

Yosef Haim Abraham

206208250

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. 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

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 $\forall x \in V \quad \|Ax\| = \|x\|$.

$$\begin{aligned} \|Ax\|^2 &= \langle Ax, Ax \rangle = (Ax)^T Ax = (x^T A^T) A x = \\ &= x^T x = \langle x, x \rangle = \|x\|^2 \end{aligned}$$

$\overbrace{A^T A = I_d}^{< \text{ from } A}$

Q. 63)

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.

$$A = \underbrace{\begin{bmatrix} 1 & 1 & 0 \\ 1 & -1 & 2 \end{bmatrix}}_{m \times n}$$

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 $AV = U\Sigma$.

$$A^\top A = \begin{bmatrix} 1 & 1 \\ 1 & -1 \\ 0 & 2 \end{bmatrix} \cdot \begin{bmatrix} 1 & 1 & 0 \\ 1 & -1 & 2 \end{bmatrix} = \begin{bmatrix} 2 & 0 & 2 \\ 0 & 2 & -2 \\ 2 & -2 & 4 \end{bmatrix} \quad : \text{ eigen}$$

$$AA^\top = \begin{bmatrix} 1 & 1 & 0 \\ 1 & -1 & 2 \\ 0 & 2 & 2 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & -1 \\ 0 & 2 \end{bmatrix} = \begin{bmatrix} 2 & 0 \\ 0 & 6 \end{bmatrix}$$

: AA^\top နဲ့ $A^\top A$ ပဲ မြတ်စွာ မျှတော်မူ မှုပါ။

$U_2 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, U_1 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$ မှာ မြတ်စွာ မြတ်စွာ မျှတော်မူ မှုပါ။ $\lambda_1 = 6, \lambda_2 = 2$

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

$$A^T A - 2I \quad \text{iff} \quad A^T A - 2I = 0 \quad (\Leftrightarrow)$$

$$\begin{bmatrix} 0 & 0 & 2 \\ 0 & 0 & -2 \\ 2 & -2 & 2 \end{bmatrix} \xrightarrow{\text{Row operations}} \begin{bmatrix} 1 & -1 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} \xrightarrow{\text{Row operations}} \begin{bmatrix} 1 & -1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix}$$

$$\begin{aligned} x - y &= 0 \\ z &= 0 \end{aligned} \quad \Leftarrow$$

$$\text{Ker}(A^T A - 2I) = \left\{ (t, t, 0) \mid t \in \mathbb{R} \right\} \quad \Leftarrow$$

$$= \text{Span}\{(1, 1, 0)\}$$

$$\therefore \text{N.P.} \quad | \text{P.R.} \quad A^T A - 6I = 0 \quad (\Leftrightarrow)$$

$$\begin{bmatrix} -4 & 0 & 2 \\ 0 & -4 & -2 \\ 2 & -2 & -2 \end{bmatrix} \xrightarrow{\text{Row operations}} \begin{bmatrix} 1 & 0 & -\frac{1}{2} \\ 0 & 1 & \frac{1}{2} \\ 1 & -1 & -1 \end{bmatrix} \xrightarrow{\text{Row operations}} \begin{bmatrix} 1 & 0 & -\frac{1}{2} \\ 0 & 1 & \frac{1}{2} \\ 0 & -1 & -\frac{1}{2} \end{bmatrix}$$

$$\xrightarrow{\text{Row operations}} \begin{bmatrix} 1 & 0 & -\frac{1}{2} \\ 0 & 1 & \frac{1}{2} \\ 0 & 0 & 0 \end{bmatrix}$$

$$y + \frac{1}{2}z = 0, \quad x - \frac{1}{2}z = 0 \quad \Leftarrow$$

$$\text{Ker}(A^T A - 6I) = \left\{ \left(\frac{t}{2}, -\frac{t}{2}, t \right) \mid t \in \mathbb{R} \right\} \quad \Leftarrow$$

$$= \text{Span}(\{(1, -1, 2)\})$$

$$\text{Defn: } A^T A = 0 \quad \text{rank } A = 1$$

$$\begin{bmatrix} 2 & 0 & 2 \\ 0 & 2 & -2 \\ 2 & -2 & 4 \end{bmatrix} \xrightarrow{\text{Row operations}} \begin{bmatrix} 2 & 0 & 2 \\ 0 & 2 & -2 \\ 0 & -2 & 2 \end{bmatrix} \xrightarrow{\text{Row operations}} \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & -1 \\ 0 & 0 & 0 \end{bmatrix}$$

$$\text{Ker}(A^T A) = \{(-t, t, t) \mid t \in \mathbb{R}\} = \text{Span}(\{(-1, 1, 1)\}) \subset$$

$$V = \begin{bmatrix} \frac{1}{\sqrt{6}} & \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{3}} \\ -\frac{1}{\sqrt{6}} & \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{3}} \\ \frac{2}{\sqrt{6}} & 0 & \frac{1}{\sqrt{3}} \end{bmatrix} \quad \text{since } V$$

$$A = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} \sqrt{6} & 0 & 0 \\ 0 & \sqrt{2} & 0 \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{6}} & \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{3}} \\ -\frac{1}{\sqrt{6}} & \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{3}} \\ \frac{2}{\sqrt{6}} & 0 & \frac{1}{\sqrt{3}} \end{bmatrix} \quad \text{is } V$$

$A \sim \text{SVD}$ can be used

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.

$$\begin{aligned}
 & \text{Def} \quad \mathbf{v} = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix}, \quad \mathbf{u} = \begin{bmatrix} u_1 \\ \vdots \\ u_m \end{bmatrix} \quad \text{Def} \\
 & \mathbf{v} \cdot \mathbf{u}^\top = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} \begin{bmatrix} u_1, \dots, u_m \end{bmatrix} = \begin{bmatrix} v_1 u_1 & v_1 u_2 & \cdots & v_1 u_m \\ v_2 u_1 & v_2 u_2 & \cdots & v_2 u_m \\ \vdots & \vdots & \ddots & \vdots \\ v_n u_1 & v_n u_2 & \cdots & v_n u_m \end{bmatrix} = \\
 & = \begin{bmatrix} | & & | \\ u_1 \cdot \mathbf{v} & \cdots & u_m \cdot \mathbf{v} \\ | & & | \end{bmatrix}
 \end{aligned}$$

$$\text{rank}(A) = 1 \quad \leftarrow \quad \text{Span } \mathbf{v} \cdot \mathbf{u}^\top = \text{Span } \{\mathbf{v}\} \quad \text{rank}$$

Q

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 .

Finne på 0022 3.8.13. 16/11 Br

$$\mathbf{x} = \sum_{i=1}^n \langle \mathbf{x}, \mathbf{u}_i \rangle \mathbf{u}_i$$

:p(1)

$$\sum_{i=1}^n \langle \mathbf{x}, \mathbf{u}_i \rangle \mathbf{u}_i = \sum_{i=1}^n \left\langle \sum_{j=1}^n a_j \mathbf{u}_j, \mathbf{u}_i \right\rangle \mathbf{u}_i =$$

$$= \sum_{i=1}^n \sum_{j=1}^n a_j \langle \mathbf{u}_j, \mathbf{u}_i \rangle \mathbf{u}_i =$$

$$= \sum_{i=1}^n a_i \|\mathbf{u}_i\|^2 \mathbf{u}_i = \sum_{i=1}^n a_i \cdot \mathbf{u}_i = \mathbf{x}$$

$$\langle \mathbf{u}_i, \mathbf{u}_j \rangle = 0 \Leftarrow i \neq j$$

$$\|\mathbf{u}_i\|^2 = 1$$

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. Hint: for any given fixed value of $n \in \mathbb{N}$ bound from above the probability of deviating more than ε .

$$\mathbb{E}\hat{\mu}_n = \frac{1}{n} \sum \mathbb{E}X_i = N \quad ! \quad \text{p.d.s.t. } \mathbb{E}X_i = N \text{ p.o.j.}, \varepsilon > 0 \quad \text{?}$$

$$\begin{aligned} P(|\hat{\mu}_n - N| \geq \varepsilon) &\leq \frac{\text{Var}(\hat{\mu}_n)}{\varepsilon^2} = \\ &= \frac{\text{Var}\left(\frac{1}{n} \sum X_i\right)}{\varepsilon^2} = \frac{1}{\varepsilon^2} \cdot \frac{1}{n^2} \cdot \sum_{i=1}^n \text{Var}(X_i) = \\ &= \frac{1}{\varepsilon^2} \cdot \frac{1}{n^2} \cdot n \cdot \text{Var}(X_n) = \frac{\text{Var}(X_n)}{\varepsilon^2 \cdot n} \xrightarrow{n \rightarrow \infty} 0 \end{aligned}$$

□ . e>3p

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.

$$\Theta = (\hat{\mu}, \hat{\Sigma}) \in \mathbb{R}^d \times \mathbb{R}_{>0}^{d \times d}$$

$$\mathcal{L}(\theta | \mathbf{x}_1, \dots, \mathbf{x}_m) = f_{\theta}(\mathbf{x}_1, \dots, \mathbf{x}_m) =$$

$\rightarrow \mathbf{x}_1, \dots, \mathbf{x}_m$

$$\begin{aligned} \downarrow \prod_{i=1}^m f_{\theta}(\mathbf{x}_i) &= \prod_{i=1}^m \frac{1}{\sqrt{(2\pi)^d \cdot |\hat{\Sigma}|}} \cdot \exp \left(-\frac{1}{2} (\mathbf{x}_i - \hat{\mu})^T \cdot \hat{\Sigma}^{-1} \cdot (\mathbf{x}_i - \hat{\mu}) \right) = \\ &= \left(\frac{1}{((2\pi)^d \cdot |\hat{\Sigma}|)^{\frac{1}{2}}} \right)^m \cdot \exp \left\{ \sum_{i=1}^m -\frac{1}{2} (\mathbf{x}_i - \hat{\mu})^T \cdot \hat{\Sigma}^{-1} \cdot (\mathbf{x}_i - \hat{\mu}) \right\} \end{aligned}$$

: $f_{\theta}(\mathbf{x})$ ergeben \log

$$\mathcal{L}^{\log}(\theta | \mathbf{x}_1, \dots, \mathbf{x}_m) = \sum_{i=1}^m -\frac{1}{2} (\mathbf{x}_i - \hat{\mu})^T \cdot \hat{\Sigma}^{-1} \cdot (\mathbf{x}_i - \hat{\mu}) - \frac{m}{2} \cdot \log((2\pi)^d \cdot |\hat{\Sigma}|)$$

□

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

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

$$\text{def } U = \begin{bmatrix} | & & | \\ u_1 & \dots & u_n \\ | & & | \end{bmatrix} \quad (\text{orth})$$

$$f(\sigma) = U \cdot \text{diag}(\sigma) \cdot U^\top \cdot x = \begin{bmatrix} | & & | \\ u_1 & \dots & u_n \\ | & & | \end{bmatrix} \cdot \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_n \end{bmatrix} \cdot \begin{bmatrix} -u_1^\top \\ \vdots \\ -u_n^\top \end{bmatrix} \cdot x =$$

$$= \begin{bmatrix} | & & | \\ \sigma_1 u_1 & \dots & \sigma_n u_n \\ | & & | \end{bmatrix} \cdot \begin{bmatrix} \langle u_1, x \rangle \\ \vdots \\ \langle u_n, x \rangle \end{bmatrix} = \langle u_1, x \rangle \cdot \sigma_1 u_1 + \dots + \langle u_n, x \rangle \cdot \sigma_n u_n$$

$$f_j(\sigma) = \sigma_j \cdot u_j^\top \cdot \langle u_j, x \rangle + \dots + \sigma_n \cdot u_n^\top \cdot \langle u_n, x \rangle \quad \Leftarrow$$

$$\frac{\partial f_j}{\partial x_i}(\sigma) = u_i^\top \cdot \langle u_j, x \rangle \quad \Leftarrow$$

$$Df|_x = \begin{bmatrix} u_1^\top \cdot \langle u_1, x \rangle & \dots & u_n^\top \cdot \langle u_n, x \rangle \\ \vdots & \ddots & \vdots \\ u_1^\top \cdot \langle u_n, x \rangle & \dots & u_n^\top \cdot \langle u_n, x \rangle \end{bmatrix} = \boxed{Df|_x}$$

$$= \begin{bmatrix} | & | \\ \langle u_1, x \rangle \cdot u_1 & \dots & \langle u_n, x \rangle \cdot u_n \\ | & | \end{bmatrix} =$$

$$= \begin{bmatrix} | & & | \\ u_1 & \dots & u_n \\ | & & | \end{bmatrix} \cdot \begin{bmatrix} \langle u_1, x \rangle & & 0 \\ 0 & \ddots & \\ & & \langle u_n, x \rangle \end{bmatrix} =$$

$$= U \cdot \text{diag} \left(\begin{bmatrix} \langle u_1, x \rangle \\ \vdots \\ \langle u_n, x \rangle \end{bmatrix} \right) = U \cdot \text{diag} ([x]_{B_U})$$

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

$$g_2(\sigma) := f(\sigma) - y \quad ! \quad g_2(\sigma) = \frac{1}{2} \|\sigma\|^2 \quad \text{!o}$$

$$\text{sic} \quad g_2 : \mathbb{R}^n \rightarrow \mathbb{R}^n \quad ! \quad g_2 : \mathbb{R}^n \rightarrow \mathbb{R}$$

$$\text{!o} \quad h(\sigma) = (g_1 \circ g_2)(\sigma)$$

$$Dh|_{\sigma} = Dg_1|_{g_2(\sigma)} \cdot Dg_2|_{\sigma}$$

$$Dg_2|_{\sigma} = 2 \cdot \frac{1}{2} \cdot \sigma^T = \sigma^T \quad \text{!o}$$

$$Dg_1|_{\sigma} = U \cdot \text{diag} \left(\begin{bmatrix} \langle u_1, x \rangle \\ \vdots \\ \langle u_n, x \rangle \end{bmatrix} \right)$$

$$Dh|_{\sigma} = (f(\sigma) - y)^T \cdot U \cdot \text{diag} \left(\begin{bmatrix} \langle u_1, x \rangle \\ \vdots \\ \langle u_n, x \rangle \end{bmatrix} \right) \quad \text{!o}$$

$$\nabla h(\sigma) = \text{diag} \left(\begin{bmatrix} \langle u_1, x \rangle \\ \vdots \\ \langle u_n, x \rangle \end{bmatrix} \right) \cdot U^T \cdot (f(\sigma) - y) \quad \Leftarrow$$

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

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

$$h(x) = \sum_{l=1}^d f_l(x) \quad (1) \quad , \quad f_i(x) = e^{x_i} \quad \because \quad f_i : \mathbb{R}^d \rightarrow \mathbb{R} \quad (2)$$

$\left\{ \begin{array}{ll} l & l \neq j \\ i & i=j \end{array} \right.$

$$\frac{\partial S_j}{\partial x_i}(x) = S_j(x)(1 - S_i(x))$$

$$\frac{\partial S_j}{\partial x_i}(x) = \frac{\partial}{\partial x_i} \frac{f_j(x)}{h(x)} = \frac{\cancel{\frac{\partial f_j(x)}{\partial x_i} \cdot h(x)} - \cancel{\frac{\partial h(x)}{\partial x_i} \cdot f_j(x)}}{h^2(x)}$$

$\stackrel{i \neq j}{=} \quad \text{if } i \neq j$

$$= -\frac{f_i(x) \cdot f_j(x)}{h^2(x)} = -\frac{f_i(x)}{h(x)} \cdot \frac{f_j(x)}{h(x)} = -S_i(x) \cdot S_j(x)$$

$$J_x(S) = \begin{bmatrix} S_1(x)(1 - S_1(x)) & \cdots & -S_d(x) \cdot S_1(x) \\ -S_1(x) \cdot S_2(x) & \ddots & \vdots \\ \vdots & \ddots & -S_d(x) \cdot S_{d-1}(x) \\ -S_1(x) \cdot S_d(x) & \cdots & S_d(x)(1 - S_d(x)) \end{bmatrix}$$

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 .

$$J_{(x,y)}(f) = \begin{bmatrix} 3x^2 - 5y & -5x - 5y^4 \end{bmatrix} \quad \text{: prem. misl}$$

$$H_{(x,y)}(f) = \begin{bmatrix} 6x & -5 \\ -5 & -20y^3 \end{bmatrix} \quad \text{: jde misl}$$

14

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 .

$i \in I$ $\forall s \in t \in [0,1] \quad ! \quad x, y \in C \quad \text{ה'ג' .9}$

$t \alpha + (1-t)\beta \in C: \quad \forall i \quad x, y \in C:$

$t x + (1-t)y \in C \quad \forall i \quad i \in I$

$x = x_1 + x_2 \quad \forall \quad x, y \in C \quad \text{ה'ג' .10}$

$y = y_1 + y_2$

$x_1, y_1 \in C_1 \quad ! \quad x_1, y_1 \in C_1 \quad !$

$t \in [0,1] \quad \text{ה'ג'}$

$$tx + (1-t)y = tx_1 + t x_2 + (1-t)y_1 + (1-t)y_2 =$$

sn

$$= tx_1 + (1-t)y_1 + tx_2 + (1-t)y_2$$

$$\begin{matrix} \cap \\ (\gamma_1)_1 C_1 \end{matrix} \quad \begin{matrix} \cap \\ (\gamma_2)_1 C_2 \end{matrix}$$

$$\textcircled{1} \quad \text{Let } \exists \quad tx + (1-t)y \in C \quad \Leftarrow$$

$$y = \lambda y' \quad ! \quad x = \lambda x' \quad \text{so} \quad y, y' \in \lambda C \quad \text{so} \quad \underline{\underline{y}}$$

$$\therefore x', y' \in C \quad \therefore$$

$$\text{so} \quad t \in [0,1] \quad \text{so}$$

$$tx + (1-t)y = \lambda t x' + \lambda (1-t)y' =$$

$$= \lambda (t x' + (1-t)y')$$

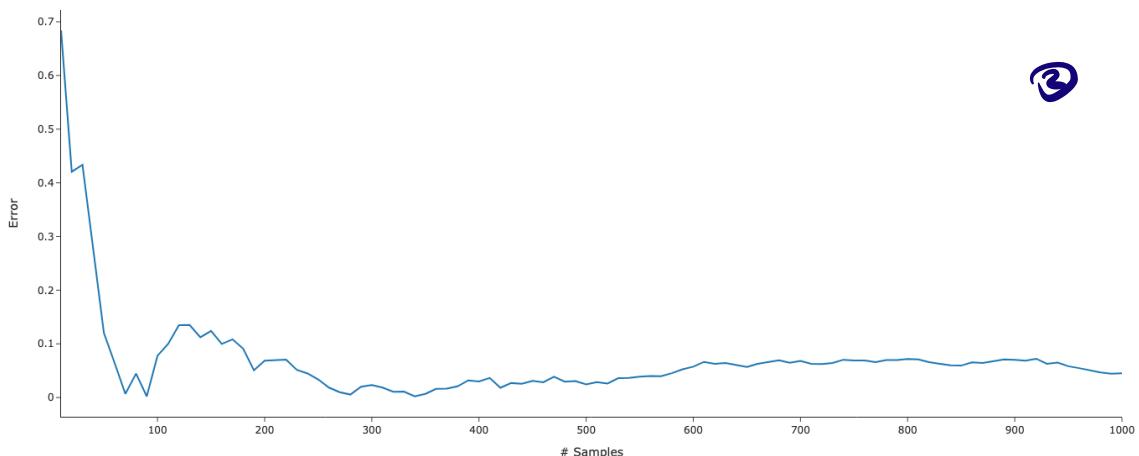
$$\begin{matrix} \cap \\ (\gamma_1)_1 C \end{matrix}$$

$$tx + (1-t)y \in \lambda C \quad \Leftarrow$$

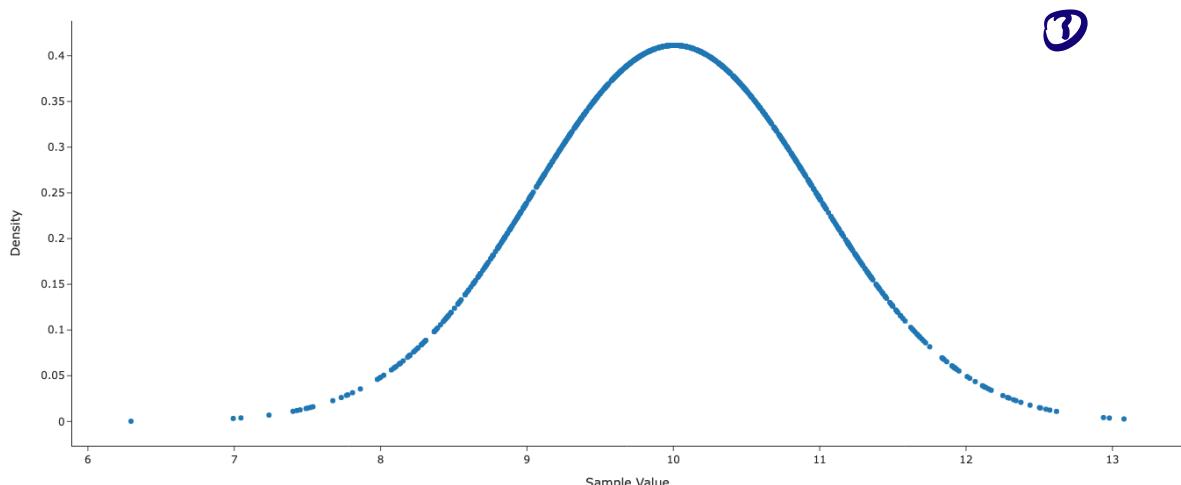
$$\textcircled{2} \quad \text{Let } \exists$$

- ✓ 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?

Difference between estimated and real expectation



Empirical PDF graph



הנרא לנו שפונקציית הספיבוב
ה empirica מושגת באמצעות פונקציית הספיבוב

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

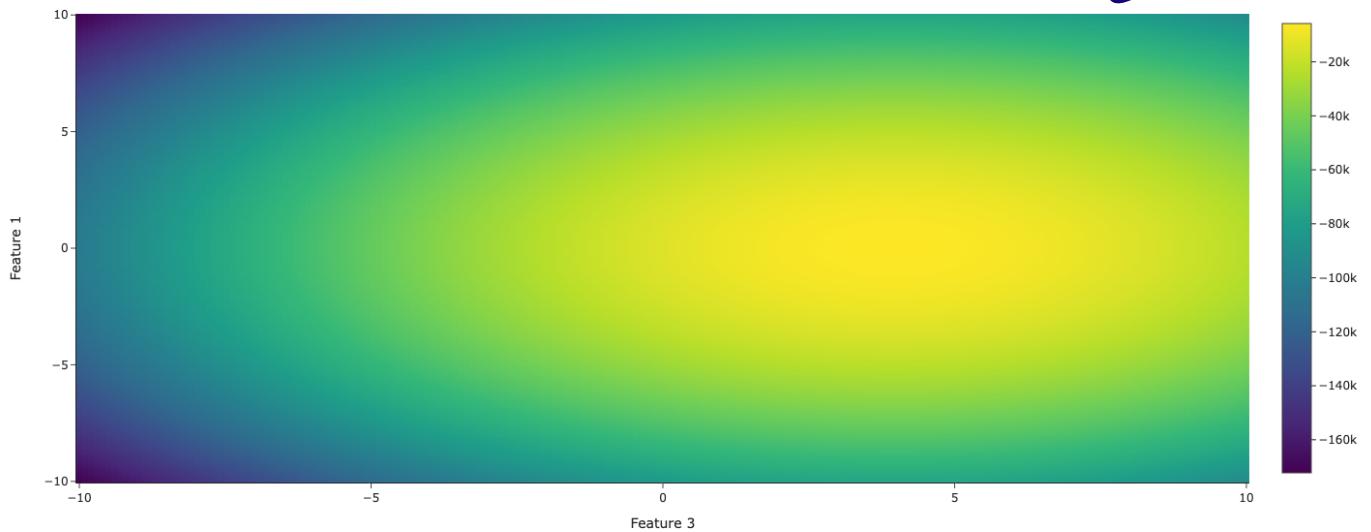
$$\mu = \begin{bmatrix} 0 \\ 0 \\ 4 \\ 0 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} 1 & 0.2 & 0 & 0.5 \\ 0.2 & 2 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0.5 & 0 & 0 & 1 \end{bmatrix}$$

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

Log likelihood of multivariate normal distribution with mean=[f1, 0, f3, 0]
Maximum likelihood is: -5806.292, Achieved at f1=-0.05, f3=3.97.

5 + 6



השאלה מבקשת לisset מושג של מודל מוגן על ידי מטריצת קוריאנס וקטורית. מטריצת קוריאנס מוגנת אם ורק אם היא חיובית-definitiva (positive definite).

השאלה מבקשת לisset מושג של מודל מוגן על ידי מטריצת קוריאנס וקטורית. מטריצת קוריאנס מוגנת אם ורק אם היא חיובית-definitiva (positive definite).