1}^k e^{x_l}}$$

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

$$S(\mathbf{x})_j = \frac{g_j(\mathbf{x})}{h(\mathbf{x})} \quad \text{where } h(\mathbf{x}) = \sum_{k=1}^d g_k(\mathbf{x}), \quad g_i(\mathbf{x}) = e^{x_i}$$

$$\frac{\partial}{\partial x_j} \frac{e^{x_i}}{\sum_{k=1}^d e^{x_k}} = S_i(1-S_j) \quad \text{if } i=j \quad \text{else} \quad \text{if } i \neq j$$

$$(J_x(S))_{i,j} = \frac{\partial}{\partial x_i} S_j = \frac{\partial}{\partial x_i} \frac{e^{x_i}}{\sum_{k=1}^d e^{x_k}} = \frac{\partial}{\partial x_i} \frac{g_i}{h} = \frac{g'_i \cdot h - h' \cdot g_i}{h^2}$$

$$\frac{\partial}{\partial x_j} g_i = \frac{\partial}{\partial x_j} e^{x_i} \underset{i \neq j}{=} 0$$

$$\frac{\partial}{\partial x_j} h = \frac{\partial}{\partial x_j} \sum_{k=1}^d g_k(\mathbf{x}) = e^{x_j} \rightarrow i \neq j$$

$$(J_x(S))_{i,j} = \frac{g'_i \cdot h - h' \cdot g_i}{h^2} = \frac{0 \cdot h - e^{x_i} \cdot g_i}{h^2} = \frac{-e^{x_j} \cdot e^{x_i}}{\left(\sum_{k=1}^d g_k(\mathbf{x})\right)^2} = \frac{-e^{x_j}}{\sum_{k=1}^d g_k(\mathbf{x})} \cdot \frac{e^{x_i}}{\sum_{k=1}^d g_k(\mathbf{x})} = -S_j \cdot S_i$$

$$(J_x(S))_{i,j} = \begin{cases} S_i(1-S_j) & i=j \\ -S_j \cdot S_i & i \neq j \end{cases} \quad \text{softmax}$$

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

.8

∴ $\text{defn. } H(f) \in \mathbb{R}^{2 \times 2} \quad | \circ f : \mathbb{R}^2 \rightarrow \mathbb{R} \quad \text{given, } f(x, y) = x^3 - 5xy - y^5$

$$[H(f)]_{1,1} = \frac{\partial^2}{\partial x \partial x} f(x, y) = \frac{\partial^2}{\partial x \partial x} x^3 - 5xy - y^5 = \frac{\partial}{\partial x} 3x^2 - 5y = 6x$$

$x \downarrow \text{if } x \neq 0$ $x \downarrow \text{if } x \neq 0$

$$[H(f)]_{1,2} = \frac{\partial^2}{\partial x \partial y} f(x, y) = \frac{\partial^2}{\partial x \partial y} x^3 - 5xy - y^5 = \frac{\partial}{\partial x} -5x - 5y^4 = -5$$

$y \downarrow \text{if } y \neq 0$ $y \downarrow \text{if } y \neq 0$

$$[H(f)]_{2,1} = \frac{\partial^2}{\partial y \partial x} f(x, y) = \frac{\partial^2}{\partial y \partial x} x^3 - 5xy - y^5 = \frac{\partial}{\partial y} 3x^2 - 5y = -5$$

$x \downarrow \text{if } x \neq 0$ $y \downarrow \text{if } y \neq 0$

$$[H(f)]_{2,2} = \frac{\partial^2}{\partial y \partial y} f(x, y) = \frac{\partial^2}{\partial y \partial y} x^3 - 5xy - y^5 = \frac{\partial}{\partial y} -5x - 5y^4 = -20y^3$$

$x \downarrow \text{if } x \neq 0$ $y \downarrow \text{if } y \neq 0$

$$H(f) = \begin{bmatrix} 6x & -5 \\ -5 & -20y^3 \end{bmatrix} : \text{Ans}$$

Prove that the intersection $C := \bigcap_{i \in I} C_i$ for $\{C_i : i \in I\}$ a collection of convex sets is convex.

.9

$\alpha \cdot v + (1-\alpha)u \in C$ $\forall u, v \in [0, 1]$ $\text{by } u, v \in C = \bigcap_{i \in I} C_i$ defn

$u, v \in C_j$ $\forall j \in I$ $\text{by } C_j \text{ is convex}$, $u, v \in C = \bigcap_{i \in I} C_i$

$(C_j \subseteq C \text{ and } \alpha \cdot v + (1-\alpha)u \in C_j \text{ for all } \alpha \in [0, 1]) \Rightarrow \alpha \cdot v + (1-\alpha)u \in C_j$

Thus C is convex.

✓

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. 10

$$\alpha \cdot v + (1-\alpha)u \in C_1 + C_2 \quad \text{for all } \alpha \in [0, 1] \quad \text{by } u, v \in C_1 + C_2 \quad \text{by } \boxed{\text{def}} \quad \boxed{\text{def}}$$

Given $v_1, v_2 \in C_1$ and $u_1, u_2 \in C_2$

$$v_1 + v_2 = v, \quad u_1 + u_2 = u \quad \text{so} \quad v_1, v_2 \in C_1 \quad \text{and} \quad u_1, u_2 \in C_2 \quad \text{so } v \in C_1, u \in C_2 \quad \text{by def}$$

$$\alpha \cdot v + (1-\alpha)u = \alpha(v_1 + v_2) + (1-\alpha)(u_1 + u_2) = \alpha v_1 + \alpha v_2 + (1-\alpha)u_1 + (1-\alpha)u_2 = \alpha v_1 + (1-\alpha)u_1 + \alpha v_2 + (1-\alpha)u_2$$

$$\alpha \cdot v + (1-\alpha)u = \alpha v_1 + (1-\alpha)u_1 + \alpha v_2 + (1-\alpha)u_2 \in C_1 + C_2, \quad \alpha v_1 + (1-\alpha)u_1 \in C_1 \quad \text{so by def } \boxed{\text{def}}$$

$$\alpha \cdot v + (1-\alpha)u = \alpha v_1 + (1-\alpha)u_1 + \alpha v_2 + (1-\alpha)u_2 \in C_1 + C_2 \quad : \alpha \in [0, 1] \quad \text{by def}, \quad \text{by def}$$

QED

Prove that the set $\lambda C := \{\lambda c : c \in C\}$ is convex, for any convex set C , and every scalar λ . ||

$$\lambda \in \mathbb{R} \quad . \quad u \cdot v + (1-u)v \in \lambda \cdot C \quad \text{forall } u \in [0,1] \quad \text{bfs} \quad u, v \in \lambda \cdot C \quad \text{bfs}$$

$\text{N}(\Gamma)$ \subset C_1 .

$$\lambda \bar{u} = u, \lambda \bar{v} = v \quad -\text{if } \exists \bar{u}, \bar{v} \in C \text{ such that } u, v \in \lambda C$$

$\vdash \forall x N \Leftarrow [0, 1] \quad \text{fs}$

$$\begin{aligned} \lambda \cdot v + (1-\lambda) \cdot u &= \lambda \cdot \lambda \bar{v} + (1-\lambda) \cdot \lambda \bar{u} = \lambda (\lambda \bar{v} + (1-\lambda) \cdot \bar{u}) \in \lambda C \\ &\downarrow \\ \lambda \cdot \bar{v} + (1-\lambda) \cdot \bar{u} &\in C \end{aligned}$$

כ' כ מאריך

$\lambda \in \mathbb{R} \cup \{\infty\}$, $\alpha \in [0, 1]$

1

• בראכן הינה לא כימי

Let $x_1, x_2, \dots \stackrel{iid}{\sim} \mathcal{P}$ 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* (find the definition in the course book, page 14, Definition 1.1.10 under "Consistency"). Hint: for any given fixed value of $n \in \mathbb{N}$ bound from above the probability of deviating more than ε . 12

$$\cdot \lim_{n \rightarrow \infty} P(|\hat{\mu}_n - \mu| > \varepsilon) = 0 \quad 0 < \varepsilon \text{ so } \underline{\text{def}} \quad \underline{\text{def}}$$

(unbiased estimator) Given by part i.) $\hat{\mu}_n$ is the sample mean. $\text{var}(x_i) = \sigma^2$ so

$$\text{as desired we have p.s. } |\hat{\mu}_n - \mu| = |\hat{\mu}_n - E(\hat{\mu}_n)| \quad \text{so} \quad E(\hat{\mu}_n) = \mu \quad \text{p.s.}$$

$$P(|\hat{\mu}_n - \mu| > \varepsilon) = P(|\hat{\mu}_n - E(\hat{\mu}_n)| > \varepsilon) \leq \frac{\text{Var}(\hat{\mu}_n)}{\varepsilon^2} \quad \text{so by def}$$

$$P(|\hat{\mu}_n - \mu| > \varepsilon) \leq \frac{\text{Var}(\hat{\mu}_n)}{\varepsilon^2} = \frac{\frac{\sigma^2}{n}}{\varepsilon^2} \quad \text{so} \quad \text{Var}(\hat{\mu}_n) = \frac{\sigma^2}{n} \quad \text{so by def}$$

$$\cdot P(|\hat{\mu}_n - \mu| > \varepsilon) \leq \frac{\sigma^2}{\varepsilon^2 \cdot n} \quad \text{so by def}$$

$$0 \leftarrow 0 \leq P(|\hat{\mu}_n - \mu| > \varepsilon) \leq \frac{\sigma^2}{\varepsilon^2 \cdot n} \xrightarrow{n \rightarrow \infty} 0 \quad \text{so by def, } \frac{\sigma^2}{\varepsilon^2 \cdot n} \xrightarrow{n \rightarrow \infty} 0$$

so by def

$$\therefore \lim_{n \rightarrow \infty} P(|\hat{\mu}_n - \mu| > \varepsilon) = 0 \quad \text{so by def}$$

■

Let $\mathbf{x}_1, \dots, \mathbf{x}_m \stackrel{iid}{\sim} \mathcal{N}(\mu, \Sigma)$ 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.

.13

(μ , σ^2) $X_1, \dots, X_m \stackrel{\text{i.i.d}}{\sim} N(\mu, \sigma^2)$: $p(x)$

$\Theta = (\mu, \Sigma)$ proj. $N(\mu, \Sigma)$ k log-likelihood \rightarrow unk bgn

• $\text{E}[\nu] = N(\mu, \Sigma)$ So גודלה גודל f גודל גודל

$$\mathcal{L}(\theta | x_1, \dots, x_m) = f_{\theta}(x_1, \dots, x_m) = \prod_{i=1}^m f_{\theta}(x_i) = \prod_{i=1}^m \frac{1}{\sqrt{(2\pi)^d \cdot |\Sigma|}} \exp\left(-\frac{1}{2}(x_i - \mu)^T \Sigma^{-1} (x_i - \mu)\right) =$$

$$\exp(x) \cdot \exp(y) = \exp(x+y) = \left(\frac{1}{\sqrt{(2\pi t)^d} |\Sigma|^{1/2}} \right)^M \cdot \exp \left(\sum_{i=1}^M \left(-\frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right) \right) =$$

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

↑ k əvənd əm log ſtos , (likelihood - function) $N(\mu, \sigma^2)$ ke ſtolkun jia is

: log-likelihood function \rightarrow

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

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

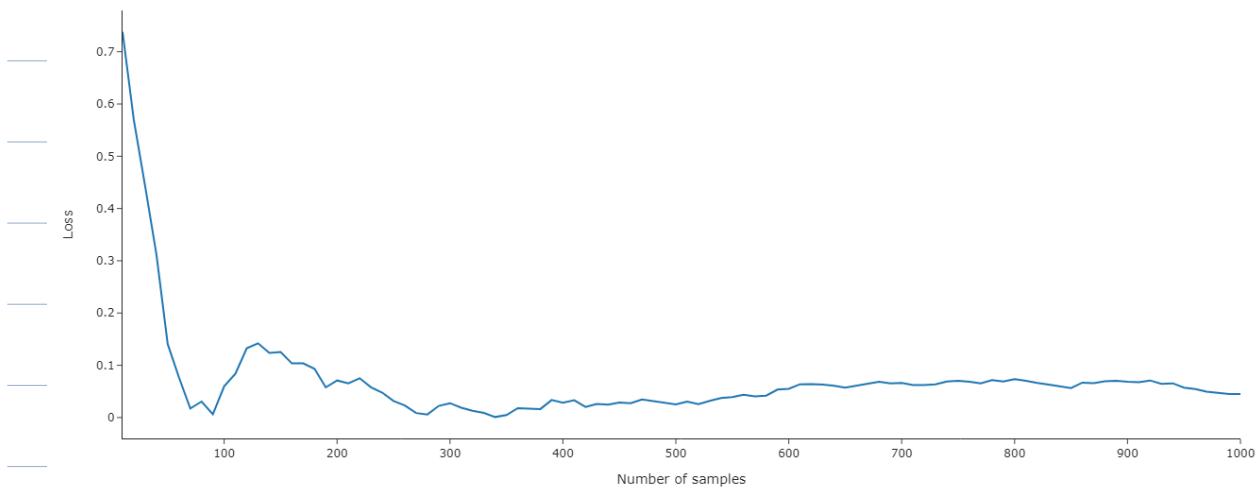
$$= -\frac{m}{2} \left(d \log(\lambda) t + \log(|\Sigma|) \right) + \sum_{i=1}^m \left(-\frac{1}{2} (x_i - \mu)^t \Sigma^{-1} (x_i - \mu) \right)$$

• Enlarge log-likelihood function - 7

loss $\beta\mathcal{D}$

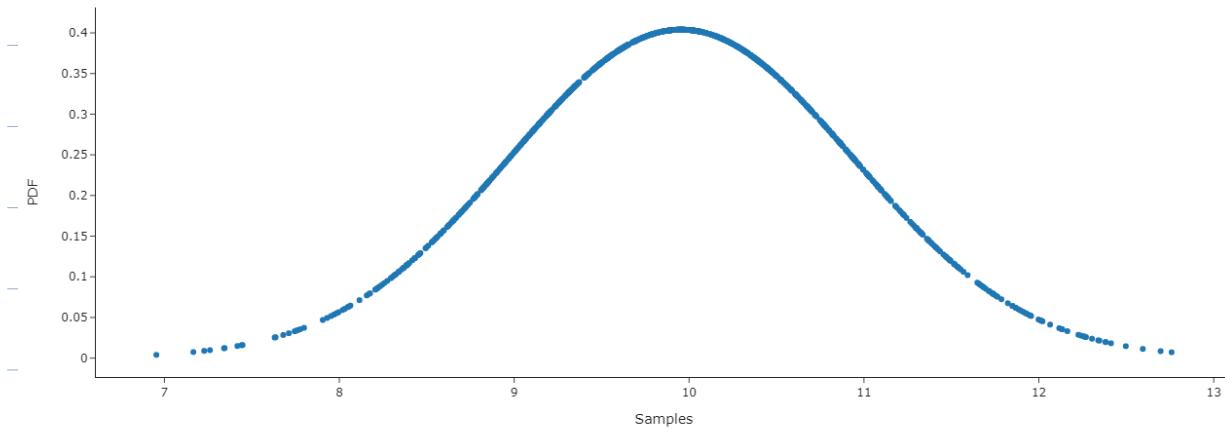
: 2 γ_{KLN} line 92

Losses as function of sample size



: 3 γ_{KLN} Scatter 92

PDF of 1000 $\sim N(10, 1)$ samples



5. Make heatmap f2

Likelihoods as function of f_3 and f_1

