

Set 3 - Finite Differences & PCA

Issued: October 26, 2022

Hand in (optional): November 9, 2022, 08:00

Question 1: Diffusion (50 points)

The diffusion of a substance can be described by the equation

$$\frac{\partial c(x, y, t)}{\partial t} = D \left(\frac{\partial^2 c(x, y, t)}{\partial x^2} + \frac{\partial^2 c(x, y, t)}{\partial y^2} \right), \quad (1)$$

where c is the concentration of the substance at position (x, y) and at time t , and D is the diffusion constant. The diffusion process happens in the domain $|x| < L/2$ and $|y| < L/2$. The concentration is zero on the boundaries of the domain for $t \geq 0$. The initial concentration is

$$c(x, y, 0) = \begin{cases} 1, & \text{if } |x| < L/4 \text{ and } |y| < L/4, \\ 0, & \text{otherwise.} \end{cases}$$

- a) Write down the 2-dimensional discretized diffusion process for eq. (1). Assume a uniform grid with spacing h and a central finite difference scheme in space and forward Euler time integration. For the forward Euler integration assume time intervals of size dt . Make sure that you annotate all variables.

In the following subquestions, you will work with the codes provided in `/hw04/code_q1/skeleton_code/`. Please have a look at the README file for further information. We recommend compiling and running the code on the Euler cluster https://scicomp.ethz.ch/wiki/Main_Page.

- b) Based on the discretization found in the previous subquestion, find the maximal timestep dt using the von Neumann stability analysis. Replace the hardcoded timestep ($t = 0.0001$) in the code with your solution.
- c) Based on the discretization found in the first subquestion, provide a cache-friendly implementation of the diffusion equation in the method `advance`. I.e. avoid copying of memory if possible and mind the access patterns of the memory. Blocking must not be implemented.
- d) Plot the total concentration as a function of time for $t \in [0, 0.5]$ using $D = 1$, $L = 2$ and $N = 100$. The concentration can be read from the file `diagnostics.dat` (column 0 and 1). Qualitatively explain the behaviour of the graph in less than 3 sentences, is this result expected?
- e) Parallelize the diffusion process (your implementation from subquestion 1c) in the method `advance` using OpenMP.
- f) Parallelize the integration of the concentration (marked with TODO in `compute_diagnostics`) and the calculation of the histogram (marked with TODO in the method `compute_histogram`) using OpenMP.

Question 2: Dimensionality reduction with PCA (50 points)

Principal Component Analysis (PCA) is a classical method to perform dimensionality reduction and uncover structures in data. The method learns an orthogonal transformation that eliminates linear correlations and tries to capture as much variance in the data as possible. You are provided with a two dimensional toy dataset with $N = 1024$ samples plotted in Figure 1. The principal components can be computed using the covariance method, by constructing the covariance matrix of the data $C \in \mathbb{R}^{D \times D}$ and identifying its eigenvalue decomposition. The covariance matrix is given by

$$C = \frac{1}{N-1} X^T X \quad (2)$$

where $X \in \mathbb{R}^{N \times D}$ is constructed by stacking the dataset. Assuming that the input data are independent (but not uncorrelated), the covariance matrix is symmetric. The eigenvalue decomposition of the symmetric matrix C reads

$$C = V \Lambda V^{-1}, \quad (3)$$

where $\Lambda = \text{diag}(\lambda_i)$ is a diagonal matrix with the eigenvalues stored in **descending** order. Due to the symmetry of the real covariance matrix, the eigenvectors are orthogonal to each other and they form an orthonormal basis, $V^{-1} = V^T$. The PCA components are the columns of the eigenvector matrix $V = [\mathbf{v}_1, \dots, \mathbf{v}_D]$. The transformed data read $\mathbf{y} = V^T \mathbf{x}$ and are linearly uncorrelated.

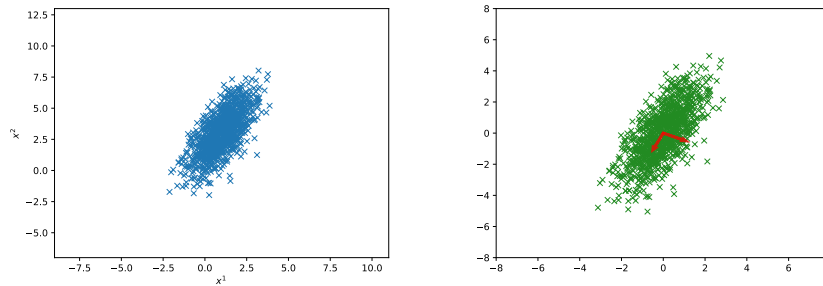


Figure 1: Toy two dimensional dataset along with its PCA components.

- a) Use an appropriate routine provided by the LAPACK software library to compute the PCA of the dataset based on the covariance method. The Intel Math Kernel Library (MKL) includes a high-performance implementation of LAPACK. In order to use the MKL on Euler, you have to load the module with `module load mkl`. After loading the module, you can include the header `#include <mkl_lapack.h>` to access the LAPACK routines.

Complete the steps in the skeleton code provided in the file `main_pca.cpp`. PCA is sensitive to the relative scaling of the input. For this exercise, **we consider relative input scaling to be relevant**, so the data should **not** be standardized (scaled with zero mean and unit variance). However, the data need to be centered. Write a program that performs the following tasks:

1. read the provided dataset (saved in a memory allocation of $N \times D$, where N is the data-set size and D the data dimension)
2. center the data

3. compute the covariance matrix of the data
4. call the `dsyev_()` routine of LAPACK to compute the eigenvalues of the matrix.

Can you recover the principal components plotted in Figure 1? Report the computed eigenvalues.

- b) In the following, we apply the PCA routine to a scenario closer to the real world. You are provided with a dataset with $N = 1280$ images of faces. Each face consists of a grayscale image $\mathbf{I} \in \mathbb{R}^{H \times W}$ that is flattened to a single data point $\mathbf{x} \in \mathbb{R}^D$, with $D = HW$, ignoring spatial structure. For this exercise $H = 50, W = 37$ so $D = 1850$. Perform PCA on this dataset and save **only** the $M = 10$ principal components (corresponding to the highest eigenvalues) in a `.txt` file. You have to save a matrix $V_r \in \mathbb{R}^{D \times M}$, obtained from the first M columns of the eigenvectors matrix $V \in \mathbb{R}^{D \times D}$ computed by PCA. Use the provided python routine to plot the principal components you computed.
- c) The components computed by PCA can be used to compress the dataset. The transformation to the compressed form reads $Z = XV_r$, where $Z \in \mathbb{R}^{N \times M}$. Use the computed $M = 10$ components to compress the data and **report** the compression ratio by measuring the number of floating point numbers needed to represent the original dataset and the compressed one (hint: for the purpose of compression you need to take into account the cost of the scaler, e.g. standardization).
- d) Try to **reconstruct** the original dataset using the compressed data, the $M = 10$ PCA components and the data scaling factors and save the result in a `.txt` file. The reconstruction can be obtained by $X = ZV_r^T$. A python routine is provided that can be used to plot the components computed by your method, plot the reconstruction and benchmark (qualitatively) against a reference python implementation of PCA. **Validate** your result with the provided python routines.