

Introduction to Machine Learning (67577)

Exercise 1

Estimation Theory & Mathematical Background

Second Semester, 2022

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
2.2	Estimation Theory	3
3	Practical Part	3
3.1	Univariate Gaussian Estimation	3
3.2	Multivariate Gaussian Estimation	3

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.
- Do not forget to answer the Moodle quiz of this assignment.

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 $\forall x \in V \ ||Ax|| = ||x||$.
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. 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 $\lambda_1 \geq \dots \geq \lambda_n$ the eigenvalues of C_0 , with the corresponding normalized eigenvectors $\mathbf{v}_1, \dots, \mathbf{v}_n$.

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: $\lim_{k \rightarrow \infty} b_k = \pm v_1$.

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

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

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

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

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

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

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

2.2 Estimation Theory

8. 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 .
9. 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 $\Sigma \in \mathbb{R}^{d \times d}$. Derive the log-likelihood function of $\mathcal{N}(\mu, \Sigma)$. 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 the IML.HUJI GitHub repository and setup the virtual environment as specified in the instructions. Write the necessary code in the files specified in the questions.

3.1 Univariate Gaussian Estimation

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

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

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

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

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 $\boldsymbol{\mu} = [f_1, 0, f_3, 0]^\top$ 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