ETH Lectures 401-0663-00L Numerical Methods for Computer Science 401-2673-00L Numerical Methods for CSE

Homework Problems and Projects

Prof. R. Hiptmair, SAM, ETH Zurich

Autumn Term 2022, Version of January 6, 2023 (C) Seminar für Angewandte Mathematik, ETH Zürich

Link to the current version of this homework collection

0.1 General Information

0.1.1 Weekly Homework Assignments

All problems will be published in this single ".pdf" file. Every week, we publish a list of problem numbers on the lecture Moodle page. These are the current assignments for that week. Details are given on the lecture Moodle page.

0.1.2 Importance of Homework

Homework assignments are **not** mandatory. However, it is **very important** that you constantly exercise with the material you learn. "Solving" the homework assignments one week before the main exam by looking at the solutions will likely result in failure.

We provide hints and solutions from the beginning. It is your responsibility to look at the solutions only if you are stuck with a difficult problem and you tried to find a solution for a sufficient amount of time.

Most homework problems come with an estimate for the time required for its solution. These times are what we expect a bright and well-prepared student to require for that problem during an exam.

0.1.3 Corrections and Grading of Assignments

Homework assignments will not be graded.

However, you may submit your handwritten solutions by handing them to your tutor directly or uploading them via the Moodle upload interface, see the lecture Moodle page for instructions. The assistants will

also be able to examine your codes written in CODEEXPERT. You may submit your solutions even if those are incomplete and/or incorrect, and, in particular, if you have found a solution different from the master solution. Please do not submit solutions you copied from others or from the published solutions.

0.1.4 Codes and templates

For each problem involving coding you will find templates on the platform CODEEXPERT. You can access the platform with your NETHZ username and password. You have to register for the group that belongs to your tutor. Please, do not register to any other group. No local setup is required on your machine and you do not need to install any software or libraries.

Templates must be used as starting point for your solutions. You can create your own .cpp files, but you are strongly recommended to avoid it. CodeExpert will automatically test your codes and compare your solutions with the expected one. All solutions we provide will be based on templates. A similar workflow will also be used during the exam. All templates come with a main() function that cannot be modified: this will call other functions that you have to edit.

More information on the usage of Code Expert can be found here.

0.1.5 Hints and Solutions

Hints and solutions are kept in separate ".pdf" files. Links are supplies in this document. You can also download the ".pdf" files with hints and solutions from the website.

Chapter 0

Introduction

Problem 0-1: Simple operations with vectors and matrices in EIGEN

The problem encourages first steps in C++ coding using the EIGEN template library. The purpose of this exercise is to learn to declare and initialize **Eigen::Matrix** objects, and to get familiar with some handy typedefs and methods.

This problem involves coding in C++ and is related to [Lecture \rightarrow Section 1.2.1].

The central data type in EIGEN is the matrix, of which vectors are just a special version. Matrix data type fall into three fundamental categories

- (I) Variable-size and resizable dense matrices declared with the template argument Eigen::Dynamic,
- (II) Fixed-size (small) matrices, whose size must be known at compile time,
- (III) Variable size sparse matrices that are stored in special formats and which will be treated in [Lecture → Section 2.7.2].

Remark. For all codes in this problem, there are many ways how to meet the specification. Please think about it and discuss alternatives.

that initializes an $n \times n$ upper triangular matrix [Lecture \rightarrow Def. 1.1.2.3], whose non-zero entries all contain the value passed in the val argument.

HIDDEN HINT 1 for (0-1.c) \rightarrow 0-1-3-0:ctr.pdf

SOLUTION for (0-1.c) $\rightarrow 0-1-3-1:s2.pdf$

(0-1.d) • (15 min.) Vectors and matrices in EIGEN can also have integer and complex entries. The corresponding types have to be tagged with the \pm and \pm and \pm and the integer and complex entries. The Coperations involving matrices/vectors with different types of entries usually entail type casting.

In the file ${\tt MatrixClass.hpp}$ supplement the missing parts of the function ${\tt casting}()$ such that it returns the real part of the Euclidean (standard) scalar product of the vectors ${\tt u}$ and ${\tt v}$. To that end use EIGEN's dot product accessible via the method ${\tt dot}()$ of Eigen::Vector.

SOLUTION for (0-1.d)
$$\rightarrow 0-1-4-0$$
:s4.pdf

(0-1.e) \boxdot (20 min.) Given $n \in \mathbb{N}$ the matrices $\mathbf{A} \in \mathbb{R}^{n,n}$ and $\mathbf{B} \in \mathbb{C}^{n,n}$ are defined as follows

$$(\mathbf{A})_{k,\ell} := egin{cases} 5 & \text{, if } k \geq \ell \text{ ,} \\ 0 & \text{, else} \end{cases}$$
 , $(\mathbf{B})_{k,\ell} := rac{k + \imath \ell}{k - \imath \ell}$, $k,\ell \in \{1,\ldots,n\}$,

where $\iota \in \mathbb{C}$ is the imaginary unit, $\iota^2 = -1$. In the file MatrixClass.hpp realize complete the C++ function

Eigen::VectorXcd arithmetics(int n);

that evaluates the expression $\mathbf{B}(\mathbf{A} - 5\mathbf{I})[1, 2, \dots, n]^{\top} \in \mathbb{C}^n$.

HIDDEN HINT 1 for (0-1.e) $\rightarrow 0-1-5-0$:art.pdf

SOLUTION for (0-1.e) $\rightarrow 0-1-5-1:s3.pdf$

End Problem 0-1, 75 min.

Problem 0-2: EIGEN block operations on matrices

EIGEN matrices offer a range of methods to access sub-matrices (blocks), instead of individual entries at a time. In this exercise we practice their use.

This problem practices certain operations in Eigen and is connected with [Lecture \rightarrow § 1.2.1.5], which you should read beforehand.

Remark. The codes given as solutions may leave ample room for improvement. You should have the ambition to do better than the master solution.

(0-2.a) (15 min.) In the file MatrixBlocks.hpp supplement the missing lines of the function

```
Eigen::MatrixXd
zero_row_col(const Eigen::MatrixXd &A, int p, int q );
```

so that it returns a matrix that arises from A by setting to zero the p-th row and q-th column.

```
HIDDEN HINT 1 for (0-2.a) \rightarrow 0-2-1-0:t1.pdf
```

```
SOLUTION for (0-2.a) \rightarrow 0-2-1-1:s1.pdf
```

(0-2.b) ☑ (20 min.) Again in the file MatrixBlocks.hpp complete the implementation of the C++ function

```
MatrixXd swap_left_right_blocks(const MatrixXd &A, int p) {
```

that returns (see [Lecture \rightarrow Section 1.1.1] for explanations concerning the notations; PYTHON users beware!)

$$\left[(\mathbf{A})_{:,p+1:m'} (\mathbf{A})_{:,1:p} \right] \quad \text{for} \quad \mathbf{A} \in \mathbb{R}^{n,m} \; , \quad m,n \in \mathbb{N} \; , \quad 1 \leq p < m \; .$$

The matrix A is passed through the argument A.

HIDDEN HINT 1 for (0-2.b) \rightarrow 0-2-2-0:t2.pdf

SOLUTION for (0-2.b)
$$\rightarrow 0-2-2-1:s2.pdf$$

(0-2.c) (20 min.) Supplement the missing lines of the function

```
MatrixXd tridiagonal(int n, double a, double b, double c);
```

in the file MatrixBlocks.hpp that generates a square tridiagonal matrix

```
\begin{bmatrix}
b & c & 0 & \dots & & & \dots & 0 \\
a & b & c & 0 & & & & \vdots \\
0 & \ddots & \ddots & \ddots & \ddots & \ddots & & & & \vdots \\
\vdots & & & & & & & & & & & \\
\vdots & & & & & & & & & & \\
& & & & & \ddots & \ddots & \ddots & 0 \\
\vdots & & & & & & & & & & & \\
0 & \dots & & & & & & & & & \\
\end{bmatrix} \in \mathbb{R}^{n,n},
```

with the value $b \in \mathbb{R}$ on the main diagonal, the value $a \in \mathbb{R}$ on the first sub-diagonal, and the value $c \in \mathbb{R}$ on the first super-diagonal.

HIDDEN HINT 1 for (0-2.c) $\rightarrow 0-2-3-0:t3.pdf$

Solution for (0-2.c) $\rightarrow 0-2-3-1:s3.pdf$

lack

End Problem 0-2, 55 min.

Problem 0-3: EIGEN's reduction operations on matrices

EIGEN's matrices have a range of methods that reduce them to a vector or even a single number, see **EIGEN** documentation. We will also learn about the class template **Eigen::Array**, see **EIGEN** documentation.

The problem practises the use of Eigen and can be regarded as a supplement to [Lecture \rightarrow Section 1.2.1].

Remark. All the coding tasks can be tackled in many different ways. The codes provided as master solution just demonstrate one way.

All functions you will be asked to complete reside in the file MatrixReduce.hpp.

that computes the average

$$\frac{1}{mn}\sum_{i=1}^{n}\sum_{j=1}^{m}(\mathbf{A})_{i,j}$$

of the entries of the matrix $\mathbf{A} \in \mathbb{R}^{n,m}$, $n,m \in \mathbb{N}$. That matrix is passed as argument A.

HIDDEN HINT 1 for (0-3.a) $\rightarrow 0-3-1-0:s1h1.pdf$

HIDDEN HINT 2 for (0-3.a) \rightarrow 0-3-1-1:tavg.pdf

SOLUTION for (0-3.a) \rightarrow 0-3-1-2:s1.pdf

(0-3.b) (20 min.) Supply the missing parts of the C++ function

```
double percent_zero(const MatrixXd &A);
```

the computes the percentage of (exact) zeros among the entries of the matrix object A.

HIDDEN HINT 1 for (0-3.b) $\rightarrow 0-3-2-0:s2h1.pdf$

HIDDEN HINT 2 for (0-3.b) \rightarrow 0-3-2-1:pzt.pdf

SOLUTION for (0-3.b) $\rightarrow 0-3-2-2:s2.pdf$

(0-3.c) (20 min.) Complete the implementation of the function

```
bool has_zero_column(const MatrixXd &A);
```

that tests whether the matrix passed in A has a column that is exactly zero.

HIDDEN HINT 1 for (0-3.c) \rightarrow 0-3-3-0:s3h1.pdf

HIDDEN HINT 2 for (0-3.c) \rightarrow 0-3-3-1:hzct.pdf

SOLUTION for (0-3.c) $\rightarrow 0-3-3-2:s3.pdf$

(0-3.d) (25 min.) Finish the implementation of the C++ function

```
Eigen::MatrixXd columns_sum_to_zero(const Eigen::MatrixXd &A);
```

that returns a matrix A_0 that agrees with the *square* matrix $A \in \mathbb{R}^{n,n}$, $n \in \mathbb{N}$, passed in A except for the diagonal. Its diagonal entries are set in way that makes all column sums vanish:

$$\sum_{i=1}^n (\mathbf{A}_0)_{i,j} = 0 \quad \forall j \in \{1,\ldots,n\} \ .$$

HIDDEN HINT 1 for (0-3.d) \rightarrow 0-3-4-0:cstzt.pdf

SOLUTION for (0-3.d) $\rightarrow 0-3-4-1:s4.pdf$

End Problem 0-3, 80 min.

Chapter 1

Computing with Matrices and Vectors

Problem 1-1: Arrow matrix × vector multiplication

Innocent looking linear algebra operation can burn considerable CPU power when implemented carelessly, see [Lecture \rightarrow Code 1.3.1.11]. In this problem we study operations involving so-called arrow matrices, that is, matrices, for which only a few rows/columns and the diagonal are populated, see the spy-plots in [Lecture \rightarrow Rem. 1.3.1.5].

This problem is related to [Lecture \rightarrow Section 1.4.3] and contains implementation in C++ based on EIGEN.

Let $n \in \mathbb{N}$, n > 0. An "arrow matrix" $\mathbf{A} \in \mathbb{R}^{n \times n}$ is constructed given two vectors $\mathbf{a} \in \mathbb{R}^n$ and $\mathbf{d} \in \mathbb{R}^n$. The matrix is then squared and multiplied with a vector $\mathbf{x} \in \mathbb{R}^n$. This is implemented in the following C++ function:

C++11-code 1.1.1: Computing A^2x for an arrow matrix A

```
void arrow matrix 2 times x(const Eigen::VectorXd &d, const Eigen::VectorXd &a,
                                const Eigen::VectorXd &x, Eigen::VectorXd &y) {
3
     assert(d.size() == a.size() && a.size() == x.size() &&
5
            "Vector size must be the same!");
     const unsigned int n = d.size();
6
     // In this lines, we extract the blocks used to construct the matrix A.
     Eigen:: VectorXd d head = d.head(n - 1);
     Eigen::VectorXd a head = a.head(n - 1);
10
     Eigen::MatrixXd d diag = d head.asDiagonal();
11
12
     Eigen:: MatrixXd A(n, n);
13
14
     // We build the matrix A using the "comma initialization": each expression
15
    // separated by a comma is a "block" of the matrix we are building. d\diag is
17
    // the top left (n-1)x(n-1) block, a\head is the top right vertical vector,
    // a head.transpose() is the bottom left horizontal vector,
18
     // d(n-1) is a single element (a 1x1 matrix) on the bottom right corner.
19
     // This is how the matrix looks like:
     // A = | D | a
21
22
            | a ^T | d(n-1) |
23
    A \ll d_{diag}, a_{head}, a_{head}. transpose(), d(n - 1);
25
     y = A * A * x;
```

```
27 }
```

Get it on ₩ GitLab (ArrowMatrix.hpp).

(1-1.a) • (10 min.) For general vectors $\mathbf{d} = [d_1, \dots, d_n]^{\top}$ and $\mathbf{a} = [a_1, \dots, a_n]^{\top}$ sketch the pattern of nonzero entries of the matrix \mathbf{A} created in the function $\mathtt{arrow_matrix_2_times_x}$ in Code 1.1.1. [Lecture \rightarrow § 1.1.2.2] shows a way to visualize that pattern.

HIDDEN HINT 1 for (1-1.a) \rightarrow 1-1-1-0:arrmatha.pdf

SOLUTION for (1-1.a)
$$\rightarrow 1-1-1-1$$
: arrws.pdf

A key concern for the developer of numerical algorithms is their computational cost in case of large problems sizes, recall [Lecture \rightarrow Section 1.4].

Timings of arrow matrix 2 times x

We measure the runtime of the function $arrow_matrix_2_times_x$ with respect to the matrix size n and plot the results in Figure 1. (Get it on \checkmark GitLab (ArrowMatrix.hpp).)

The plot shows timings for arrow_matrix_2_times_x (log-log plot).

Machine details: Intel(R) Core(TM) i7-6700K CPU @ 4.00GHz. Compiled with: gcc 6.2.1 (flags: -O3). ▷

Give a detailed description of the behavior of the measured runtimes and an explanation for it.

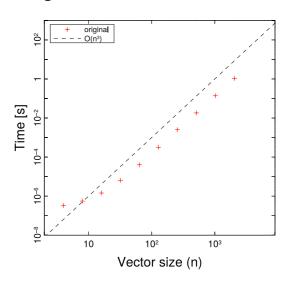


Fig. 1

Solution for (1-1.b)
$$\rightarrow 1-1-2-0$$
:arrwc.pdf

Write an efficient C++ function

that computes the same product as in Code 1.1.1, but with optimal asymptotic complexity with respect to n. Compute the complexity of your code and explain why you obtain such result.

Here d passes the vector $[d_1, \ldots, d_n]^{\top}$ and a passes the vector $[a_1, \ldots, a_n]^{\top}$.

SOLUTION for (1-1.c) $\rightarrow 1-1-3-0$:arrwi.pdf

What is the asymptotic complexity of your algorithm from sub-problem (1-1.c) (with respect to matrix size n)?

SOLUTION for (1-1.d) \rightarrow 1-1-4-0:arrwce.pdf

1. Computing with Matrices and Vectors, 1.1. Arrow matrix×vector multiplication

10

Compare the runtime of your implementation and the implementation given in Code 1.1.1 for $n = 2^5, \dots, 2^{12}$. Beware, for large n (n > 2048) the computations may take a long time.

Output the times in seconds, using 3 decimal digits in scientific notation. You can use **std::setw**, **std::**precision(**int**) and **std::**scientific from the standard library to output formatted text (include iomanip). For example:

Remark: run your code using optimization flags (-O3). With CMake, you can achieve this using -DCMAKE_BUILD_TYPE=Release as a CMake option.

§1.1.6 (Measuring runtimes in a C++ code) In order to measure runtimes you have two options: either you use std::chrono or use the Timer class.

How to use Timer

If you want to time a code, include timer.h, create a new Timer object, as demonstrated in the following code.

C++11-code 1.1.7: Usage of Timer Timer t; t.start(); // HERE CODE TO TIME t.stop(); 6 // Now you can get the time passed between start and stop using t.duration(); // You can start() and stop() again, a number of times // Ideally: repeat experiment many times and use min() to obtain the // fastest run 10 t.min(); 11 // You can also obtain the mean(): 12 t.mean();

All times will be outputted in seconds. Timer is simply a wrapper to std::chrono.

Advanced user: how to use std::chrono

Find the documentation of chrono on © C++ reference.

First include the chrono STL header file (note: this requires C++11). The chrono header provides the function:

```
C++11-code 1.1.8: now() function

std::chrono::high_resolution_clock::time_point
std::chrono::high_resolution_clock::now();
```

that returns the current time using the highest possible precision offered by the machine. For simplicity, we rename the return type:

C++11-code 1.1.9: now() function using time_point_t = std::chrono::high_resolution_clock::time_point; // Declare a starting point and a termination time time_point_t start, end;

The difference between two objects of type time_point_t (e.g. end - start) is an object of type std::chrono::high_resolution_clock::duration. In order to convert the chrono's duration type (which, in principle, can be anything: seconds, milliseconds, ...), to a fixed duration (say nanoseconds, std::chrono::nanoseconds), use the duration cast:

```
using duration_t = std::chrono::nanoseconds;
duration_t elapsed = std::chrono::duration_cast<duration_t >(end - start);
```

To obtain the actual number (as integer) used to represent elasped, use elapsed.count ().

Note: the data type used to represent std::chorno::seconds, std::chorno::milliseconds, etc. is of integer type. Thus, you cannot obtain fractions of this units. Make sure you use a sufficiently "refined" unit of measure (e.g. std::chrono::nanoseconds).

Solution for (1-1.e) \rightarrow 1-1-5-0:arrwc.pdf

End Problem 1-1, 65 min.

Problem 1-2: Gram-Schmidt orthonormalization with EIGEN

In this exercise we rely on EIGEN to implement a fundamental algorithm from linear algebra, Gram-Schmidt orthonormalization, see [NIS02] or [Lecture \rightarrow § 1.5.1.1], in order to practice the use of EIGEN's basic vector and matrix operations.

This problem relies on [Lecture \rightarrow Ex. 0.3.5.29], involves C++ implementation with EIGEN.

In [Lecture \rightarrow Ex. 0.3.5.29] and [Lecture \rightarrow Code 0.3.5.30] you can find an implementation of the famous Gram-Schmidt orthonormalization algorithm from linear algebra. That code makes use of a rather simple (and inefficient) implementation of vector and matrix classes. In this exercise we want to use EIGEN to implement Gram-Schmidt orthonormalization.

(1-2.a) lacktriangledown (15 min.) Determine the asymptotic complexity of the function <code>gramschmidt()</code> from [Lecture \rightarrow Code 0.3.5.30] in terms of the matrix dimension n, if the argument matrix A has size $n \times n$ and no premature termination occurs.

SOLUTION for (1-2.a)
$$\rightarrow$$
 1-2-1-0:grsi.pdf

```
Eigen::MatrixXd gram_schmidt(const Eigen::MatrixXd &A);
```

The output vectors should be returned as the columns of a matrix.

```
HIDDEN HINT 1 for (1-2.b) \rightarrow 1-2-2-0:gss1h1.pdf
Solution for (1-2.b) \rightarrow 1-2-2-1:grsi.pdf
```

Implement a function

```
bool testGramSchmidt(unsigned int n);
```

that tests your implementation by applying the function gram_schmidt to the matrix

$$\mathbf{A} \in \mathbb{R}^{n,n}$$
 , $(\mathbf{A})_{i,j} := i + 2j$, $1 \le i, j \le n$.

then checks the orthonormality of the columns of the output matrix, returning **true**, if orthonormality is confirmed.

```
HIDDEN HINT 1 for (1-2.c) \rightarrow 1-2-3-0:sp2mn.pdf

HIDDEN HINT 2 for (1-2.c) \rightarrow 1-2-3-1:sp2h2.pdf

HIDDEN HINT 3 for (1-2.c) \rightarrow 1-2-3-2:sp2Q2.pdf

SOLUTION for (1-2.c) \rightarrow 1-2-3-3:grst.pdf
```

End Problem 1-2, 80 min.

Problem 1-3: Kronecker product

In [Lecture \rightarrow Def. 1.4.3.7] we learned about the so-called Kronecker product. In this problem we revisit the discussion of [Lecture \rightarrow Ex. 1.4.3.8].

This problem is connected with [Lecture \rightarrow Section 1.4.3].

Definition [Lecture → Def. 1.4.3.7]. Kronecker product

The Kronecker product $\mathbf{A} \otimes \mathbf{B}$ of two matrices $\mathbf{A} \in \mathbb{K}^{m,n}$ and $\mathbf{B} \in \mathbb{K}^{l,k}$, $m,n,l,k \in \mathbb{N}$, is the $(ml) \times (nk)$ -matrix

$$\mathbf{A} \otimes \mathbf{B} := \begin{bmatrix} (\mathbf{A})_{11} \mathbf{B} & (\mathbf{A})_{1,2} \mathbf{B} & \dots & \dots & (\mathbf{A})_{1,n} \mathbf{B} \\ (\mathbf{A})_{2,1} \mathbf{B} & (\mathbf{A})_{2,2} \mathbf{B} & & & \vdots \\ \vdots & & \vdots & & & \vdots \\ (\mathbf{A})_{m,1} \mathbf{B} & (\mathbf{A})_{m,2} \mathbf{B} & \dots & \dots & (\mathbf{A})_{m,n} \mathbf{B} \end{bmatrix} \in \mathbb{K}^{ml,nk} .$$

(1-3.a) \bigcirc (10 min.) Compute the Kronecker product $\mathbf{C} = \mathbf{A} \otimes \mathbf{B}$ of the 2×2 matrices

$$\mathbf{A} = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \quad , \quad \mathbf{B} = \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix} . \tag{1.3.1}$$

SOLUTION for (1-3.a) $\rightarrow 1-3-1-0$: cmpm.pdf

(1-3.b) (20 min.)

Based on EIGEN, in the file kron.hpp implement a C++ function

```
Eigen::MatrixXd kron(const Eigen::MatrixXd & A, const
    Eigen::MatrixXd & B);
```

that computes the Kronecker product of the argument matrices A and B and returns the resulting matrix.

HIDDEN HINT 1 for (1-3.b) $\rightarrow 1-3-2-0$: kps2h1.pdf

SOLUTION for (1-3.b) $\rightarrow 1-3-2-1$: cmpk.pdf

What is the asymptotic complexity ([Lecture \rightarrow Def. 1.4.1.1]) in terms of the problem size parameter $n \rightarrow \infty$ of your implementation of kron from (1-3.b), when we pass $n \times n$ square matrices as arguments.

You may use the *Landau* symbol from [Lecture \rightarrow Def. 1.4.1.2] to state your answer.

SOLUTION for (1-3.c) $\rightarrow 1-3-3-0$: cmpk1.pdf

(1-3.d) (30 min.)

In general, computing the entire matrix is unnecessary. Devise an *efficient* implementation of an EIGEN-based C++ function (in the file kron.hpp)

```
Eigen::VectorXd kron_mult(
  const Eigen::MatrixXd & A, const Eigen::MatrixXd & B,
  const Eigen::VectorXd & x);
```

for the computation of $\mathbf{y} = (\mathbf{A} \otimes \mathbf{B})\mathbf{x}$ for square matrices $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n,n}$. Do not use reshaping through the Map method of Eigen's matrix classes; use access to sub-vectors by the segment method instead. The meaning of the arguments of kron_mult should be self-explanatory.

```
HIDDEN HINT 1 for (1-3.d) \rightarrow 1-3-4-0:h1.pdf
```

```
SOLUTION for (1-3.d) \rightarrow 1-3-4-1: cmpk.pdf
```

(1-3.e) \odot (30 min.) Devise another efficient implementation of a C++ code for the computation of $y=(A\otimes B)x$, A, $B\in\mathbb{R}^{n,n}$. This time use EIGEN's "matrix reshaping" functions. The function to implement in the file kron.hpp is

```
Eigen::VectorXd kron_reshape(
  const Eigen::MatrixXd &A, const Eigen::MatrixXd &B,
  const Eigen::VectorXd &x);
```

As its counterpart kron_mult () from Sub-problem (1-3.d) it is supposed to compute $y = (A \otimes B)x$.

```
HIDDEN HINT 1 for (1-3.e) \rightarrow 1-3-5-0:kph3.pdf
```

HIDDEN HINT 2 for (1-3.e) $\rightarrow 1-3-5-1$:ss3rs.pdf

```
SOLUTION for (1-3.e) \rightarrow 1-3-5-2:cmpk.pdf
```

Report the measurements in seconds with scientific notation using 3 decimal digits.

```
SOLUTION for (1-3.f) \rightarrow 1-3-6-0: cmpk2.pdf
```

End Problem 1-3, 130 min.

Problem 1-4: Fast matrix multiplication with EIGEN

[Lecture \to Rem. 1.4.2.2] presents Strassen's algorithm that can achieve the multiplication of two dense square matrices of size $n=2^k$, $k\in\mathbb{N}$, with an asymptotic complexity better than $O(n^3)$ (which is the complexity of a loop-based matrix multiplication). This problem centers around the implementation of this algorithm in EIGEN.

(1-4.a) Using EIGEN, implement a recursive function

```
Eigen::MatrixXd strassenMatMult(const Eigen::MatrixXd &A, const
    Eigen::MatrixXd &B);
```

that uses Strassen's algorithm ([Lecture \rightarrow Rem. 1.4.2.2]) to multiply the two square matrices **A** and **B** of size $n = 2^k$, $k \in \mathbb{N}$, and returns the result as output.

```
SOLUTION for (1-4.a) \rightarrow 1-4-1-0:fastmi.pdf
```

(1-4.b) Validate the correctness of your code by comparing the result with EIGEN's built-in matrix multiplication.

```
HIDDEN HINT 1 for (1-4.b) \rightarrow 1-4-2-0:hs1.pdf
```

```
SOLUTION for (1-4.b) \rightarrow 1-4-2-1: fastmv.pdf
```

(1-4.c) • Measure the runtime of your function strassenMatMult for random matrices of sizes 2^k , k = 3, ..., 6, and compare with the matrix multiplication offered by the *-operator of EIGEN.

You can use the class Timer or the STL class std::chrono as explained in Problem 1-1, § 1.1.6.

Recommendation: Use the optimization capabilities of your C++ compiler. With CMake, this can be done by appending $-DCMAKE_BUILD_TYPE=Release$ to the CMake invocation. Please note that this also disables various integrity checks.

Warning: For big values of n, the computations may take some time. Consider reducing n while debugging.

SOLUTION for (1-4.c)
$$\rightarrow$$
 1-4-3-0:fastmt.pdf

End Problem 1-4

Problem 1-5: Householder reflections

This problem is a supplement to [Lecture \rightarrow Section 1.5.1] and is related to Gram-Schmidt orthogonalization, see [Lecture \rightarrow Code 1.5.1.4].

For solving this problem, it is useful (but not necessary) to remember the QR-decomposition of a matrix from linear algebra [NIS02], see also [Lecture \rightarrow Thm. 3.3.3.4]. This problem is meant to practise some (advanced) linear algebra skills. It also addresses issues of asymptotic complexity, see [Lecture \rightarrow Section 1.4.1]. Involves implementation based on EIGEN.

(1-5.a) \square (20 min.) Let $\mathbf{v} \in \mathbb{R}^n$, $n \in \mathbb{N}$, n > 0. Consider the following algorithm given as a pseudocode:

Algorithm 1.5.1: Householder reflection

$$\mathbf{w} \leftarrow \mathbf{v}/\|\mathbf{v}\|_{2}$$

$$\mathbf{u} \leftarrow \mathbf{w}$$

$$\mathbf{u}_{1} \leftarrow \mathbf{u}_{1} + 1$$

$$\mathbf{q} \leftarrow \mathbf{u}/\|\mathbf{u}\|_{2}$$

$$\mathbf{X} \leftarrow \mathbf{I} - 2\mathbf{q}\mathbf{q}^{\top}$$

$$\mathbf{Z} \leftarrow ((\mathbf{X})_{:2:n})$$

Using the facilities of EIGEN write a C++ function with signature

```
void houserefl(const Eigen::VectorXd &v, Eigen::MatrixXd &Z);
```

that implements the algorithm from Pseudocode 1.5.1.

HIDDEN HINT 1 for (1-5.a) \rightarrow 1-5-1-0:hrfH1.pdf

SOLUTION for (1-5.a) $\rightarrow 1-5-1-1$:impl.pdf

(1-5.b) (15 min.) Show that the matrix X, defined in Code 1.5.1, satisfies:

$$\mathbf{X}^{\mathsf{T}}\mathbf{X} = \mathbf{I}_{n}$$

gather* Here I_n is the identity matrix of size n.

HIDDEN HINT 1 for (1-5.b) \rightarrow 1-5-2-0:norm.pdf

SOLUTION for (1-5.b)
$$\rightarrow 1-5-2-1$$
: orth.pdf

(1-5.c) (15 min.) Show that the first column of X, in Code 1.5.1, is a multiple of the vector v.

HIDDEN HINT 1 for (1-5.c) $\rightarrow 1-5-3-0$:mulh.pdf

SOLUTION for (1-5.c)
$$\rightarrow$$
 1-5-3-1:muls.pdf

What property does the set of columns of the matrix Z have? What is the purpose of the function houserefl?

SOLUTION for (1-5.d)
$$\rightarrow$$
 1-5-4-0:propertyimpl.pdf

What is the asymptotic complexity of the function houserefl as the length n of the input vector \mathbf{v} tends to ∞ ?

Specify it in leading order using the Landau symbol O.

SOLUTION for (1-5.e) \rightarrow 1-5-5-0:impl.pdf

Δ

End Problem 1-5, 70 min.

Problem 1-6: Matrix powers

This problems studies a (moderately) efficient way to compute large integer powers of matrices in EIGEN.

This problem practises EIGEN and also revisits the concept of asyymptotic complexity, hence is based on [Lecture \rightarrow Section 1.2.1] and [Lecture \rightarrow Section 1.4].

This problems deals with computing integer powers of square matrices, defined as:

$$\mathbf{A}^k := \overbrace{\mathbf{A} \cdot \ldots \cdot \mathbf{A}}^{k ext{ times}}, \quad k \in \mathbb{N}$$

for $\mathbf{A} \in \mathbb{C}^{n \times n}$ and $n \in \mathbb{N}$, n > 0.

(1-6.a) ☐ (30 min.) Implement a C++ function

```
Eigen::MatrixXcd matPow(Eigen::MatrixXcd & A, unsigned int k);
```

that, using only basic linear algebra operations (including matrix-vector or matrix-matrix multiplications), computes the k^{th} power of the $n \times n$ matrix **A** as efficiently as possible.

Here, MatrixXcd is a *complex matrix*. A complex matrix is similar to MatrixXd, and differs only by the fact that if contains elements of \mathbb{C} .

HIDDEN HINT 1 for (1-6.a) \rightarrow 1-6-1-0:.pdf

Remark, see also [Lecture \rightarrow § 1.2.1.1]: a Eigen::MatrixXcd is defined internally (in Eigen), as:

```
using MatrixXcd = Eigen::Matrix<std::complex, Eigen::Dynamic,
    Eigen::Dynamic>;
```

where Dynamic is an enum type with value -1. A generic matrix can be defined in EIGEN as

```
Eigen::Matrix<T, rows, cols> M;
```

where rows and cols are integers (number of rows resp. columns) and T is the underlying type. If rows (resp. cols) is -1 (or Dynamic), the Matrix will have a variable number of rows (resp. columns), see [Lecture $\rightarrow \S$ 1.2.1.1].

Remark: Matrix multiplication in EIGEN is not affected by aliasing issues, cf. \blacksquare EIGEN documentation. Therefore you can safely write A = A * A.

Remark: For code validation the EIGEN implementation of Matrix power (MatrixXcd::pow(unsigned int)) can be used writing:

#include <unsupported/Eigen/MatrixFunctions>

```
Solution for (1-6.a) \rightarrow 1-6-1-1:powi.pdf
```

(1-6.b) • (20 min.) Find (in leading order and state using the Landau symbol O) the asymptotic complexity in terms of $k \to \infty$ (and $n \to \infty$) of your implementation of matPow () taking into account that in EIGEN a matrix-matrix multiplication requires a $O(n^3)$ computational effort for $n \to \infty$.

```
SOLUTION for (1-6.b) \rightarrow 1-6-2-0: powc.pdf
```

 k. Compare this with the function \mathtt{matPow} () from Sub-problem (1-6.a). For n>0, use the complex $n\times n$ matrix $\mathbf A$ defined by

$$(\mathbf{A})_{j,k} := \frac{1}{\sqrt{n}} \exp\left(\frac{2\pi \imath jk}{n}\right), \quad i,j \in \{1,\ldots,n\},$$

to test the two functions (1 is the complex imaginary unit).

You should use the class Timer or the STL library std::chrono as explained in Problem 1-1.

Compute and output the runtime of the computation of the power for a matrix of size N=3. The runtime should be the minimum runtime of 10 trial runs for $k=2,\ldots,2^{31}$. Output the time is seconds, using 3 decimal digits in scientific notation. Do so in the C++ function

void tabulateRuntime();

Solution for (1-6.c) $\rightarrow 1-6-3-0$:powr.pdf

End Problem 1-6, 80 min.

Problem 1-7: Structured matrix-vector product

In [Lecture \rightarrow Ex. 1.4.3.5] we saw how the particular structure of a matrix can be exploited to compute a matrix-vector product with substantially reduced computational effort. This problem presents a similar case.

Involves a moderate amount of coding in C++ based on EIGEN.

Let $n \in \mathbb{N}$, n > 0 and consider the real $n \times n$ matrix **A** defined by

$$(\mathbf{A})_{i,j} = a_{i,j} = \min\{i,j\}, \ i,j = 1,\dots, n. \tag{1.7.1}$$

The matrix-vector product $\mathbf{y} = \mathbf{A}\mathbf{x}$ can be implemented in C++ as

```
C++11-code 1.7.2: Computing y = Ax for A from (1.7.1)

Eigen::VectorXd one = Eigen::VectorXd::Ones(n);

Eigen::VectorXd linsp = Eigen::VectorXd::LinSpaced(n, 1, n);

y = ((one * linsp.transpose()).cwiseMin(linsp * one.transpose())) * x;

Get it on ✓ GitLab (multAmin.hpp).
```

(1-7.a) • (10 min.) What is the asymptotic complexity (for $n \to \infty$) of the evaluation of the C++ code displayed above, with respect to the problem size parameter n?

```
SOLUTION for (1-7.a) \rightarrow 1-7-1-0:strc.pdf
```

(1-7.b) ☐ (30 min.) Write an efficient C++ function

```
void multAmin(const Eigen::VectorXd &x, Eigen::VectorXd &y);
```

that computes the same multiplication as Code 1.7.2 but with a better asymptotic complexity with respect to n.

You can test your implementation by comparing the returned values with the ones obtained with Code 1.7.2.

```
HIDDEN HINT 1 for (1-7.b) \rightarrow 1-7-2-0:smv:h1.pdf
SOLUTION for (1-7.b) \rightarrow 1-7-2-1:stri.pdf
```

(1-7.c) \bigcirc (10 min.) What is the asymptotic complexity (in terms of problem size parameter n) of your function multAmin?

```
Solution for (1-7.c) \rightarrow 1-7-3-0:cmpk.pdf
```

(1-7.d) \bigcirc (20 min.) Compare the runtimes of your implementation and the implementation given in Code 1.7.2 for $n=2^4,\ldots,2^{10}$.

You can use the class Timer or the STL class std::chrono as explained in Problem 1-1, § 1.1.6.

Report the measurements in seconds with scientific notation using 3 decimal digits.

```
SOLUTION for (1-7.d) \rightarrow 1-7-4-0:cmpk.pdf
```

```
C++11-code 1.7.6: Initializing B

Eigen:: MatrixXd B = Eigen:: MatrixXd:: Zero(n, n);

for (unsigned int i = 0; i < n; ++i) {
```

```
B(i, i) = 2;

if (i < n - 1) B(i + 1, i) = -1.0;

if (i > 0) B(i - 1, i) = -1.0;

B(n - 1, n - 1) = 1.0;

Get it on \checkmark GitLab (multAmin.hpp).
```

Sketch the matrix **B** created by these lines.

```
SOLUTION for (1-7.e) \rightarrow 1-7-5-0:cmpk.pdf (1-7.f) \odot (15 min.) With \bf A from (1.7.1) and \bf B from (1-7.e) write a short EIGEN-based C++ function Eigen::MatrixXd multABunitv();
```

that, for n = 10, computes \mathbf{ABe}_j , \mathbf{e}_j the j-th unit vector in \mathbb{R}^n , $j = 1, \dots, n$, and returns the computed vectors as the columns of a matrix in addition to printing that matrix.

Formulate a conjecture based on you observations.

```
SOLUTION for (1-7.f) \rightarrow 1-7-6-0:cmpkinv.pdf
```

End Problem 1-7, 95 min.

Problem 1-8: Avoiding cancellation

In [Lecture \rightarrow Section 1.5.4] we saw that the so-called *cancellation phenomenon* is a major cause of numerical instability, *cf.* [Lecture \rightarrow § 1.5.4.5]. Cancellation is the massive amplification of *relative errors* when subtracting two real numbers in floating point representation of about the same value. Fortunately, expressions vulnerable to cancellation can often be recast in a mathematically equivalent form that is no longer affected by cancellation, see [Lecture \rightarrow § 1.5.4.16]. There, we studied several examples, and this problem gives some more.

Involves simple implementation in C++.

(1-8.a) (10 min.) We consider the function

$$f_1(x_0, h) := \sin(x_0 + h) - \sin(x_0)$$
 (1.8.1)

Derive an *equivalent* expression $f_2(x_0, h) = f_1(x_0, h)$, which no longer involves the difference of return values of trigonometric functions. In other words, f_1 and f_2 give the same values, in exact arithmetic, for any given argument values x_0 and h.

HIDDEN HINT 1 for (1-8.a) \rightarrow 1-8-1-0:A1h1.pdf

Solution for (1-8.a)
$$\rightarrow$$
 1-8-1-1:sA1.pdf

Suggest a formula that avoids cancellation errors for computing the difference quotient approximation

$$f'(x) \approx \frac{f(x_0 + h) - f(x_0)}{h}$$
 (1.8.3)

of the derivative of $f(x) := \sin(x)$ at $x = x_0$.

Write a C++ function

void sinederv();

that implements

- 1. the difference quotient (1.8.3) directly, and
- 2. your new formula

and computes the resulting approximation of f'(1.2), for $h = 1 \cdot 10^{-20}, 1 \cdot 10^{-19}, \dots, 1$. Tabulate the relative errors of the results using $\cos(1.2)$ as exact value.

Explain the observed behaviour of the error.

HIDDEN HINT 1 for (1-8.b) \rightarrow 1-8-2-0:A2h1.pdf

Solution for (1-8.b)
$$\rightarrow$$
 1-8-2-1:A2s.pdf

(1-8.c) \Box (15 min.) Rewrite function $f(x) := \ln(x - \sqrt{x^2 - 1})$, x > 1, into a mathematically equivalent expression that is more suitable for numerical evaluation for any x > 1. Explain why, and provide a numerical example, which highlights the superiority of your new formula.

HIDDEN HINT 1 for (1-8.c) \rightarrow 1-8-3-0:logcnch.pdf

SOLUTION for (1-8.c)
$$\rightarrow$$
 1-8-3-1:canclog.pdf

(1-8.d) \Box (15 min.) For the following expressions, state the numerical difficulties that might affect a straightforward implementation for certain values of x (which ones?), and rewrite the formulas in a way that is more suitable for numerical computation.

1.
$$\sqrt{x + \frac{1}{x}} - \sqrt{x - \frac{1}{x}}$$
, for $x > 1$.

2.
$$\sqrt{\frac{1}{a^2} + \frac{1}{b^2}}$$
, for $a, b > 0$.

Solution for (1-8.d) \rightarrow 1-8-4-0:cancsqr.pdf

End Problem 1-8, 60 min.

Problem 1-9: Complexity of a C++ function

C++11-code 1.9.1: Code for getit.

In this problem we recall a concept from linear algebra: the diagonalization of a square matrix. Unless you can still remember what this means, please look up the chapter on "eigenvalues" in your linear algebra lecture notes. This problem also has a subtle relationship with Problem 1-6.

We consider the C++ function defined in getit.hpp (cf. Code 1.9.1).

W * (V.array().pow(k) * (W.partialPivLu().solve(cx)).array()).matrix();

Get it on ₩ GitLab (getit.hpp).

return ret.real();

Eigen::VectorXcd ret =

10

11

12

14 15 }

See the EIGEN documentation for eigenvalue decompositions in Eigen.

The operator array() for $v \in \mathbb{R}^n$ transforms the vector to an EIGEN array, on which operations are performed componentwise. The function matrix() is the "inverse" of array(), transforming an array to a matrix, on which vector/matrix operations are performed.

The code W.partialPivLU().solve(cx) performs a LU decomposition and solves the system W*x = cx for x. See the **EIGEN** documentation.

The MatrixBase::cast<T>() method applied to a matrix, will perform a componentwise cast of the matrix to a matrix of type T. See the EIGEN documentation.

(1-9.a) What is the output of getit, when A is a diagonalizable $n \times n$ matrix, $x \in \mathbb{R}^n$ and $k \in \mathbb{N}$?

```
SOLUTION for (1-9.a) \rightarrow 1-9-1-0:eigpoww.pdf
```

(1-9.b) $\ \ \ \ \ \ \$ Fix $k \in \mathbb{N}$. Discuss (in detail) the asymptotic complexity of getit for $n \to \infty$.

You may use the fact that computing eigenvalues and eigenvectors of an $n \times n$ matrix has asymptotic complexity $O(n^3)$ for $n \to \infty$.

```
Solution for (1-9.b) \rightarrow 1-9-2-0:eigpowc.pdf
```

End Problem 1-9

Problem 1-10: Approximating the Hyperbolic Sine

In this problem, we study how Taylor expansions can be used to avoid cancellation errors in the approximation of the hyperbolic sine, cf. the discussion in [Lecture \rightarrow Ex. 1.5.4.26].

The problem involves a little implementation in C++. The background is provided in [Lecture \rightarrow Section 1.5.4].

The hyperbolic sine $\sinh()$ is the function $\sinh(x) = \frac{1}{2}(e^x - e^{-x}), x \in \mathbb{R}$. The following code implements a C++ function providing $\sinh(x)$ for $x \in \mathbb{R}$.

C++ code 1.10.1: Implementation of sinh

```
double sinh_unstable(double x) {
   const double t = std::exp(x);
   return .5 * (t - 1. / t);
}
```

Get it on ₩ GitLab (sinh.hpp).

(1-10.a) \odot (30 min.) Explain why the function given in Code 1.10.1 may not give a good approximation of the hyperbolic sine for small values of x, and compute the relative error

$$\epsilon_{rel} := \frac{|\mathrm{sinh_unstable}(x) - \mathrm{sinh}(x)|}{|\mathrm{sinh}(x)|}$$

with $x = 10^{-k}$, k = 1, 2, ..., 10. To that end implement a C++ function **void sinhError**();

that tabulates those relative errors. As "exact value" use the result of the C++ built-in function std::sinh (include cmath).

SOLUTION for (1-10.a)
$$\rightarrow$$
 1-10-1-0:sinhex.pdf

(1-10.b) \odot (15 min.) Write the Taylor expansion with m terms around x=0 of the function $x\mapsto e^x$. Specify the remainder term as in [Lecture \to Ex. 1.5.4.26].

SOLUTION for (1-10.b)
$$\rightarrow$$
 1-10-2-0:sinhwrt.pdf

(1-10.c) \odot (15 min.) Prove that, for every $x \ge 0$ the following inequality holds true:

$$\sinh x \ge x \quad \forall x \ge 0 \ . \tag{1.10.4}$$

```
SOLUTION for (1-10.c) \rightarrow 1-10-3-0:sinhleq.pdf
```

Based on **Taylor expansion**, find an approximation for $\sinh(x)$, for every $0 \le x \le 10^{-3}$, so that the relative approximation error ϵ_{rel} is smaller than 10^{-15} . Follow the considerations of [Lecture \rightarrow Ex. 1.5.4.26].

SOLUTION for (1-10.d)
$$\rightarrow$$
 1-10-4-0:sinhty.pdf

End Problem 1-10, 90 min.

Problem 1-11: Complex roots

Universally, operations with complex numbers can be reduced to manipulating their real and imaginary parts. This problem studies how he complex square root can be realized in this way and then has to deal with potentially perilous cancellation.

This problem is related to [Lecture \rightarrow Section 1.5.4] and involves a little implementation in C++.

Given a complex number w=u+iv, $u,v\in\mathbb{R}$ with $v\geq 0$, its root (main branch) $\sqrt{w}=x+\imath y$ can be defined by the following real and imaginary parts

$$x := \sqrt{(\sqrt{u^2 + v^2} + u)/2}, \qquad (1.11.1)$$

$$y := \sqrt{(\sqrt{u^2 + v^2} - u)/2}, \tag{1.11.2}$$

where $\sqrt{}$ stands for the standard root function defined on \mathbb{R}_0^+ .

(1-11.a) \Box (15 min.) Show that in fact $(x + y)^2 = w$.

SOLUTION for (1-11.a)
$$\rightarrow 1-11-1-0$$
:sols0cr.pdf

(1-11.b) • (15 min.) For what $w \in \mathbb{C}$ will the direct implementation of (1.11.1) and (1.11.2) be vulnerable to **cancellation**?

SOLUTION for (1-11.b)
$$\rightarrow$$
 1-11-2-0:root1.pdf

(1-11.c) ⊡ (5 min.) [depends on Sub-problem (1-11.a)]

Compute xy as an expression of u and v.

SOLUTION for (1-11.c)
$$\rightarrow$$
 1-11-3-0:root2.pdf

(1-11.d) (30 min.) Implement a function

```
std::complex<double> myroot( std::complex<double> w );
```

that computes the root of w as given by (1.11.1) and (1.11.2) without the risk of cancellation. You may use only real arithmetic: for instance, you may not apply the function std::sqrt with *complex* arguments.

SOLUTION for (1-11.d)
$$\rightarrow$$
 1-11-4-0:root3.pdf

End Problem 1-11, 65 min.

Problem 1-12: Symmetric Gauss-Seidel iteration

This problem investigates a so-called stationary linear vector iteration from the implementation point of view.

Implementation in C++ based on Eigen is requested, [Lecture \rightarrow Section 1.2.1] required

For a square matrix $\mathbf{A} \in \mathbb{R}^{n,n}$, define $\mathbf{D}_{\mathbf{A}}$, $\mathbf{L}_{\mathbf{A}}$, $\mathbf{U}_{\mathbf{A}} \in \mathbb{R}^{n,n}$ by:

$$(\mathbf{D}_{\mathbf{A}})_{i,j} := \begin{cases} (\mathbf{A})_{i,j}, & i = j, \\ 0, & i \neq j \end{cases}, (\mathbf{L}_{\mathbf{A}})_{i,j} := \begin{cases} (\mathbf{A})_{i,j}, & i > j, \\ 0, & i \leq j \end{cases}, (\mathbf{U}_{\mathbf{A}})_{i,j} := \begin{cases} (\mathbf{A})_{i,j}, & i < j, \\ 0, & i \geq j \end{cases}.$$
(1.12.1)

The symmetric Gauss-Seidel iteration associated with the linear system of equations $\mathbf{A}\mathbf{x} = \mathbf{b}$ is defined as

$$\mathbf{x}^{(k+1)} = (\mathbf{U_A} + \mathbf{D_A})^{-1}\mathbf{b} - (\mathbf{U_A} + \mathbf{D_A})^{-1}\mathbf{L_A}(\mathbf{L_A} + \mathbf{D_A})^{-1}(\mathbf{b} - \mathbf{U_A}\mathbf{x}^{(k)}). \tag{1.12.2}$$

(1-12.a) (15 min.) Give a necessary and sufficient condition on A, such that the iteration (1.12.2) is well-defined.

SOLUTION for (1-12.a)
$$\rightarrow$$
 1-12-1-0:gsitw.pdf

(1-12.b) \odot (20 min.) Assume that (1.12.2) is well-defined. Show that a fixed point x of the iteration (1.12.2) is a solution of the linear system of equations Ax = b.

SOLUTION for (1-12.b)
$$\rightarrow$$
 1-12-2-0:sinhex.pdf

(1-12.c) (20 min.) Implement a C++ function

solving the linear system Ax = b using the iterative scheme (1.12.2). To that end, apply the iterative scheme to an initial guess $x^{(0)}$ passed trough x. The approximated solution given by the final iterate is then stored in x as an output.

Use a correction based termination criterion with relative tolerance rtol based on the Euclidean vector norm.

SOLUTION for (1-12.c)
$$\rightarrow$$
 1-12-3-0:gsitf.pdf

(1-12.d) \odot (15 min.) Test your implementation (n=9, rtol = 10e-8) with the linear system given by

$$\mathbf{A} = \begin{bmatrix} 3 & 1 & 0 & \dots & 0 \\ 2 & \ddots & \ddots & \ddots & \vdots \\ 0 & \ddots & \ddots & \ddots & 0 \\ \vdots & \ddots & \ddots & \ddots & 1 \\ 0 & \dots & 0 & 2 & 3 \end{bmatrix} \in \mathbb{R}^{n,n}, \mathbf{b} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \in \mathbb{R}^{n}$$

$$(1.12.11)$$

output the l^2 norm of the residual $\mathbf{b} - \mathbf{A}\widetilde{\mathbf{x}}$ of the approximated solution \mathbf{x} . Use \mathbf{b} as your starting vector for the iteration.

To that end, implement a C++ function

double testGSIt(unsigned int n);

SOLUTION for (1-12.d) \rightarrow 1-12-4-0:gsitim.pdf

•

(1-12.e) (15 min.) Using the same matrix \mathbf{A} and the same r.h.s. vector \mathbf{b} as above in Subproblem (1-12.d), we have tabulated the quantity $\|\mathbf{x}^{(k)} - \mathbf{A}^{-1}\mathbf{b}\|_2$ for $k = 1, \dots, 20$.

\boldsymbol{k}	$\ \mathbf{x}^{(k)} - \mathbf{A}^{-1}\mathbf{b}\ _2$	k	$\ \mathbf{x}^{(k)} - \mathbf{A}^{-1}\mathbf{b}\ _2$
1	0.172631	11	0.00341563
2	0.0623049	12	0.00234255
3	0.0464605	13	0.00160173
4	0.0360226	14	0.00109288
5	0.0272568	15	0.000744595
6	0.0200715	16	0.000506784
7	0.0144525	17	0.000344682
8	0.0102313	18	0.000234316
9	0.00715417	19	0.000159234
10	0.00495865	20	0.000108185

Describe qualitatively and quantitatively the convergence of the iterative scheme with respect to the number of iterations k.

SOLUTION for (1-12.e) \rightarrow 1-12-5-0:gsitc.pdf

End Problem 1-12, 85 min.

Problem 1-13: Structured matrix-vector product II

In this problem, you'll implement a specialized matrix-vector multiplication function, exploiting the special structure of the matrix to obtain a speed-up compared to ordinary matrix-vector multiplication.

This problem practises elementary numerical linear algebra and the use of EIGEN.

Consider the following Matrix-vector multiplication:

$$\mathbf{x} = \mathbf{A}\mathbf{y} \tag{1.13.1}$$

with

$$\mathbf{A} = \begin{bmatrix} a_1 & & & a_n \\ & \ddots & & \ddots \\ & & \ddots & & \\ & & \ddots & & \ddots \\ & & & a_1 & & a_n \end{bmatrix} \in \mathbb{R}^{n,n}, \qquad (1.13.2)$$

that is

$$a_{ij} = egin{cases} a_i & ext{, if } i=j \ a_j & ext{. if } i=n+1-j \ 0 & ext{, else.} \end{cases}$$

template < class Vector>

void xmatmult(const Vector& a, const Vector& y, Vector& x);

that efficiently evaluates (1.13.1) for A as described above.

HIDDEN HINT 1 for (1-13.a) \rightarrow 1-13-1-0:template.pdf

SOLUTION for (1-13.a)
$$\rightarrow$$
 1-13-1-1:xmatmulta.pdf

(1-13.b) \odot (10 min.) What is the complexity of your function with respect to the problem size n?

SOLUTION for (1-13.b)
$$\rightarrow$$
 1-13-2-0:xmatmultb.pdf

Tabulate the runtimes for matrix sizes $n = 2^k$, $k \in \{1, ..., 14\}$.

HIDDEN HINT 1 for (1-13.c) \rightarrow 1-13-3-0:xm3h1.pdf

HIDDEN HINT 2 for (1-13.c) \rightarrow 1-13-3-1:initialize.pdf

SOLUTION for (1-13.c) \rightarrow 1-13-3-2:xmatmultc.pdf

End Problem 1-13, 45 min.

Problem 1-14: Avoiding Cancellation by Rewriting Formulas

This problem studies two situations where careful implementation is necessary to avoid severe error amplification due to cancellation:

- 1. Blindly applying the well-known analytic formula for the roots of a quadratic polynomial can result in an implementation vulnerable to cancellation.
- 2. A code unstable for certain arguments will result from straightforwardly implementing the simple formula for the inverse function of the hyperbolic sine.

This problem revisits [Lecture \rightarrow Ex. 1.5.4.17].

```
(1-14.a) oxdots (4 min.) The task is to compute all real roots of the real quadratic polynomial t\mapsto t^2+\alpha t+\beta for given \alpha,\beta\in\mathbb{R}.
```

The following C++ code is meant to implement a function

```
Eigen::Vector2d zerosquadpol(double a, double b);
```

that computes the zeros of $t \mapsto t^2 + \alpha t + \beta$ in a stable way. The arguments a and b pass the coefficients α and β of the polynomial. Supplement the missing parts.

C++ code 1.14.1: Stable computation of real roots of a quadratic polynomial Eigen::Vector2d zerosquadpol(double a, double b) { Eigen::Vector2d z(2); double D = a * a - 4 * b; 3 if (D < 0) { throw "no real zeros"; }</pre> else { 5 double wD = std::sqrt(D); 6 >= 0) { 7 8 double t = Z << 9 10 } else { double t = 11 Z << 12 13 14 return z; 15 16

```
HIDDEN HINT 1 for (1-14.a) \rightarrow 1-14-1-0:srh1.pdf
HIDDEN HINT 2 for (1-14.a) \rightarrow 1-14-1-1:srh2.pdf
HIDDEN HINT 3 for (1-14.a) \rightarrow 1-14-1-2:srh3.pdf
SOLUTION for (1-14.a) \rightarrow 1-14-1-3:srp1.pdf
```

(1-14.b) \square (3 min.) The inverse function for the hyperbolic sine $\sinh(x) = \frac{1}{2}(e^x - e^{-x})$ is $\operatorname{arsinh}(x) := \log(x + \sqrt{x^2 + 1})$, $x \in \mathbb{R}$. (1.14.4)

The following C++ code is to supply a *cancellation-free* implementation of $\operatorname{arsinh}(x)$ for all $x \in \mathbb{R}$:

HIDDEN HINT 1 for (1-14.b) \rightarrow 1-14-2-0:strbh1.pdf

Solution for (1-14.b) \rightarrow 1-14-2-1:srhpb.pdf

End Problem 1-14, 7 min.

Chapter 2

Direct Methods for Linear Systems of Equations

Problem 2-1: Solving a small linear system

In this problem, you will use EIGEN to solve a simple linear algebra exercise.

This problem is meant to practise the solution of linear systems of equations in EIGEN, see [Lecture \rightarrow § 2.5.0.4].

We have the following information: Daniel buys 3 mangoes, 1 kiwi, 7 lychees and 2 bananas. Peter buys 2 mangoes and 1 pineapple. Julia needs 3 pomegranates and 1 mango for her fruit juice, whereas Ralf needs 5 kiwis and 1 banana for his fruit salad. Marcus buys 1 pineapple and 2 bananas. Manfred buys 20 lychees and 1 kiwi. Daniel pays CHF 11.10, Peter pays CHF 17.00 (and so spends all his pocket money), Julia pays CHF 6.10, Ralf pays CHF 5.25, Marcus pays CHF 12.50 and Manfred pays CHF 7.00.

(2-1.a) \odot Study [Lecture \rightarrow § 2.5.0.8] and $\stackrel{\blacksquare}{\blacksquare}$ EIGEN documentation to learn how to solve linear systems using direct elimination solvers of EIGEN.

(2-1.b) Write a C++ function

```
Eigen::VectorXd fruitPrice();
```

that returns the prices of the fruits in the following order:

Mango - Kiwi - Lychee - Banana - Pomegranate - Pineapple

SOLUTION for (2-1.b) $\rightarrow 2-1-2-0$:prb:sola.pdf

End Problem 2-1

Problem 2-2: Structured linear systems with pivoting

This problem deals with a block structured system and ways to efficiently implement the solution of such system. For this problem, you should have understood Gauss elimination with partial pivoting for the solution of a linear system, see [Lecture \rightarrow Section 2.3.3].

This problem is connected with [Lecture \rightarrow Section 2.3.2], [Lecture \rightarrow Section 2.5], and involves C++ implementation using Eigen.

We consider a block partitioned linear system of equations

$$\mathbf{A}\mathbf{x} = \mathbf{b}$$
, $\mathbf{A} = \begin{bmatrix} \mathbf{D}_1 & \mathbf{C} \\ \mathbf{C} & \mathbf{D}_2 \end{bmatrix} \in \mathbb{R}^{2n,2n}$, (2.2.1)

where all the $n \times n$ -matrices D_1 , D_2 and C are diagonal. Hence, the matrix A can be described through three n-vectors \mathbf{d}_1 , \mathbf{c} and \mathbf{d}_2 , which provide the diagonals of the matrix blocks. These vectors will be passed as arguments d_1 , c, and d_2 to the C++ codes below.

(2-2.a) (20 min.) Which permutation of rows and columns converts the matrix into a tridiagonal matrix?

```
HIDDEN HINT 1 for (2-2.a) \rightarrow 2-2-1-0:blse1h1.pdf
```

SOLUTION for (2-2.a) \rightarrow 2-2-1-1:blockrnt.pdf

(2-2.b) ☐ (20 min.) In the file blocklsepiv.hpp write an efficient C++ function

```
Eigen::VectorXd multA(
  const Eigen::VectorXd & d1, const Eigen::VectorXd & d2,
  const Eigen::VectorXd & c, const Eigen::VectorXd & x);
```

that returns y := Ax. The argument x passes a column vector $\mathbf{x} \in \mathbb{R}^{2n}$ and the other arguments the vectors defining the matrix as in (2.2.1).

SOLUTION for (2-2.b)
$$\rightarrow$$
 2-2-2-0:blockim.pdf

(2-2.c) (15 min.) Compute the LU-factors of A as given in (2.2.1), where you may assume that they exist.

```
HIDDEN HINT 1 for (2-2.c) \rightarrow 2-2-3-0:hblocklu.pdf
```

SOLUTION for (2-2.c)
$$\rightarrow$$
 2-2-3-1:blocklu.pdf

When does an LU-decomposition of A from (2.2.1) according to [Lecture \rightarrow Def. 2.3.2.3] exist?

SOLUTION for (2-2.d)
$$\rightarrow$$
 2-2-4-0:plse2a.pdf

(2-2.e) (15 min.) [depends on Sub-problem (2-2.c)]

Give an example of a matrix A with the structure as in (2.2.1)

- that possesses an LU-decomposition exists, but fails to be regular,
- for which an LU-decomposition does not exist, though it is invertible.

HIDDEN HINT 1 for (2-2.e) \rightarrow 2-2-5-0:plseh2b.pdf

SOLUTION for (2-2.e) \rightarrow 2-2-5-1:plse2ps.pdf

2. Direct Methods for Linear Systems of Equations, 2.2. Structured linear systems with pivoting

(2-2.f) (45 min.) [depends on Sub-problem (2-2.c) and Sub-problem (2-2.a)]

Write an *efficient* C++ function

```
Eigen::VectorXd solveA(
  const Eigen::VectorXd & d1, const Eigen::VectorXd & d2,
  const Eigen::VectorXd & c, const Eigen::VectorXd & b);
```

that solves $\mathbf{A}\mathbf{x} = \mathbf{b}$ with Gaussian elimination. The argument b passes the right-hand-side vector \mathbf{b} , whereas the other arguments provide information about the matrix \mathbf{A} as defined in (2.2.1).

- You must not use EIGEN's built-in linear solvers for the full matrix.
- Test for "near singularity" of the matrix.

```
HIDDEN HINT 1 for (2-2.f) \rightarrow 2-2-6-0:hsmalllse.pdf
SOLUTION for (2-2.f) \rightarrow 2-2-6-1:blockim.pdf
```

Remark. Probably the use of EIGEN's *sparse elimination solvers* as presented in [Lecture \rightarrow Section 2.7.3] also yields an efficient implementation, because that solver should detect the special decoupling of unknowns encoded in A.

(2-2.g) \Box (10 min.) State the **asymptotic complexity** of your implementation of solveA in term of the problem size parameter $n \to \infty$.

```
Solution for (2-2.g) \rightarrow 2-2-7-0:blockco.pdf
```

(2-2.h) Determine the asymptotic complexity of solveA in a numerical experiment through measuring runtimes.

As test case use $\mathbf{d}_1 := \mathbf{1}^\top = \begin{bmatrix} 1, \dots, n \end{bmatrix}^\top = -\mathbf{d}_2$, $\mathbf{c} = \mathbf{1}$ and $\mathbf{b} = \begin{bmatrix} \mathbf{d}_1 \\ \mathbf{d}_1 \end{bmatrix}$. Tabulate and plot the runtimes (in seconds, 3 decimal digit, scientific notation) for $n = 2^\ell$, $\ell = 2, \dots, 20$.

```
SOLUTION for (2-2.h) \rightarrow 2-2-8-0:blockrnt.pdf
```

End Problem 2-2, 140 min.

Problem 2-3: Block-wise LU-decomposition

Using block matrix operations and the LU-decomposition as described in [Lecture \rightarrow Rem. 2.3.2.19], you should implement a function to solve a LSE for a matrix of particular structure.

Moderate Eigen-based implementation in C++. Connected with [Lecture → Section 2.6]

Let the matrix $\mathbf{A} \in \mathbb{R}^{n+1,n+1}$ be block-partitioned according to

$$\mathbf{A} = \begin{bmatrix} \mathbf{R} & \mathbf{v} \\ \mathbf{u}^\top & 0 \end{bmatrix}$$

where $\mathbf{v}, \mathbf{u} \in \mathbb{R}^n$ and $\mathbf{R} \in \mathbb{R}^{n,n}$ is *upper triangular* and regular.

(2-3.a) \odot (15 min.) Refresh yourself on the definition of the LU-decomposition of a matrix [Lecture \rightarrow Def. 2.3.2.3] and study [Lecture \rightarrow Rem. 2.3.2.19].

(2-3.b) \bigcirc (10 min.) Find the LU-decomposition of **A**.

SOLUTION for (2-3.b) \rightarrow 2-3-2-0:bllua.pdfHere it is important that R was assumed to be invertible.

Show that **A** is regular, if and only if $\mathbf{u}^{\top} \mathbf{R}^{-1} \mathbf{v} \neq 0$.

HIDDEN HINT 1 for (2-3.c) $\rightarrow 2-3-3-0$:bh1.pdf

Solution for (2-3.c) \rightarrow 2-3-3-1:bllub.pdf

In the file blockLU.hpp write a C++ function

```
Eigen::VectorXd solve_LSE(const Eigen::MatrixXd& R, const
    Eigen::VectorXd& v, const Eigen::VectorXd& u, const
    Eigen::VectorXd& b);
```

that solves $\mathbf{A}\mathbf{x} = \mathbf{b}$ for \mathbf{A} described above and returns \mathbf{x} . Make use of the block-LU-decomposition you derived in Sub-problem (2-3.b) and only use elementary operations, in particular no built-in solver classes of EIGEN.

HIDDEN HINT 1 for (2-3.d) $\rightarrow 2-3-4-0$: <tag>.pdf

Solution for (2-3.d) \rightarrow 2-3-4-1:blluc.pdf

End Problem 2-3. 50 min.

Problem 2-4: Cholesky and QR decomposition

This problem is about the Cholesky and QR decomposition and the relationship between them. The Cholesky decomposition is explained in [Lecture \rightarrow § 2.8.0.13]. The QR decomposition is explained in [Lecture \rightarrow Section 3.3.3]. Please study these topics before tackling this problem.

This problem relies on both [Lecture \rightarrow § 2.8.0.13] and [Lecture \rightarrow Section 3.3.3]

(2-4.a) \square Show that, for every matrix $\mathbf{A} \in \mathbb{R}^{m,n}$ such that $\operatorname{rank}(\mathbf{A}) = n$, the product matrix $\mathbf{A}^{\top}\mathbf{A}$ admits a Cholesky decomposition.

```
HIDDEN HINT 1 for (2-4.a) \rightarrow 2-4-1-0:CholeskyQR1h.pdf
Solution for (2-4.a) \rightarrow 2-4-1-1:CholeskyQR1s.pdf
```

(2-4.b) The following C++ functions with EIGEN are given:

C++11-code 2.4.1: QR-decomposition via Cholesky decomposition

```
void CholeskyQR(const Eigen::MatrixXd& A, Eigen::MatrixXd& R,
                   Eigen::MatrixXd& Q) {
     Eigen::MatrixXd AtA = A.transpose() * A;
    Eigen::LLT<Eigen::MatrixXd> L = AtA. IIt ();
5
    R = L.matrixL().transpose();
6
    Q = R. transpose()
7
             .triangularView < Eigen :: Lower > ()
8
             . solve (A. transpose ())
             .transpose();
10
    // .triangularView() template member only accesses the triangular
11
     // of a dense matrix and allows to easily solve linear problem
12
13
```

Get it on ₩ GitLab (choleskyQR.hpp).

C++11-code 2.4.2: Standard way to do QR-decomposition in EIGEN

```
void DirectQR(const Eigen::MatrixXd& A, Eigen::MatrixXd& R,
                Eigen::MatrixXd& Q) {
3
    const size_t m = A.rows();
    const size t n = A.cols();
5
6
    Eigen:: HouseholderQR<Eigen:: MatrixXd> QR = A. householderQr();
    Q = QR.householderQ() * Eigen::MatrixXd::Identity(m, std::min(m, n));
8
    R = Eigen:: MatrixXd::Identity(std::min(m, n), m) *
        QR. matrixQR().triangularView < Eigen::Upper > ();
10
    // If A: m x n, then Q: m x m and R: m x n.
11
    // If m > n, however, the extra columns of Q and extra rows of R are
12
        not
    // needed.
                 Matlab returns this "economy-size" format calling
13
        "gr(A, 0)", which
    // does not compute these extra entries. With the code above, Eigen
        is smart
    // enough to not compute the discarded vectors.
15
16
Get it on ₩ GitLab (choleskyQR.hpp).
```

Prove that, for every matrix A, CholeskyQR and DirectQR produce the same output matrices Q and R, if there were no roundoff errors (Such functions are called algebraically equivalent).

```
HIDDEN HINT 1 for (2-4.b) \rightarrow 2-4-2-0:CholeskyQR2h.pdf
SOLUTION for (2-4.b) \rightarrow 2-4-2-1:CholeskyQR2s.pdf
```

•

(2-4.c) $oxed{oxed{oxed{2}}}$ Let EPS denote the machine precision. Why does the function <code>CholeskyQR</code> from Subproblem (2-4.b) fail to return the correct result for $old A = \begin{bmatrix} 1 & 1 \\ \frac{1}{2} \text{EPS} & 0 \\ 0 & \frac{1}{2} \text{EPS} \end{bmatrix}$?

HIDDEN HINT 1 for (2-4.c) \rightarrow 2-4-3-0:cencel.pdf

HIDDEN HINT 2 for (2-4.c) \rightarrow 2-4-3-1: CholeskyQR3h.pdf

Solution for (2-4.c) \rightarrow 2-4-3-2:CholeskyQR3s.pdf

End Problem 2-4

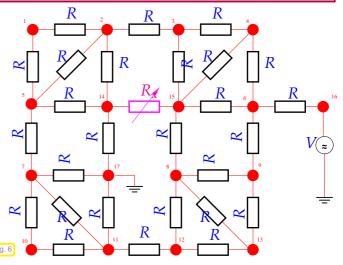
Problem 2-5: Resistance to impedance map

In [Lecture \rightarrow § 2.6.0.12], we learned about the Sherman-Morrison-Woodbury update formula [Lecture \rightarrow Lemma 2.6.0.21], which allows the efficient solution of a linear system of equations after a low-rank update according to [Lecture \rightarrow Eq. (2.6.0.16)], provided that the setup phase of an elimination (\rightarrow [Lecture \rightarrow § 2.3.2.15]) solver has already been done for the system matrix.

This problem is based on [Lecture \rightarrow Ex. 2.1.0.3] and [Lecture \rightarrow Ex. 2.6.0.24] and heavily draws on techniques introduced in [Lecture \rightarrow Section 2.6].

In this problem, we examine the concrete application from [Lecture \rightarrow Ex. 2.6.0.24], where the update formula is key to efficient implementation. This application is the computation of the impedance of the circuit drawn in Fig. 6 as a function of a variable resistance of a *single* circuit element, the resistor between nodes $\sharp 14$ and $\sharp 15$.

With the exception of that resistor, the circuit contains only identical linear resistors with the resistance R, that is, the relationship between branch currents and voltages is $I = R^{-1}U$ throughout. Excitation is provided by a voltage V imposed at node $\sharp 16$.



Here we consider DC operation (stationary setting), that is, all branch currents, voltages, and nodal potentials are real-valued. Moreover, all currents are given in 1A, all voltages in 1V, and all resistances in $1\Omega = \frac{V}{A}$. Thus we can drop the physical units in the following considerations and pretend that we deal with naked numbers even when computing physical quantities.

(2-5.a)
☐ (15 min.) Study [Lecture → Ex. 2.1.0.3] that explains how to compute voltages and currents in a linear circuit by means of nodal analysis. Understand how this leads to a linear system of equations for the unknown nodal potentials.

(2-5.b) ② (20 min.) Use nodal analysis to derive the 15×17 linear system of equations $\widehat{\mathbf{A}}(R_x)\mathbf{x} = \mathbf{0}$ satisfied by the nodal potentials U_i of the circuit from Fig. 6. Here, \mathbf{x} stands for the vector of unknown nodal potentials: $U_i = (\mathbf{x})_i$ is that at node $\sharp i, i \in \{1, \ldots, 17\}$. All resistors except for the controlled one (R_x) , colored magenta in Fig. 6) have the same resistance R. Use the numbering of nodes indicated in Fig. 6.

Solution for (2-5.b)
$$\rightarrow$$
 2-5-2-0:.pdf

The nodal potential of node $\sharp 16$ is fixed to the value V by connecting it to a grounded voltage source. In addition. node $\sharp 17$ is grounded, that is, its nodal potential is set to zero.

Taking into account these constraints we arrive at a 15×15 linear system of equations $\mathbf{A}(R_x)\mathbf{x} = \mathbf{b}$ for the unknown nodal potentials U_1, \dots, U_{15} , which are collected in the vector $\mathbf{x} := [U_1, \dots, U_{15}]^{\top}$. Determine the system matrix $\mathbf{A} \in \mathbb{R}^{15,15}$ and the right-hand side vector $\mathbf{b} \in \mathbb{R}^{15}$.

SOLUTION for (2-5.c)
$$\rightarrow$$
 2-5-3-0:circmat.pdf

(2-5.d) **(20** min.) [depends on Sub-problem (2-5.c)]

Provide an efficient implementation of the constructor and of the evaluation operator of the C++ class

class NodalPotentials {

```
public:
NodalPotentials() = delete;
NodalPotentials(const NodalPotentials&) = default;
NodalPotentials(NodalPotentials&) = default;
NodalPotentials& operator=(const NodalPotentials&) = default;
NodalPotentials& operator=(NodalPotentials&) = default;
~NodalPotentials() = default;

NodalPotentials(double R, double Rx);
Eigen::VectorXd operator () (double V) const;
private:
...
};
```

This class is meant to simulate the circuit from Fig. 6. Its constructor accepts the resistances R and R_x (in Ω) as arguments. The evaluation operator should return the vector $[U_1, \ldots, U_{15}]^{\top} \in \mathbb{R}^{15}$ of nodal voltages when the current source excites the circuit with a voltage V.

```
HIDDEN HINT 1 for (2-5.d) \rightarrow 2-5-4-0:h2al.pdf

HIDDEN HINT 2 for (2-5.d) \rightarrow 2-5-4-1:h2a2.pdf

SOLUTION for (2-5.d) \rightarrow 2-5-4-2:s2a.pdf
```

Characterize the change in the circuit matrix $\mathbf{A}(R_x)$ derived in Sub-problem (2-5.c) induced by a change in the value of R_x as a low-rank modification (\rightarrow [Lecture \rightarrow § 2.6.0.12]) of the circuit matrix $\mathbf{A}(0)$ of the form

$$\mathbf{A}(R_x) = \mathbf{A}(0) + \mathbf{u}\mathbf{u}^{\top}$$
 with suitable $\mathbf{u} \in \mathbb{R}^{15}$. (2.5.6)

```
HIDDEN HINT 1 for (2-5.e) \rightarrow 2-5-5-0:circh.pdf
```

SOLUTION for (2-5.e) \rightarrow 2-5-5-1:circsmw.pdf

(2-5.f) **(15 min.)** [depends on Sub-problem (2-5.e)]

How can the solution of $\mathbf{A}(R_x) = \mathbf{b}$ be computed relying on

- a solver for linear systems of equations with system matrix A(0),
- simple vector operations with asymptotic cost proportional to the vector length?

```
SOLUTION for (2-5.f) \rightarrow 2-5-6-0:sc3a.pdf
```

Based on the EIGEN library, complete the implementation of the C++ class

```
class ImpedanceMap {
  public:
  ImpedanceMap() = delete;
  ImpedanceMap(const ImpedanceMap&) = default;
  ImpedanceMap(ImpedanceMap&&) = default;
  ImpedanceMap& operator=(const ImpedanceMap&) = default;
  ImpedanceMap& operator=(ImpedanceMap&&) = default;
  result of the const ImpedanceMap& operator = (ImpedanceMap&&) = default;
  result of the const ImpedanceMap& operator = (ImpedanceMap&&) = default;
  result of the const ImpedanceMap& operator = (ImpedanceMap&&) = default;
  result of the const ImpedanceMap& operator = (ImpedanceMap&&) = default;
  result of the const ImpedanceMap& operator = (ImpedanceMap&&) = default;
  result of the const ImpedanceMap& operator = (ImpedanceMap&&) = default;
  result of the const ImpedanceMap&& operator = (ImpedanceMap&&) = default;
  result of the const ImpedanceMap&& operator = (ImpedanceMap&&) = default;
  result of the const ImpedanceMap&& operator = (ImpedanceMap&&) = default;
  result of the const ImpedanceMap&& operator = (ImpedanceMap&&) = default;
  result of the const ImpedanceMap&& operator = (ImpedanceMap&&) = default;
  result of the const ImpedanceMap&& operator = (ImpedanceMap&&) = default;
  result of the const ImpedanceMap&& operator = (ImpedanceMap&&) = default;
  result of the const ImpedanceMap&& operator = (ImpedanceMap&&) = default;
  result of the const ImpedanceMap&& operator = (ImpedanceMap&&) = default;
  result of the const ImpedanceMap&& operator = (ImpedanceMap&&) = default;
  result of the const ImpedanceMap&& operator = (ImpedanceMap&&) = default;
  result of the const ImpedanceMap&& operator = (ImpedanceMap&&) = default;
  result of the const ImpedanceMap&& operator = (ImpedanceMap&&) = default;
  result of the const ImpedanceMap&& operator = (ImpedanceMap&&) = default;
  result of the const ImpedanceMap&& operator = (ImpedanceMap&&) = default;
  result of the const ImpedanceMap&& operator = (ImpedanceMap&&) = default;
  result of the const ImpedanceMap&& operator = (ImpedanceMap&&) = default;
  result of the const ImpedanceMap&& oper
```

```
ImpedanceMap(double R, double V);
double operator () (double Rx) const;
private:
....
};
```

whose evaluation operator operator () returns the impedance R_* of the circuit from Fig. 6, when supplied with a concrete value for R_x . The impedance is defined as

$$R_* := rac{V}{I_V}$$
 , $I_V :=$ current through voltage source.

This function should be implemented efficiently using the result of Sub-problem (2-5.f). The setup phase for the direct solver should be carried out in the constructor, *cf.* [Lecture \rightarrow § 2.3.2.15], [Lecture \rightarrow Rem. 2.3.2.16].

```
HIDDEN HINT 1 for (2-5.g) \rightarrow 2-5-7-0:h40.pdf
HIDDEN HINT 2 for (2-5.g) \rightarrow 2-5-7-1:sp4h1.pdf
Solution for (2-5.g) \rightarrow 2-5-7-2:sc3a.pdf
```

End Problem 2-5, 150 min.

•

Problem 2-6: Banded matrix

Banded matrices are an important class of structured matrices; see [Lecture \rightarrow Section 2.7.5] and [Lecture \rightarrow Def. 2.7.5.1]. We will study ways to exploit during computations the knowledge that only some (sub)diagonals of a matrix are nonzero.

Connected with [Lecture \rightarrow Section 1.4.3] and [Lecture \rightarrow Section 2.7.5]

For $n \in \mathbb{N}$, consider the following matrix $(a_i, b_i \in \mathbb{R})$:

$$\mathbf{A} := \begin{bmatrix} 2 & a_1 & 0 & \dots & \dots & 0 \\ 0 & 2 & a_2 & 0 & \dots & \dots & 0 \\ b_1 & 0 & \ddots & \ddots & \ddots & \ddots & \vdots \\ 0 & b_2 & \ddots & \ddots & \ddots & \ddots & \vdots \\ \vdots & 0 & \ddots & \ddots & \ddots & \ddots & 0 \\ \vdots & \vdots & \ddots & \ddots & \ddots & \ddots & a_{n-1} \\ 0 & 0 & \dots & 0 & b_{n-2} & 0 & 2 \end{bmatrix} \in \mathbb{R}^{n,n}$$
 (2.6.1)

This matrix is an instance of a banded matrix.

(2-6.a) \Box (15 min.) In the file efficientbandmult.hpp implement an efficient C++ function multAx for the computation of y = Ax:

```
Eigen::VectorXd multAx(
   const Eigen::VectorXd & a, const VectorXd & b,
   const Eigen::VectorXd & x);
```

Eigen::VectorXd a and Eigen::VectorXd b contain the nonzero entries along the subdiagonals of matrix \mathbf{A} from (2.6.1). The function is supposed to return $\mathbf{y} = \mathbf{A}\mathbf{x}$, when Eigen::VectorXd x supplies are the vector $\mathbf{x} \in \mathbb{R}^n$.

```
SOLUTION for (2-6.a) \rightarrow 2-6-1-0:effbandimpl.pdf
```

(2-6.b) \odot (30 min.) Show that **A** is invertible if $a_i, b_i \in [0, 1]$.

```
HIDDEN HINT 1 for (2-6.b) \rightarrow 2-6-2-0:effbandinvh.pdf
```

```
SOLUTION for (2-6.b) \rightarrow 2-6-2-1:effbandinv.pdf
```

(2-6.c) (20 min.) Assume that $b_i = 0$ for all i = 1, ..., n - 2. Implement an efficient C++ function solving $\mathbf{A}\mathbf{x} = \mathbf{r}$:

Eigen::VectorXd a contains the nonzero entries along the upper subdiagonal of matrix A as defined in (2.6.1). r provides the right-hand-side vector of the linear system. The function should return the solution vector.

```
HIDDEN HINT 1 for (2-6.c) \rightarrow 2-6-3-0:effbandimplh.pdf
SOLUTION for (2-6.c) \rightarrow 2-6-3-1:effbandimplsol.pdf
```

(2-6.d) \odot (60 min.) For general $a_i, b_i \in [0, 1]$, implement an *efficient* C++ function that computes the solution of the linear system of equations $\mathbf{A}\mathbf{x} = \mathbf{r}$ by Gaussian elimination:

```
Eigen::VectorXd solvelseA(
  const Eigen::VectorXd & a, const Eigen::VectorXd & b,
  const Eigen::VectorXd & r);
```

The vectors a and b pass information about the matrix A as in Sub-problem (2-6.c), whereas r supplies the right-hand-side vector R. The return value is the solution x.

You must not use any built-in linear solvers of EIGEN.

```
HIDDEN HINT 1 for (2-6.d) \rightarrow 2-6-4-0:effbandgaussh.pdf Solution for (2-6.d) \rightarrow 2-6-4-1:effbandhausssol.pdf
```

What is the asymptotic complexity of your implementation of solvelseA for $n \to \infty$?

```
HIDDEN HINT 1 for (2-6.e) \rightarrow 2-6-5-0:effbandcompl.pdf
```

```
SOLUTION for (2-6.e) \rightarrow 2-6-5-1:effbandcompls.pdf
```

(2-6.f) Implement a C++ function

```
Eigen::VectorXd solvelseASparse(
  const Eigen::VectorXd & a, const Eigen::VectorXd & b,
  const Eigen::VectorXd & r);
```

with exactly the same specification as solvelseA() in Sub-problem (2-6.d), but this time using EIGEN's sparse elimination solver as explained in [Lecture \rightarrow Section 2.7.3].

```
HIDDEN HINT 1 for (2-6.f) \rightarrow 2-6-6-0:effbandsprs.pdf
```

```
Solution for (2-6.f) \rightarrow 2-6-6-1:effbandsprssol.pdf
```

End Problem 2-6, 135 min.

Problem 2-7: Properties of Householder reflections

[Lecture \rightarrow Section 3.3.3] describes an orthogonal transformation that can be used to map a vector onto another vector of the same length: the Householder transformation [Lecture \rightarrow Eq. (3.3.3.12)]. In this problem we examine properties of the "Householder matrices" and construct some of them.

Simple implementation in C++/EIGEN is requested

Householder transformation are described by square matrices of the form

$$\mathbf{H} = \mathbf{I} - 2\mathbf{v}\mathbf{v}^H$$
 with $\mathbf{v}^H\mathbf{v} = 1$, $\mathbf{v} \in \mathbb{C}^n$. (2.7.1)

We study a few important properties of this class of matrices and delve into the geometric meaning of Householder reflections.

(2-7.a) Prove the following properties:

- i) HH = I.
- ii) $\mathbf{H}\mathbf{H}^H = \mathbf{I}$.
- iii) $|det(\mathbf{H})| = 1$.
- iv) For every $\mathbf{s} \in \mathbb{C}^n$ perpendicular to \mathbf{v} (i.e. $\mathbf{v}^H \mathbf{s} = 0$): $\mathbf{H} \mathbf{s} = \mathbf{s}$.
- v) For every $z \in \mathbb{C}^n$ collinear to v (i.e. z = cv): Hz = -z.

Solution for (2-7.a)
$$\rightarrow$$
 2-7-1-0:householder1.pdf

(2-7.b) .

In real space the transformation \mathbf{H} can be understood as a reflection across a hyperplane perpendicular to \mathbf{v} . Compute the matrix \mathbf{H} and sketch the corresponding transformation for $\mathbf{v} = \frac{1}{\sqrt{10}}[1,3]^T$, n=2. What is the line of reflection?

SOLUTION for (2-7.b)
$$\rightarrow$$
 2-7-2-0:householder2.pdf

(2-7.c) .

The matrix

$$\mathbf{C} = \begin{bmatrix} 1 & 4 \\ 2 & 3 \end{bmatrix}$$

should be transformed to an upper triangular matrix using the Householder transformation (HC = R).

Compute:

- the vector v that defines H (is it unique?);
- · the matrix H;
- the images of the columns of C under the transformation H.

Solution for (2-7.c)
$$\rightarrow$$
 2-7-3-0:householder3.pdf

(2-7.d) Implement a function

void applyHouseholder(Eigen::VectorXd &x, const Eigen::VectorXd &v)

that efficiently computes $\mathbf{x} \leftarrow \mathbf{H}\mathbf{x}, \mathbf{x} \in \mathbb{R}^n$, for the Householder matrix ($\mathbf{v} \in \mathbb{R}^n$)

$$\mathbf{H} = \mathbf{I} - 2\mathbf{w}\mathbf{w}^{\top}$$
 , $\mathbf{w} := \frac{\mathbf{v}}{\|\mathbf{v}\|_2}$.

What is the asymptotic complexity of your implementation for $n \to \infty$?

SOLUTION for (2-7.d)
$$\rightarrow$$
 2-7-4-0:householder2.pdf

(2-7.e) lacksquare As we learned in [Lecture o Rem. 3.3.3.21], the sequence of Householder transformations $\mathbf{H}_1,\ldots,\mathbf{H}_{n-1}$ effecting the conversion of an $n\times n$ matrix into upper triangular form (o [Lecture o § 3.3.3.11]) is represented by the sequence of their defining *unit vectors* $\mathbf{v}_i \in \mathbb{R}^n$, also following the convention that $(\mathbf{v}_i)_i \geq 0$: $\mathbf{H}_i := \mathbf{I} - 2\mathbf{v}_i\mathbf{v}_i^{\top}$. The vectors \mathbf{v}_i are stored in the *strictly lower triangular* part of another $n\times n$ matrix $\mathbf{V} \in \mathbb{R}^{n,n}$. More precisely, we have

$$\mathbf{v}_i = [0, \dots, 0, (\mathbf{v}_i)_i, (\mathbf{V})_{i+1:n,i}]^{\top} \in \mathbb{R}^n, \quad i = 1, \dots, n-1.$$

Write a C++/EIGEN function

that computes $\mathbf{x} \leftarrow \mathbf{H}_{n-1}^{-1} \dots \mathbf{H}_1^{-1} \mathbf{x}$ based on Householder vectors stored in \mathbf{V} as described above.

```
HIDDEN HINT 1 for (2-7.e) \rightarrow 2-7-5-0:house:mb.pdf
```

HIDDEN HINT 2 for (2-7.e) \rightarrow 2-7-5-1:house:h1.pdf

Solution for (2-7.e) \rightarrow 2-7-5-2:householder2.pdf

End Problem 2-7

Problem 2-8: Lyapunov equation

Any linear system of equations with a finite number of unknowns can be written in the "canonical form" $\mathbf{A}\mathbf{x} = \mathbf{b}$ with a system matrix \mathbf{A} and a right hand side vector \mathbf{b} . However, the linear system may be given in a different form and it may not be obvious how to extract the system matrix. We propose an intriguing example and also present an important *matrix equation*, the so-called Lyapunov equation.

Related to [Lecture \rightarrow Section 2.5] and [Lecture \rightarrow Section 2.7.2]

Given the square matrix $\mathbf{A} \in \mathbb{R}^{n,n}$, we consider the equation

$$\mathbf{AX} + \mathbf{XA}^{\top} = \mathbf{I} \,, \tag{2.8.1}$$

where the matrix $\mathbf{X} \in \mathbb{R}^{n,n}$ is the unknown.

(2-8.a) \Box (10 min.) Show that for a fixed matrix $\mathbf{A} \in \mathbb{R}^{n,n}$ the following mapping is *linear*:

$$L: \left\{ \begin{array}{ccc} \mathbb{R}^{n,n} & \to & \mathbb{R}^{n,n} \\ \mathbf{X} & \mapsto & \mathbf{A}\mathbf{X} + \mathbf{X}\mathbf{A}^{\top} \end{array} \right.$$

HIDDEN HINT 1 for (2-8.a) \rightarrow 2-8-1-0:ly1h.pdf

Solution for (2-8.a)
$$\rightarrow$$
 2-8-1-1:ly1s.pdf

In the following let $\operatorname{vec}(\mathbf{M}) \in \mathbb{R}^{n^2}$ denote the column vector obtained by storing the internal coefficient array of a matrix $\mathbf{M} \in \mathbb{R}^{n,n}$ (\rightarrow [Lecture \rightarrow Eq. (1.2.3.5)]) in column major format, i.e. a data array with n^2 components [Lecture \rightarrow Rem. 1.2.3.4]. In PYTHON, MATLAB, $\operatorname{vec}(\mathbf{M})$ would be the column vector obtained by reshape (M, n*n, 1), see [Lecture \rightarrow Rem. 1.2.3.6] for the implementation based on EIGEN.

Owing to the observation made in Sub-problem (2-8.a), (2.8.1) is equivalent to the following linear system of equations:

$$\mathbf{C}\operatorname{vec}(\mathbf{X}) = \mathbf{b} \tag{2.8.2}$$

The system matrix is $\mathbf{C} \in \mathbb{R}^{n^2,n^2}$ and the right-hand side vector $\mathbf{b} \in \mathbb{R}^{n^2}$.

(2-8.b) \odot (10 min.) Recall the notion of *sparse matrix*: see [Lecture \rightarrow Section 2.7] and, in particular, [Lecture \rightarrow Notion 2.7.0.1] and [Lecture \rightarrow Def. 2.7.0.3].

(2-8.c) \Box (15 min.) Determine C and b from Eq. (2.8.2) for n=2 and

$$\mathbf{A} = \begin{bmatrix} 2 & 1 \\ -1 & 3 \end{bmatrix}$$

HIDDEN HINT 1 for (2-8.c) \rightarrow 2-8-3-0:ly3h.pdf

SOLUTION for (2-8.c)
$$\rightarrow$$
 2-8-3-1:1y3s.pdf

Solution for (2-8.d)
$$\rightarrow$$
 2-8-4-0:ly4s.pdf

In the file lyapunov.hpp implement a C++ function that builds the EIGEN matrix C from A:

```
Eigen::SparseMatrix<double> buildC(const Eigen::MatrixXd& A)
```

Make sure that the initialisation is done efficiently using an intermediate triplet format, see [Lecture \rightarrow Section 2.7.2].

HIDDEN HINT 1 for (2-8.e) \rightarrow 2-8-5-0:arrmatha.pdf

SOLUTION for (2-8.e) \rightarrow 2-8-5-1:arrwc.pdf

(2-8.f) (20 min.) In the file lyapunov.hpp a C++ function that returns the solution of (2.8.1), i.e. the $n \times n$ -matrix X if $A \in \mathbb{R}^{n,n}$ is passed in the argument A.

```
void solveLyapunov(const Eigen::MatrixXd& A, Eigen::MatrixXd& X)
```

Use a suitable sparse elimination solver provided by EIGEN, see [Lecture \rightarrow Section 2.7.3].

HIDDEN HINT 1 for (2-8.f) \rightarrow 2-8-6-0:ly6h.pdf

SOLUTION for (2-8.f)
$$\rightarrow$$
 2-8-6-1:ly6s.pdf

bool testLyapunov();

that is meant to validate your implementation of buildC and solveLyapunov for n = 5 and:

$$\mathbf{A} = \begin{bmatrix} 10 & 2 & 3 & 4 & 5 \\ 6 & 20 & 8 & 9 & 1 \\ 1 & 2 & 30 & 4 & 5 \\ 6 & 7 & 8 & 20 & 0 \\ 1 & 2 & 3 & 4 & 10 \end{bmatrix}.$$

SOLUTION for (2-8.g) $\rightarrow 2-8-7-0:1y7s.pdf$

SOLUTION for (2-8.h)
$$\rightarrow$$
 2-8-8-0:1y8s.pdf

End Problem 2-8, 135 min.

Problem 2-9: Partitioned Matrix

Based on the block view of matrix multiplication presented in [Lecture \rightarrow § 1.3.1.13], we looked at *block elimination* for the solution of block partitioned linear systems of equations in [Lecture \rightarrow § 2.6.0.2]. Also of interest are [Lecture \rightarrow Rem. 2.3.2.19] and [Lecture \rightarrow Rem. 2.3.2.17] where LU-factorisation is viewed from a block perspective. Closely related to this problem is [Lecture \rightarrow Ex. 2.6.0.5], which you should study again as warm-up to this problem.

This problem practises block Gaussian elimination and is related to [Lecture → Section 2.6]

Let the matrix $\mathbf{A} \in \mathbb{R}^{n+1,n+1}$ be block partitioned according to

$$\mathbf{A} = \begin{bmatrix} \mathbf{R} & \mathbf{v} \\ \mathbf{u}^T & 0 \end{bmatrix}, \tag{2.9.1}$$

where $\mathbf{v} \in \mathbb{R}^n$, $\mathbf{u} \in \mathbb{R}^n$, and $\mathbf{R} \in \mathbb{R}^{n \times n}$ is upper triangular and regular.

(2-9.a) (10 min.) Give a necessary and sufficient condition for the triangular matrix **R** to be invertible.

HIDDEN HINT 1 for (2-9.a) \rightarrow 2-9-1-0:s1.pdf

SOLUTION for (2-9.a)
$$\rightarrow$$
 2-9-1-1:partinv.pdf

(2-9.b) \odot (15 min.) Determine expressions for the sub-vector $\mathbf{z} \in \mathbb{R}^n$ and the last component $\xi \in \mathbb{R}$ of the solution vector of the linear system of equations

$$\begin{bmatrix} \mathbf{R} & \mathbf{v} \\ \mathbf{u}^T & 0 \end{bmatrix} \begin{bmatrix} \mathbf{z} \\ \boldsymbol{\xi} \end{bmatrix} = \begin{bmatrix} \mathbf{b} \\ \boldsymbol{\beta} \end{bmatrix}$$
 (2.9.2)

for arbitrary $\mathbf{b} \in \mathbb{R}^n$, $\beta \in \mathbb{R}$.

Use block Gaussian elimination as presented in [Lecture \rightarrow § 2.6.0.2].

SOLUTION for (2-9.b)
$$\rightarrow$$
 2-9-2-0:partgauss.pdf

(2-9.c) (10 min.) [depends on Sub-problem (2-9.b)]

Show that **A** is regular if and only if $\mathbf{u}^T \mathbf{R}^{-1} \mathbf{v} \neq 0$.

SOLUTION for (2-9.c)
$$\rightarrow$$
 2-9-3-0:partareg.pdf

Implement the C++ function

```
void solvelse(
  const Eigen::MatrixXd & R,
  const Eigen::VectorXd & v,
  const Eigen::VectorXd & u,
  const Eigen::VectorXd & bb,
  Eigen::VectorXd & x);
```

for the efficient computation of the solution of $\mathbf{A}\mathbf{x} = \mathbf{b}$ (with \mathbf{A} as in (2.9.1)), which is to be returned in the argument x. Perform size checks on input matrices and vectors.

Use the decomposition from Sub-problem (2-9.b). You can rely on the triangularView() function to inform EIGEN of the triangular structure of \mathbf{R} , see [Lecture \rightarrow Code 1.2.1.6].

SOLUTION for (2-9.d)
$$\rightarrow$$
 2-9-4-0:partimpl.pdf

```
(2-9.e) (20 min.) [depends on Sub-problem (2-9.d)]
```

Code a C++ function

```
bool testSolveLSE(const Eigen::MatrixXd & R,
const Eigen::VectorXd & v, const Eigen::VectorXd & u,
const Eigen::VectorXd & b, Eigen::VectorXd & x);
```

that tests your implementation of solvelse by comparison with a standard LU-solver provided by EIGEN. Return **true**, if the results "are the same".

```
HIDDEN HINT 1 for (2-9.e) \rightarrow 2-9-5-0:s3h1.pdf

HIDDEN HINT 2 for (2-9.e) \rightarrow 2-9-5-1:s3h2.pdf

SOLUTION for (2-9.e) \rightarrow 2-9-5-2:parttest.pdf

(2-9.f) (10 min.) [depends on Sub-problem (2-9.d)]
```

What is the asymptotic complexity of your implementation of solvelse () in terms of problem size parameter $n \to \infty$?

```
SOLUTION for (2-9.f) \rightarrow 2-9-6-0: partcmp.pdf
```

End Problem 2-9, 90 min.

Problem 2-10: Rank-one perturbations

This exercise centers around an application of the Sherman-Morrison-Woodbury formula (see [Lecture \rightarrow Lemma 2.6.0.21]). Please carefully revise [Lecture \rightarrow § 2.6.0.12] in the lecture notes.

Linked to [Lecture \rightarrow Rem. 2.5.0.10] and [Lecture \rightarrow § 2.6.0.12], related problem Problem 2-5

Consider the following pseudocode:

```
Algorithm 2.10.1: Code rankoneinvit Require: \mathbf{d} \in \mathbb{R}^n, tol \in \mathbb{R}, tol > 0 Ensure: l_{min} \mathbf{ev} \leftarrow \mathbf{d}; l_{\min} \leftarrow 0; l_{\text{new}} \leftarrow \min |\mathbf{d}|; while |l_{\text{new}} - l_{\min}| > tol \cdot l_{\min} \, \mathbf{do} l_{\min} \leftarrow l_{\text{new}}; \mathbf{M} \leftarrow \operatorname{diag}(\mathbf{d}) + \mathbf{ev} \cdot \mathbf{ev}^{\top}; \mathbf{ev} \leftarrow \mathbf{M}^{-1} \mathbf{ev}; \mathbf{ev} \leftarrow \mathbf{ev}/|\mathbf{ev}|; l_{\text{new}} \leftarrow \mathbf{ev}^{\top} \mathbf{Mev}; \mathbf{end} \, \mathbf{while} return l_{\text{new}}
```

Here diag designates a mapping $\mathbb{R}^n \mapsto \mathbb{R}^{n,n}$ that converts a vector into a diagonal matrix with the vector components as diagonal entries.

We assume that this code does something meaningful in a larger context that we do not care about. It might be taken from an old code that we have to upgrade.

(2-10.a) (30 min.) Directly port the function rankoneinvit (Algorithm 2.10.1) to C++ using EIGEN. To that end create a function

```
double rankoneinvit(const Eigen::VectorXd &d, const double &tol);
```

The C++ code should perform exactly the same computations. Recall that in EIGEN the asDiagonal () method converts a vector into a diagonal matrix, see EIGEN documentation. EIGEN also offers the normalize () method for vectors.

```
HIDDEN HINT 1 for (2-10.a) \rightarrow 2-10-1-0:s1h1.pdf
```

```
SOLUTION for (2-10.a) \rightarrow 2-10-1-1:smwps.pdf
```

```
HIDDEN HINT 1 for (2-10.b) \rightarrow 2-10-2-0:smwch.pdf
```

```
SOLUTION for (2-10.b) \rightarrow 2-10-2-1:smwcs.pdf
```

Implement an EIGEN-based C++ function

```
double rankoneinvit_fast(const Eigen::VectorXd &d, const double
    &tol);
```

that is *algebraically equivalent* to rankoneinvit () from, but enjoys a way better asymptotic complexity by avoiding any invocation of Gaussian elimination/LU-factorization.

HIDDEN HINT 1 for (2-10.c) \rightarrow 2-10-3-0:smweh.pdf

Solution for (2-10.c) \rightarrow 2-10-3-1:smwes.pdf

What is the asymptotic complexity of the execution of the loop body in your implementation of rankoneinvit_fast() from Sub-problem (2-10.c)?

SOLUTION for (2-10.d) \rightarrow 2-10-4-0:smwecs.pdf

•

(2-10.e) (30 min.) Tabulate the runtimes of the two C++ implementations with different vector sizes $n=2^p$, with $p=2,3,\ldots,8$. As test vector d use Eigen::VectorXd::LinSpaced(n,1,2). Estimate the laws of the different asymptotic complexities of the two implementations using the measured runtimes.

SOLUTION for (2-10.e) \rightarrow 2-10-5-0:smwcds.pdf

 \blacksquare

End Problem 2-10, 125 min.

Problem 2-11: Sequential linear systems

Consider a sequence of linear systems when all the linear systems share the same matrix \mathbf{A} : $\mathbf{A} \times \mathbf{b}_i$. The computational cost of this problem may be reduced by performing an LU decomposition of \mathbf{A} only once and reusing it for the different systems.

Related to [Lecture \rightarrow Rem. 2.5.0.10]

Consider the following pseudocode with input data $\mathbf{A} \in \mathbb{R}^{n \times n}$ and $\mathbf{b} \in \mathbb{R}^n$:

```
Algorithm 2.11.1: Code solvepermb Require: Matrix \mathbf{A} \in \mathbb{R}^{n \times n}, vector \mathbf{b} \in \mathbb{R}^n Ensure: Matrix \mathbf{X} = [\mathbf{X}_1, \dots, \mathbf{X}_n] \in \mathbb{R}^{n \times n}, \mathbf{X}_i \in \mathbb{R}^n while not all cyclic permutations of \mathbf{b} tested yet do \mathbf{b} \leftarrow [b_n, b_1, b_2, \dots, b_{n-1}]^\top \mathbf{X}_i = \mathbf{A}^{-1}\mathbf{b} end while
```

(2-11.a) • What is the worst case asymptotic complexity of the function solvepermb for $n \to \infty$?

HIDDEN HINT 1 for (2-11.a) \rightarrow 2-11-1-0:permch.pdf

SOLUTION for (2-11.a) \rightarrow 2-11-1-1:permcs.pdf

(2-11.b) • Port the function solvepermb (Code 2.11.1) to C++ using EIGEN.

HIDDEN HINT 1 for (2-11.b) \rightarrow 2-11-2-0:permph.pdf

SOLUTION for (2-11.b) \rightarrow 2-11-2-1:permps.pdf

(2-11.c) \odot Design an efficient C++ implementation of function solvepermb using EIGEN with asymptotic complexity $O(n^3)$.

HIDDEN HINT 1 for (2-11.c) \rightarrow 2-11-3-0:permeh.pdf

SOLUTION for (2-11.c) \rightarrow 2-11-3-1:permes.pdf

End Problem 2-11

Problem 2-12: Structured linear systems

In this problem we come across the example of a structured matrix, for which a linear system can be solved very efficiently, though this is not obvious.

Related to |Lecture → Section 2.6|

Consider the linear system Ax = b, where the $n \times n$ matrix A has the following structure:

$$\mathbf{A} = \begin{bmatrix} a_1 & 0 & 0 & \dots & 0 \\ a_1 & a_2 & 0 & \ddots & \vdots \\ a_1 & a_2 & a_3 & \ddots & \\ \vdots & \vdots & & \ddots & 0 \\ a_1 & a_2 & a_3 & \dots & a_n \end{bmatrix}$$

(2-12.a) \odot Give necessary and sufficient conditions for the vector $\mathbf{a}=(a_1,\ldots,a_n)^\top\in\mathbb{R}^n$ such that the matrix **A** is non-singular.

```
HIDDEN HINT 1 for (2-12.a) \rightarrow 2-12-1-0:structuredLSE1h.pdf
```

SOLUTION for (2-12.a)
$$\rightarrow$$
 2-12-1-1:structuredLSE1s.pdf

(2-12.b) Write a C++ function that builds the matrix A given the vector a:

```
Eigen::MatrixXd buildA(const Eigen::VectorXd & a);
```

You can use EIGEN classes for your code.

```
HIDDEN HINT 1 for (2-12.b) \rightarrow 2-12-2-0:structuredLSE2h.pdf
```

```
SOLUTION for (2-12.b) \rightarrow 2-12-2-1: structuredLSE2s.pdf
```

(2-12.c) Implement a function in C++ which solves the linear system Ax = b by means of "structure" oblivious" Gaussian elimination, for which EIGEN's dedicated classes and methods should be used.

```
void solveA(const Eigen::VectorXd & a, const Eigen::VectorXd & b,
  Eigen::VectorXd & x);
```

The input is composed of vectors **a** and **b**.

Run this function with **a** and **b** made of random elements and explore the cases $n = 2^k$, $k = 4, \dots, 12$. What is the asymptotic complexity in n of this naive implementation?

```
HIDDEN HINT 1 for (2-12.c) \rightarrow 2-12-3-0:structuredLSE3h.pdf
```

```
SOLUTION for (2-12.c) \rightarrow 2-12-3-1:structuredLSE3s.pdf
```

(2-12.d) \Box Given a generic vector $\mathbf{a} \in \mathbb{R}^n$, compute symbolically (i.e. by hand) the inverse of matrix \mathbf{A} and the solution x of Ax = b.

HIDDEN HINT 1 for (2-12.d) \rightarrow 2-12-4-0:structuredLSE4h.pdf

```
SOLUTION for (2-12.d) \rightarrow 2-12-4-1:structuredLSE4s.pdf
```

(2-12.e) Implement a function in C++ which efficiently solves the linear system Ax = b, using EIGEN classes. The input is consists of vectors **a** and **b**.

```
HIDDEN HINT 1 for (2-12.e) \rightarrow 2-12-5-0:structuredLSE5h.pdf
```

Solution for (2-12.e)
$$\rightarrow$$
 2-12-5-1:structuredLSE5s.pdf

(2-12.f) • Compare the timing of the naive implementation of (2-12.c) with the efficient solution of (2-12.e).

HIDDEN HINT 1 for (2-12.f) \rightarrow 2-12-6-0:structuredLSE6h.pdf

Solution for (2-12.f) \rightarrow 2-12-6-1:structuredLSE6s.pdf

End Problem 2-12

Problem 2-13: Conversion of Triplet (COO) Matrix Format to CRS Format

This exercise is about sparse matrices. Make sure you are prepared on the subject by reading [Lecture \rightarrow Section 2.7] in the lecture notes. In particular, read through the various *sparse storage formats* discussed in class (cf. [Lecture \rightarrow Section 2.7.1]).

Focus of the problem is implementation in C++, using EIGEN in one part.

Remark. The ultimate goal is to devise a function that converts a matrix given in **triplet (or COOrdinate) format** (COO, see [Lecture \rightarrow § 2.7.1.1]) to the **compressed row storage format** (CRS, see [Lecture \rightarrow § 2.7.1.4]). You do not need to rely on EIGEN to solve this problem.

Recall that the COO format stores a collection of triplets (i, j, v), with $i, j \in \mathbb{N}$, $i, j \geq 0$ (the indices) and $v \in \mathbb{R}$ (the value in cell (i, j)). Multiple triplets corresponding to the same cell (i, j) are allowed, meaning that multiple values with the same indices (i, j) should be summed up when fully expressing the matrix.

In C++ the data for a matrix stored in COO format and with entry type **SCALAR** can be stored in the following data structure

```
template <typename SCALAR>
struct TripletMatrix {
   std::size_t n_rows {0}; // Number of rows
   std::size_t n_cols {0}; // Number of columns
   std::vector<std::tuple<std::size_t, std::size_t, SCALAR>> triplets;
};
```

Also remember that the CRS format uses three vectors:

- 1. val, which stores the values of the nonzero cells, one row after the other.
- 2. col ind, which stores the column indices of the nonzero cells, one row after the other.
- 3. row_ptr, which stores the index of the entry in val/col_ind containing the first element of each row.

The case of rows composed of only zero elements requires special consideration. The convention is that row_ptr stores two consecutive entries with the same values, meaning that in vectors val and col_ind the next row begins at the same index of the previous row (which implies that the previous row must be empty). This convention should be taken into account when e.g. iterating through row_ptr . However, to simplify things, in this problem you may always consider *matrices without all-zero rows*.

A suitable C++ data structure for storing a matrix with entries of type **SCALAR** in CRS format is the following

```
template <typename SCALAR>
struct CRSMatrix {
   std::size_t n_rows {0}; // Number of rows
   std::size_t n_cols {0}; // Number of columns
   std::vector <SCALAR> val; // Value array
   std::vector <std::size_t > col_ind; // Column indices
   std::vector <std::size_t > row_ptr; // Row pointers
};

(2-13.a) ③ (30 min.) Implement the C++ functions
   template <typename SCALAR>
```

```
Eigen::Matrix<SCALAR, Eigen::Dynamic, Eigen::Dynamic>
    densify(const TripletMatrix<SCALAR> &M);
template <typename SCALAR>
Eigen::Matrix<SCALAR, Eigen::Dynamic, Eigen::Dynamic>
    densify(const CRSMatrix<SCALAR> &M);
```

that convert the argument matrices into EIGEN's dense matrix storage type.

```
SOLUTION for (2-13.a) \rightarrow 2-13-1-0:triplettoCRS3s.pdf
```

(2-13.b) ③ (60 min.) Write a C++ function that converts a matrix **T** in COO format to a matrix in CRS format and returns the latter:

```
template <typename SCALAR>
CRSMatrix<SCALAR> tripletToCRS(const TripletMatrix<SCALAR>& T);
```

Your implementation should have (almost, maybe up to a logarithm) optimal complexity, that is, its computational effort should scale linearly with the number of triplets in the COO format.

```
HIDDEN HINT 1 for (2-13.b) \rightarrow 2-13-2-0:triplettoCRS4h.pdf
SOLUTION for (2-13.b) \rightarrow 2-13-2-1:triplettoCRS4s.pdf
```

(2-13.c) • What is the worst-case complexity of your function tripletToCRS ((2-13.b)) in the number of triplets?

```
HIDDEN HINT 1 for (2-13.c) \rightarrow 2-13-3-0:triplettoCRS5h.pdf
SOLUTION for (2-13.c) \rightarrow 2-13-3-1:triplettoCRS5s.pdf
```

(2-13.d) (20 min.) [depends on Sub-problem (2-13.a)]

Test the correctness of your implementation of tripletToCRS() by checking whether the function densify() from Sub-problem (2-13.a) yield the same matrix when applied before or after the call to tripletToCRS().

For testing use a **TripletMatrix<double>** object containing an $n \times n$ -matrix defined by 5n triplets with random indices and value 1. To generate random indices use **std:** rand()% n. Implement the test as a C++ function

```
bool testTripletToCRS(std::size_t n);
```

which accepts the matrix dimension as arguments and returns true, if the test is passed.

```
HIDDEN HINT 1 for (2-13.d) \rightarrow 2-13-4-0:6h1.pdf
SOLUTION for (2-13.d) \rightarrow 2-13-4-1:triplettoCRS6s.pdf
```

End Problem 2-13, 110 min.

Problem 2-14: Sparse matrices in CCS format

Sparse matrices must be stored in special formats that avoid storing the many zero entries. Sometimes one has to manipulate sparse matrices on that basic level. Therefore it is important to be familiar with some details of the storage formats. This exercise focuses on the CCS format.

Related to [Lecture \rightarrow Section 2.7.1]

In [Lecture \rightarrow § 2.7.1.4] the so-called *CRS format* (Compressed Row Storage) was introduced. It uses three arrays (val, col_ind and row_ptr) to store a sparse matrix.

EIGEN can also handle the *CCS format* (Compressed Column Storage), which is like CRS for the transposed matrix. It relies on the arrays val, row_ind and col_ptr.

```
void CCS(const Eigen::MatrixXd & A, Eigen::VectorXd & val,
Eigen::VectorXd & row_ind, Eigen::VectorXd & col_ptr);
```

The CCS data are returned through the vectors val, row_ind, and col_ptr.

Prohibited. While you can use EIGEN classes such as Eigen::MatrixXd, do not use EIGEN methods to directly access val, row_ind and col_ptr. In other words, you must not use EIGEN's class Eigen::SparseMatrix and simply run makeCompressed().

You may assume that A does not have a a column with all elements equal to 0 (Otherwise, it may become problematic to update col_ptr.).

In this problem you may test for exact equality with 0.0, because zero matrix entries are not supposed to be results of floating point computations.

```
SOLUTION for (2-14.a) \rightarrow 2-14-1-0: smwps.pdf
```

(2-14.b) (5 min.) [depends on Sub-problem (2-14.a)]

What is the computational complexity of your CCS function:

- with respect to the matrix size n, where $\mathbf{A} \in \mathbb{R}^{n \times n}$?
- with respect to nnz, which denotes the number of nonzero elements in A?

```
Solution for (2-14.b) \rightarrow 2-14-2-0: smwps.pdf
```

End Problem 2-14, 35 min.

Problem 2-15: Ellpack sparse matrix format

The number of processing cores of high-performance computers has seen an exponential growth. However, the increase in *memory bandwidth*, which is the rate at which data can be read from or written into memory by a processor, has been slower. Hence, memory bandwidth is a bottleneck for several algorithms.

Several papers present storage formats to improve the performance of sparse matrix-vector multiplications. One example is the **Ellpack format**, where the number of nonzero entries per row is bounded and shorter rows are padded with zeros.

Related to [Lecture \rightarrow Section 2.7.1], which should be studied before tackling this problem.

The C++ class **EllpackMat** whose definition is listed below implements the Ellpack sparse matrix format for a generic matrix $\mathbf{A} \in \mathbb{R}^{m,n}$:

C++ code 2.15.1: Declaration of EllpackMat class

```
class EllpackMat {
   public:
3
    EllpackMat(const Triplets &triplets, index_t m, index_t n);
    double operator()(index_t i, index_t j) const;
    void mvmult(const Vector &x, Vector &y) const;
    void mtvmult(const Vector &x, Vector &y) const;
    index t get maxcols() const;
   private:
10
    std::vector<double> val; //< Vector containing values</pre>
11
    // corresponding to entries in 'col'
12
    std::vector<index_t> col; //< Vector containing column</pre>
13
    // indices of the entries in 'val'.
    // The position of a cell in 'val' and 'col'
15
    // is determined by its row number and original position in 'triplets'
16
    index_t maxcols; //< Number of non-empty columns</pre>
17
    index_t m, n; //< Number of rows, number of columns</pre>
18
  };
19
```

Get it on ₩ GitLab (ellpack.hpp).

In the code above the following typedefs are used

```
using Triplet = Eigen::Triplet<double>;
using Triplets = std::vector<Triplet>;
```

Further, the data member maxcols is defined as the maximum number of non-zero entries in a row of the matrix:

```
\max cols := \max \{ \# [(i,j) \mid \mathbf{A}_{i,j} \neq 0, j = 1,...,n], i = 1,...,m \}.
```

The entry read-access **operator** () method, which takes the entries $(\mathbf{A})_{i,j}$ row and column indices i and j as arguments, is implemented as follows. It completely "defines" the Ellpack sparse matrix format.

C++11-code 2.15.2: Definition of entry-access operator operator ()

```
double EllpackMat::operator()(index_t i, index_t j) const {
   assert(0 <= i && i < m && 0 <= j && j < n && "Index out of bounds!");

for (index_t l = i * maxcols; l < (i + 1) * maxcols; ++l) {
   if (col.at(l) == j) return val.at(l);
```

It is well defined for i = 0, ..., m-1 and for j = 0, ..., n-1 (C++ indexing!).

(2-15.a) ☐ (60 min.) Implement an efficient constructor of EllpackMat that builds an object of type EllpackMat from data forming a matrix in triplet format. In other words, implement the constructor with the following signature:

```
EllpackMat(const Triplets & triplets, index_t m, index_t n);
```

triplets is a std::vector of EIGEN triplets (see [Lecture \rightarrow Section 2.7.2]). The arguments m and n represent the number of rows and columns of the matrix A.

- The data in the triplet vector must be compatible with the matrix size provided.
- to the constructor. No assumption is made on the ordering of the triplets.
- Values belonging to multiple occurrences of the same index pair are to be summed up.



Taking care of repeated index pairs requires sophisticated bookkeeping.

```
HIDDEN HINT 1 for (2-15.a) \rightarrow 2-15-1-0:ellpack1h.pdf
HIDDEN HINT 2 for (2-15.a) \rightarrow 2-15-1-1:ellp1h1.pdf
HIDDEN HINT 3 for (2-15.a) \rightarrow 2-15-1-2:sellh2.pdf
SOLUTION for (2-15.a) \rightarrow 2-15-1-3:ellpack1s.pdf
```

```
void mvmult(const Eigen::VectorXd &x, Eigen::VectorXd &y) const;
```

A is the matrix represented by \star this. The implementation must run with the optimal complexity of O(nnz(A)).

```
HIDDEN HINT 1 for (2-15.b) \rightarrow 2-15-2-0:ellpack2h.pdf
```

```
SOLUTION for (2-15.b) \rightarrow 2-15-2-1:ellpack2s.pdf
```

```
HIDDEN HINT 1 for (2-15.c) \rightarrow 2-15-3-0:ellpack3h.pdf
```

Solution for (2-15.c)
$$\rightarrow$$
 2-15-3-1:ellpack3s.pdf

(2-15.d) (10 min.) [depends on Sub-problem (2-15.a)]

Discuss the asymptotic complexity of your implementation of the constructor of **EllpackMat** in terms of growing number of triplets, $\sharp \text{triplets} \to \infty$.

SOLUTION for (2-15.d)
$$\rightarrow$$
 2-15-4-0:sellx.pdf

End Problem 2-15, 120 min.

Problem 2-16: Grid functions

This exercise deals with construction of sparse matrices from triplet format, see [Lecture \rightarrow Section 2.7.2]. This task is relevant for applications like image processing and the numerical solution of partial differential equations.

Connected with [Lecture \rightarrow Section 2.7]

Consider the following matrix $S \in \mathbb{R}^{3,3}$:

$$\mathbf{S} := \begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$$

Denote by $s_{i,j}$ the entries of **S**. Consider the following *linear* operator from $\mathbb{R}^{n,m}$ to $\mathbb{R}^{n,m}$:

$$L: \left\{ \begin{array}{ccc} \mathbb{R}^{n,m} & \to & \mathbb{R}^{n,m} \\ \mathbf{X} & \mapsto & L(\mathbf{X}): & (L(\mathbf{X}))_{i,j} := \sum\limits_{k,l=1}^{3} s_{k,l}(\mathbf{X})_{i+k-2,j+l-2}, \ i \in \{1,\ldots,n\}, \ m \in \{1,\ldots,m\} \end{array} \right.$$

where we adopt the convention that in the sum non-existent matrix entries are set to zero.

(2-16.a) \Box (10 min.) Show that L is a linear operator.

SOLUTION for (2-16.a)
$$\rightarrow$$
 2-16-1-0:grfl.pdf

(2-16.b) • The space $\mathbb{R}^{n,m}$, as a vector-space, is isomorphic to $\mathbb{R}^{n \cdot m}$, via the mapping $\mathbf{X} \mapsto \text{vec}(\mathbf{X})$, see [Lecture \to Eq. (1.2.3.5)].

From linear algebra, we know that any linear operator over a finite dimensional (real) vector space can be represented by a (real) matrix multiplication. We define the matrix $\mathbf{A} \in \mathbb{R}^{nm,nm}$ s.t.

$$\mathbf{A} \cdot \text{vec}(\mathbf{X}) = \text{vec}(L(\mathbf{X})).$$

Write down the matrix **A** for n = m = 3.

Solution for (2-16.b)
$$\rightarrow 2-16-2-0$$
:gfrwa.pdf

(2-16.c) (15 min.) Write a C++ function

which, given a matrix $\mathbf{X} \in \mathbb{R}^{n,m}$ and a function $f : \mathbb{N}^2 \to \mathbb{R}$, fills the matrix \mathbf{X} according to the rule $(\mathbf{X})_{i,j} := f(i,j), i \in \{1,\ldots,n\}, j \in \{1,\ldots,m\}.$

Additionally, inside the auxiliary function $define_f()$ define a lambda function f, with signature **double** (Eigen::Index, Eigen::Index), implementing:

$$f: \left\{ \begin{array}{ccc} \mathbb{N}^2 & \to & \mathbb{R}, \\ (i,j) & \mapsto & \begin{cases} 1 & i > n/4 \wedge i < 3n/4 \wedge j > m/4 \wedge j < 3m/4 \,, \\ 0 & \text{otherwise}. \end{cases} \right.$$

SOLUTION for (2-16.c) \rightarrow 2-16-3-0:gfrwa.pdf

```
(2-16.d) (30 min.) [depends on Sub-problem (2-16.c)]
```

Implement a function

which, given the stencil matrix $S \in \mathbb{R}^{3,3}$ and the size of the matrix X (passed through Xn, Xm), builds and returns the sparse matrix A in CRS format. The matrix A can be built using an intermediate triplet/COO format.

```
HIDDEN HINT 1 for (2-16.d) \rightarrow 2-16-4-0:2h1.pdf
```

```
SOLUTION for (2-16.d) \rightarrow 2-16-4-1:gfrwb.pdf
```

(2-16.e) (20 min.) Implement a function

which, given a matrix X and the sparse matrix A, returns the matrix Y s.t. vec(Y) := A vec(X). You can use objects of type Eigen::Map to "reshape" a matrix into a column vector.

```
SOLUTION for (2-16.e) \rightarrow 2-16-5-0:gfrwm.pdf
```

(2-16.f) : (20 min.)

Implement a function

which, given a matrix \mathbf{Y} and the sparse matrix \mathbf{A} , returns the matrix \mathbf{X} s.t. $\operatorname{vec}(\mathbf{Y}) := \mathbf{A} \operatorname{vec}(\mathbf{X})$. You can use objects of type Eigen: :Map to "reshape" a matrix into a column vector. Objects of this type are read and write compatible.

```
SOLUTION for (2-16.f) \rightarrow 2-16-6-0:gfrwi.pdf
```

End Problem 2-16, 95 min.

Problem 2-17: Efficient sparse matrix-matrix multiplication in COO format

The *triplet (or COOrdinate) list format* works very well for defining a matrix or adding/modifying its elements. This is however not so true for matrix operations such as matrix-matrix multiplications, unlike the CSC/CSR storage.

This problem is about (sparse) matrix-matrix multiplication in COO format. We will consider the following items:

- 1. The worst case of sparse matrix-matrix multiplication, when one wants to use sparse matrix storage formats.
- 2. The most efficient way to tackle this problem with the COO format.
- 3. The asymptotic complexity of such algorithm (which is, by definition, the complexity under the worst-case scenario).

Related to [Lecture \rightarrow Section 2.7.1]

For simplicity, in the following we will only deal with *binary* matrices. Binary matrices are only made of 0 or 1. At the same time, we will allow for duplicates in our implementations of the triplet format (see [Lecture \rightarrow § 2.7.1.1]). The matrix entry associated to the pair (i, j) is defined to be the sum of all triplets corresponding to it.

For the subproblems involving coding, we define types Trip = Triplet < double > and TripVec = std::vector < trip>.

(2-17.a) lacktriangle Consider multiplications between sparse matrices: AB = C. Is the product C guaranteed to be sparse? Make an example of two sparse matrices which, when multiplied, return a dense matrix.

```
HIDDEN HINT 1 for (2-17.a) \rightarrow 2-17-1-0:matmatCOO1h.pdf
```

SOLUTION for (2-17.a) \rightarrow 2-17-1-1:matmatCOO1s.pdf

(2-17.b) \Box Implement a C++ function which returns the input matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ in COO format:

```
TripVec Mat2COO(const Eigen::MatrixXd &A);
```

You can use EIGEN classes.

```
SOLUTION for (2-17.b) \rightarrow 2-17-2-0: matmatCOO2s.pdf
```

(2-17.c) Implement a C++ function which computes the product between two matrices in COO format:

```
TripVec COOprod_naive(const TripVec &A, const TripVec &B);
```

You can use EIGEN classes.

```
HIDDEN HINT 1 for (2-17.c) \rightarrow 2-17-3-0:matmatCOO3h.pdf
```

```
SOLUTION for (2-17.c) \rightarrow 2-17-3-1:matmatC003s.pdf
```

(2-17.d) • What is the asymptotic complexity of your naive implementation in (2-17.c)?

HIDDEN HINT 1 for (2-17.d) \rightarrow 2-17-4-0:matmatCOO4h.pdf

```
SOLUTION for (2-17.d) \rightarrow 2-17-4-1:matmatC004s.pdf
```

```
TripVec COOprod_effic(TripVec &A, TripVec &B);
```

You can use EIGEN classes.

Can you do better than (2-17.c)?

HIDDEN HINT 1 for (2-17.e) \rightarrow 2-17-5-0:matmatCOO5h.pdf

SOLUTION for (2-17.e) \rightarrow 2-17-5-1:matmatC005s.pdf

•

(2-17.f) • What is the asymptotic complexity of your efficient implementation in (2-17.e)?

HIDDEN HINT 1 for (2-17.f) $\rightarrow 2-17-6-0$: matmatCOO6h.pdf

SOLUTION for (2-17.f) \rightarrow 2-17-6-1:matmatCOO6s.pdf

•

(2-17.g) \odot Compare the timing of COOprod_naive (2-17.c) and COOprod_effic (2-17.e) for random matrices with different dimensions n. Perform the comparison twice: first only for products between sparse matrices, then for any kind of matrix.

HIDDEN HINT 1 for (2-17.g) \rightarrow 2-17-7-0:matmatCOO7h.pdf

SOLUTION for (2-17.g) $\rightarrow 2-17-7-1$: matmatCOO7s.pdf

End Problem 2-17

Problem 2-18: Gauss-Seidel Iteation for a Sparse Linear System of Equations

This problem studies are particular iterative method defined by a matrix and a vector, the so-called Gauss-Seidel iteration. We consider its algorithmic aspects when the matrix is given in a raw CRS format.

Familiarity with [Lecture \rightarrow Section 2.7.1] is esssential for this problem. Some concepts from [Lecture \rightarrow Section 8.2.2] will also be used. Implementation in C++ based on Eigen will be reauested.

Given the *regular* square matrix $A \in \mathbb{R}^{n,n}$ and the vector $\mathbf{b} \in \mathbb{R}^n$, $n \in \mathbb{N}$, the **Gauss-Seidel iteration** $\mathbf{x}^{(k+1)} := \mathbf{\Phi}(\mathbf{x}^{(k)}), \, \mathbf{\Phi} : \mathbb{R}^n \to \mathbb{R}^n$, is defined as

$$\left(\mathbf{x}^{(k+1)}\right)_{i} = (\mathbf{A})_{i,i}^{-1} \left((\mathbf{b})_{i} - \sum_{j=1}^{i-1} (\mathbf{A})_{i,j} \left(\mathbf{x}^{(k+1)}\right)_{j} - \sum_{j=i+1}^{n} (\mathbf{A})_{i,j} \left(\mathbf{x}^{(k)}\right)_{j}\right), \quad i = 1, 2, \dots, n \text{ (in this order!) }.$$
(2.18.1)

Here sums for which the lower index is larger than the upper are meant to evaluate to zero.

The matrix A is available in a raw compressed row-storage (CRS) format, encapsulated in the following data type:

```
struct CRSMatrix {
 unsigned int m;
                                      // number of rows
 unsigned int n;
                                      // number of columns
  std::vector<double> val;
                                      // value array
  std::vector<unsigned int> col_ind; // same length as value array
  std::vector<unsigned int> row_ptr; // length m+1, row_ptr[m] ==
    val.size()
};
```

Refer to [Lecture \rightarrow § 2.7.1.4] for further explanations on the meaning of the fields of **CRSMatrix**.

(2-18.a) (10 min.) For what square matrices A and vectors b is the Gauss-Seidel iteration (2.18.1) well-defined?

```
Iteration (2.18.1) well-defined \forall \mathbf{A} \in \mathbb{R}^{n,n}, \mathbf{b} \in \mathbb{R}^n:
```

```
(2-18.b)  (45 min.)
```

SOLUTION for (2-18.a) \rightarrow 2-18-1-0:gscrss1.pdf

In the file gaussseidelcrs.hpp implement the C++ function

```
bool GaussSeidelstep crs(
  const CRSMatrix &A, const Eigen::VectorXd &b,
 Eigen::VectorXd &x);
```

that carries out a single step of the Gauss-Seidel iteration according to (2.18.1). The vector $\mathbf{x}^{(k)}$ is passed through x and it should return $\mathbf{x}^{(k+1)}$ also in x. **false** is returned in case the iteration is not well-defined, true otherwise.

```
SOLUTION for (2-18.b) \rightarrow 2-18-2-0: gscrss2.pdf
```

For given regular $A \in \mathbb{R}^{n,n}$ and $b \in \mathbb{R}^n$ assume that the Gauss-Seidel **(2-18.c)** (20 min.) iteration (2.18.1) is well-defined. Characterize the set of all **fixed points** of the Gauss-Seidel iteration (2.18.1) by a compact formula:

```
\{\mathbf{x}^* \in \mathbb{R}^n : \mathbf{x}^* \text{ is fixed point of (2.18.1)}\} =
```

A fixed point of an iteration $\mathbf{x}^{(k+1)} = F(\mathbf{x}^{(k)})$ is a vector \mathbf{x}^* such that $\mathbf{x}^* = F(\mathbf{x}^*)$.

```
SOLUTION for (2-18.c) \rightarrow 2-18-3-0: gss3.pdf
```

Develop a C++ function

```
bool GaussSeidel_iteration(
  const CRSMatrix &A, const Eigen::VectorXd &b, Eigen::VectorXd &x,
  double atol = 1.0E-8, double rtol = 1.0E-6,
  unsigned int maxit = 100);
```

which carries out Gauss-Seidel iterations for the data $A \in \mathbb{R}^{n,n}$, $b \in \mathbb{R}^n$ (arguments A and b), and with initial guess $\mathbf{x}^{(0)}$ passed in x. The same argument is used to return the final iterate. The iteration is to terminate after maxit steps. If this happens, false should be returned, true otherwise.

The function should rely on the correction-based termination criterion [Lecture \rightarrow Eq. (8.5.3.2)] with relative and absolute tolerance passed in rtol and atol, respectively and based on the Euclidean vector norm.

```
Solution for (2-18.d) \rightarrow 2-18-4-0:gss33.pdf
```

We apply the Gauss-Seidel iteration for

$$\mathbf{A} = \begin{bmatrix} 10 & 0 & 0 & 0 & -2 & 0 \\ 3 & 9 & 0 & 0 & 0 & 3 \\ 0 & 2 & 8 & 2 & 0 & 0 \\ 3 & 0 & 1 & 10 & 5 & 0 \\ 0 & 2 & 0 & 3 & 13 & 2 \\ 0 & 4 & 0 & 0 & 2 & 11 \end{bmatrix} , \mathbf{b} = \mathbf{0}.$$

and initial guess $\mathbf{x}^{(0)} = [1, 2, 3, 4, 5, 6]^{\mathsf{T}}$. The following table displays the progress of the iteration:

k	$\left\ \mathbf{x}^{(k)}\right\ _2$	$\left\ \mathbf{x}^{(k)}-\mathbf{x}^{(k-1)}\right\ $	$\left\ \mathbf{x}^{(k)}\right\ _2:\left\ \mathbf{x}^{(k-1)}\right\ _2$
1	3.863865e+00	1.127304e+01	
2	8.297690e-01	3.751824e+00	2.147510e-01
3	6.432366e-02	7.857529e-01	7.751996e-02
4	1.476504e-02	5.984439e-02	2.295429e-01
5	4.068577e-03	1.098126e-02	2.755548e-01
6	9.204751e-04	3.155445e-03	2.262401e-01
7	1.880926e-04	7.327086e-04	2.043429e-01
8	3.636497e-05	1.517466e-04	1.933355e-01
9	6.803531e-06	2.956273e-05	1.870902e-01
10	1.246867e-06	5.556761e-06	1.832676e-01
11	2.254491e-07	1.021425e-06	1.808125e-01
12	4.039721e-08	1.850525e-07	1.791855e-01
13	7.194116e-09	3.320315e-08	1.780844e-01
14	1.275722e-09	5.918398e-09	1.773286e-01

Describe qualitatively and quantitatively the observed asymptotic convergence of the iteration for this example.

Solution for (2-18.e) \rightarrow 2-18-5-0:gss5.pdf

End Problem 2-18, 110 min.

Problem 2-19: A special Sylvester Equation

A Sylvester equations is a linear matrix equation of the form

$$\mathbf{AX} + \mathbf{XB} = \mathbf{C}$$
 , $\mathbf{A}, \mathbf{B}, \mathbf{C} \in \mathbb{R}^{n,n}$,

which yields a linear system of equations for the entries of the unknown matrix $X \in \mathbb{R}^{n,n}$. In this exercise we consider a very special specimen of Sylvester equation and recast it for the sake of efficient direct solution with EIGEN.

Related to various topics from class and also Problem 2-8.

(2-19.a)
$$\odot$$
 (5 min.) Refresh your knowledge about s.p.d. matrices, see [Lecture \rightarrow Def. 1.1.2.6] and [Lecture \rightarrow Lemma 1.1.2.7].

(2-19.b) \odot (5 min.) Recall the "vectorization" $\operatorname{vec}(\mathbf{X}) \in \mathbb{R}^{n^2}$ of a matrix $\mathbf{X} \in \mathbb{R}^{n,n}$ from [Lecture \rightarrow Rem. 1.2.3.4] and jot down $\operatorname{vec}(\mathbf{X})$ for

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{bmatrix} \in \mathbb{R}^{2,2} .$$

SOLUTION for (2-19.b)
$$\rightarrow 2-19-2-0$$
:sla.pdf

(2-19.c) \Box (15 min.) For $\mathbf{S}, \mathbf{T} \in \mathbb{R}^{2,2}$ find the matrix $\mathbf{C} \in \mathbb{R}^{4,4}$ such that

$$\operatorname{vec}(\mathbf{SXT}) = \mathbf{C}\operatorname{vec}(\mathbf{X}) \quad \forall \mathbf{X} \in \mathbb{R}^{2,2}$$
,

that is, express the entries of C in terms of the entries $(S)_{i,j}$ and $(T)_{i,j}$ of S and T.

SOLUTION for (2-19.c)
$$\rightarrow 2-19-3-0:slb.pdf$$

$$vec(SXT) = C vec(X)$$
.

HIDDEN HINT 1 for (2-19.d) $\rightarrow 2-19-4-0:s1h1.pdf$

SOLUTION for (2-19.d)
$$\rightarrow 2-19-4-1$$
:s2.pdf

Throughout the remainder of this problem let $A \in \mathbb{R}^{n,n}$, $n \in \mathbb{N}$, stand for a *symmetric, positive definite* (s.p.d.) matrix, whose rows and columns contain at most five non-zero entries. Hence, for large n the matrix A can be considered *sparse* in the sense of [Lecture \rightarrow Notion 2.7.0.1].

Then consider the special Sylvester equation:

seek
$$X \in \mathbb{R}^{n,n}$$
: $XA^{-1} + AX = I$, (2.19.4)

and the equivalent matrix equation:

seek
$$\mathbf{X} \in \mathbb{R}^{n,n}$$
: : $\mathbf{X} + \mathbf{A}\mathbf{X}\mathbf{A} = \mathbf{A}$. (2.19.5)

(2-19.e) \odot (5 min.) Both (2.19.4) and (2.19.5) can be recast as an $n^2 \times n^2$ linear system of equations for the unknown vector $\text{vec}(\mathbf{X})$. Based on the result of Sub-problem (2-19.d) express their coefficient matrices by means of the Kronecker product [Lecture \rightarrow Def. 1.4.3.7].

Solution for (2-19.e)
$$\rightarrow 2-19-5-0$$
:s3a.pdf

(2-19.f)
☐ (10 min.) Based on the result of Sub-problem (2-19.e) and in light of the observation made in [Lecture → Ex. 2.7.4.4] give a sharp upper bound for the number of non-zero entries of the coefficient matrices of those linear systems for (2.19.4) and (2.19.5).

```
HIDDEN HINT 1 for (2-19.f) \rightarrow 2-19-6-0:s3h1.pdf
```

Solution for (2-19.f) $\rightarrow 2-19-6-1:s4.pdf$

(2-19.g) (15 min.) Show that both (2.19.4) and (2.19.5) give rise to linear systems of equations for the entries of **X** with *s.p.d.* coefficient matrices.

For your argument you can appeal to the following result:

Theorem 2.19.9. Eigenvalues of Kronecker-product matrices

Assume that both $S, T \in \mathbb{R}^{n,n}$ can be diagonalized and write $\sigma(S), \sigma(T) \subset \mathbb{C}$ for the sets of their eigenvalues. Then the set of eigenvalues of the Kronecker product is given by

$$\sigma(\mathbf{S} \otimes \mathbf{T}) = \{ \lambda \mu : \lambda \in \sigma(\mathbf{S}), \mu \in \sigma(\mathbf{T}) \}$$
.

HIDDEN HINT 1 for (2-19.g) \rightarrow 2-19-7-0:s4h1.pdf

HIDDEN HINT 2 for (2-19.g) $\rightarrow 2-19-7-1:s4h2.pdf$

SOLUTION for (2-19.g) $\rightarrow 2-19-7-2:s4.pdf$

_

```
template <typename Vector>
```

```
Eigen::SparseMatrix<double> solveDiagSylvesterEq(const Vector
   &diagA);
```

that solves (2.19.4) with the diagonal of $\bf A$ passed through the vector argument diagA. The result $\bf X$ is returned as a sparse matrix. The type **Vector** must supply a size() method and component access through **operator** [].

```
SOLUTION for (2-19.h) \rightarrow 2-19-8-0:s5.pdf
```

```
Eigen::SparseMatrix<double> sparseKron(const
    Eigen::SparseMatrix<double> &M);
```

that computes the Kronecker product $\mathbf{M} \otimes \mathbf{M}$ for the square sparse matrix $\mathbf{M} \in \mathbb{R}^{n,n}$ and returns the result in sparse-matrix format.

For the sake of efficient initialization you have to use the reserve () function and determine the maximum number of non-zero entries per column for $\bf A$ in advance. To do this use the member functions of **SparseMatrix** that allow direct access to the data vectors of the CCS format, *cf.* [Lecture \rightarrow § 2.7.1.4]:

- const double *valuePtr() const: returns a pointer to the array of values of non-zero matrix entries.
- const StorageIndex* innerIndexPtr() const: returns a pointer to the array of row indices for every non-zero matrix entry,
- const StorageIndex* outerIndexPtr() const: returns a pointer to the vector of starting positions of columns.

These access methods also permit you to run through all non-zero entries of A and retrieve their locations.

```
SOLUTION for (2-19.i) \rightarrow 2-19-9-0:s6.pdf

(2-19.j) ① (20 min.) [depends on Sub-problem (2-19.i), Sub-problem (2-19.e)]

Write a C++ function

Eigen::MatrixXd solveSpecialSylvesterEq(
const Eigen::SparseMatrix<double> &A);

that solves (2.19.4) and returns the solution X.
```

End Problem 2-19, 155 min.

SOLUTION for (2-19.j) $\rightarrow 2-19-10-0:s7.pdf$

Problem 2-20: Matrix-Vector Product in CRS format

Sparse matrices have to be stored in special formats in order to save memory and inform algorithms about the position of non-zero entries. This problem will examine the CRS format

This problems is related to [Lecture \rightarrow § 2.7.1.4] and assumes familiarity with C++.

The following data structure is used to store an $m \times n$ -matrix in compressed row-storage (CRS) format:

This format is defined by the relationship ("C++ indexing")

```
	ext{val}[k] = a_{ij} \Leftrightarrow \left\{ egin{array}{ll} 	ext{col\_ind}[k] = j \ , \ 	ext{row\_ptr}[i] \leq k < 	ext{row\_ptr}[i+1] \ , \end{array} 
ight. \quad 0 \leq k \leq 	ext{val.size()} \ , \ 	ext{for } i \in \{0, \ldots, m-1\}, j \in \{0, \ldots, n-1\}. \end{array} 
ight.
```

(2-20.a)
☐ (40 min.) The function crsmv is supposed to return the product of a matrix in CRS format passed through M and of a vector given as argument x. Supplement the missing parts of the following listing code by writing valid C++ code into the boxes.

```
template <typename VECTORTYPE_I, typename VECTORTYPE_II>
VECTORTYPE_I crsmv(const CRSMatrix &M,
                    const VECTORTYPE_II &x) {
  assert((x.size() == M.
         && "Size mismatch between x and M");
  VECTORTYPE_I y (M.m);
  for (int k = 0; k <
                                 ; ++k) {
    y[k] = 0;
    for (int j = M.
                              ; j <
                                                ; ++j) {
                                                                 ];
      У[
                                            * x [
    }
  return y;
}
```

SOLUTION for (2-20.a) \rightarrow 2-20-1-0:crsmv1.pdf

End Problem 2-20, 40 min.

Problem 2-21: Asymptotic Complexity of Numerical Linear Algebra Operations

Most numerical algorithms rely on basic operations from numerical linear algebra. In order to gauge the computational cost of these algorithms, precise knowledge about the asymptotic computational effort involved in these operations is required. This problem asks you to analyze a given algorithm in this respect.

This problem draws on [Lecture \rightarrow Section 1.4.2] and [Lecture \rightarrow Section 2.5]. No C++ implementation is requested.

The function

```
std::pair<double, Eigen::VectorXd>
invit(Eigen::MatrixXd& A, double tol);
```

implements a so-called inverse power iteration. It takes a general (invertible) matrix $\mathbf{A} \in \mathbb{R}^{n,n}$, $n \in \mathbb{N}$, as argument \mathbb{A} .

The C++ code for invit () is listed in Code 2.21.1.

C++ code 2.21.1: Function implementing inverse power iteration

```
std::pair<double, Eigen::VectorXd> invit(Eigen::MatrixXd& A,
                                              double tol = 1.0E-6) {
3
     assert((A.cols() == A.rows()) && "A must be square!");
     const unsigned int n = A.cols();
5
     Eigen::VectorXd y_new = Eigen::VectorXd::Random(n);
6
     Eigen::VectorXd y_old(n);
7
     double I_new = y_new.transpose() * A * y_new;
     double | old;
     const auto A_lu_dec = A.lu(); //
10
     y_new = y_new / y_new.norm();
11
12
     do {
       l_old = l_new;
13
       y_old = y_new; //
14
       y_new = A_lu_dec.solve(y_old); //
15
       y_new /= y_new.norm(); //
16
       I_new = y_new.transpose() * A * y_new; //
17
     } while (std::abs(l_new - l_old) > tol * std::abs(l_new));
18
     return {l_new, y_new};
19
20
```

(2-21.a) (1 min.) What is the asymptotic complexity of the operations in Line 10,

```
const auto A_lu_dec = A.lu();
```

of Code 2.21.1?

Cost(Line 10 of Code 2.21.1) =
$$O(n \to \infty)$$

SOLUTION for (2-21.a) \rightarrow 2-21-1-0:cits1.pdf

(2-21.b) (1 min.) What is the asymptotic complexity of the operations in Line 14, $y_old = y_new$;

of Code 2.21.1?

Cost(Line 14 of Code 2.21.1) =
$$O(n \rightarrow \infty)$$

SOLUTION for (2-21.b) \rightarrow 2-21-2-0:cits2.pdf

(2-21.c) (1 min.) What is the asymptotic complexity of the operations in Line 15, $y_new = A_lu_dec.solve(y_old)$;

of Code 2.21.1?

Cost(Line 15 of Code 2.21.1) =
$$O(n \rightarrow \infty)$$

SOLUTION for (2-21.c) \rightarrow 2-21-3-0:cits3.pdf

(2-21.d) (1 min.) What is the asymptotic complexity of the operations in Line 16, y_new /= y_new.norm();

of Code 2.21.1?

Cost(Line 16 of Code 2.21.1) =
$$O(n \to \infty)$$

SOLUTION for (2-21.d) \rightarrow 2-21-4-0:cits4.pdf

(2-21.e)
○ (1 min.) What is the asymptotic complexity of the operations in Line 17,

l_new = y_new.transpose() * A * y_new;

of Code 2.21.1?

Cost(Line 17 of Code 2.21.1) =
$$O(n \rightarrow \infty)$$

SOLUTION for (2-21.e) \rightarrow 2-21-5-0:cits5.pdf

End Problem 2-21, 5 min.

Chapter 3

Direct Methods for Linear Least Squares Problems

Problem 3-1: Hidden linear regression

This problem is about a hidden linear regression: in fact, at first glance it will seem that we are dealing with a nonlinear system of equations. However, we will be able to reduce the system to a linear form.

Linear regression was introduced in [Lecture \rightarrow Ex. 3.0.1.1] and [Lecture \rightarrow Ex. 3.0.1.4]. You have to be familiar with normal equations, see [Lecture \rightarrow Section 3.1.2] and, in particular, [Lecture \rightarrow Ex. 3.1.2.4]. This problems also addresses the solution of linear least-squares problems in EIGEN, see [Lecture \rightarrow Code 3.2.0.1] and [Lecture \rightarrow Code 3.3.4.2].

We consider the function $f(t) = \alpha e^{\beta t}$ with unknown parameters $\alpha > 0$, $\beta \in \mathbb{R}$. Given are measurements (t_i, y_i) , i = 1, ..., n, $1 < n \in \mathbb{N}$. We want to determine $1 < n \in \mathbb{R}$ such that $1 < n \in \mathbb{R}$

$$f(t_i) = y_i$$
, with $y_i > 0$, $i = 1, ..., n$. (3.1.1)

(3-1.a) □ (10 min.) Which overdetermined system of *nonlinear* equations directly stems from (3.1.1)?

SOLUTION for (3-1.a)
$$\rightarrow$$
 3-1-1-0:AdaptedLinReg1s.pdf

(3-1.b) (15 min.) In which overdetermined system of *linear* equations can you transform the nonlinear system derived in (3-1.a)?

HIDDEN HINT 1 for (3-1.b) \rightarrow 3-1-2-0:AdaptedLinReg2h.pdf

SOLUTION for (3-1.b)
$$\rightarrow$$
 3-1-2-1:AdaptedLinReg2s.pdf

(3-1.c) (15 min.) Determine the **normal equations** for the overdetermined linear system of (3-1.b).

HIDDEN HINT 1 for (3-1.c) \rightarrow 3-1-3-0:AdaptedLinReg3h.pdf

```
SOLUTION for (3-1.c) \rightarrow 3-1-3-1:AdaptedLinReg3s.pdf
```

(3-1.d) \square (20 min.) Consider the data $(t_i, y_i) \in \mathbb{R}^2$, i = 1, ..., m, stored in two vectors t and y of equal length. In the file adaptedling, hpp implement a C++ function

```
Eigen::VectorXd linReg(const Eigen::VectorXd &t, const
    Eigen::VectorXd &y);
```

that uses the *method of normal equations* discussed in [Lecture \rightarrow Section 3.2] to solve the 1D linear regression problem as described in [Lecture \rightarrow Ex. 3.0.1.1] in least-squares sense. This function should return the estimated parameters α and β in a vector of length 2.

```
SOLUTION for (3-1.d) \rightarrow 3-1-4-0: AdaptedLinReg4s.pdf (3-1.e) (15 min.) [depends on Sub-problem (3-1.d), Sub-problem (3-1.a)]
```

Now consider the data (t_i, y_i) stored in two vectors t and f of equal length. Implement a C++ function

```
Eigen::VectorXd expFit(const Eigen::VectorXd &t, const
    Eigen::VectorXd &y);
```

that returns a least squares estimate of α and β , given the nonlinear relationship (3.1.1) between t and f:

```
HIDDEN HINT 1 for (3-1.e) \rightarrow 3-1-5-0:AdaptedLinReg5h.pdf
SOLUTION for (3-1.e) \rightarrow 3-1-5-1:AdaptedLinReg5s.pdf
```

End Problem 3-1, 75 min.

Problem 3-2: Estimating a Tridiagonal Matrix

In [Lecture \rightarrow Section 3.1] we learned that to determine the least squares solution of an overdetermined linear system of equations $\mathbf{A}\mathbf{x} = \mathbf{b}$ we minimize the residual norm $\|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2$ w.r.t. \mathbf{x} . However, we also face a linear least squares problem when minimizing the residual norm w.r.t. the entries of \mathbf{A} . This is the unusual situation considered in this problem.

This is an exercise in converting a problem statement into a proper linear least squares problem in the form of an overdetermined linear system of equations as in [Lecture \rightarrow Section 3.0.1] and then solving it through the normal equations, see [Lecture \rightarrow Section 3.2].

The following least-squares problem arises in the estimation of material parameters in one-dimensional diffusion problems: Let two vectors $\mathbf{z}, \mathbf{c} \in \mathbb{R}^n$, $n > 2 \in \mathbb{N}$ of measured values be given. The two real numbers α^* and β^* are defined as:

$$(\alpha^*, \beta^*) = \underset{\alpha, \beta \in \mathbb{R}}{\operatorname{argmin}} \|\mathbf{T}_{\alpha, \beta} \mathbf{z} - \mathbf{c}\|_{2}, \qquad (3.2.1)$$

where $\mathbf{T}_{\alpha,\beta} \in \mathbb{R}^{n \times n}$ is the following tridiagonal matrix

$$\mathbf{T}_{\alpha,\beta} = \begin{bmatrix} \alpha & \beta & 0 & \dots & 0 \\ \beta & \alpha & \beta & \ddots & \vdots \\ 0 & \beta & \ddots & \ddots & 0 \\ \vdots & \ddots & \ddots & \alpha & \beta \\ 0 & \dots & 0 & \beta & \alpha \end{bmatrix} . \tag{3.2.2}$$

(3-2.a)
☐ (20 min.) Reformulate (3.2.1) as a linear least squares problem in standard form according to [Lecture → Def. 3.1.1.1]

$$\mathbf{x}^* = \underset{\mathbf{x} \in \mathbb{R}^k}{\operatorname{argmin}} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2$$
 (3.2.3)

For appropriate choice of $m, k \in \mathbb{N}$ define suitable $\mathbf{A} \in \mathbb{R}^{m,k}$, $\mathbf{x} \in \mathbb{R}^k$ and $\mathbf{b} \in \mathbb{R}^m$, with $m, k \in \mathbb{N}$.

HIDDEN HINT 1 for (3-2.a) \rightarrow 3-2-1-0:TridiagLeastSquares1h.pdf

SOLUTION for (3-2.a)
$$\rightarrow$$
 3-2-1-1:TridiagLeastSquares1s.pdf

Explicitly state the normal equations [Lecture \rightarrow Thm. 3.1.2.1] associated with the linear least-squares problem (3.2.3).

Solution for (3-2.b)
$$\rightarrow$$
 3-2-2-0:s1b.pdf

Write a C++ function

that computes the optimal (in least-squares sense) parameters α^* and β^* according to (3.2.1) from data vectors z and c (i.e. z and c from (3.2.1)). Use the *method of normal equations*.

HIDDEN HINT 1 for (3-2.c) \rightarrow 3-2-3-0:TridiagLeastSquares2h.pdf

After the matrix A has been initialized, the solution just copies [Lecture \rightarrow Code 3.2.0.1].

SOLUTION for (3-2.c) \rightarrow 3-2-3-1:TridiagLeastSquares2s.pdf

End Problem 3-2, 65 min.

Problem 3-3: Linear Data Fitting

Abstract linear data fitting and its connection to the linear least squares problem [Lecture \rightarrow Eq. (3.1.3.7)] is discussed in [Lecture \rightarrow Ex. 3.0.1.1]. In this problem we practise these techniques for a concrete example.

This problem relies on the methods and EIGEN facilities introduced in [Lecture \rightarrow Section 3.2] and [Lecture \rightarrow Section 3.3.4].

For certain points in time t_i we have measured data f_i for a physical quantity f whose time-dependence $t \mapsto f(t)$ is only known qualitatively, that is, up to a few unknown parameters, see [Lecture \to Ex. 3.0.1.1]:

f(t) is assumed to be a linear combination

$$f(t) = \sum_{j=1}^{4} \gamma_j \phi_j(t)$$
 (3.3.1)

of the functions

$$\phi_1(t) = \frac{1}{t}, \qquad \phi_2(t) = \frac{1}{t^2}, \qquad \phi_3(t) = e^{-(t-1)}, \qquad \phi_4(t) = e^{-2(t-1)}.$$
 (3.3.2)

We aim for calculating the coefficients γ_i , j = 1, 2, 3, 4, such that

$$\sum_{i=1}^{10} (f(t_i) - f_i)^2 \longrightarrow \min . \tag{3.3.3}$$

(3-3.a) (30 min.) In the file linfit.hpp write a C++ function

that computes a least squares estimate for the coefficients γ_i by related to the data (3.3.2) and (3.3.3) by means of solving the *normal equations*.

SOLUTION for (3-3.a)
$$\rightarrow$$
 3-3-1-0:ldf1.pdf

```
Eigen::VectorXd data_fit_qr(const Eigen::VectorXd &t, const
    Eigen::VectorXd &f);
```

that uses *orthogonal transformation techniques* as introduced in [Lecture \rightarrow Section 3.3.4] for the same task as in the previous sub-problem.

SOLUTION for (3-3.b)
$$\rightarrow$$
 3-3-2-0:ldf2.pdf

that takes the estimated parameters γ_i , i=1,2,3,4, as argument vector gamma, a vector t_vec of times τ_i , and returns the vector of values $f(\tau_i)$.

Solution for (3-3.c)
$$\rightarrow$$
 3-3-3-0:ldf3.pdf

(3-3.d) (15 min.) Using EIGEN write a C++ function

Eigen::VectorXd fitting_error(const Eigen::VectorXd &gamma);

that returns the vector

$$\left[(f(t_i) - f_i)^2 \right]_{i=1}^{10} \in \mathbb{R}^{10}$$

where f is defined according to (3.3.1) using the values for the γ_j -coefficients passed in the gamma argument.

Solution for (3-3.d) \rightarrow 3-3-4-0:s5.pdf

•

End Problem 3-3, 90 min.

Problem 3-4: QR-Factorization of Tridiagonal Matrices

In [Lecture \rightarrow Rem. 3.3.3.26] we saw that a QR-factorization of a tri-diagonal matrix features an upper triangular factor that is tridiagonal again. In this exercise we exploit this fact for the design of efficient algorithms, in particular for the stable solution of tridiagonal linear systems of equations.

This exercise assume familiarity with the concept of QR-decomposition/factorization [Lecture \rightarrow Section 3.3.3] and, in particular, Givens rotations [Lecture \rightarrow § 3.3.3.15] and their compressed storage [Lecture \rightarrow Rem. 3.3.3.21].

Throughout this problem, in order to be compatible with some software, for generic real-valued square tridiagonal matrices we are supposed to use the data structure

```
#include <Eigen/Dense>
struct TriDiagonalMatrix {
    Eigen::Index n; // Matrix size: n × n
    Eigen::VectorXd d; // n-vector of diagonal entries
    Eigen::VectorXd 1; // n-1-vector, entries of first lower diagonal
    Eigen::VectorXd u; // n-1-vector, entries of first upper diagonal
    // Simple constructor
    TriDiagonalMatrix(const Eigen::VectorXd &d,
    const Eigen::VectorXd &l,
    const Eigen::VectorXd &u)
    : n(d.size()), d(d), l(l), u(u) {
        assert((n-1) == l.size() && (n-1) == u.size());
    }
};
```

Note that in full-fledged C++ implementations, this would be a C++ class with four different constructors and a destructor explicitly defined, private data members, and access member functions. The rudimentary definition above is proposed solely for the sake of simplicity.

Let us write $\mathbf{d} \in \mathbb{R}^n$, $\mathbf{l} \in \mathbb{R}^{n-1}$, and $\mathbf{u} \in \mathbb{R}^{n-1}$ for the vectors stored in d, 1, u. Then the matrix $\mathbf{A} \in \mathbb{R}^{n,n}$ stored in an **TriDiagonalMatrix** object is defined by

$$(\mathbf{A})_{i,j} := \begin{cases} (\mathbf{d})_i & \text{for } i = j \text{,} \\ (\mathbf{u})_i & \text{for } j = i+1 \text{,} \\ (\mathbf{l})_j & \text{for } i = j+1 \text{,} \end{cases} \\ (0 & \text{else,} \end{cases}$$

$$\Leftrightarrow \mathbf{A} = \begin{bmatrix} d_1 & u_1 \\ l_1 & d_2 & \ddots \\ \vdots & \ddots & \ddots & u_{n-1} \\ l_{n-1} & d_n \end{bmatrix} \quad \mathbf{d} = \begin{bmatrix} d_1, \dots, d_n \end{bmatrix}^\top \in \mathbb{R}^n \text{,} \\ \mathbf{u} = \begin{bmatrix} u_1, \dots, u_{n-1} \end{bmatrix}^\top \in \mathbb{R}^{n-1} \text{,} \\ \mathbf{l} = \begin{bmatrix} l_1, \dots, l_{n-1} \end{bmatrix}^\top \in \mathbb{R}^{n-1} \text{.} \end{cases}$$

(3-4.a) \square (30 min.) Prove that for a *tridiagonal* invertible matrix $\mathbf{A} \in \mathbb{R}^{n,n}$ a QR-factorization $\mathbf{A} = \mathbf{QR}, \mathbf{Q} \in \mathbb{R}^{n,n}$ orthogonal, $\mathbf{R} \in \mathbb{R}^{n,n}$ upper triangular with positive diagonal entries, satisfies

$$(\mathbf{R})_{i,j} = 0$$
 for $j > i + 2$ and $(\mathbf{Q})_{i,j} = 0$ for $i > j + 1$. (3.4.1)

HIDDEN HINT 1 for (3-4.a) \rightarrow 3-4-1-0:s1h1.pdf

```
SOLUTION for (3-4.a) \rightarrow 3-4-1-1:s1.pdf
```

To store a QR-factorization A = QR of a real-valued tridiagonal matrix $A \in \mathbb{R}^{n,n}$ we use the (again rudimentary) data structure

```
class TriDiagonalQR {
   public:
    explicit TriDiagonalQR(const TriDiagonalMatrix &A);

   template <typename VecType>
    Eigen::VectorXd applyQT(const VecType &x) const;
   template <typename VecType>
   Eigen::VectorXd solve(const VecType &b) const;
   // For debugging purposes: extract factors as dense matrices
   std::pair<Eigen::MatrixXd,Eigen::MatrixXd> getQRFactors(void)
        const;

   private:
   Eigen::Index n; // size of the square matrix
   Eigen::Matrix B; // Three non-zero upper diagonals of R
   Eigen::VectorXd rho; // Encoded Givens rotations
};
```

Here the non-zero entries of **R** are stored in an $n \times 3$ -matrix **B** adopting the convention

$$(\mathbf{R})_{i,j} := \begin{cases} (\mathbf{B})_{i,1} & \text{, if } i = j \text{,} \\ (\mathbf{R})_{i,j} = (\mathbf{B})_{i,2} & \text{, if } j = i+1 \text{,} \\ (\mathbf{R})_{i,j} = (\mathbf{B})_{i,3} & \text{, if } j = i+2 \text{,} \\ 0 & \text{else,} \end{cases}$$
 $i,j \in \{1,\ldots,n\}$, (3.4.6)

that is,
$$\mathbf{R} = \begin{bmatrix} x_1 & y_1 & z_1 & & & \\ & x_2 & y_2 & z_2 & & \\ & & x_3 & y_3 & z_3 & & \\ & & \ddots & \ddots & \ddots & \\ & & & x_{n-2} & y_{n-2} & z_{n-2} \\ & & & & x_n & y_{n-1} \\ & & & & x_n & y_{n-1} \end{bmatrix} \Rightarrow \mathbf{B} = \begin{bmatrix} x_1 & y_1 & z_1 \\ \vdots & \vdots & \vdots \\ x_{n-2} & y_{n-2} & z_{n-2} \\ x_{n-1} & y_{n-1} & 0 \\ x_n & 0 & 0 \end{bmatrix}.$$

Note that three entries of **B** will never be used.

The Q-factor is stored in compressed format through the underlying n-1 Givens rotations whose target rows are clear a priori, see [Lecture \rightarrow Rem. 3.3.3.26]. Every Givens rotation is stored through a single number ρ following the convention described in [Lecture \rightarrow Rem. 3.3.3.21]. These n-1 numbers are collected in the data member rho. Its j-th entry corresponds to the Givens rotation acting on rows j & j+1.

```
(3-4.b) ☐ (20 min.) Implement a C++ function

std::tuple<double, double, double>

compGivensRotation (Eigen::Vector2d a);
```

that computes a Givens rotation mapping the vector $\mathbf{a} \in \mathbb{R}^2$ onto a multiple of the first Cartesian coordinate vector $\mathbf{e}_1 \in \mathbb{R}^2$. Use the algorithm implemented in [Lecture \to Code 3.3.3.17]. The

function is to return the triplet (ρ, γ, σ) , where $\mathbf{G} = \begin{bmatrix} \gamma & \sigma \\ -\sigma & \gamma \end{bmatrix}$ is the transformation matrix such that $\mathbf{G}^{\top}\mathbf{a} = (\pm \|\mathbf{a}\|, 0)^{\top}$, and $\rho \in \mathbb{R}$ encodes it according to the rule set

$$\rho := \begin{cases} 1 & \text{, if } \gamma = 0 \text{ ,} \\ \frac{1}{2}\operatorname{sign}(\gamma)\sigma & \text{, if } |\sigma| < |\gamma| \text{ ,} \\ 2\operatorname{sign}(\sigma)/\gamma & \text{, if } |\sigma| \geq |\gamma| \text{ ,} \end{cases} \qquad \text{[Lecture} \to \text{Eq. (3.3.3.22)]}$$

Solution for (3-4.b) \rightarrow 3-4-2-0:s2.pdf

```
(3-4.c) \square (45 min.) [depends on [Lecture \rightarrow Rem. 3.3.3.26]]
```

Write an *efficient* code for the constructor of **TriDiagonalQR**, which initializes the data members, that is, computes the QR-factorization of the tridiagonal matrix $A \in \mathbb{R}^{n,n}$ passed to it in A.

HIDDEN HINT 1 for (3-4.c) \rightarrow 3-4-3-0:h2a.pdf

HIDDEN HINT 2 for (3-4.c) \rightarrow 3-4-3-1:h2b.pdf

SOLUTION for (3-4.c) \rightarrow 3-4-3-2:s2a.pdf

The type **VecType** must describe a vector with real-valued entries supplying component access through **operator** [] **const** and a size () method.

```
HIDDEN HINT 1 for (3-4.d) \rightarrow 3-4-4-0:sp3h1.pdf
```

SOLUTION for (3-4.d) \rightarrow 3-4-4-1:s3.pdf

Devise an efficient implementation of the method solve () of **TriDiagonalQR** that accepts a vector $\mathbf{b} \in \mathbb{R}^n$ as arguments and solves the linear systems of equations $\mathbf{A}\mathbf{x} = \mathbf{b}$, $\mathbf{A} \in \mathbb{R}^{n,n}$ the tridiagonal matrix stored in the current **TriDiagonalQR** object, returning the solution $\mathbf{x} \in \mathbb{R}^n$. A **std::runtime_error** exception should be thrown in case the matrix \mathbf{A} is not invertible.

HIDDEN HINT 1 for (3-4.e) \rightarrow 3-4-5-0:h41.pdf

Solution for (3-4.e)
$$\rightarrow$$
 3-4-5-1:s4.pdf

(3-4.f) ☑ (10 min.) [depends on Sub-problem (3-4.b), Sub-problem (3-4.d), Sub-problem (3-4.e)]

What is the asymptotic computational effort for your implementations of the constructor of **TriDiago-nalQR** and of the methods applyQT() and solve() in terms of $n \to \infty$?

```
SOLUTION for (3-4.f) \rightarrow 3-4-6-0:s5.pdf
```

(3-4.g) \odot (30 min.) The following templated C++ function implements (in a non-optimal way) a so-called **inverse iteration** [Lecture \rightarrow Section 9.3.2].

```
template <typename MatrixType>
unsigned int invit(const MatrixType &A,Eigen::VectorXd &x,
double ToL = 1E-6, unsigned int maxit = 100) {
   x.normalize();
   Eigen::VectorXd x_old(x.size());
   int it = 0;
   do {
```

```
x_old = x;
    x = A.lu().solve(x_old);
    x.normalize();
}
while (((x-x_old).norm() > TOL) && (it++ < maxit));
return it;
}
```

Implement a specialization (C++ reference) of this function for MatrixType = TriDiagonalMatrix.

```
HIDDEN HINT 1 for (3-4.g) \rightarrow 3-4-7-0:s6h1.pdf Solution for (3-4.g) \rightarrow 3-4-7-1:.pdf
```

End Problem 3-4, 195 min.

Problem 3-5: Conditional minimum with SVD

A Singular Value Decomposition (SVD, see [Lecture \rightarrow Def. 3.4.1.3]) is the generalized version of the diagonalization process for any $m \times n$ matrix. We will see that the diagonal and orthogonal matrices arising from an SVD have certain properties such that it is immediate to identify the minimum of a least squares functional for that matrix. In other words, computing the SVD of a matrix is equivalent to solving its corresponding least square problem.

The solution of this problem is discussed in [Lecture \rightarrow Section 3.4.4.1]; the algorithm is implemented in [Lecture \rightarrow Code 3.4.4.2].

Let $\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^{\top}$ be the singular value decomposition of $\mathbf{A} \in \mathbb{C}^{m,n}$, $m \geq n$. Consider the following problem:

$$\mathbf{x}^* = \underset{\|\mathbf{x}\|_2 = 1}{\operatorname{argmin}} \{ \|\mathbf{A}\mathbf{x}\|_2 \} \tag{3.5.1}$$

(3-5.a) **③** (45 min.) Show that the solution of (3.5.1) is:

$$\mathbf{x}^* = \mathbf{v}_n \coloneqq \mathbf{V}_{:,n} \tag{3.5.2}$$

Show that the following holds:

$$\|\mathbf{A}\mathbf{v}_n\|_2 = \sigma_n \coloneqq \mathbf{\Sigma}_{n.n}. \tag{3.5.3}$$

HIDDEN HINT 1 for (3-5.a) \rightarrow 3-5-1-0:MinimumSVD1h.pdf

SOLUTION for (3-5.a) \rightarrow 3-5-1-1:MinimumSVD1s.pdf

End Problem 3-5, 45 min.

Problem 3-6: Sparse Approximate Inverse (SPAI [GRH97])

This problem studies the least squares aspects of the SPAI method, a technique used in the numerical solution of partial differential equations. We encounter an "exotic" sparse matrix technique and rather non-standard least squares problems. Please note that the matrices to which SPAI techniques are applied will usually be huge and extremely sparse, say, of dimension $10^7 \times 10^7$ with only 10^8 non-zero entries. Therefore sparse matrix techniques must be applied.

Requires only familiarity with [Lecture \rightarrow Section 3.1] and [Lecture \rightarrow Section 3.2].



Sub-problem (3-6.d) is connected with [Lecture \rightarrow Section 10.3]. In case this section had not (yet) been covered in the course that sub-problem should be skipped.

Let $A \in \mathbb{R}^{N,N}$, $N \in \mathbb{N}$, be a regular sparse matrix with at most $n \ll N$ non-zero entries per row and column. We define the space of matrices with the same pattern (on non-zero entries) as A:

$$\mathcal{P}(\mathbf{A}) := \{ \mathbf{X} \in \mathbb{R}^{N,N} : (\mathbf{A})_{ij} = 0 \Rightarrow (\mathbf{X})_{ij} = 0 \}.$$
 (3.6.1)

The "primitive" SPAI (sparse approximate inverse) B of A is defined as

$$\mathbf{B} := \underset{\mathbf{X} \in \mathcal{P}(\mathbf{A})}{\operatorname{argmin}} \|\mathbf{I} - \mathbf{A}\mathbf{X}\|_{F}, \qquad (3.6.2)$$

where $\|\cdot\|_F$ stands for the Frobenius norm [Lecture \to Def. 3.4.4.17]. The solution of (3.6.2) can be used as a so-called preconditioner for the acceleration of iterative methods for the solution of linear systems of equations, see [Lecture \to Chapter 10]. An extended "self-learning" variant of the SPAI method is presented in [GRH97].

(3-6.a) \odot (20 min.) Show that the columns of **B** can be computed independently of each other by solving linear least squares problems. In the statement of these linear least squares problems write \mathbf{b}_i for the columns of **B**.

```
HIDDEN HINT 1 for (3-6.a) \rightarrow 3-6-1-0:spaih1.pdf
```

SOLUTION for (3-6.a)
$$\rightarrow$$
 3-6-1-1:spail.pdf

(3-6.b) (3-6.b) (60 min.) [depends on Sub-problem (3-6.a)]

Implement an efficient C++ function

```
Eigen::SparseMatrix<double> spai(const
    Eigen::SparseMatrix<double>& A);
```

for the computation of $\bf B$ according to (3.6.2). You may rely on the normal equations associated with the linear least squares problems for the columns of $\bf B$ or you may simply invoke the least squares solver of EIGEN, see [Lecture \rightarrow Code 3.4.3.14].

HIDDEN HINT 1 for (3-6.b) \rightarrow 3-6-2-0:spai2h.pdf

SOLUTION for (3-6.b)
$$\rightarrow$$
 3-6-2-1:spai2s.pdf

(3-6.c) (15 min.) [depends on Sub-problem (3-6.b)]

What is the total asymptotic computational effort of your implementation of spai() in terms of the problem size parameters N and n.

Solution for (3-6.c) \rightarrow 3-6-3-0:SPAI3.pdf

For $n \in \mathbb{N}$, consider a sparse s.p.d. [Lecture \to Def. 1.1.2.6] matrix **A** given by the following C++-code

```
C++11-code 3.6.5: Initialization of matrix A
   Eigen::SparseMatrix<double> init A (unsigned int n) {
2
     Eigen::SparseMatrix<double> L(n, n);
3
     Eigen::SparseMatrix<double> R(n, n);
4
     L.reserve(n);
     R. reserve (n + 2 * (n - 1));
     L.insert(0, 0) = std :: exp(1. / n);
     R.insert(0, 0) = 2.;
     for (unsigned int i = 1; i < n; ++i) {
       L.insert(i, i) = std::exp((i + 1.) / n);
10
       R.insert(i, i) = 2.;
11
       R.insert(i - 1, i) = -1.;
12
       R.insert(i, i - 1) = -1.;
13
14
     Eigen::SparseMatrix<double> A = kroneckerProduct(L, R);
15
     A.makeCompressed();
     return A;
17
  ( }
18
Get it on ₩ GitLab (spai.hpp).
```

Solve the linear system of equations $\mathbf{A}x = \mathbf{b}$ using EIGEN's **ConjugateGradient** (EIGEN documentation) iterative solver with initial guess $\mathbf{0}$ and default tolerance. Do this without a preconditioner and using the preconditioner $\frac{1}{2}(\mathbf{B} + \mathbf{B}^T)$, where \mathbf{B} solves (3.6.2). Use $\mathbf{b} = [1, 1, \dots, 1]^T$. Write a C++ function

```
std::vector<std:pair<unsigned int, unsigned int>>
testSPAIPrecCG(unsigned int L);
```

that returns the number of CG iterations (with and without the SPAI preconditioner) versus the system size $N = n^2$ for $n = 2^{1,...,L}$.

```
HIDDEN HINT 1 for (3-6.d) \rightarrow 3-6-4-0:spaih4.pdf
HIDDEN HINT 2 for (3-6.d) \rightarrow 3-6-4-1:spaih4.pdf
Solution for (3-6.d) \rightarrow 3-6-4-2:spai4.pdf
```

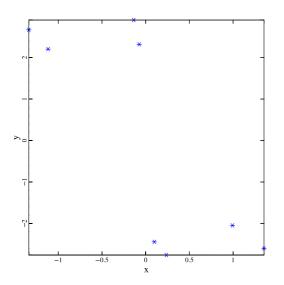
End Problem 3-6, 140 min.

Problem 3-7: Shape identification

This problem deals with pattern recognition, a central problem in image processing (e.g. in aerial photography, self-driving cars, and recognition of pictures of cats). We will deal with a very simplistic problem.

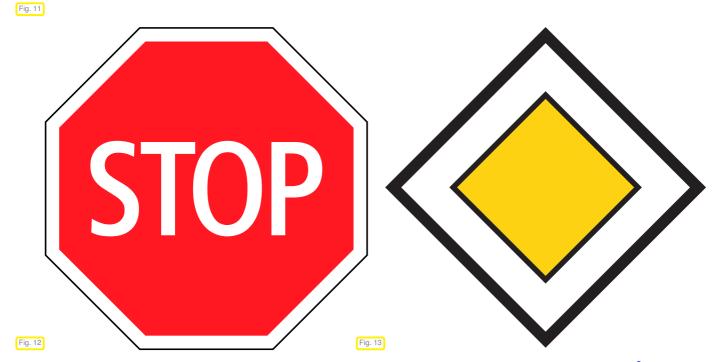
This problem practises the solution of a linear least squares problem based on the normal-equation method, see [Lecture \rightarrow Section 3.2].

A routine within the program of a self driving-car is tasked with the job of identifying road signs. The task is the following: given a collection of points $\mathbf{p}_i \in \mathbb{R}^2, i = 1, \dots, n$, we have to decide whether those point represent a stop sign, or a "priority road sign". In this exercise we consider n = 8.

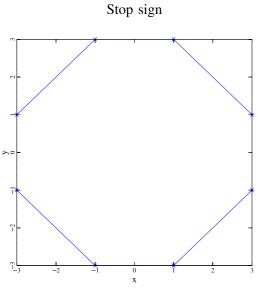


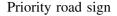
A set of possible input points \mathbf{p}^i : corners detected in the photo.

Which traffic sign had been in the focus of the camera?



The shape of the sign can be represented by a 2×8 matrix with the 8 coordinates in \mathbb{R}^2 defining the shape of the sign. We assume that the stop sign resp. the priority road sign are defined by the "model" points (known a priori) $\mathbf{x}_{stop}^i \in \mathbb{R}^2$ resp. $\mathbf{x}_{priority}^i \in \mathbb{R}^2$ for $i = 1, \dots, 8$.





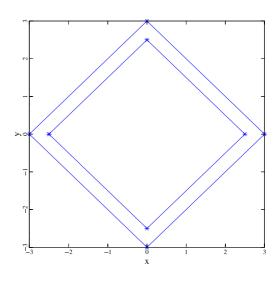


Fig. 15

Fig. 14

The 8-points \mathbf{x}_{stop}^i defining the model of a stop sign. The 8-points $\mathbf{x}_{priority}^i$ defining the model of a priority road sign.

However, in a real case scenario, one can imagine that the photographed shape is not exactly congruent to the one specified by the "model" points. Instead, one can assume that the points on the photo $(\mathbf{p}^i \in \mathbb{R}^2)$ are the result of a *linear transformation* of the original points $\mathbf{x}^i \in \mathbb{R}^2$ for $i = 1, \ldots, 8$; i.e. we can assume that there exists a matrix $\mathbf{A} \in \mathbb{R}^{2,2}$, s.t.

$$\mathbf{p}^i = \mathbf{A}\mathbf{x}^i, i = 1, \dots, n . \tag{3.7.1}$$

We do not know whether $\mathbf{x}^i = \mathbf{x}^i_{stop}$ or $\mathbf{x}^i = \mathbf{x}^i_{priority}$.

Moreover, we have some error in the data, i.e. our points do not represent exactly one of the linearly transformed shapes (it is plausible to imagine that there is some measurement error, and that the photographed shape is not exactly the same as our "model" shape), i.e. (3.7.1) is satisfied only "approximately".

With this problem, we will try to use the least square method to find the matrix **A** and to *classify* our points, i.e. to decide whether they represent a stop or a priority road sign.

$$\mathbf{B}\mathbf{v}=\mathbf{w}$$
,

whose least-squares solution \mathbf{v}^* will allow to determine the "best" linear transformation \mathbf{A} (in the least-squares sense). How is the vector of unknowns related to \mathbf{A} ?

SOLUTION for (3-7.a)
$$\rightarrow$$
 3-7-1-0: ShapeIdent1.pdf

(3-7.b) (10 min.) When does the matrix B have full rank? Give a geometric interpretation of this condition.

Solution for (3-7.b)
$$\rightarrow$$
 3-7-2-0:si2s.pdf

In the file shape_ident.hpp implement a C++ function

that returns the matrix **B**. Pass the vectors \mathbf{x}^i in the columns of a $2 \times n$ EIGEN matrix X.

SOLUTION for (3-7.c)
$$\rightarrow$$
 3-7-3-0:si3s.pdf

In the file shape_ident.hpp implement a C++ function

that computes the matrix A by solving the least squares problem based on the normal equations, see [Lecture \rightarrow Section 3.2]. It should return the Euclidean norm of the residual of the least square solution. Pass the vectors \mathbf{x}^i and the vectors \mathbf{p}^i as a $2 \times n$ EIGEN matrix X resp. P.

```
SOLUTION for (3-7.d) \rightarrow 3-7-4-0:si4s.pdf
```

Explain how the norm of the residual of the least square solution can be used to identify the shape defined by the points \mathbf{p}^{i} . Implement a function

that identifies the shape (either Stop or Priority sign) of the input points \mathbf{p}^i . The function returns an enum that classifies the shape of points specified by P. Return the linear transformation in the matrix A. The "model points" \mathbf{x}^i_{stop} resp. $\mathbf{x}^i_{priority}$ are passed through Xstop resp. Xpriority.

```
SOLUTION for (3-7.e) \rightarrow 3-7-5-0:si5s.pdf
```

End Problem 3-7, 105 min.

Problem 3-8: Low rank approximation of matrices

As explained in the course, large $m \times n$ matrices of low rank $k \ll \min\{m,n\}$ can be stored using only k(m+n) real numbers when using their SVD factorization/representation as sum of tensor products [Lecture \rightarrow Eq. (3.4.1.9)]. Thus, low rank matrices are of considerable interest for *matrix compression* (see [Lecture \rightarrow Ex. 3.4.4.24]). Unfortunately, adding two low rank matrices usually leads to an increase of the rank and entails "recompression" by computing a low-rank best approximation of the sum. This problems demonstrates an efficient approach to recompression.

Study [Lecture \rightarrow Section 3.4.4.2] to prepare for this exercise. The algorithms have to be implemented based on EIGEN and they rely on the SVD [Lecture \rightarrow Section 3.4.2] and QR factorization [Lecture \rightarrow Section 3.3.3.4].

(3-8.a) \Box (10 min.) Show that for a matrix $\mathbf{X} \in \mathbb{R}^{m,n}$ the following statements are equivalent:

```
(i) rank(\mathbf{X}) = k
```

```
(ii) \mathbf{X} = \mathbf{A}\mathbf{B}^{\top} for matrices \mathbf{A} \in \mathbb{R}^{m,k}, \mathbf{B} \in \mathbb{R}^{n,k}, k \leq \min\{m,n\}, both of full rank.
```

HIDDEN HINT 1 for (3-8.a) \rightarrow 3-8-1-0:LowRankRep1h.pdf

```
SOLUTION for (3-8.a) \rightarrow 3-8-1-1:LowRankRep1s.pdf
```

In the file lowrankrep.hpp write an EIGEN-based C++ function

```
std::pair<Eigen::MatrixXd, Eigen::MatrixXd>
factorize_X_AB(const Eigen::MatrixXd& X, unsigned int k);
```

that factorizes the matrix $\mathbf{X} \in \mathbb{R}^{m,n}$ with $\mathrm{rank}(\mathbf{X}) = k$ into $\mathbf{A}\mathbf{B}^{\top}$, as in Sub-problem (3-8.a), and returns the two matrices $\mathbf{A} \in \mathbb{R}^{m,k}$ and $\mathbf{B} \in \mathbb{R}^{n,k}$ in this order.

The function should issue a warning in case rank(X) = k is in doubt.

```
HIDDEN HINT 1 for (3-8.b) \rightarrow 3-8-2-0:LowRankRep2h.pdf
```

```
SOLUTION for (3-8.b) \rightarrow 3-8-2-1:LowRankRep2s.pdf
```

```
\mathbf{A} = \mathbf{Q}_A \mathbf{R}_A , \mathbf{Q}_A \in \mathbb{R}^{m,k} with \mathbf{Q}_A^{\top} \mathbf{Q}_A = \mathbf{I} , \mathbf{R}_A \in \mathbb{R}^{k,k} upper triangular , \mathbf{B} = \mathbf{Q}_B \mathbf{R}_B , \mathbf{Q}_B \in \mathbb{R}^{n,k} with \mathbf{Q}_B^{\top} \mathbf{Q}_B = \mathbf{I} , \mathbf{R}_B \in \mathbb{R}^{k,k} upper triangular ,
```

devise a way, how the economical singular value decomposition $\mathbf{A}\mathbf{B}^{\top} = \mathbf{U}\Sigma\mathbf{V}$ [Lecture \rightarrow Eq. (3.4.1.5)] of $\mathbf{A}\mathbf{B}^{\top}$ can be computed using only the singular value decomposition of a $k \times k$ -matrix.

```
SOLUTION for (3-8.c) \rightarrow 3-8-3-0:s4.pdf
```

(3-8.d) **(30 min.)** [depends on Sub-problem (3-8.c)]

```
Given \mathbf{A} \in \mathbb{R}^{m,k}, \mathbf{B} \in \mathbb{R}^{n,k}, with k \ll m, n, write an efficient C++ function
```

```
std::tuple<Eigen::MatrixXd, Eigen::MatrixXd, Eigen::MatrixXd>
svd_AB(const Eigen::MatrixXd& A, const Eigen::MatrixXd& B);
```

that calculates a singular value decomposition (SVD) [Lecture \to Thm. 3.4.1.1] of the product $\mathbf{A}\mathbf{B}^{\top} = \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^{\top}$, where $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{n,k}$ have orthogonal columns and $\boldsymbol{\Sigma} \in \mathbb{R}^{k,k}$ is a diagonal matrix. The SVD-factors are returned as a triple $(\mathbf{U}, \boldsymbol{\Sigma}, \mathbf{V})$.

HIDDEN HINT 1 for (3-8.d) \rightarrow 3-8-4-0:LowRankRep3h.pdf

SOLUTION for (3-8.d) \rightarrow 3-8-4-1:LowRankRep3s.pdf

(3-8.e) (10 min.) [depends on Sub-problem (3-8.d)]

What is the asymptotic computational cost of the function svd_AB () for (fixed) small k and $m = n \rightarrow \infty$? Discuss the effort required by the different steps of your algorithm.

SOLUTION for (3-8.e) \rightarrow 3-8-5-0:LowRankRep4s.pdf

•

(3-8.f) \square (10 min.) If $X, Y \in \mathbb{R}^{m,n}$ satisfy $\operatorname{rank}(Y) = \operatorname{rank}(X) = k$, show that $\operatorname{rank}(Y + X) \leq 2k$.

SOLUTION for (3-8.f) \rightarrow 3-8-6-0:LowRankRep5s.pdf

A

(3-8.g) \square (10 min.) Consider \mathbf{A}_X , $\mathbf{A}_Y \in \mathbb{R}^{m,k}$, \mathbf{B}_X , $\mathbf{B}_Y \in \mathbb{R}^{n,k}$, $\mathbf{X} = \mathbf{A}_X \mathbf{B}_X^\top$ and $\mathbf{Y} = \mathbf{A}_Y \mathbf{B}_Y^\top$. Find a factorization of the sum $\mathbf{X} + \mathbf{Y}$ as $\mathbf{X} + \mathbf{Y} = \mathbf{A}\mathbf{B}^\top$, with $\mathbf{A} \in \mathbb{R}^{m,2k}$ and $\mathbf{B} \in \mathbb{R}^{n,2k}$.

SOLUTION for (3-8.g) \rightarrow 3-8-7-0:LowRankRep6s.pdf

•

Rely on the previous subproblems to write an efficient C++ function

that returns two matrix factors $\mathbf{A}_Z \in \mathbb{R}^{m,k}$ and $\mathbf{B}_Z \in \mathbb{R}^{n,k}$ whose product $\mathbf{Z} = \mathbf{A}_Z \mathbf{B}_Z^\top$ yields the rank-k best approximation of the matrix sum $\mathbf{X} + \mathbf{Y} = \mathbf{A}_X \mathbf{B}_X^\top + \mathbf{A}_Y \mathbf{B}_Y^\top$:

$$\mathbf{Z} := \underset{\substack{\mathbf{M} \in \mathbb{R}^{m,n} \\ \operatorname{rank}(\mathbf{M}) \leq k}}{\operatorname{argmin}} \left\| \mathbf{A}_X \mathbf{B}_X^\top + \mathbf{A}_Y \mathbf{B}_Y^\top - \mathbf{M} \right\|_F$$

Here the matrices \mathbf{A}_X , $\mathbf{A}_Y \in \mathbb{R}^{m,k}$ and \mathbf{B}_X , $\mathbf{B}_Y \in \mathbb{R}^{n,k}$ are given and passed to the function through the arguments Ax, Ay, Bx, and By. You may only consider the case where 2k < m, n.

HIDDEN HINT 1 for (3-8.h) \rightarrow 3-8-8-0:LowRankRep7h.pdf

SOLUTION for (3-8.h) \rightarrow 3-8-8-1:LowRankRep7s.pdf

•

(**3-8.i**) **②** (10 min.) [depends on Sub-problem (3-8.h)]

What is the asymptotic computational cost of the function rank_k_approx for a small k and $m = n \to \infty$?

SOLUTION for (3-8.i) \rightarrow 3-8-9-0:LowRankRep8s.pdf

End Problem 3-8, 170 min.

Problem 3-9: Matrix Fitting by Constrained Least Squares

In this problem we look at a particular *constrained* linear least squares problem, as they have been discussed in [Lecture \rightarrow Section 3.6]. In particular, we will study the augmented normal equation approach to solve such type of least squares problem, see [Lecture \rightarrow Section 3.6.1]. In this context we refresh our understanding of the Lagrange multiplier technique.

Related to [Lecture \rightarrow Section 3.6]; involves implementation with EIGEN.

Consider the following problem of finding the minimal-norm matrix fitting a linear system of equations with given solution and right-hand side vector:

Given
$$n \in \mathbb{N}$$
, $\mathbf{z} \in \mathbb{R}^n$, $\mathbf{g} \in \mathbb{R}^n$, find $\mathbf{M}^* = \underset{\mathbf{M} \in \mathbb{R}^{n,n}, \ \mathbf{Mz} = \mathbf{g}}{\operatorname{argmin}} \|\mathbf{M}\|_F$, (3.9.1)

where $\|\cdot\|_F$ denotes the Frobenius norm of a matrix as introduced in [Lecture \rightarrow Def. 3.4.4.17].

Definition [Lecture \rightarrow Def. 3.4.4.17]. Frobenius norm

The Frobenius norm of $A \in \mathbb{K}^{m,n}$ is defined as

$$\|\mathbf{A}\|_F^2 := \sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2$$
.

$$\mathbf{x}^* = \underset{\mathbf{x} \in \mathbb{R}^N, \ \mathbf{C}\mathbf{x} = \mathbf{d}}{\operatorname{argmin}} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2 , \qquad (3.9.2)$$

for suitable matrices A, C and vectors b and d, see [Lecture \rightarrow Eq. (3.6.0.1)]. These matrices and vectors have to be specified based on z and g.

SOLUTION for (3-9.a)
$$\rightarrow$$
 3-9-1-0:lsqfrobe.pdf

State a necessary and sufficient condition for the matrix C found in Sub-problem (3-9.a) to possess full rank.

Solution for (3-9.b)
$$\rightarrow$$
 3-9-2-0:lsqf1.pdf

State the **augmented normal equations** [Lecture \rightarrow Eq. (3.6.1.6)] corresponding to the constrained linear least squares problem found in Sub-problem (3-9.a) and give necessary and sufficient conditions on **g** and **z** that ensure existence and uniqueness of solutions.

HIDDEN HINT 1 for (3-9.c)
$$\rightarrow$$
 3-9-3-0:lsqfh1.pdf

SOLUTION for (3-9.c)
$$\rightarrow$$
 3-9-3-1:.pdf

Write a C++ function

```
Eigen::MatrixXd min_frob(const Eigen::VectorXd &z, const
    Eigen::VectorXd &g);
```

that computes the solution \mathbf{M}^* of the minimization problem (3.9.1) given the vectors $\mathbf{z}, \mathbf{g} \in \mathbb{R}^n$. Use the augmented normal equations [Lecture \to Eq. (3.6.1.6)] derived in ((3-9.c)).

It is convenient to use EIGEN's built-in Kronecker product class **Eigen::KroneckerProduct** (#include <unsupported/Eigen/KroneckerProduct> required).

SOLUTION for (3-9.d)
$$\rightarrow$$
 3-9-4-0:lsqfrobi.pdf

Using the test vectors $\mathbf{z} = [1, 2, ..., n]^{\top} \in \mathbb{R}^n$ and $\mathbf{g} = [1, -1, 1, ..., 1, -1]^{\top} \in \mathbb{R}^n$, implement a C++ function

bool testMformula(unsigned int n);

that checks in a numerical experiment whether

$$\mathbf{M} = \frac{\mathbf{g}\mathbf{z}^{\top}}{\|\mathbf{z}\|_{2}^{2}} \tag{3.9.7}$$

is a solution to (3.9.1). To that end call the function min_frob implemented in Sub-problem (3-9.d), conduct the test and return **true** if it succeeds.

HIDDEN HINT 1 for (3-9.e) \rightarrow 3-9-5-0:lsqfrobh31.pdf

SOLUTION for (3-9.e) \rightarrow 3-9-5-1:lsqfrobc.pdf

Now, give a rigorous proof that (3.9.7) gives a solution of (3.9.1), provided that $z \neq 0$.

Solution for (3-9.f)
$$\rightarrow$$
 3-9-6-0:lsqfrx.pdf

So far, we have simply applied the formulas from [Lecture \rightarrow Section 3.6.1] to (3.9.1) and its reformulation as a constrained linear least squares problem. Now we take a closer look at the method of Lagrangian multipliers, which was used to derive these equations in class.

$$\mathbf{M}^* = \underset{\mathbf{M} \in \mathbb{R}^{n,n}}{\operatorname{argmin}} \left\{ \sup_{\mathbf{m} \in \mathbb{R}^n} [L(\mathbf{M}, \mathbf{m})] \right\}. \tag{3.9.10}$$

The matrix \mathbf{M}^* should provide the solution of (3.9.1).

HIDDEN HINT 1 for (3-9.g) \rightarrow 3-9-7-0:LSQFrobLagr1h.pdf

SOLUTION for (3-9.g)
$$\rightarrow$$
 3-9-7-1:LSQFrobLagr1s.pdf

(3-9.h) \odot (15 min.) Find the derivative $\operatorname{grad} \Phi \in \mathbb{R}^{n,n}$ of the function Φ defined as:

$$\Phi: \mathbb{R}^{n,n} \to \mathbb{R}, \quad \Phi(\mathbf{X}) = \|\mathbf{X}\|_F^2$$
 (3.9.12)

HIDDEN HINT 1 for (3-9.h) \rightarrow 3-9-8-0:LSQFrobLagr2h.pdf

SOLUTION for (3-9.h)
$$\rightarrow$$
 3-9-8-1:LSQFrobLagr2s.pdf

3. Direct Methods for Linear Least Squares Problems, 3.9. Matrix Fitting by Constrained Least Squares

Use the result of Sub-problem (3-9.h) to derive saddle point equations for $\operatorname{grad} L(\mathbf{M}, \mathbf{m}) = 0$, for the L obtained in Sub-problem (3-9.g).

HIDDEN HINT 1 for (3-9.i) \rightarrow 3-9-9-0:LSQFrobLagr3h.pdf

SOLUTION for (3-9.i)
$$\rightarrow$$
 3-9-9-1:LSQFrobLagr3s.pdf

Solve the saddle point equations obtained in Sub-problem (3-9.i) in symbolic form.

HIDDEN HINT 1 for (3-9.j) \rightarrow 3-9-10-0:LSQFrobLagr4h.pdf

SOLUTION for (3-9.j) \rightarrow 3-9-10-1:LSQFrobLagr4s.pdf

End Problem 3-9, 205 min.

Problem 3-10: Fitting of (relative) Point Locations from Distances

Given all pairwise distance of points on a line, this redundant information can be used to determine their positions through least-squares fitting. Pertinent algorithms are discussed in this problem.

This problem addresses the algorithmic details of the least-squares solution of the sparse overdetermined linear system of equations described in [Lecture \rightarrow Ex. 3.0.1.9]. The focus is on normal equation methods as introduced in [Lecture \rightarrow Section 3.2].

The numbers $x_1, \ldots, x_n \in \mathbb{R}$ describe the location of *n* points lying on a line that is identified with the real axis \mathbb{R} . We know measured values for the $m:=\binom{n}{2}=\frac{1}{2}n(n-1)$ signed distances $d_{i,j}:=x_j-x_i\in\mathbb{R}$, $1 \le i < j \le n$, from which we have to determine the x_i up to a shift. To fix the shift, we set $x_n := 0$.

This amounts to solving the overdetermined linear system of equations

$$\begin{cases} x_j - x_i = d_{i,j}, & 1 \le i < j \le n - 1, \\ -x_i = d_{i,n}, & i \in \{1, \dots, n - 1\}. \end{cases}$$
 (3.10.1)

Collecting the unknown positions in the vector $\mathbf{x} := [x_1, \dots, x_{n-1}]^{\top} \in \mathbb{R}^{n-1}$, the equations (3.10.1) can be recast as an $m \times (n-1)$ linear system of equations $\mathbf{A}\mathbf{x} = \mathbf{b}$, where $\mathbf{b} \in \mathbb{R}^m$ contains each of the distances $d_{i,j}$, $1 \le i < j \le n$, exactly once.

(3-10.a) (30 min.) How many real numbers and integers have to be stored at least in order to represent the matrix A, regarded as a matrix with real entries, in

- COO/triplet format and
- CRS format?

```
SOLUTION for (3-10.a) \rightarrow 3-10-1-0:s1.pdf
```

(3-10.b) :: (25 min.) Based on EIGEN implement (in the file disfitting.hpp) an efficient C++ function for the initialization of a sparse-matrix data structure for the matrix A:

```
Eigen::SparseMatrix<double> initA(unsigned int n);
```

It goes without saying that the parameter n passes the number n of points.

Solution for (3-10.b)
$$\rightarrow$$
 3-10-2-0:s4.pdf

(**3-10.c**) (45 min.) depends on Sub-problem (3-10.b)

According to [Lecture \rightarrow § 3.2.0.7] the **extended normal equations** for the generic overdetermined linear system $\mathbf{M}\mathbf{x} = \mathbf{c}$, $\mathbf{M} \in \mathbb{R}^{m,\ell}$, $\mathbf{c} \in \mathbb{R}^m$, $m \ge \ell$, boil down to the linear system of equations

$$\begin{bmatrix} -\mathbf{I} & \mathbf{M} \\ \mathbf{M}^{\top} & \mathbf{O} \end{bmatrix} \begin{bmatrix} \mathbf{r} \\ \mathbf{x} \end{bmatrix} = \begin{bmatrix} \mathbf{c} \\ \mathbf{0} \end{bmatrix}. \tag{3.10.5}$$

Implement an efficient C++ function

```
Eigen::VectorXd solveExtendedNormalEquations(
  const Eigen::MatrixXd &D);
```

that uses a sparse direct solver provided by EIGEN to solve the extended normal equations for the position fitting problem (3.10.1) and returns its least-squares solution $\mathbf{x} \in \mathbb{R}^{n-1}$. The argument D passes an $n \times n$ -matrix **D** whose strict upper triangular part contains the distances

$$[\mathbf{D}]_{i,j} := d_{i,j} , \quad 1 \le i < j \le n .$$

HIDDEN HINT 1 for (3-10.c) \rightarrow 3-10-3-0:h3tt.pdf

Solution for (3-10.c)
$$\rightarrow$$
 3-10-3-1:s3.pdf

(**3-10.d**) **□** (20 min.) [depends on Sub-problem (3-10.b)]

The coefficient matrix $\mathbf{M} \in \mathbb{R}^{n-1,n-1}$ of the **normal equations** [Lecture \rightarrow Eq. (3.1.2.2)] associated with the overdetermined linear system of equations (3.10.1) is of the form

$$\mathbf{M} = \mathbf{D} - \mathbf{z}\mathbf{z}^{\top}$$
 with a diagonal matrix $\mathbf{D} \in \mathbb{R}^{n-1,n-1}$, and a column vector $\mathbf{z} \in \mathbb{R}^{n-1}$. (3.10.8)

Compute **D** and **z**, and find a simple expression for the inverse M^{-1} .

HIDDEN HINT 1 for (3-10.d) \rightarrow 3-10-4-0:smw.pdf

SOLUTION for (3-10.d)
$$\rightarrow$$
 3-10-4-1:s4.pdf

[depends on Sub-problem (3-10.b) and Sub-problem (3-10.d)]

Write an efficient C++ function

```
Eigen::VectorXd solveNormalEquations()
  const Eigen::MatrixXd &D);
```

that returns the least-squares solution of (3.10.1) by solving the corresponding normal equations. The signature is the same as for solveExtendedNormalEquations () from Sub-problem (3-10.c).

HIDDEN HINT 1 for (3-10.e) \rightarrow 3-10-5-0:h4iA.pdf

Solution for (3-10.e)
$$\rightarrow$$
 3-10-5-1:s5.pdf

(3-10.f) (15 min.) [depends on Sub-problem (3-10.c) and Sub-problem (3-10.e)]

Write a C++ function

```
bool testNormalEquations(const MatrixXd &D);
```

that calls both solveExtendedNormalEquations() and solveNormalEquations() and returns true, if the results agree.

HIDDEN HINT 1 for (3-10.f) \rightarrow 3-10-6-0:dfh5a1.pdf

Solution for (3-10.f)
$$\rightarrow$$
 3-10-6-1:df5as.pdf

What is the asymptotic complexity of your implementation of solveNormalEquations() in Subproblem (3-10.e) for $n \to \infty$?

Solution for (3-10.g)
$$\rightarrow$$
 3-10-7-0:s6.pdf

End Problem 3-10, 190 min.

Problem 3-11: A Class Representing Low-Rank Matrices

Low-rank matrices play a central role in modern algorithms for data compression and the efficient handling of non-local operations. In this problem we devise a dedicated data type for efficiently storing and manipulating low-rank matrices.

This problem is connected with [Lecture \rightarrow Section 1.4.3] and [Lecture \rightarrow Section 3.4.4]. It is closely related to Problem 3-8.

Recall the following result from elementary linear algebra:

Theorem 3.11.1. Representation of low-rank matrices

Every matrix $\mathbf{M} \in \mathbb{R}^{m,n}$, $m,n \in \mathbb{N}$, with rank $r := \operatorname{rank}(\mathbf{M}) > 0$ can be factored as

$$\mathbf{M} = \mathbf{A}\mathbf{B}^{ op} \quad \textit{with} \quad egin{matrix} \mathbf{A} \in \mathbb{R}^{m,r}, \ \mathbf{B} \in \mathbb{R}^{n,r}. \end{cases}$$

This factorization is used in the design of the following EIGEN-based C++ class (File matrixlowrank.hpp) meant to store real matrices $\mathbf{M} \in \mathbb{R}^{m,n}$ with maximal rank $r < \min\{m,n\}$ and to provide special operations on them:

```
class MatrixLowRank {
public:
 MatrixLowRank (unsigned int m, unsigned int n, unsigned int r);
 MatrixLowRank (const Eigen::MatrixXd &A, const Eigen::MatrixXd &B);
 Eigen::Index rows() const { return _m; };
 Eigen::Index cols() const { return _n; };
 Eigen::Index rank() const { return _r; };
 Eigen::MatrixXd operator*(const Eigen::MatrixXd &) const;
 MatrixLowRank & operator *= (const Eigen::MatrixXd &);
 MatrixLowRank &addTo(const MatrixLowRank &, double tol = 1E-6);
private:
 unsigned int _m; // num. of rows
 unsigned int _n;
                     // num. of columns
 unsigned int _r;  // maximal rank, =1 for zero matrix
 Eigen::MatrixXd _A; // factor matrix A
 Eigen::MatrixXd _B; // factor matrix B
};
```

The convention for the case M = 0 is that we set r = 1 and the low-rank factors to zero.

```
(3-11.a)  (15 min.)
```

```
Give a proof of Thm. 3.11.1.
```

```
HIDDEN HINT 1 for (3-11.a) \rightarrow 3-11-1-0:h1.pdf
```

```
SOLUTION for (3-11.a) \rightarrow 3-11-1-1:s1.pdf
```

(3-11.b) ☑ (15 min.) In the file matrixlowrank.hpp implement the method

```
Eigen::MatrixXd MatrixLowRank::Operator *
  (const Eigen::MatrixXd &X) const;
```

which is supposed to return the product of the $m \times n$ low-rank matrix stored in the MatrixLowRank object (left factor) with a generic dense matrix $\mathbf{X} \in \mathbb{R}^{n,k}$ (right factor). Take care, that your code is efficient.

SOLUTION for (3-11.b)
$$\rightarrow$$
 3-11-2-0:s3.pdf

(3-11.c) (5 min.) [depends on Sub-problem (3-11.b)]

What is the asymptotic complexity of your implementation of MatrixLowRank::operator * from Sub-problem (3-11.b) in terms of $m, n, k \to \infty$, $r := \operatorname{rank}(\mathbf{M})$ fixed. Here, the size of the argument matrix is $n \times k$, $k \in \mathbb{N}$, and the **MatrixLowRank** object stores $\mathbf{M} \in \mathbb{R}^{m,n}$.

Solution for (3-11.c)
$$\rightarrow$$
 3-11-3-0:s3.pdf

(3-11.d) (10 min.) Provide an efficient implementation of the method (File matrixlowrank.hpp)

```
MatrixLowRank & operator *= (const Eigen::MatrixXd &X);
```

for the *in-situ* (*right*) *multiplication* of a low-rank $m \times n$ matrix stored in a **MatrixLowRank** object with a generic matrix $X \in \mathbb{R}^{n,k}$. "In-situ" means that after the operation the MatrixLowRank object stores the result of the matrix multiplication. The "internal rank" _r of the MatrixLowRank object should not change.

SOLUTION for (3-11.d)
$$\rightarrow 3-11-4-0:s4.pdf$$

(3-11.e) (3-11.e) (40 min.) With an eye on efficiency, in the file matrixlowrank.hpp complete the implementation of the method

```
MatrixLowRank &MatrixLowRank::addTo(
  const MatrixLowRank &X, double atol, double rtol);
```

that replaces the matrix $\mathbf{M} \in \mathbb{R}^{m,n}$ stored in the MatrixLowRank object with $\operatorname{trunc_{tol}}(\mathbf{M} + \mathbf{X})$, where

- $X \in \mathbb{R}^{m,n}$ is the matrix passed through the argument X,
- and, for $\mathbf{Y} \in \mathbb{R}^{m,n}$,

$$\mathsf{trunc_{tol}}(\mathbf{Y}) := \mathsf{argmin} \left\{ egin{align*} & \|\mathbf{Y} - \mathbf{R}\|_2 \leq \mathsf{rtol} \|\mathbf{Y}\|_2 \ , \ & \mathsf{or} \ & \|\mathbf{Y} - \mathbf{R}\|_2 \leq \mathsf{atol} \ \end{array}
ight\},$$

where $\|\cdot\|_2$ designates the Euclidean matrix norm. The "internal rank" $_r$ must be updated.

Your implementation may assume that

$$rank(\mathbf{M}) + rank(\mathbf{X}) \le min\{m, n\}$$
.

```
HIDDEN HINT 1 for (3-11.e) \rightarrow 3-11-5-0:h5a.pdf
HIDDEN HINT 2 for (3-11.e) \rightarrow 3-11-5-1:h5t.pdf
SOLUTION for (3-11.e) \rightarrow 3-11-5-2:s4.pdf
```

End Problem 3-11, 85 min.

Problem 3-12: Polar Decomposition of Matrices

In class we have seen various product decompositions/factorizations of matrices like the LU-decomposition [Lecture \rightarrow Section 2.3.2], the QR-decomposition [Lecture \rightarrow Section 3.3.3], and the singular-value decomposition (SVD) [Lecture \rightarrow Section 3.4]. In this problems we study another factorization, which is sometimes used in numerical methods, though it is not as important as the decompositions listed before.

This problem is related to [Lecture \rightarrow Section 3.3.3.2] and [Lecture \rightarrow Section 3.4.2] and requires familiarity with the Eigen-based C++ implementation discussed in those sections.

Theorem 3.12.1. Polar decomposition

```
For every matrix \mathbf{X} \in \mathbb{R}^{m,n}, m \ge n, there is a matrix \mathbf{Q} \in \mathbb{R}^{m,n} with orthonormal columns, \mathbf{Q}^{\top}\mathbf{Q} = \mathbf{I}_n, and a symmetric positive semi-definite [Lecture \rightarrow Def. 1.1.2.6] matrix \mathbf{M} \in \mathbb{R}^{n,n} such that \mathbf{X} = \mathbf{Q}\mathbf{M}.
```

The matrix factorization postulated in Thm. 3.12.1 is called **polar decomposition** of **X**.

```
(3-12.a) \odot (30 min.) Give a proof of Thm. 3.12.1. 
HIDDEN HINT 1 for (3-12.a) \to 3-12-1-0:h1p.pdf 
SOLUTION for (3-12.a) \to 3-12-1-1:s1.pdf
```

Based on the data types of EIGEN, the polar decomposition of a "tall/slim" real matrix is to be implemented as the following C++ class.

```
class PolarDecomposition {
  public:
    explicit PolarDecomposition(const Eigen::MatrixXd &X) { initialize(X); }
  PolarDecomposition(const Eigen::MatrixXd &A, const Eigen::MatrixXd &B);
  PolarDecomposition(const PolarDecomposition &) = default;
    ~PolarDecomposition() = default;

    // Left multiplication of M with the Q-factor of the polar decomposition
    void applyQ(Eigen::MatrixXd &Y) const { Y.applyOnTheLeft(Q_); }
    // Left multiplication of M with the M-factor of the polar decomposition
    void applyM(Eigen::MatrixXd &Y) const { Y.applyOnTheLeft(M_); }

    private:
    void initialize(const Eigen::MatrixXd &X);
    Eigen::MatrixXd Q_; // factor Q
    Eigen::MatrixXd M_; // factor M
};
```

The following specification of the class if given:

· The constructor

```
PolarDecomposition(const Eigen::MatrixXd &X);
```

should compute the polar decomposition factors Q and M according to Thm. 3.12.1 of the matrix passed in X and store them in the data members Q and M.

· The constructor

```
PolarDecomposition(const Eigen::MatrixXd &A, const
    Eigen::Matrix &B);
```

is supposed the initialize the data members Q and M with the polar decomposition factors $Q \in \mathbb{R}^{m,n}$ and $M \in \mathbb{R}^{n,n}$ of the matrix $AB^{\top} \in \mathbb{R}^{m,n}$, where the matrices $A \in \mathbb{R}^{m,k}$ and $B \in \mathbb{R}^{n,k}$ are passed through the arguments A and B.

• The methods applyQ() and applyM() realize the operations

$$Y \leftarrow QY$$
 , $Y \leftarrow MY$,

where Q and M are the factors of the polar decomposition stored in the **PolarDecomposition** object.

(3-12.b) (15 min.) Regardless of the implementation, what is the *minimal* asymptotic computational cost, that is, a *sharp lower bound* of the asymptotic computational effort, for a call of the second constructor

```
PolarDecomposition(const Eigen::MatrixXd &A, const Eigen::Matrix
    &B);
```

of a **PolarDecomposition** object for a $m \times n$ -matrix, $m \ge n$, for $m, n \to \infty$, and small fixed k?

Minimal asymptotic cost = O(for $m,n o \infty$.

SOLUTION for (3-12.b) $\rightarrow 3-12-2-0$: .pdf

In the file polardecomposition.hpp implement the method

```
void PolarDecomposition::initialize(const Eigen::MatrixXd &X);
```

that sets the data members $_Q$ and $_M$ of the class **PolarDecomposition**. These data members store the factors Q and M of the polar decomposition of the argument matrix X according to Thm. 3.12.1.

```
HIDDEN HINT 1 for (3-12.c) \rightarrow 3-12-3-0:pds2h.pdf
```

SOLUTION for (3-12.c) $\rightarrow 3-12-3-1:s2.pdf$

Write a code for the constructor

which is supposed to initialize the data members Q and M with the polar decomposition factors $Q \in \mathbb{R}^{m,n}$ and $M \in \mathbb{R}^{n,n}$ of the matrix $X := AB^{\top} \in \mathbb{R}^{m,n}$, where the matrices $A \in \mathbb{R}^{m,k}$ and $B \in \mathbb{R}^{n,k}$ are passed through the arguments A and A. We assume that A is small and fixed and A is small and A is small

```
HIDDEN HINT 1 for (3-12.d) \rightarrow 3-12-4-0:pcs3h1.pdf
```

SOLUTION for (3-12.d) $\rightarrow 3-12-4-1:s2.pdf$

End Problem 3-12, 195 min.

Problem 3-13: QR-Iteration

We have learned about the QR-decomposition/factorization as a tool for solving linear systems of equations [Lecture \rightarrow Rem. 3.3.4.3] and linear least-squares problems [Lecture \rightarrow Section 3.3.4]. Yet, this matrix factorization has many more important applications. For instance, it is a crucial ingredient for the so-called QR-algorithm, an iteration for solving algebraic eigenvalue problems [GOL13]. This problem will examine some algorithmic aspects of that algorithm.

You are supposed to know the contents of [Lecture \rightarrow Section 3.3.3] and implementation in Eigen.

We recall the QR-decomposition of matrices:

Theorem [Lecture → Thm. 3.3.3.4]. QR-decomposition

For any matrix $\mathbf{A} \in \mathbb{K}^{n,k}$ with $\operatorname{rank}(\mathbf{A}) = k$ there exists

(i) a unique Matrix $\mathbf{Q}_0 \in \mathbb{R}^{n,k}$ that satisfies $\mathbf{Q}_0^H \mathbf{Q}_0 = \mathbf{I}_k$, and a unique upper triangular Matrix $\mathbf{R}_0 \in \mathbb{K}^{k,k}$ with $(\mathbf{R})_{i,i} > 0$, $i \in \{1, \dots, k\}$, such that

$$\mathbf{A} = \mathbf{Q}_0 \cdot \mathbf{R}_0$$
 ("economical" QR-decomposition),

(ii) a unitary Matrix $\mathbf{Q} \in \mathbb{K}^{n,n}$ and a unique upper triangular $\mathbf{R} \in \mathbb{K}^{n,k}$ with $(\mathbf{R})_{i,i} > 0$, $i \in \{1, \ldots, n\}$, such that

$$A = Q \cdot R$$
 (full QR-decomposition).

If $\mathbb{K} = \mathbb{R}$ all matrices will be real and \mathbb{Q} is then orthogonal.

(3-13.a) (30 min.) Prove the following fact:

Lemma 3.13.1. Lower bandwidth of Q-factor

Let T = QR be a QR-decomposition of a regular tridiagonal matrix $T \in \mathbb{R}^{n,n}$, $n \in \mathbb{N}$. Then Q has lower bandwidth 1, that is,

$$i,j \in \{1,\ldots,n\}$$
, $i > j+1 \Rightarrow (\mathbf{Q})_{i,j} = 0$.

HIDDEN HINT 1 for (3-13.a) $\rightarrow 3-13-1-0:s1h1.pdf$

SOLUTION for (3-13.a) \rightarrow 3-13-1-1:s1.pdf

Lemma 3.13.3. QR-based similarity transform

If T = QR is the QR-decomposition of a regular, symmetric, and tridiagonal matrix $T \in \mathbb{R}^{n,n}$, $n \in \mathbb{N}$, then T' := RQ is also symmetric and tridiagonal.

HIDDEN HINT 1 for (3-13.b) $\rightarrow 3-13-2-0:s2h1.pdf$

SOLUTION for (3-13.b) $\rightarrow 3-13-2-1:s2.pdf$

(3-13.c) \boxtimes (120 min.) A *symmetric* and *tridiagonal* matrix $\mathbf{T} \in \mathbb{R}^{n,n}$ can be encoded by two

vectors $\mathbf{d} = \mathbf{d}(\mathbf{T}) \in \mathbb{R}^n$ and $\mathbf{u} = \mathbf{u}(\mathbf{T}) \in \mathbb{R}^{n-1}$ when setting

In the file <code>qriteration.hpp</code> implement a C++ function

that takes the vectors $\mathbf{d}(\mathbf{T}) \in \mathbb{R}^n$ and $\mathbf{u}(\mathbf{T}) \in \mathbb{R}^{n-1}$ describing a symmetric, tridiagonal matrix $\mathbf{T} \in \mathbb{R}^{n,n}$ as arguments, and

- returns the tuple $(\mathbf{d}(\mathbf{T}'), \mathbf{u}(\mathbf{T}'))$ of the defining vectors $\mathbf{d}(\mathbf{T}')$ and $\mathbf{u}(\mathbf{T}')$ of the symmetric, tridiagonal matrix $\mathbf{T}' := \mathbf{Q}^{\top} \mathbf{T} \mathbf{Q}$, where \mathbf{Q} is the Q-factor of the QR-decomposition $\mathbf{T} = \mathbf{Q} \mathbf{R}$ of \mathbf{T} , cf. Lemma 3.13.1,
- and features an asymptotic complexity of O(n) for $n \to \infty$.

```
HIDDEN HINT 1 for (3-13.c) \rightarrow 3-13-3-0:s3hgiv.pdf
HIDDEN HINT 2 for (3-13.c) \rightarrow 3-13-3-1:s3hbd.pdf
Solution for (3-13.c) \rightarrow 3-13-3-2:.pdf
```

End Problem 3-13, 180 min.

Problem 3-14: Economical Singular-Value Decomposition

For most applications requiring the singular-value decomposition (SVD) of matrix, its economical (thin) variant is sufficient. This problems examines some of its features.

This problems is based on the contents of [Lecture \rightarrow Section 3.4.1].

Given a matrix $\mathbf{A} \in \mathbb{R}^{m,n}$, $m,n \in \mathbb{N}$, $\mathbf{A} \neq \mathbf{O}$, we write

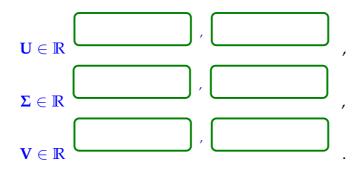
 $\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^{\top}$

for its **economical**

(thin) SVD.

(**3-14.a**) (12 min.)

Fill in the correct matrix sizes



HIDDEN HINT 1 for (3-14.a) \rightarrow 3-14-1-0:ecosvdh1.pdf

SOLUTION for (3-14.a) \rightarrow 3-14-1-1:ecogrs1.pdf

(3-14.b) \odot (24 min.) $m, n \in \mathbb{N}$.

Decide, whether the following statements are true for *every* $\mathbf{A} \in \mathbb{R}^{m,n} \setminus \{\mathbf{O}\}$,

1. The matrix **U** is an orthogonal matrix.





2. The set of all columns of **U** is an orthonormal basis of $\mathcal{R}(\mathbf{A})$.





3. V has orthogonal rows.





4. $\mathbf{V}^{\mathsf{T}}\mathbf{V} = \mathbf{I}$.



5. $nnz(\Sigma) := \sharp \{(i,j) : (\Sigma)_{i,j} \neq 0\} \leq n.$



O FALSE

6. $\mathbf{U}\mathbf{V}^{\top} = \mathbf{I}$.



O FALSE

SOLUTION for (3-14.b) \rightarrow 3-14-2-0:ecoqrs2.pdf

(3-14.c) \odot (12 min.) What is the asymptotic computational effort for computing the economical SVD of a dense matrix $\mathbf{A} \in \mathbb{R}^{m,n}, m, n \in \mathbb{N}$?

 $\operatorname{cost}(\operatorname{economical}\,\mathsf{SVD}\,\mathsf{of}\,\mathbf{A}\in\mathbb{R}^{m,n})=O($

for $m, n \to \infty$.

SOLUTION for (3-14.c) \rightarrow 3-14-3-0:ecoqrs3.pdf

End Problem 3-14, 48 min.

Problem 3-15: Low-rank Approximation of Concatenated Matrices

Low-rank approximation is a widely used classical technique for the compression of data sets. In this problem we study the recompression of two already compressed data sets after they have been merged. In linaer-algebra terms this boils down to computing the economical singular-value decomposition of a matrix arising from concatenating two low-rank factorized matrices.

This problem requires the concepts, techniques and results presented in [Lecture \rightarrow Section 3.4.2] and [Lecture \rightarrow Section 3.4.4.2].

Given are two low-rank matrices in factorized form

$$\mathbf{X}_1 := \mathbf{A}_1 \mathbf{B}_1^{\top}$$
 $\mathbf{A}_1, \mathbf{A}_2 \in \mathbb{R}^{m,k}$, $\mathbf{X}_2 := \mathbf{A}_2 \mathbf{B}_2^{\top}$, $\mathbf{B}_1, \mathbf{B}_2 \in \mathbb{R}^{n,k}$, (3.15.1)

for $m, n \in \mathbb{N}$, $k \leq \min\{m, n\}$. We "concatenate" the two matrices X_1 and X_2 into

$$\mathbf{X} := [\mathbf{X}_1 \ \mathbf{X}_2] \in \mathbb{R}^{m,2n} \ .$$
 (3.15.2)

Our ultimate goal is the efficient low-rank compression of X, given A_1 , A_2 , B_1 , B_2 .

Give a *sharp bound* for the maximal rank of the matrix **X** defined in (3.15.2)?

$$rank(\mathbf{X}) \leq$$

SOLUTION for (3-15.a) \rightarrow 3-15-1-0:cmlrcs1.pdf

(3-15.b) (90 min.) In the sequel we assume $2k \leq \min\{m, n\}$. In the file concatmatlrc.hpp implement a C++ function

```
std::tuple<Eigen::MatrixXd, Eigen::VectorXd, Eigen::MatrixXd>
 eco_svd_concat(
    const Eigen::MatrixXd &A1, const Eigen::MatrixXd &B1,
    const Eigen::MatrixXd &A2, const Eigen::MatrixXd &B2);
```

which accepts four matrices \mathbf{A}_1 , \mathbf{B}_1 , \mathbf{A}_2 , \mathbf{B}_2 with sizes \mathbf{A}_1 , $\mathbf{A}_2 \in \mathbb{R}^{m,k}$ and \mathbf{B}_1 , $\mathbf{B}_2 \in \mathbb{R}^{n,k}$, see (3.15.1), as arguments. The function is supposed to compute the factors

(i)
$$\mathbf{U} \in \mathbb{R}^{m,2k}$$
 with $\mathbf{U}^{\top}\mathbf{U} = \mathbf{I}_{2k}$, (3.15.3a)

(ii)
$$\Sigma = \operatorname{diag}(\sigma_1, \sigma_2, \dots, \sigma_{2k}) \in \mathbb{R}^{2k,2k}$$
 with $\sigma_1 \ge \sigma_2 \ge \dots \ge \sigma_{2k} \ge 0$, (3.15.3b)

(iii)
$$\mathbf{V} \in \mathbb{R}^{2n,2k}$$
 with $\mathbf{V}^{\top}\mathbf{V} = \mathbf{I}_{2k}$, (3.15.3c)

of the economical/thin SVD $X = U\Sigma V^{\top}$ of the concatenation $X, X := [A_1B_1^{\top}A_2B_2^{\top}]$ as defined in (3.15.2). The function should return the tuple $(\mathbf{U}, [\sigma_1, \dots, \sigma_{2k}]^{\top}, \mathbf{V})$.

The asymptotic complexity of your implementation for $m, n \to \infty$ and for small, fixed k, $2k \leq \min\{m, n\}$, must be *optimal*.

HIDDEN HINT 1 for (3-15.b) \rightarrow 3-15-2-0:clrm2h2.pdf

SOLUTION for (3-15.b) \rightarrow 3-15-2-1:cmlrcs2.pdf

```
(3-15.c) (30 min.) Complete the C++ function
```

```
bool test_eco_svd_concat(
    const Eigen::MatrixXd &A1, const Eigen::MatrixXd &B1,
    const Eigen::MatrixXd &A2, const Eigen::MatrixXd &B2);
```

that verifies that for the given argument matrices the implemented function $eco_svd_concat()$ complies with the specifications, that is, that its return values provide a factorization of X from (3.15.2) and meet the requirements (3.15.3). Return true, if this is the case.

```
HIDDEN HINT 1 for (3-15.c) \rightarrow 3-15-3-0:lrcmh3.pdf
```

SOLUTION for (3-15.c)
$$\rightarrow$$
 3-15-3-1:lrcms3.pdf

(**3-15.d**) **(**30 min.)

Based on $eco_svd_concat()$ as specified in Sub-problem (3-15.b) in the file concatmatlrc.hpp implement a C++ function

```
std::pair<Eigen::MatrixXd, Eigen::MatrixXd>
    concat_low_rank_best(
    const Eigen::MatrixXd &A1, const Eigen::MatrixXd &B1,
    const Eigen::MatrixXd &A2, const Eigen::MatrixXd &B2,
    double tol);
```

that computes a matrix $\mathbf{M} \in \mathbb{R}^{m,2n}$ of *minimal rank* $r \in \mathbb{N}$ such that

$$\|\mathbf{X} - \mathbf{M}\|_{2} \le \text{tol} \cdot \|\mathbf{X}\|_{2}$$
, (3.15.11)

where X is defined in (3.15.2), $\|\cdot\|_2$ stands for the Euclidean matrix norm, and the tolerance tol \in]0,1[is passed as an argument.

The matrix M should be returned in low-rank factorized form

$$\mathbf{M} = \mathbf{A}_M \mathbf{B}_M^{\top}$$
 , $\mathbf{A}_M \in \mathbb{R}^{m,r}$, $\mathbf{B}_M \in \mathbb{R}^{2n,r}$, (3.15.12)

as tuple $(\mathbf{A}_M, \mathbf{B}_M)$.

SOLUTION for (3-15.d)
$$\rightarrow$$
 3-15-4-0:crm4s.pdf

End Problem 3-15, 165 min.

Problem 3-16: Compact Storage Format for QR-Factorization

It is essential for the efficiency of orthogonalization techniques that orthogonal matrices, which commonly occur as factors in matrix products. are stored as products of elementary orthogonal transformations like Householder reflections or Givens rotations, see [Lecture \rightarrow Rem. 3.3.3.21]. This approach is adopted in all numerical libraries including EIGEN. This problem addresses some related algorithmic issues.

This problem assumes familiarity with orthogonal matrices [Lecture \rightarrow Section 3.3.2] and Householder reflections [Lecture \rightarrow § 3.3.3.11].

A legacy FORTRAN-77 numerical library routine returns a raw array (type **double** \star) of n^2 **double** floating-point numbers, $n \in \mathbb{N}$, which represents an $n \times n$ -matrix $\mathbf{A} \in \mathbb{R}^{n,n}$ according to the following specification:

- (I) The array encodes an $n \times n$ matrix $\mathbf{M} \in \mathbb{R}^{n,n}$ in column-major format [Lecture \rightarrow Section 1.2.3].
- (II) The column vectors $\mathbf{u}_i \in \mathbb{R}^n$, $j = 1, \dots, n-1$, defined as¹

are the building blocks for the product matrix

$$\mathbf{Q} = (\mathbf{I}_n - 2\mathbf{u}_1\mathbf{u}_1^{\top})(\mathbf{I}_n - 2\mathbf{u}_2\mathbf{u}_2^{\top}) \cdot \cdots \cdot (\mathbf{I}_n - 2\mathbf{u}_{n-1}\mathbf{u}_{n-1}^{\top}).$$
 (3.16.2)

(III) The matrix **A** is given as $\mathbf{A} = \mathbf{Q}\mathbf{R}$, where $\mathbf{R} \in \mathbb{R}^{n,n}$ is the *upper triangular* part of **M**:

Your task is to write an EIGEN-based wrapper class for the raw data that offers certain operations with the matrix **A**. The class definition is a follows:

C++ code 3.16.4: Class definition of CompactStorageQR

```
class CompactStorageQR {
    public:
3
      * \brief Construct a new CompactStorageQR object
5
6
      * \param data raw array which has to be persistent
8
      * \param n number of rows/columns of encoded matrix
     CompactStorageQR(const double *data, unsigned int n) : n_(n), M_(data, n, n) {
10
       for (unsigned int k = 0; k < n_{-} - 1; ++k) {
11
         double sn{0};
12
         for (unsigned int j = k + 1; j < n; ++j) {
13
           const double t\{M_{(j, k)}\};
14
           sn += t * t;
15
         }
16
```

¹A sum whose lower summation bound is larger than the upper summation bound is understood to evaluate to zero.

```
if (sn > 1.0) {
17
           throw std::runtime_error(
18
                "CompactStorageQR: Illegal subdiagonal column norm!");
19
20
       }
21
     }
22
23
24
      * \brief Determinant of the matrix
25
      * \return double the determinant
27
28
     double det() const;
29
30
31
      * \brief Right multiplication with another matrix
32
33
      * \param X matrix to multiply
      * \return Eigen::MatrixXd product
35
     Eigen::MatrixXd matmult(const Eigen::MatrixXd &X) const;
37
38
39
      * \brief Solution of linear systems of equations
40
41
      * \param B right hand side
      * \return Eigen::MatrixXd solution matrix
43
44
     Eigen::MatrixXd solve(const Eigen::MatrixXd &B) const;
45
    private:
47
     unsigned int n;
                                              // Matrix dimensions
48
     Eigen::Map<const Eigen::MatrixXd> M_; // Raw data wrapped into a matrix
49
   };
Get it on ₩ GitLab (compactstoragegr.hpp).
```

const Eigen::MatrixXd &X) const;

which returns the matrix product $\mathbf{A} \cdot \mathbf{X}$ of the matrix $\mathbf{A} \in \mathbb{R}^{n,n}$ stored in the **CompactStorageQR** object and of a matrix $\mathbf{X} \in \mathbb{R}^{n,k}$, $k \in \mathbb{N}$, passed in \mathbf{X} . Your implementation must not incur asymptotic computational cost worse than $O(n^2k)$ for $n, k \to \infty$.

```
SOLUTION for (3-16.a) \rightarrow 3-16-1-0:.pdf
```

(3-16.b) ☑ (30 min.) In the file compactstorageqr.hpp supply the missing parts of the implementation of

```
Eigen::MatrixXd CompactStorageQR::solve(const Eigen::MatrixXd &B)
    const;
```

which is to compute $\mathbf{A}^{-1}\mathbf{B}$ for $\mathbf{A} \in \mathbb{R}^{n,n}$ stored in the current **CompactStorageQR** object and $\mathbf{B} \in \mathbb{R}^{n,k}$ contained in the argument \mathbf{B} . Again, the asymptotic complexity of your implementation must be $O(n^2k)$ for $n,k \to \infty$.

```
HIDDEN HINT 1 for (3-16.b) \rightarrow 3-16-2-0:csqrs2h1.pdf
SOLUTION for (3-16.b) \rightarrow 3-16-2-1:.pdf
```

which gives the **determinant** det A of the matrix $A \in \mathbb{R}^{n,n}$ stored in the current **CompactStorageQR** object.

HIDDEN HINT 1 for (3-16.c) \rightarrow 3-16-3-0:cqrdeth1.pdf Solution for (3-16.c) \rightarrow 3-16-3-1:csqr3.pdf

End Problem 3-16, 105 min.

Problem 3-17: Extended Normal Equations

In [Lecture \rightarrow § 3.2.0.7] we learned that sparsity of the system matrix of a large overdetermined linear system of equations, which is lost when forming the regular normal equations, can be preserved by switching to the so-called extended normal equations. This problem reviews those and examines the initialization of the corresponding system matrix.

The problem is based on [Lecture \rightarrow Section 3.1.2], [Lecture \rightarrow § 3.2.0.7], and requires familiarity with [Lecture \rightarrow Section 2.7.2].

Throughout this problem we consider a large *sparse* overdetermined linear system of equations $\mathbf{A}\mathbf{x} = \mathbf{b}$, with $\mathbf{A} \in \mathbb{R}^{m,n}$, $\mathbf{b} \in \mathbb{R}^m$ given, $m \gg n \gg 1$. The matrix \mathbf{A} is assumed to have full rank.

(3-17.a) \odot (3 min.) By introducing the residual $\mathbf{r} := \mathbf{A}\mathbf{x} - \mathbf{b}$ as additional auxiliary variable the normal equations [Lecture \rightarrow Eq. (3.1.2.2)] for $\mathbf{A}\mathbf{x} = \mathbf{b}$ can be converted into the **extended normal equations**

$$\begin{bmatrix} \mathbf{r}_e \\ \mathbf{x}_e \end{bmatrix} = \begin{bmatrix} \mathbf{r}_e \\ \mathbf{0} \end{bmatrix}. \tag{3.17.1}$$

Complete the blocks of the system matrix and of the right-hand-side vector such that the solution of (3.17.1) provides the (unique) least-squares solution of $\mathbf{A}\mathbf{x} = \mathbf{b}$ in \mathbf{x}_e and the residual $\mathbf{r} := \mathbf{A}\mathbf{x} - \mathbf{b}$ in \mathbf{r}_e .

SOLUTION for (3-17.a) \rightarrow 3-17-1-0:extneqs1.pdf

What is the number of non-zero entries of the system matrix A_e of the extended normal equation (3.17.1) in terms of the number nnz(A) of non-zero entries of A?

$$nnz(\mathbf{A}_e) =$$

SOLUTION for (3-17.b) \rightarrow 3-17-2-0:extneqs2.pdf

The function

```
TripletMatrix
  extNeqSysMatCOO(
    unsigned int m, unsigned int n,
    const TripletMatrix &A_coo);
```

creates the matrix \mathbf{A}_{ℓ} of the extended normal equations (3.17.1) in triplet (COO) format when given the matrix $\mathbf{A} \in \mathbb{R}^{m,n}$ in triplet format through the argument \mathbb{A} , and its dimensions through \mathbb{R} and \mathbb{R} . The following typedefs have been used:

```
using Triplet = Eigen::Triplet<double>;
using TripletMatrix = std::vector<Triplet>;
```

Complete the implementation of extNeqSysMatCOO () given in Code 3.17.3.

C++ code 3.17.3: Initialization of A_{ℓ} in triplet format using Triplet = Eigen::Triplet <double>; using TripletMatrix = std::vector<Triplet>; TripletMatrix extNeqSysMatCOO(unsigned int m, unsigned int n, const TripletMatrix &A_coo) { TripletMatrix A_ext_coo{}; 5 for (int j = 0; $j < (int)m; ++j) {$ 6 A_ext_coo.push_back(Triplet { }); 8 9 for (const auto &t : A_coo) { 10 A_ext_coo.push_back(11 Triplet { }); 12 A_ext_coo.push_back(13 Triplet { }); 14 15 return A_ext_coo; 16 17 }

```
HIDDEN HINT 1 for (3-17.c) \rightarrow 3-17-3-0:extneqh1.pdf
SOLUTION for (3-17.c) \rightarrow 3-17-3-1:extneqs3.pdf
```

End Problem 3-17, 7 min.

.

Chapter 4

Filtering Algorithms

Problem 4-1: Autofocus with FFT

In this problem, we will use 2D frequency analysis to find the "best focused" image among a collection of out-of-focus photos. To that end, we will implement an algorithm based on 2D DFT as introduced in [Lecture \rightarrow Section 4.2.5].

This problem takes for granted that you know about the 2D discrete Fourier transform and its connection with discrete convolutions, see [Lecture \rightarrow Section 4.2.2] and [Lecture \rightarrow Section 4.2.5].

Let a grey-scale image consisting of $n \times m$ pixels be given as a matrix $P \in \mathbb{R}^{n,m}$ as in [Lecture \rightarrow Ex. 4.2.5.16]; each element of the matrix indicates the gray-value of the pixel as a number between 0 and v_{max} . This image is regarded as the "perfect image".

If a camera is poorly focused due to an inadequate arrangement of the lenses, it will record a blurred image $\mathbf{B} \in \mathbb{R}^{n,m}$. As before [Lecture \to Eq. (4.2.5.17)], the blurring operation can be modeled through the 2D (periodic) discrete convolution and can be undone provided that the point spread function (PSF) is known. However, in this problem we must not assume any knowledge of the PSF.

Assume that the only data available to us is the "black-box" C++ function (declared in autofocus.hpp):

```
Eigen::MatrixXd set_focus(double f);
```

which returns the potentially blurred image $\mathbf{B}(f)$, when the focus parameter f is set to a particular value. The actual operation of $\mathtt{set_focus}$ is utterly obscure. The focus parameter can be read as the distance of the lenses of a camera that can be changed through turning on and off a stepper motor. The following sequence of images illustrates the impact of setting different focus parameters:





Fig. 17





Fig. 18

Fig. 19

The problem of *autofocusing* is to determine the value of the focus parameter f, which yields an image $\mathbf{B}(f)$ with the least blur. The idea of the autofocusing algorithm is the following: *The less the image is marred by blur, the larger its "high frequency content"*. The latter can found out based on the discrete Fourier transform, more precisely, its 2D version.

To translate the autofocus idea into an algorithm we have to give a quantitative meaning to the notion of "high frequency content" of a (blurred) image \mathbb{C} , which is done by looking at the second moments of its 2D discrete Fourier transform:

$$V(\mathbf{C}) = \sum_{k_1=0}^{n-1} \sum_{k_2=0}^{m-1} \left(\left(\frac{n}{2} - \left| k_1 - \frac{n}{2} \right| \right)^2 + \left(\frac{m}{2} - \left| k_2 - \frac{m}{2} \right| \right)^2 \right) |\hat{\mathbf{C}}_{k_1, k_2}|^2.$$
 (4.1.1)

Here, $\hat{C} \in \mathbb{C}^{n,m}$ stand for the 2D DFT of the image C. Hence, the autofocusing policy can be rephrased as

find $f \in \mathbb{R}$ such that $V(\mathbf{B}(f))$ becomes maximal.

(4-1.a) (20 min.) This sub-problem supplies us with a tool to write an image file from a C++ code. Write a C++ function

```
void save_image(double focus);
```

that saves the image $\mathbf{B}(f)$ returned by $\mathtt{set_focus}()$ for f=0,1,2,3 in Portable Graymap (PGM) format to the file $./\mathtt{cx_out/image_focus} < f>.png$, where < f> is a string containing the argument focus.

For this task you can use objects of type PGMObject (found in "pgm.hpp"). You can find example of usages of this class in "examples/pgm_example.cpp".

```
Solution for (4-1.a) \rightarrow 4-1-1-0:render.pdf
```

The file fft.hpp supplies a few FFT-based convenience functions. Among them is fft2r, which performs the 2D discrete Fourier transform [Lecture \rightarrow Eq. (4.2.5.4)] of a real-values matrix, *cf.* [Lecture \rightarrow Code 4.2.5.6], which results in a complex matrix however.

```
C++ code 4.1.3: 2D DFT of a Eigen::MatrixXd

Eigen::MatrixXcd fftr(const Eigen::MatrixXd& X) {
    const long m = X.rows(), n = X.cols();

Eigen::MatrixXcd Y(m, n);

Eigen::FFT<double> fft;
```

```
for (long j = 0; j < n; ++j) {
    Eigen::VectorXd Xj = X.col(j);
    Y.col(j) = fft.fwd(Xj);
}

return Y;
}

Eigen::MatrixXcd fft2r(const Eigen::MatrixXd& X) {
    return fft(fftr(X).transpose()).transpose();
}</pre>
```

(4-1.b) (30 min.) The two-dimensional discrete Fourier transform can be performed using the C++ function fft2r(). Write a C++ function

```
void plot_freq(double focus);
```

that creates 2D color plots of the (modulus of the) 2D DFTs of the images obtained in sub-problem (4-1.a). Clamp the data between 0 and 8000, this means when a value exceeds 8000, then set it to 8000.

```
HIDDEN HINT 1 for (4-1.b) \rightarrow 4-1-2-0:fft2h1.pdf

SOLUTION for (4-1.b) \rightarrow 4-1-2-1:plotfft2.pdf

(4-1.c) \bigcirc (15 min.) [depends on Sub-problem (4-1.b)]
```

In the plots generated in Sub-problem (4-1.b), mark (or explain) the regions corresponding to the high and low frequencies.

```
HIDDEN HINT 1 for (4-1.c) \rightarrow 4-1-3-0:afh1.pdf
SOLUTION for (4-1.c) \rightarrow 4-1-3-1:explain.pdf
```

```
double high_frequency_content(const Eigen::MatrixXd & C);
```

that returns the value $V(\mathbf{C})$ for a given matrix \mathbf{C} . Write a C++ function

```
void plotV();
```

that uses MATPLOTLIBCPP (function plot ()) and plots the function $V(\mathbf{B}(f))$ in terms of the focus parameter f. Use 100 equidistantly spaced sampling points in the interval [0,5].

```
Solution for (4-1.d) \rightarrow 4-1-4-0:plotV.pdf
```

(4-1.e) (60 min.) Write an *efficient* (you want the auto-focus procedure to take not longer than a few milliseconds!) C++ function

```
double autofocus();
```

that numerically determines optimal focus parameter $f_0 \in [0,5]$ for which the 2nd moment $V(\mathbf{B}(f))$ is maximized. What is the resulting focus parameter f_0 ?

Two particular challenges have to be tackled in the implementation of autofocus ():

(I) To locate the change of slope of $V(\mathbf{B}(f))$ in [0,5] you should use the **bisection algorithm** [Lecture \to Section 8.4.1] for finding a zero of the derivative $f \mapsto \frac{dV(\mathbf{B}(f))}{df}$, which will, hopefully, mark the location of the maximum of V(f).

The bisection algorithm for finding a zero of a continuous function $\varphi:[a,b]\to\mathbb{R}$ with $\varphi(a)\cdot\varphi(b)<0$ is rather simple and lucidly demonstrated in [Lecture \to Code 8.4.1.2]. Study that sample code first.

(II) Of course, the derivative is not available, so that we have to approximate it crudely by means of a difference quotient

$$\frac{dV(f)}{df} \approx \frac{V(\mathbf{B}(f+\delta f)) - V(\mathbf{B}(f-\delta f))}{2\delta f}.$$

You can assume that the smallest step of the auto-focus mechanism is 0.05 and so the natural choice for δf is $\delta f = 0.05$. This maximal resolution also tells you a priori how many bisection steps should be performed.

Exploit the structure of the function $V(\mathbf{B}(f))$ that you observed in sub-problem (4-1.d). Use as few evaluations of $V(\mathbf{B}(f))$ as possible!

SOLUTION for (4-1.e) \rightarrow 4-1-5-0:bisect.pdf

End Problem 4-1, 155 min.

Problem 4-2: FFT and least squares

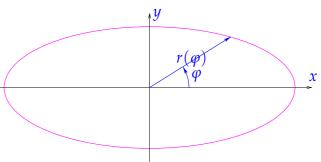
This problem deals with an application of fast Fourier transformation in the context of a least squares (data fitting) problems.

You should be familiar with both DFT [Lecture \rightarrow Section 4.2], the idea of the least-squares solution of fitting problems [Lecture \rightarrow Section 3.1.1], and the normal equation method [Lecture \rightarrow Section 3.2].

The orbit of a planet around the sun is a closed, planar curve $\mathcal{C} \subset \mathbb{R}^2$, which can be described by a so-called polar representation, closely related to polar coordinates:

$$\mathcal{C}:=\left\{egin{bmatrix} d(arphi)\cosarphi\ d(arphi)\sinarphi \end{bmatrix}:\ 0\leqarphi<2\pi
ight\}$$
 ,

where $d:[0,2\pi] \to \mathbb{R}^+$ is a given function of the polar angle φ providing the distance of a point on the curve from the origin in the direction of that angle.



We put the sun in the origin of the coordinate system. For equispaced polar angles (\rightarrow Fig. 25) $\varphi_j = 2\pi \frac{j}{n}$, $j = 0, \ldots, n-1$, $n \in \mathbb{N}$ the distance of the planet from the sun is measured and found to be $d_j \in \mathbb{R}$ $j = 0, \ldots, n-1$. We now discuss a numerical method providing the planet's path in polar representation.

Write $\mathcal{P}_m^{\mathbb{C}}$, $m \in \mathbb{N}$, for the space of real trigonometric polynomials of degree $\leq m$. These are 2π -periodic functions $\mathbb{R} \to \mathbb{R}$ of the form

$$p(t) = \sum_{k=0}^{m} c_k \cos(kt) , \quad c_k \in \mathbb{R} .$$
 (4.2.1)

The idea is to approximate the distance $d(\varphi)$ of the planet from the sun as a function of the angle by a trigonometric polynomial $p^* \in \mathcal{P}_m^{\mathbb{C}}$, which minimized the sum of squares of "distance mismatches":

$$p^* = \underset{p \in \mathcal{P}_{c}^{C}}{\operatorname{argmin}} \sum_{j=0}^{n-1} |p(\varphi_j) - d_j|^2.$$
 (4.2.2)

Throughout, we assume that the number of distance data is more than twice as big as the degree of p^* : 2m < n.

(4-2.a) (30 min.) Recast this task as a standard linear least squares problem

$$\mathbf{x}^* = \underset{\mathbf{x} \in \mathbb{R}^?}{\operatorname{argmin}} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2, \qquad (4.2.3)$$

with a suitable matrix $\mathbf{A} \in \mathbb{R}^{?,?}$ and vectors $\mathbf{b}, \mathbf{x} \in \mathbb{R}^?$.

SOLUTION for (4-2.a)
$$\rightarrow$$
 4-2-1-0:relsq.pdf

State the **normal equations** [Lecture \rightarrow Thm. 3.1.2.1] for the linear least-squares problem found in Sub-problem (4-2.a) in explicit form.

HIDDEN HINT 1 for (4-2.b) \rightarrow 4-2-2-0:spnh1.pdf

SOLUTION for (4-2.b)
$$\rightarrow$$
 4-2-2-1:normleg.pdf

(4-2.c) (20 min.) [depends on Sub-problem (4-2.a), Sub-problem (4-2.b)]

Write a short EIGEN-based C++ function

bool testNormEqMatrix(unsigned int n, unsigned int m);

that tests, whether the formula for the matrix product $\mathbf{A}^{\top}\mathbf{A}$ that you have found in Sub-problem (4-2.b) is correct.

To that end, form the dense matrix A as an object of type **Eigen::MatrixXd**, compute $A^{\top}A$ and compare with the matrix given by your formula. Of course, the comparison has to be done in a numerically sound way following the advice of [Lecture \rightarrow Rem. 1.5.3.15].

SOLUTION for (4-2.c)
$$\rightarrow$$
 4-2-3-0:ts.pdf

(4-2.d) □ (45 min.) [depends on Sub-problem (4-2.b)]

Write a C++ function

```
Eigen::VectorXd find_c(const Eigen::VectorXd &d, unsigned int m);
```

that computes the coefficients c_k , $k=0,\ldots,m$, of p^* in the form (4.2.1) for arbitrary inputs $d_j \in \mathbb{R}$, $j=0,\ldots,n-1$. These are passed as vector d.

We aim for an efficient implementation: the function must not have asymptotic complexity worse than $O(n \log n)$ in terms of the problem size parameters $m, n \to \infty$.

HIDDEN HINT 1 for (4-2.d) \rightarrow 4-2-4-0:imph1.pdf

SOLUTION for (4-2.d)
$$\rightarrow$$
 4-2-4-1:implem.pdf

(4-2.e) (60 min.) [depends on Sub-problem (4-2.d)]

Since Kepler we know that the orbits of planets around the sun are ellipses. Choose n = 100 and sample the distance data form the polar form of an ellipse:

$$d(\varphi) = \frac{1}{\sqrt{1 - c^2 \cos^2 \varphi}}, \quad \varphi \in [0, 2\pi], \quad c = 0.8.$$
 (4.2.11)

- 1. Using your implementation of find_c compute and tabulate the coefficients c_k of p^* for m=1,2,3.
- 2. Using MATPLOTLIBCPP's plot () function create a plot of the ellipse given by (4.2.11).
- 3. Into the same plot insert the curves described by the trigonometric polynomials p^* of degrees m=1,2,3.

This should be accomplished by a C++ function

HIDDEN HINT 1 for (4-2.e) \rightarrow 4-2-5-0:sxh1.pdf

SOLUTION for (4-2.e) $\rightarrow 4-2-5-1$:sx.pdf

End Problem 4-2, 185 min.

Problem 4-3: Multiplication and division of polynomials based on DFT

In [Lecture \rightarrow § 4.1.3.6] we saw that the formula of discrete convolution [Lecture \rightarrow Def. 4.1.3.3] also gives the coefficients of products of polynomials. Thus, in light of the realization of discrete convolutions by means of DFT [Lecture → Section 4.2.2], the Fast Fourier Transform (FFT) becomes relevant for efficiently multiplying polynomials of very high degree.

This problem requires knowledge about DFT and discrete convolutions, see [Lecture \rightarrow Section 4.1.3] and [Lecture → Section 4.2.2]. For the connection of DFT, FFT, and polynomial multiplication also see the YouTube Video.

Polynomials. We call a polynomial in x of degree $N \in \mathbb{N}$ a function $p : \mathbb{K} \to \mathbb{K}$ of the form

$$p(x) = \sum_{j=0}^{N} \gamma_j x^j$$
 with coefficients $\gamma_j \in \mathbb{K}$, (4.3.1)

also called monomial representation. Under pointwise addition and multiplication the polynomials of degree N form a vector space \mathcal{P}_N of dimension N+1. Associating the vector $\mathbf{c}:=[\gamma_0,\ldots,\gamma_N]\in\mathbb{K}^{N+1}$ to every $p\in\mathcal{P}_N$ as given in (4.3.1) induces an isomorphism of the vector spaces \mathcal{P}_N and \mathbb{K}^{N+1} . In other words, we can identify $p \in \mathcal{P}_N$ with its coefficient vector \mathbf{c} .

Thus, when we have to represent polynomials $\in \mathcal{P}_N$ in a C++ code we can do this via objects of type Eigen::VectorXd that contain their coefficient vectors. This convention is adopted throughout this problem.

Let two large numbers $m, n \gg 1$ and two polynomials of degrees m-1 and n-1, respectively, in monomial representation

$$u(x) = \sum_{j=0}^{m-1} \alpha_j x^j, \quad v(x) = \sum_{j=0}^{m-1} \beta_j x^j, \quad \alpha_j, \beta_j \in \mathbb{C} ,$$
 (4.3.2)

be given.

The tasks that we are going to tackle in this problem are the following.

- 1. Polynomial multiplication: Efficiently compute the coefficients of the polynomial p = uv (point wise product, a polynomial of degree $\leq m + n - 2$).
- 2. Polynomial division: Given the polynomials p and u in terms of their coefficients, the degree of p at least as bis as that of u, find a polynomial v, if it exists, such that p = uv.

(4-3.a) (25 min.) Write a C++ function

```
Eigen::VectorXd polyMult_naive(const Eigen::VectorXd & u, const
  Eigen::VectorXd & v);
```

that "naively", i.e. using a simple loop-based implementation, computes and returns the vector of coefficients of the polynomial p = uv. The coefficient vectors for the polynomials u and v are passed in uand v.

Also determine the asymptotic complexity of your algorithm for $m, n \to \infty$ (separately).

HIDDEN HINT 1 for (4-3.a) \rightarrow 4-3-1-0:PolyDiv1h.pdf

SOLUTION for (4-3.a) $\rightarrow 4-3-1-1$: PolyDiv1s.pdf

(4-3.b) (30 min.) Write a C++ function

that *efficiently* computes the vector of monomial coefficients of the product polynomial p = uv of the polynomials u and v, passed through their respective coefficient vectors. To that end use the DFT functionality of Eigen, see [Lecture \rightarrow Code 4.2.1.22].

```
HIDDEN HINT 1 for (4-3.b) \rightarrow 4-3-2-0:PolyDiv2h.pdf

SOLUTION for (4-3.b) \rightarrow 4-3-2-1:PolyDiv2s.pdf
```

(4-3.c) • (10 min.) Determine the complexity of your implementation of algorithm in (4-3.b) making use of the fact that EIGEN's DFT relies on the FFT algorithm [Lecture \rightarrow Section 4.3] and thus can be carried out with an asymptotic effort $O(n \log n)$ for vector length $n \rightarrow \infty$.

```
HIDDEN HINT 1 for (4-3.c) \rightarrow 4-3-3-0:PolyDiv3h.pdf
Solution for (4-3.c) \rightarrow 4-3-3-1:PolyDiv3s.pdf
```

(4-3.d) \odot (20 min.) Give an interpretation of the entries of the vector $\mathbf{DFT}_n\mathbf{u}$ obtained by applying a discrete Fourier transform to the coefficient vector \mathbf{u} of a polynomial u of degree n-1? What is an equivalent operation that can be performed for the polynomial u which leads to the same result?

```
HIDDEN HINT 1 for (4-3.d) \rightarrow 4-3-4-0:PolyDiv4h.pdf
Solution for (4-3.d) \rightarrow 4-3-4-1:PolyDiv4s.pdf
```

The following sub-problems deal with polynomial division. Polynomial division with remainder can be carried out using Euclid's algorithm and is captured by the following theorem:

Theorem 4.3.6. Polynomial devision with remainder

For every $p \in \mathcal{P}_N$ and $u \in \mathcal{P}_m$, $m \leq N$, there exist unique polynomials $q \in \mathcal{P}_{N-m}$ and $r \in \mathcal{P}_{m-1}$ such that

$$p = uq + r$$
.

Using the notations of the theorem, we are interested in determining, whether the polynomial p can be divided by u exactly, that is, whether $r \equiv 0$. In this case we also want to compute the quotient polynomial q.

Write a C++ function

```
Eigen::VectorXd polyDiv(const Eigen::VectorXd & p, const
    Eigen::VectorXd & u);
```

that takes the coefficient vectors p and u of two polynomials p and u, respectively, as arguments, where the degree of p must be at least as large as that of u. The function is to return the coefficient vector of a q, such that p = uv, if it exists. If not, a vector of length zero should be returned.

```
HIDDEN HINT 1 for (4-3.e) \rightarrow 4-3-5-0:PolyDiv6h.pdf
HIDDEN HINT 2 for (4-3.e) \rightarrow 4-3-5-1:pdxh2.pdf
SOLUTION for (4-3.e) \rightarrow 4-3-5-2:PolyDiv6s.pdf
```

End Problem 4-3, 130 min.

Problem 4-4: Solving triangular Toeplitz systems

In [Lecture \rightarrow Section 4.5] we learned about Toeplitz matrices, the class of matrices with constant diagonals [Lecture \rightarrow Def. 4.5.1.7]. Obviously, $m \times n$ Toeplitz matrices are *data-sparse* in the sense that it takes only m+n-1 steps to encode them completely. Therefore, operations with Toeplitz matrices can often be done with asymptotic computational cost significantly lower than that of the same operations for a generic matrix of the same size. In [Lecture \rightarrow Section 4.5.2] we have seen this for FFT-based matrix \times vector multiplication for Toeplitz matrices.

In this problem we study an efficient FFT-based algorithm for solving triangular linear systems of equations whose coefficient matrix is Toeplitz. Such linear systems are faced, for instance, when inverting a finite, linear, time-invariant, causal channel (LT-FIR) as introduced in [Lecture \rightarrow Section 4.1].

In [Lecture \to Section 4.1.2] we found that the output operator of a causal finite linear time-invariant filter with impulse response $\mathbf{h} = [h_0, \dots, h_{n-1}]^{\top} \in \mathbb{R}^n$ can be obtained by discrete convolution [Lecture \to Eq. (4.1.2.4)]; see also [Lecture \to Def. 4.1.3.3]. Given an input signal $\mathbf{x} := [x_0, \dots, x_{n-1}]^{\top} \in \mathbb{R}^n$, the first n components y_0, \dots, y_{n-1} of the output are obtained by:

$$\begin{bmatrix} y_0 \\ \vdots \\ \vdots \\ y_{n-1} \end{bmatrix} = \begin{bmatrix} h_0 & 0 & \cdots & \cdots & 0 \\ h_1 & h_0 & 0 & & \vdots \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \vdots & & \ddots & \ddots & \ddots & \vdots \\ \vdots & & & \ddots & \ddots & 0 \\ h_{n-1} & \cdots & & & h_1 & h_0 \end{bmatrix} \begin{bmatrix} x_0 \\ \vdots \\ \vdots \\ x_{n-1} \end{bmatrix}$$

$$=: \mathbf{H}_n \in \mathbb{R}^{n,n}$$

$$(4.4.1)$$

Note that the matrix $\mathbf{H}_n \in \mathbb{R}^{n,n}$ is a lower triangular matrix with constant diagonals (linear algebra indexing)

$$(\mathbf{H})_{l,j} = egin{cases} h_{l-j} & ext{if } l \geq j \ 0 & ext{else}, \end{cases} \quad l,j \in \{1,\ldots,n\}$$

This matrix is clearly not circulant: see [Lecture \rightarrow Def. 4.1.4.12]. At the same time, it still belongs to the class of *Toeplitz matrices*, see [Lecture \rightarrow Def. 4.5.1.7], which generalizes the class of circulant matrices. (Recall [Lecture \rightarrow Def. 4.5.1.7]: $\mathbf{T} \in \mathbb{K}^{m,n}$ is a Toeplitz matrix if there is a sequence/array $\mathbf{u} = (u_{-m+1}, \ldots, u_{n-1}) \in \mathbb{K}^{m+n-1}$ such that $t_{ij} = u_{j-i}$ for $1 \le i \le m$ and $1 \le j \le n$.)

If you want to recover the input signal from the output (*deconvolution*), you will face the following crucial task for digital signal processing: *solve a lower triangular LSE with a Toeplitz system matrix*. Of course, forward elimination, see [Lecture \rightarrow Section 2.3], can achieve this with asymptotic complexity $\mathcal{O}(n^2)$: see [Lecture \rightarrow Eq. (2.3.1.6)]. However, since a Toeplitz matrix merely has an information content of $\mathcal{O}(n)$ numbers, there might be a more efficient way – which ought to be explored by you in this problem.

(4-4.a) \Box (15 min.) Given a Toeplitz matrix $\mathbf{T} \in \mathbb{R}^{n,n}$, find another matrix $\mathbf{S} \in \mathbb{R}^{n,n}$ such that the following block matrix is a *circulant* matrix:

$$\mathbf{C} = egin{bmatrix} \mathbf{T} & \mathbf{S} \\ \mathbf{S} & \mathbf{T} \end{bmatrix} \in \mathbb{R}^{2n,2n} \ .$$

HIDDEN HINT 1 for (4-4.a) \rightarrow 4-4-1-0:Toeplitz1h.pdf

SOLUTION for (4-4.a) \rightarrow 4-4-1-1:Toeplitz1s.pdf

•

Write the following C++ function:

This function should take as input column vectors $\mathbf{c} \in \mathbb{R}^m$ and $\mathbf{r} \in \mathbb{R}^n$ and return a Toeplitz matrix $\mathbf{T} \in \mathbb{R}^{m,n}$ [Lecture \to Def. 4.5.1.7] such that (linear-algebra indexing)

$$(\mathbf{T})_{i,j} = \begin{cases} (c)_{i-j+1} & \text{, if } i \geq j \text{,} \\ (r)_{j-i+1} & \text{, if } j > i \text{,} \end{cases} \quad 1 \leq i \leq m, \ 1 \leq j \leq n \ .$$

SOLUTION for (4-4.b) $\rightarrow 4-4-2-0$:tplcpp.pdf

Show that the following two C++ functions realize the same linear mapping $x \mapsto y$, when supplied with the same arguments.

```
C++11-code 4.4.4: Function toepmult()
```

```
Eigen::VectorXd toepmult(const Eigen::VectorXd &c, const Eigen::VectorXd &r,
                            const Eigen::VectorXd &x) {
     assert(c.size() == r.size() && c.size() == x.size() &&
            "c, r, x have different lengths!");
5
     const unsigned int n = c.size();
6
     // Complex arithmetic required here
     Eigen::VectorXcd cr tmp = c.cast<std::complex<double>>();
     Eigen::VectorXcd x_tmp = x.cast<std::complex<double>>();
10
     // Prepare vector encoding circulant matrix {f C}
     cr_tmp.conservativeResize(2 * n);
13
     cr_tmp.tail(n) = Eigen::VectorXcd::Zero(n);
14
     cr_{tmp.tail}(n-1).real() = r.tail(n-1).reverse();
15
    // Zero padding
17
    x_tmp.conservativeResize(2 * n);
18
     x_tmp.tail(n) = Eigen::VectorXcd::Zero(n);
19
20
    // Periodic discrete convolution from [Lecture \rightarrow Code 4.2.2.4]
21
     Eigen::VectorXd y = pconvfft(cr_tmp, x_tmp).real();
22
    y.conservativeResize(n);
23
24
     return y;
25
```

Get it on ₩ GitLab (toeplitz.hpp).

The function pconvfft is defined in [Lecture \rightarrow Section 4.2.2], [Lecture \rightarrow Code 4.2.2.4], and toeplitz () is the simple C++ function from Sub-problem (4-4.b).

```
HIDDEN HINT 1 for (4-4.c) \rightarrow 4-4-3-0: Toeplitz1h.pdf
```

```
SOLUTION for (4-4.c) \rightarrow 4-4-3-1:Toeplitz2s.pdf
```

(4-4.d) • (10 min.) What is the asymptotic complexity of toepmatmult () (\rightarrow Code 4.4.3) and toepmult () (\rightarrow Code 4.4.4)?

```
SOLUTION for (4-4.d) \rightarrow 4-4-4-0:Toeplitz3s.pdf
```

(4-4.e)

(45 min.) Explain in detail what the following C++ functions are meant for and why they are algebraically equivalent (i.e. equivalent when ignoring roundoff errors), provided that h.size() = 2^1.

C++11-code 4.4.6: Function ttmatsolve Eigen::VectorXd ttmatsolve(const Eigen::VectorXd &h, const Eigen::VectorXd &y) { assert(h.size() == y.size() && "h and y have different lengths!"); 3 const unsigned int n = h.size(); 5 **Eigen::VectorXd** h_tmp = **Eigen::VectorXd::Zero**(n); 6 $h_{tmp}(0) = h(0);$ 8 Eigen::MatrixXd T = toeplitz(h, h_tmp); 9 10 Eigen::VectorXd x = T.fullPivLu().solve(y); 11 12 return x; 13 14 | } Get it on ₩ GitLab (toeplitz.hpp).

C++11-code 4.4.7: Function ttrecsolve

```
Eigen::VectorXd ttrecsolve(const Eigen::VectorXd &h, const Eigen::VectorXd &y,
                               int | | {
3
     assert(h.size() == y.size() && "h and y have different lengths!");
4
     // Result vector
    Eigen:: VectorXd x;
6
    // Trivial case of asn 1x1 LSE
    if (| == 0) {
       x.resize(1);
       x(0) = y(0) / h(0);
10
     } else {
11
       int n = std::pow(2, 1);
12
       int m = n / 2;
13
14
       // Check matching length of vectors
15
       assert(h.size() == n && y.size() == n &&
16
              "h and y have length different from 2^1!");
17
18
       Eigen:: VectorXd x1 = ttrecsolve(h.head(m), y.head(m), l - 1);
19
       Eigen::VectorXd y2 =
20
21
           y.segment(m, m) -
           toepmult(h.segment(m, m), h.segment(1, m).reverse(), x1);
22
       Eigen:: VectorXd x2 = ttrecsolve(h.head(m), y2, I - 1);
       x.resize(n);
25
```

```
x. head (m) = x1;

x. tail (m) = x2;

}

return x;

}

Get it on 	GitLab (toeplitz.hpp).
```

HIDDEN HINT 1 for (4-4.e) \rightarrow 4-4-5-0:Toeplitz4h.pdf

SOLUTION for (4-4.e) \rightarrow 4-4-5-1:Toeplitz4s.pdf

(4-4.f) • (5 min.) Why is the algorithm implemented in ttrecsolve() called a "divide & conquer" method?

```
SOLUTION for (4-4.f) \rightarrow 4-4-6-0: Toeplitz5s.pdf
```

```
h = Eigen::VectorXd::LinSpaced(n,1,n).cwiseInverse();
```

and $n = 2^l$, l = 3, ..., 9. Plot the timing results for both functions in doubly logarithmic scale. Describe your observations.

```
HIDDEN HINT 1 for (4-4.g) \rightarrow 4-4-7-0:tp6h1.pdf
```

```
Solution for (4-4.g) \rightarrow 4-4-7-1:Toeplitz6s.pdf
```

(4-4.h) \odot (20 min.) Derive the asymptotic complexity of both functions ttmatsolve() and ttrecsolve() in terms of the problem size parameter $n \to \infty$.

```
Solution for (4-4.h) \rightarrow 4-4-8-0:Toeplitz7s.pdf
```

(4-4.i) \odot (45 min.) For the case that n = h.size() is **not** a power of 2, implement a *wrapper function*

```
Eigen::VectorXd ttsolve(const Eigen::VectorXd &h, const
    Eigen::VectorXd &y);
```

for ttrecsolve that, in the absence of roundoff errors, would behave exactly like ttmatsolve, i.e. ttsolve(h,y) == ttmatsolve(h,y). The term "wrapper function" indicates that ttsolve() is supposed to call ttrecsolve() after preparing suitable input vectors.

```
HIDDEN HINT 1 for (4-4.i) \rightarrow 4-4-9-0:tplhext.pdf
```

```
SOLUTION for (4-4.i) \rightarrow 4-4-9-1:Toeplitz8s.pdf
```

End Problem 4-4, 220 min.

Problem 4-5: Implementing Discrete Convolution

[Lecture \rightarrow Section 4.2.2] discussed in detail the implementation of discrete convolutions of signals stored in vectors by means of the discrete Fourier transform (DFT). This problem revisits these techniques and adds a new twist to implementation.

This problem assumes solid knowledge of [Lecture \rightarrow Section 4.2.2].

From [Lecture \rightarrow Section 4.1.3] we recall the definitions

Definition [Lecture → Def. 4.1.3.3]. Discrete convolution

Given $\mathbf{x} = [x_0, \dots, x_{m-1}]^{\top} \in \mathbb{K}^m$, $\mathbf{h} = [h_0, \dots, h_{n-1}]^{\top} \in \mathbb{K}^n$ their **discrete convolution** (DCONV) $\mathbf{h} * \mathbf{x}$ is the vector $\mathbf{y} = [y_0, \dots, y_{m+n-2}]^{\top} \in \mathbb{K}^{m+n-1}$ with components

$$y_k = \sum_{j=0}^{m-1} h_{k-j} x_j$$
, $k = 0, ..., m+n-2$, [Lecture \to Eq. (4.1.3.4)]

where we have adopted the convention $h_j := 0$ for j < 0 or $j \ge n$.

and

Definition *cf.* [Lecture \rightarrow Def. 4.1.4.7]. Discrete periodic convolution (of vectors)

The discrete periodic convolution $p *_n x$ of two vectors $p, x \in \mathbb{K}^n$ is the vector $y \in \mathbb{K}^n$ with components (C++ indexing)

$$(\mathbf{y})_k := \sum_{j=0}^{n-1} (\mathbf{p})_{(k-j) \bmod n} (\mathbf{x})_j, \quad k \in \{0, \dots, n-1\}.$$

$$((k-j) \mod n \in \{0, ..., n-1\} \triangleq (k-j) \text{ %n in C++ syntax.})$$

$$\mathbf{h} * \mathbf{x} = \left(\widetilde{\mathbf{h}} *_{N} \widetilde{\mathbf{x}}\right)_{?:?}. \tag{4.5.1}$$

C++ indexing has to be used and the notational convention is $(\mathbf{v})_{r:s} = [(\mathbf{v})_r, \dots, (\mathbf{v})_s] \in \mathbb{K}^{s-r+1}$ for $r, s \in \mathbb{N}, r \leq s$. Note the different discrete convolutions on the left-hand and right-hand side of (4.5.1).

SOLUTION for (4-5.a)
$$\rightarrow$$
 4-5-1-0:ptdcs1s.pdf

The discrete Fourier transform in Eigen relies on the library KissFFT, which offers a very efficient implementation of DFT with asymptotic computational cost $O(n \log n)$ for vectors of length $n := 2^{\ell}$, $\ell \in \mathbb{N}$, but may not be fast for vectors of other lengths. To avoid degraded performance of a numerical code, the following rule is imposed throughout this problem:

EIGEN's built-in DFT implementations may be invoked only for vectors of length 2^ℓ for some $\ell \in \mathbb{N}!$

(4-5.b) (60 min.) Based on EIGEN, in the file discreteconvolution.hpp implement an efficient function

that computes the discrete periodic convolution $\mathbf{p} *_n \mathbf{x}$ as defined in [Lecture \rightarrow Def. 4.1.4.7] of two real vectors $\mathbf{p}, \mathbf{x} \in \mathbb{R}^n$ of the same length.

```
HIDDEN HINT 1 for (4-5.b) \rightarrow 4-5-2-0:dcv2h1.pdf Solution for (4-5.b) \rightarrow 4-5-2-1:dcsol1.pdf
```

End Problem 4-5, 75 min.

Problem 4-6: Asymptotic Cost of EIGEN-Based Functions

This short problem is about reading off the asymptotic computational cost of some C++ functions performing numerical linear algebra tasks based on EIGEN. It relies on information provided in [Lecture \rightarrow Section 1.4.2], [Lecture \rightarrow § 2.5.0.4], [Lecture \rightarrow § 3.3.3.37], and [Lecture \rightarrow Section 4.3].

Just reading of C++ code is required.

(4-6.a) • (8 min.) The function slvtriag() listed in Code 4.6.1 expects two vector arguments of equal length $n \in \mathbb{N}$. What is its asymptotic computational complexity for $n \to \infty$?

$$\operatorname{cost}(\operatorname{slvtriag}) = O($$
) for $n \to \infty$.

```
C++ code 4.6.1: Function slvtriag()

Eigen::VectorXd slvtriag(const Eigen::VectorXd &u, const Eigen::VectorXd &v) {
   assert((u.size() == v.size()) && "Size mismatch");
   const Eigen::MatrixXd A{u * v.transpose()};
   return A.triangularView < Eigen::Upper > ().solve(u);
}
```

SOLUTION for (4-6.a) \rightarrow 4-6-1-0:s1.pdf

(4-6.b) \odot (8 min.) What is the asymptotic complexity of the C++ function slvsymrom() listed as Code 4.6.2 in terms of the length n of its argument vector \mathbf{u} ?

```
\operatorname{cost}(\operatorname{slvsymrom}) = O( \hspace{1cm} ) \hspace{1cm} \text{for} \hspace{0.2cm} n 	o \infty \, .
```

A sharp bound is expected.

```
C++ code 4.6.2: Function slvsymrom()

Eigen::VectorXd slvsymrom(const Eigen::VectorXd &u) {
    const unsigned int n = u.size();
    return (Eigen::MatrixXd::Identity(n, n) + u * u.transpose()).lu().solve(u);
}
```

SOLUTION for (4-6.b) \rightarrow 4-6-2-0:s2.pdf

(4-6.c) Code 4.6.3 lists the C++ function getcompbas(), whose argument A is a densely populated matrix $A \in \mathbb{R}^{n,k}$, k < n. Give a sharp asymptotic bound for the asymptotic computational cost of getcompbas() for $n \to \infty$ assuming k to be small and fixed.

```
\mathrm{cost}(\mathtt{getcompbas}) = O( ) for n \to \infty.
```

```
C++ code 4.6.3: Function getcompbas()

Eigen::MatrixXd getcompbas(const Eigen::MatrixXd &A) {
   const int n = A.rows();
   const int k = A.cols();
   Eigen::HouseholderQR<Eigen::MatrixXd> qr(A);
```

```
return qr.householderQ() *

(Eigen::MatrixXd(n, n - k) << Eigen::MatrixXd::Zero(k, n - k),

Eigen::MatrixXd::Identity(n - k, n - k))

finished();
```

HIDDEN HINT 1 for (4-6.c) \rightarrow 4-6-3-0:asch31.pdf

SOLUTION for (4-6.c) \rightarrow 4-6-3-1:s3.pdf

(4-6.d) (12 min.) Code 4.6.4 shows the implementation of the C++ function fft2_square(), which requires a complex $n \times n$ -matrix $\mathbf{Y} \in \mathbb{C}^{n,n}$ as argument Y. What is a sharp bound on the asymptotic complexity of fft2_square() as $n \to \infty$?

```
\operatorname{cost}(\operatorname{fft2\_square}) = O( ) for n \to \infty.
```

```
C++ code 4.6.4: Function fft2_square()
  Eigen::MatrixXcd fft2_square(const Eigen::MatrixXcd &Y) {
     const unsigned int n = Y.rows();
3
     assert((n == Y.cols()) && "Matrix Y must be square");
     Eigen:: MatrixXcd tmp(n, n);
5
     Eigen::FFT<double> fft;
     for (int k = 0; k < n; ++k) {
       const Eigen::VectorXcd tv(Y.row(k));
8
       tmp.row(k) = fft.fwd(tv).transpose();
9
10
     for (int k = 0; k < n; ++k) {
11
       const Eigen::VectorXcd tv(tmp.col(k));
12
       tmp.col(k) = fft.fwd(tv);
13
14
15
     return tmp;
  }
16
```

SOLUTION for (4-6.d) $\rightarrow 4-6-4-0:s4.pdf$

End Problem 4-6, 40 min.

Problem 4-7: Modified Cosine Transform

Many so-called trigonometric transformations can be reduced to the discrete Fourier transform (DFT). One such example is studied in this problem.

Assumes familiarity with the discrete Fourier transform, C++, and the complex exponential

For a given vector $\mathbf{a} = [a_0, \dots, a_{n-1}]^{\top} \in \mathbb{R}^n$, $n \in \mathbb{N}$, we want to compute another vector $\mathbf{f} := [f_0, \dots, f_{n-1}]^{\top} \in \mathbb{R}^n$ according to the formula

$$f_k := \sum_{j=0}^{n-1} a_j \cos\left(\pi \frac{jk}{n}\right), \quad k = 0, \dots, n-1.$$
 (4.7.1)

Using the identity $\cos x = \frac{1}{2}(\exp(-\iota x) + \exp(\iota x))$, we find the alternative formula

$$f_{k} = \sum_{\ell=0}^{n-1} \frac{1}{2} a_{\ell} \exp\left(-2\pi \imath \frac{k\ell}{2n}\right) + \exp\left(\pi \imath \frac{k(n-1)}{n}\right) \cdot \sum_{\ell=0}^{n-1} \frac{1}{2} a_{n-1-\ell} \exp\left(-2\pi \imath \frac{k\ell}{2n}\right)$$
(4.7.2)

for k = 0, ..., n - 1.

Recall the definition of the discrete Fourier transform (DFT) as multiplication of a vector with the Fourier matrix \mathbf{F}_n :

Definition [Lecture → Def. 4.2.1.18].

The linear map $\mathsf{DFT}_n:\mathbb{C}^n\mapsto\mathbb{C}^n$, $\mathsf{DFT}_n(\mathbf{y}):=\mathbf{F}_n\mathbf{y}$, $\mathbf{y}\in\mathbb{C}^n$, is called **discrete Fourier transform** (DFT), i.e. for $[c_0,\ldots,c_{n-1}]:=\mathsf{DFT}_n(\mathbf{y})$

$$c_k := \sum_{j=0}^{n-1} y_j \exp\left(-2\pi \imath \frac{kj}{n}\right)$$
 , $k = 0, \dots, n-1$. [Lecture \to Eq. (4.2.1.19)]

(4-7.a) (10 min.) Taking the cue from (4.7.2) we can write ("C++ indexing" throughout)

$$\mathbf{f} = (\mathsf{DFT}_{2n}(\mathbf{x}))_{0:n-1} + \mathbf{D}(\mathsf{DFT}_{2n}(\mathbf{y}))_{0:n-1},$$
 (4.7.3)

with suitable vectors $\mathbf{x}, \mathbf{y} \in \mathbb{C}^{2n}$ and a diagonal matrix $\mathbf{D} \in \mathbb{C}^{n,n}$. Characterize the vectors \mathbf{x}, \mathbf{y} (in terms of \mathbf{a}) and the matrix \mathbf{D} .

$$(\mathbf{x})_{\ell} = \left\{ \begin{array}{c} & \text{for } \ell \in \{ \\ & \end{array} \right\}, \\ & \text{for } \ell \in \{ \\ \end{array} \right\}, \\ & (\mathbf{D})_{k,k} = \left\{ \begin{array}{c} & \text{for } \ell \in \{ \\ & \end{array} \right\}, \\ & k = 0, \dots, n-1 \ . \\ \end{array}$$

SOLUTION for (4-7.a) \rightarrow 4-7-1-0:mcts1.pdf

(4-7.b) (10 min.) The C++ function modcostrf() realizes the mapping $a \rightarrow f$ according to (4.7.1). Supplement the missing parts of the following listing by writing valid C++ code in the boxes.

```
Eigen::VectorXd modcostrf(const Eigen::VectorXd &a) {
   using Comp = std::complex<double>;
   unsigned int n = a.size();
```

```
Eigen::VectorXd f(n);
  Eigen::FFT<double> fft;
  Eigen::VectorXcd tmp(
                                  );
  tmp <<
                          , Eigen::VectorXcd::Zero(n);
  Eigen::VectorXcd v1 = fft.
                                               ;
  tmp << 0.5 * a.reverse(),</pre>
                                                                  ;
  Eigen::VectorXcd v2 = fft.
  Comp fac = std::exp(((n - 1) * M_PI * Comp(0, 1)) / (double)n);
  Comp d(1.0, 0.0);
  for (int k = 0; k < n; ++k) {
    f[k] = (v1[
                                                        ]).real();
                                          v2 [
                            ;
  }
  return f;
}
```

The constant M_PI is the number π .

```
HIDDEN HINT 1 for (4-7.b) \rightarrow 4-7-2-0:mcth1.pdf
```

HIDDEN HINT 2 for (4-7.b) \rightarrow 4-7-2-1:mcth2.pdf

SOLUTION for (4-7.b) \rightarrow 4-7-2-2:.pdf

End Problem 4-7, 20 min.

Problem 4-8: Discrete Fourier Transform and Applications

The discrete Fourier transform (DFT) is a fundamental tool is signal theory and processing. This problem studies some of its properties and the connection with periodic convolution.

This problem is based upon [Lecture \rightarrow Section 4.2.1] and [Lecture \rightarrow Section 4.2.2].

Recall the discrete Fourier transform (DFT)

Definition [Lecture → Def. 4.2.1.18]. discrete Fourier transform

The linear map $\mathsf{DFT}_n:\mathbb{C}^n\mapsto\mathbb{C}^n$, $\mathsf{DFT}_n(\mathbf{y}):=\mathbf{F}_n\mathbf{y},\mathbf{y}\in\mathbb{C}^n$, is called **discrete Fourier transform** (DFT), i.e. for $[c_0,\ldots,c_{n-1}]:=\mathsf{DFT}_n(\mathbf{y})$

$$c_k = \sum_{j=0}^{n-1} y_j \, \omega_n^{kj} = \sum_{j=0}^{n-1} y_j \, \exp\left(-2\pi \imath \frac{kj}{n}\right)$$
 , $k = 0, \dots, n-1$. [Lecture \to Eq. (4.2.1.19)]

Here \mathbf{F}_n is the Fourier matrix

$$\mathbf{F}_n = \left[\exp\left(-2\pi \imath \frac{\ell j}{n}\right)\right]_{\ell,j=0}^n \in \mathbb{C}^{n,n}, \quad n \in \mathbb{N}.$$

about which we known

7 }

Lemma [Lecture → Lemma 4.2.1.14]. Properties of Fourier matrices

The scaled Fourier-matrix $\frac{1}{\sqrt{n}}\mathbf{F}_n$ is unitary (\rightarrow [Lecture \rightarrow Def. 6.3.1.2]): $\mathbf{F}_n^{-1} = \frac{1}{n}\mathbf{F}_n^{\mathrm{H}} = \frac{1}{n}\overline{\mathbf{F}}_n$.

In Eigen the discrete Fourier transform of a vector is available through the helper class **Eigen::FFT** and its member function fwd (). The inverse DFT can be carried out by calling the member function inv (), see [Lecture \rightarrow Code 4.2.1.22].

(4-8.a) (10 min.) Code 4.8.1 displays the incomplete listing of a C++ function

```
Eigen::VectorXcd inverseDFT(const Eigen::VectorXcd &c);
```

that implements <code>Eigen::FFT::inv()</code> only based on calls to <code>Eigen::FFT::fwd()</code>. Fill in the missing code parts so that the implementation conforms with the specification.


```
HIDDEN HINT 1 for (4-8.a) \rightarrow 4-8-1-0:fftmhal.pdf
```

SOLUTION for (4-8.a) \rightarrow 4-8-1-1:fftmas.pdf

(4-8.b) : (15 min.) Based on

Lemma [Lecture → Lemma 4.2.1.16]. Diagonalization of circulant matrices

For any circulant matrix $C \in \mathbb{K}^{n,n}$, $c_{ij} = u_{i-j}$, $(u_k)_{k \in \mathbb{Z}}$ an n-periodic sequence, holds true

$$\mathbf{C}\overline{\mathbf{F}}_n = \overline{\mathbf{F}}_n \operatorname{diag}(d_0,\ldots,d_{n-1})$$
, $[d_0,\ldots,d_{n-1}]^{\top} = \mathbf{F}_n[u_0,\ldots,u_{n-1}]^{\top}$.

the C++ function

```
Eigen::VectorXcd circmatvec(const Eigen::VectorXcd &u,
const Eigen::VectorXcd &x);
```

realizes the left-multiplication of the n-vector passed in \times with a circulant $n \times n$ -matrix defined by the n-vector u (periodically extended to an n-periodic sequence (u_n)).

Complete the listing given in Code 4.8.3.

SOLUTION for (4-8.b) \rightarrow 4-8-2-0:fftm1.pdf

End Problem 4-8, 25 min.

Chapter 5

Data Interpolation and Data Fitting in 1D

Problem 5-1: Evaluating the derivatives of interpolating polynomials

In [Lecture \to Ex. 5.1.0.8] we learned about the importance of data interpolation for obtaining functor representations of constitutive relationships $t \mapsto f(t)$. Numerical methods like Newton's method [Lecture \to Code 8.4.2.2] often require information about the derivative f' as well. Therfore, we need efficient ways to evaluate the derivatives of interpolants. In this problem we discuss this issue for polynomial interpolation for (i) monomial representation and (ii) "update-friendly" point evaluation. We generalize the Horner scheme [Lecture \to § 5.2.1.5] and the Aitken-Neville algorithm [Lecture \to § 5.2.3.8].

This problem is connected with [Lecture \rightarrow Section 5.2.3], simple coding in C++/EIGEN is requested.

We first deal with polynomials in monomial representation [Lecture \rightarrow Rem. 5.2.1.4].

(5-1.a) • (20 min.) Using the Horner scheme [Lecture \rightarrow § 5.2.1.5], in the file <code>eval_deriv.hpp</code> write an efficient C++ implementation of a template function which returns the pair (p(x), p'(x)), where p is a polynomial with coefficients in c:

```
std::pair<double, double > evaldp(const CoeffVec &c, double x);
```

Vector c contains the coefficient of the polynomial in the monomial basis, following the PYTHON/MATLAB convention (leading coefficient in c [0]), that is,

$$p(t) = c[0]t^{k} + c[1]t^{k-1} + \dots + c[k-1]t + c[k], \quad t \in \mathbb{R}.$$
 (5.1.1)

As indicated above, the type **CoeffVec** must provide component access via **operator** [] (**int**) **const** and a size() method.

```
SOLUTION for (5-1.a) \rightarrow 5-1-1-0: Horner1s.pdf
```

(5-1.b) • (15 min.) For the sake of testing, write a naive C++ implementation of the evaluation requested in (5-1.a)

```
std::pair<double, double > evaldp_naive(const CoeffVec &c, double x);
```

which returns the same pair (p(x), p'(x)). However, this time p(x) and p'(x) should be calculated by means of simply summing the contributions of the monomials $t \mapsto t^k$.

SOLUTION for (5-1.b)
$$\rightarrow 5-1-2-0$$
: Horner2s.pdf

What are the asymptotic complexities of the two functions you implemented in (5-1.a) and (5-1.b) in terms of raising the polynomial degree?

```
SOLUTION for (5-1.c) \rightarrow 5-1-3-0:Horner3s.pdf
```

Implement a function

```
bool polyTestTime(const unsigned int d);
```

that

- (i) validates your implementation of the two functions $evaldp()/evaldp_naive()$ from (5-1.a)/(5-1.b) by comparing their results for x=1.23 and a polynomial of degree n with monomial coefficients $c_0=1, c_1=2, \ldots, c_n=n$ and $n=2,4,8,\ldots,2^d$. Return **false** in case of a discrepancy.
- (ii) measures the runtimes of the two functions for every degree and tabulates them. To that end use the **Timer** class defined in timer.h:

```
Timer t;
t.start(); /* do something */ t.stop();
double time = t.duration();
```

HIDDEN HINT 1 for (5-1.d) \rightarrow 5-1-4-0:h4h1.pdf

```
SOLUTION for (5-1.d) \rightarrow 5-1-4-1: Horner4s.pdf
```

Now we revisit polynomial interpolation. In [Lecture \rightarrow Section 5.2.3.2] we learned about an efficient and "update-friendly" scheme for evaluating Lagrange interpolants at a *single or a few* points. This so-called **Aitken-Neville algorithm**, see [Lecture \rightarrow Code 5.2.3.10], can be extended to return the derivative value of the polynomial interpolant as well. This will be explored next.

(5-1.e) (30 min.) In the file eval_deriv.hpp write an efficient C++ function

that returns the row vector $[p'(x_1), \ldots, p'(x_m)]^{\top}$, when the argument x passes $[x_1, \ldots, x_m]$, $m \in \mathbb{N}$ small. Here, p' denotes the *derivative* of the polynomial $p \in \mathcal{P}_n$ interpolating the data points (t_i, y_i) , $i = 0, \ldots, n$, for pairwise different $t_i \in \mathbb{R}$ and data values $y_i \in \mathbb{R}$.

To that end differentiate the recursion formula [Lecture \rightarrow Eq. (5.2.3.9)] and devise an algorithm in the spirit of the Aitken-Neville algorithm implemented in [Lecture \rightarrow Code 5.2.3.10].

```
HIDDEN HINT 1 for (5-1.e) \rightarrow 5-1-5-0:hornxh1.pdf
SOLUTION for (5-1.e) \rightarrow 5-1-5-1:dipoleval.pdf
```

(5-1.f) (30 min.) For validation purposes devise an alternative, less efficient, implementation of dipoleval()

based on the following algorithm:

1. Use the function

```
Eigen::VectorXd polyfit(const Eigen::VectorXd& t,
                        const Eigen::VectorXd& y,
                        unsigned int degree);
```

(supplied in polyfit.hpp) to compute the monomial coefficients of the Lagrange interpolant of the data points.

- 2. Compute the monomial coefficients of the derivative.
- 3. Use the function

```
void polyval(const Eigen::VectorXd& p,
             const Eigen::VectorXd& x,
             Eigen::VectorXd& y);
```

(defined in polyval.hpp) to evaluate the derivative at a number of points. This function evaluates the polynomial with monomial coefficients passed in p at all points in x, cf. [Lecture \rightarrow Code 5.2.1.7.

Use dipoleval_alt () to verify the correctness of your implementation of dipoleval () by computing the derivative of the degree-10 interpolating polynomial for the values $(t_i, sin(t_i))$, where t_i are equidistant points in [0,3]. This test should be conducted within the function

```
bool testDipolEval(void);
```

Return **false** in case the test fails, **true** otherwise.

```
SOLUTION for (5-1.f) \rightarrow 5-1-6-0:test.pdf
```

```
(5-1.g) • (20 min.)
                      Based on your version of dipoleval () write a C++ function
    void plotPolyInterpolant(const std::string &filename);
```

that saves plots the derivative of the polynomial interpolant for the test data as specified in Subproblem (5-1.f) in the file <filename>.png. To that end evaluate the polynomial in 100 equidistant points in the interval [0,3] and, based on the obtained value, use the plot () function of MATPLOTLIBCPP to draw the polynomials. Do not forget to add axis labels "t" and "p'(t)", and a ti-

```
HIDDEN HINT 1 for (5-1.q) \rightarrow 5-1-7-0:hx1.pdf
SOLUTION for (5-1.g) \rightarrow 5-1-7-1:sx.pdf
```

End Problem 5-1, 135 min.

Problem 5-2: Lagrange Polynomials

This short problem studies a particular aspect of the **Lagrange polynomials** introduced in [Lecture \rightarrow § 5.2.2.3].

This exercise just requires basic knowledge of [Lecture → Section 5.2.2]

As in [Lecture \to Eq. (5.2.2.4)] we denote by L_i the i-th Lagrange polynomial for given nodes $t_j \in \mathbb{R}$, $j=0,\ldots,n,\ t_i\neq t_j$, if $i\neq j$. Then the polynomial Lagrange interpolant p through the data points $(t_i,y_i)_{i=0}^n$ has the representation

$$p(x) = \sum_{i=0}^{n} y_i L_i(x) \quad \text{with} \quad L_i(x) := \prod_{\substack{j=0 \ j \neq i}}^{n} \frac{x - t_j}{t_i - t_j}.$$
 (5.2.1)

$$\sum_{i=0}^{n} L_i(t) = 1 \quad \forall t \in \mathbb{R} . \tag{5.2.2}$$

HIDDEN HINT 1 for (5-2.a) $\rightarrow 5-2-1-0$:glgh1.pdf

SOLUTION for (5-2.a)
$$\rightarrow 5-2-1-1$$
:glgs1.pdf

(5-2.b) : (15 min.) Show that

$$\sum_{i=0}^{n} L_i(0) t_i^j = \begin{cases} 1 & \text{for } j = 0, \\ 0 & \text{for } j = 1, \dots, n. \end{cases}$$
 (5.2.3)

(Here, read t_i^j as t_i raised to the power j.)

HIDDEN HINT 1 for (5-2.b) $\rightarrow 5-2-2-0$:glgh1.pdf

SOLUTION for (5-2.b)
$$\rightarrow$$
 5-2-2-1:glgs1.pdf

(5-2.c) \odot (15 min.) Show that p(x) can also be written as

$$p(x) = \omega(x) \sum_{i=0}^{n} \frac{y_i}{(x - t_i)\omega'(t_i)}$$
 with $\omega(t) := \prod_{j=0}^{n} (t - t_j)$ (5.2.4)

Here ω' is the derivative of ω .

HIDDEN HINT 1 for (5-2.c) \rightarrow 5-2-3-0:LagrangePoly1h.pdf

SOLUTION for (5-2.c) \rightarrow 5-2-3-1:LagrangePoly1s.pdf

End Problem 5-2, 40 min.

Problem 5-3: Generalized Lagrange polynomials for Hermite interpolation

[Lecture \rightarrow Rem. 5.2.2.14] addresses a particular generalization of the Lagrange polynomial interpolation considered in [Lecture \rightarrow Section 5.2.2]: Hermite interpolation, where beside the usual "nodal" interpolation conditions also the derivative of the interpolating polynomials is prescribed in the nodes, see [Lecture \rightarrow Eq. (5.2.2.15)] (for $l_i = 1$ throughout). The associated generalized Lagrange polynomials were defined in [Lecture \rightarrow Def. 5.2.2.17] and in this problem we aim to find concrete formulas for them in the spirit of [Lecture \rightarrow Eq. (5.2.2.4)].

This is a purely theoretical problem.

Let n+1 different nodes t_i , $i=0,\ldots,n, n\in\mathbb{N}$, be given. For numbers $y_i,c_i\in\mathbb{R}, j=0,\ldots,n$, Hermite interpolation seeks a polynomial $p \in \mathcal{P}_{2n+1}$ satisfying $p(t_i) = y_i$ and $p'(t_i) = c_i$, $i = 0, \ldots, n$. Here p'designates the derivative of p.

The **generalized Lagrange polynomials** for Hermite interpolation L_i , $K_i \in \mathcal{P}_{2n+1}$ satisfy

$$L_i(t_j) = \delta_{ij}$$
 , $L_i'(t_j) = 0$, (5.3.1a)

$$K_i(t_j) = 0$$
, $K'_i(t_j) = \delta_{ij}$, (5.3.1b)

for $i, j \in \{0, ..., n\}$.

Explicitly state the linear system of equations for the basis expansion coefficients of the Hermite interpolant with respect to the *monomial basis* $\{x \mapsto x^k, k = 0, \dots, 2n + 1\}$ of \mathcal{P}_{2n+1} .

HIDDEN HINT 1 for (5-3.a) $\rightarrow 5-3-1-0$: gagh0.pdf

SOLUTION for (5-3.a)
$$\rightarrow$$
 5-3-1-1:glgs1.pdf

Give concrete formulas for K_0 , K_1 , L_0 , L_1 for n = 1, $t_0 = -1$, $t_1 = 1$.

HIDDEN HINT 1 for (5-3.b) $\rightarrow 5-3-2-0$:glgh1.pdf

SOLUTION for (5-3.b) \rightarrow 5-3-2-1:glgs1.pdfNote that $L_0(t) = L_1(-t)$, $K_0(t) = K_1(-t)$, which reflects the symmetry of the interval [-1,1].

- **(5-3.c)** (10 min.) Find the general formula describing all polynomials in *t*
 - 1. of degree 2 that have a double zero in t = a, $a \in \mathbb{R}$,
 - 2. of degree n > 2 that have a double zero in t = a, $a \in \mathbb{R}$.

A polynomial p is said to have a double zero in t=a, if p(a)=0 and p'(a)=0.

SOLUTION for (5-3.c)
$$\rightarrow$$
 5-3-3-0:glgs1.pdf

Let $g_1, \ldots, g_m, m \in \mathbb{N}$, be functions in $C^1(I)$ (space of continuously differentiable **(5-3.d)** (15 min.) functions). Find a formula for the derivative of $h(t) = \prod_{k=1}^{m} g_k(t)$.

SOLUTION for (5-3.d) $\rightarrow 5-3-4-0$:glgs1.pdf

depends on Sub-problem (5-3.c), Sub-problem (5-3.d)

Find general formulas for the polynomials L_i , K_i as defined through (5.3.1). Of course, these formulas will involve the nodes t_i .

```
HIDDEN HINT 1 for (5-3.e) \rightarrow 5-3-5-0:glghl.pdf
HIDDEN HINT 2 for (5-3.e) \rightarrow 5-3-5-1:glghx.pdf
Solution for (5-3.e) \rightarrow 5-3-5-2:glgs1.pdf
```

End Problem 5-3, 85 min.

Problem 5-4: Piecewise linear interpolation

[Lecture \rightarrow Ex. 5.1.0.15] introduced piecewise linear interpolation as a simple linear interpolation scheme. It finds an interpolant in the space spanned by the so-called tent functions, which are cardinal basis functions. Formulas are given in [Lecture \rightarrow Eq. (5.1.0.17)].

Study [Lecture \rightarrow Ex. 5.1.0.15] before tackling this problem. Problem involves coding in C++.

representing the piecewise linear interpolant of the data points (t_i, y_i) passed to the constructor through two vectors t and y of equal length. In addition the class should provide the evaluation operator **double operator** (**double** x) **const**, which returns the value of the piecewise linear interpolant at $x \in \mathbb{R}$. Pay attention to efficient implementation.

Note that you must not assume that the nodes passed in t are sorted.

```
SOLUTION for (5-4.a) \rightarrow 5-4-1-0:impl.pdf
```

(5-4.b) \odot (10 min.) Discuss the asymptotic effort involved in the constructor and the evaluation operator of the class **LinearInterpolant** as the number n of data points becomes large.

```
Solution for (5-4.b) \rightarrow 5-4-2-0:scost.pdf
```

End Problem 5-4, 30 min.

Problem 5-5: Cardinal basis for trigonometric interpolation

In [Lecture \rightarrow Section 5.6] we learned about interpolation into the space \mathcal{P}_{2n}^T or trigonometric polynomials from [Lecture \rightarrow Def. 5.6.1.1]. In this task we investigate the cardinal basis functions for this interpolation operator and will also obtain a result about its stability.

Requires everything from [Lecture → Section 5.6] and involves coding in C++.

The (ordered) real trigonometric basis of \mathcal{P}_{2n}^T is

$$\mathcal{B}_{Re} := \left\{ C_0 := t \mapsto 1, S_1 := t \mapsto \sin(2\pi t), C_1 := t \mapsto \cos(2\pi t), S_2 := t \mapsto \sin(4\pi t), \qquad (5.5.1) \right.$$

$$C_2 := t \mapsto \cos(4\pi t), \dots, S_n := t \mapsto \sin(2n\pi t), C_n := t \mapsto \cos(2n\pi t) \right\}.$$

Another (ordered) basis of \mathcal{P}_{2n}^T is

$$\mathcal{B}_{\exp} := \{ B_k := t \mapsto \exp(2\pi k \iota t) : k = -n, \dots, n \} .$$
 (5.5.2)

(Strictly speaking, we have to take the real part of these functions.) Both bases have 2n + 1 elements, which agrees with the dimension of \mathcal{P}_{2n}^T [Lecture \rightarrow Cor. 5.6.1.6].

(5-5.a) \odot (15 min.) Give the matrix $\mathbf{S} \in \mathbb{C}^{2n+1,2n+1}$ that effects the transformation of a coefficient representation with respect to \mathcal{B}_{Re} into a coefficient representation with respect to \mathcal{B}_{exp} .

HIDDEN HINT 1 for (5-5.a) $\rightarrow 5-5-1-0$:trh1.pdf

SOLUTION for (5-5.a)
$$\rightarrow$$
 5-5-1-1:trisol1.pdf

(5-5.b) $oldsymbol{oldsymbol{arphi}}$ (10 min.) The shift operator $S_{ au}:C^0(\mathbb{R})\to C^0(\mathbb{R})$ acts on a function according to $S_{ au}f(t)=f(t- au)$, for $au\in\mathbb{R}$. If $\mathbf{c}\in\mathbb{C}^{2n+1}$ is the coefficient vector of $p\in\mathcal{P}_{2n}^T$ with respect to \mathcal{B}_{\exp} , what is the coefficient vector $\tilde{\mathbf{c}}\in\mathbb{C}^{2n+1}$ of $S_{ au}p$?

HIDDEN HINT 1 for (5-5.b) $\rightarrow 5-5-2-0$:tch2.pdf

SOLUTION for (5-5.b)
$$\rightarrow$$
 5-5-2-1:trisol2.pdf

Now we consider the set of interpolation nodes

$$\mathcal{T}_n := \left\{ t_j := \frac{j+1/2}{2n+1} : j = 0, \dots, 2n \right\}.$$
 (5.5.3)

The **cardinal basis functions** $b_0, \ldots, b_{2n} \in \mathcal{P}_{2n}^T$ of \mathcal{P}_{2n}^T with respect to \mathcal{T}_n are defined by

$$b_j(t_k) = \delta_{kj}$$
 [Kronecker symbol], $j, k = 0, ..., 2n$ (5.5.4)

(5-5.c) • (30 min.) Use the function trigpolyvalequid, see [Lecture \rightarrow Code 5.6.3.9] to plot the cardinal basis function (implement plot_basis in trigipcardfn.hpp) $t \mapsto b_0(t)$ in [0,1] for n=5. To that end evaluate $b_0(t)$ in 1000 equidistant points $\frac{k}{1000}$, $k=0,\ldots,999$.

Solution for (5-5.c)
$$\rightarrow$$
 5-5-3-0:plotsol.pdf

(5-5.d) • Show that

$$b_{k+1} = S_{\delta}b_k \quad \text{with} \quad \delta := \frac{1}{2n+1} \ .$$
 (5.5.6)

HIDDEN HINT 1 for (5-5.d) $\rightarrow 5-5-4-0$:h1.pdf

SOLUTION for (5-5.d)
$$\rightarrow$$
 5-5-4-1:trisol1.pdf

(5-5.e) \odot (15 min.) Describe the entries of the "interpolation matrix", that is, the system matrix of [Lecture \rightarrow Eq. (5.1.0.23)], for the trigonometric interpolation problem with node set \mathcal{T}_n and with respect to the basis \mathcal{B}_{exp} of \mathcal{P}_{2n}^T .

HIDDEN HINT 1 for (5-5.e) $\rightarrow 5-5-5-0$:tklh1.pdf

HIDDEN HINT 2 for (5-5.e) $\rightarrow 5-5-5-1$:tmp.pdf

SOLUTION for (5-5.e) \rightarrow 5-5-5-2:trisol1.pdf

What is the inverse of the interpolation matrix found in the previous problem?

HIDDEN HINT 1 for (5-5.f) $\rightarrow 5-5-6-0$:h1.pdf

SOLUTION for (5-5.f) $\rightarrow 5-5-6-1$:trisol1.pdf

HIDDEN HINT 1 for (5-5.g) $\rightarrow 5-5-7-0:h1.pdf$

SOLUTION for (5-5.g)
$$\rightarrow$$
 5-5-7-1:trisol1.pdf

(5-5.h) (30 min.) Write a function

double trigIpL(std::size_t n);

that approximately computes the **Lebesgue constant** $\lambda(n)$, see [Lecture \rightarrow § 5.2.4.8], for trigonometric interpolation based on the node set \mathcal{T}_n . The Lebesgue constant is available through the formula

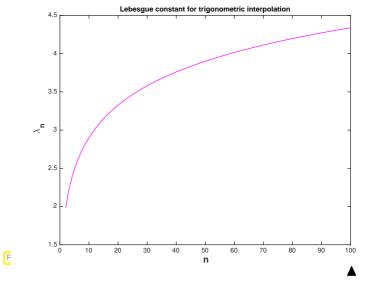
$$\lambda(n) = \sup_{\mathbf{y} \in \mathbb{C}^{2n+1} \setminus \{\mathbf{0}\}} \frac{\|\mathbf{I}_n \mathbf{y}\|_{L^{\infty}([0,1])}}{\|\mathbf{y}\|_{\infty}} = \sup_{t \in [0,1]} \sum_{k=0}^{2n} |b_k(t)|,$$
 (5.5.11)

where $I_n: \mathbb{C}^{2n+1} \to C^0(\mathbb{R})$ is the trigonometric interpolation operator.

Use sampling in 10^4 equidistant points to obtain a good guess for the supremum norm of a function on [0,1].

Note. Efficiency of the implementation need not be a concern in this numerical experiment.

SOLUTION for (5-5.h) $\rightarrow 5-5-8-0$:trisol1.pdf



Graph of $n \mapsto \lambda(n)$

(5-5.i) \odot (5 min.) The following table gives the values $\lambda(n)$ for $n=2^k, k=2,3,\ldots,8$.

2^k	$\lambda(2^k)$
4	2.7
8	3.07
16	3.46
32	3.89
64	4.32
128	4.75
256	5.18
512	5.62

From these numbers predict the asymptotic behavior of $\lambda(n)$ for $n \to \infty$. Is it $O(\log n)$, $O(n \log n)$, or $O(n^2)$?

SOLUTION for (5-5.i) $\rightarrow 5-5-9-0$:xtrsol1.pdf

End Problem 5-5, 165 min.

Problem 5-6: Quadratic Splines

[Lecture \to Def. 5.4.1.1] introduces spline spaces $\mathcal{S}_{d,\mathcal{M}}$ of any degree $d \in \mathbb{N}_0$ on a node set $\mathcal{M} \subset \mathbb{R}$. [Lecture \to Section 5.4.2] discusses interpolation by means of cubic splines, which is the most important case. In this problem we practise spline interpolation for quadratic splines in order to understand the general principles. We will see a situation, where knot sets and node sets have to be different.

Familiarity with [Lecture \rightarrow Section 5.4.2] is required. This exercise involves substantial implementation in C++ and EIGEN.

Consider a 1-periodic function $f: \mathbb{R} \to \mathbb{R}$, that is, f(t+1) = f(t) for all $t \in \mathbb{R}$, and a set of knots

$$\mathcal{M} := \{0 = t_0 < t_1 < t_2 < \dots < t_{n-1} < t_n = 1\} \subset [0,1]$$
.

We want to approximate f using a 1-periodic quadratic spline function $s \in \mathcal{S}_{2,\mathcal{M}}$, which interpolates f in the midpoints of the intervals $[t_{i-1}, t_i]$, $j = 0, \ldots, n$:

$$s\left(\frac{1}{2}(t_{j-1}+t_j)\right)=y_j:=f\left(\frac{1}{2}(t_{j-1}+t_j)\right), \quad j=1,\ldots,n.$$
 (5.6.1)

In analogy to the local representation of a cubic spline function according to [Lecture \rightarrow Eq. (5.4.2.6)], we locally in each knot interval parameterize a quadratic spline function $s \in \mathcal{S}_{2,\mathcal{M}}$ according to

$$s|_{[t_{j-1},t_j]}(t) = d_j \tau^2 + c_j 4 \tau (1-\tau) + d_{j-1} (1-\tau)^2, \quad \tau := \frac{t-t_{j-1}}{t_j-t_{j-1}}, \quad j=1,\ldots,n, \quad (5.6.2)$$

with $c_j, d_k \in \mathbb{R}$, $j = 1, \ldots, n$, $k = 0, \ldots, n$.

Remark. In this problem we see a case of spline interpolation for which the **knots** t_j , j = 0, ..., n, do not coincide with the interpolation **nodes**, which are the first coordinates of the data points to be interpolated.

HIDDEN HINT 1 for (5-6.a) $\rightarrow 5-6-1-0:s0h1.pdf$

Solution for (5-6.a)
$$\rightarrow 5-6-1-1:s0.pdf$$

(5-6.b) \odot (10 min.) What is the dimension of the subspace of *1-periodic* spline functions in $\mathcal{S}_{2.\mathcal{M}}$?

HIDDEN HINT 1 for (5-6.b) $\rightarrow 5-6-2-0:s1h1.pdf$

SOLUTION for (5-6.b)
$$\rightarrow 5-6-2-1:s1.pdf$$

What kind of continuity is already guaranteed by the mere use of the representation (5.6.2)?

SOLUTION for (5-6.c)
$$\rightarrow 5-6-3-0$$
:s2.pdf

HIDDEN HINT 1 for (5-6.d) $\rightarrow 5-6-4-0$:s3h1.pdf

Solution for (5-6.d)
$$\rightarrow$$
 5-6-4-1:s3.pdf

Show that the system matrix found in Sub-problem (5-6.d) can be written as a rank-1 modification of a tridiagonal matrix.

```
SOLUTION for (5-6.e) \rightarrow 5-6-5-0:s3a.pdf
```

```
(5-6.f) (15 min.) [depends on Sub-problem (5-6.d) and Sub-problem (5-6.e)]
```

Outline an efficient algorithm for the solution of the linear system of equations found in Subproblem (5-6.d). What is the asymptotic complexity of the algorithm in terms of the number n of data values for $n \to \infty$

```
SOLUTION for (5-6.f) \rightarrow 5-6-6-0:s6.pdf
```

```
(5-6.g) (25 min.) [depends on Sub-problem (5-6.d)]
```

Implement a C++ function

that returns the coefficient vector $[c_1,\ldots,c_n]^{\top}\in\mathbb{R}^n$ when supplied a sorted knot vector \mathbf{t} (of length n-1, because $t_0=0$ and $t_n=1$ will be taken for granted), an n-vector \mathbf{y} of data values y_j , $j=1,\ldots,n$, that specify the function values of the quadratic spline s at the midpoints of the knot intervals.

You can immediately use a suitable sparse direct solver provided by EIGEN, which is "smart" enough to solve the linear system (??) with optimal computational effort. You need not employ the techniques of [Lecture \rightarrow § 2.6.0.12].

```
SOLUTION for (5-6.g) \rightarrow 5-6-7-0:s3b.pdf
```

Realize a C++ function

that determines the coefficients d_j , $j=0,\ldots,n$, in the local representation (5.6.2) assuming that the values c_j , $j=1,\ldots,n$, have already been computed. The argument t is to pass the variable knot positions t_j , $j=1,\ldots,n-1$.

```
SOLUTION for (5-6.h) \rightarrow 5-6-8-0:3c.pdf
```

Based on compute_c () from Sub-problem (5-6.g) implement an efficient C++ function

which takes as input a (sorted) knot vector t (of length n-1, because $t_0=0$ and $t_n=1$ will be taken for granted), an n-vector y containing the values $y_j \in \mathbb{R}$ of a function f at the midpoints $\frac{1}{2}(t_{j-1}+t_j)$, $j=1,\ldots,n$, and a sorted N-vector x of evaluation points in [0,1].

The function is to return the values of the interpolating quadratic spline s at the positions x.

```
SOLUTION for (5-6.i) \rightarrow 5-6-9-0:s4.pdf
```

void plotquadspline(const std::string &filename);

that saves a plot (in <filename>.png) of the interpolating 1-periodic quadratic spline s for $f(t) := \exp(\sin(2\pi t))$, n=10, and $\mathcal{M} = \left\{\frac{j}{n}\right\}_{j=0}^{10}$, that is, the spline s is to fulfill $s(\frac{1}{2}(t_j+t_{j-1})) = f(\frac{1}{2}(t_j+t_{j-1}),\ j=1,\ldots,10.$ To that end evaluate the spline in 200 equidistant points in the interval [0,1] and, based on the obtained value, use the plot () function of MATPLOTLIBCPP to draw it. Do not forget to add axis labels "t" and "s(t)", and a suitable title.

SOLUTION for (5-6.j) $\rightarrow 5-6-10-0:s5.pdf$

(5-6.k) (30 min.) [depends on Sub-problem (5-6.i)]

Write a C++ function

std::vector<double> qsp_error(unsigned int q);

that approximately computes the $L^{\infty}([0,1])$ -norm of the error of the approximation of $f(t) := \exp(\sin(2\pi t))$ by its 1-periodic quadratic spline interpolant s_n based on the equidistant knot set $\mathcal{M}_n = \{j/n\}_{j=0}^n$:

$$E_n := \|f - s_n\|_{L^{\infty}([0,1])} := \sup_{x \in [0,1]} |f(x) - s_n(x)|.$$

To approximate the $L^{\infty}([0,1])$ -norm you can evaluate the difference $|f(x) - s_n(x)|$ at $N \gg n$ equidistant points in [0,1] and then take the maximum; choose N=10.000.

The function is supposed to return the error norms E_n for $n=2,4,8,\ldots,2^q$, where q is passed in the argument q. Describe quantitatively and qualitatively the convergence of the error norms in terms of $n\to\infty$.

HIDDEN HINT 1 for (5-6.k) $\rightarrow 5-6-11-0:s7h1.pdf$

SOLUTION for (5-6.k) \rightarrow 5-6-11-1:s7.pdf

End Problem 5-6, 200 min.

Problem 5-7: Linear monotonicity-preserving C^1 -interpolation scheme

In Sections [Lecture \rightarrow Section 5.3.2], [Lecture \rightarrow Section 5.3.2], [Lecture \rightarrow Section 5.3.3], different shape preserving operator are introduced. In the first case (piecewise linear functions) the interpolating function is continuous but not C^1 , in the two remaining cases (Hermite polynomials and splines) the interpolating operator is not linear. In all these examples the monotonicity is preserved, i.e., monotonic data yield monotonic interpolants. This problem shows that monotonicity preserving interpolators that are both C^1 and linear present an unpleasant and unexpected feature.

A purely theoretical problem. As preparation study [Lecture \rightarrow § 5.1.0.21] and [Lecture \rightarrow Def. 5.1.0.25] together with [Lecture \rightarrow Section 5.3.3.2].

The goal of this problem is to prove the following theorem.

Theorem 5.7.1. Linear monotonicity preserving interpolation operator

Given the node set $\mathcal{T} = \{t_0, \dots, t_n\} \subset [t_0, t_n]$, with $t_0 < t_1 < \dots < t_n$, consider an interpolation operator $I_{\mathcal{T}} : \mathbb{R}^{n+1} \to C^0([t_0, t_n])$, for which, by definition, we have $(I_{\mathcal{T}}(\mathbf{y}))(t_j) = y_j$ for every $j = 0, \dots, n$ and every vector of data values $\mathbf{y} = [y_0, \dots, y_n]^{\top} \in \mathbb{R}^{n+1}$. Assume that the operator satisfies

- i) smoothness: $I_{\mathcal{T}}(\mathbf{y})$ belongs to $C^1([t_0,t_n])$ for any $\mathbf{y} \in \mathbb{R}^{n+1}$,
- ii) linearity: $I_{\mathcal{T}}(\alpha \mathbf{y} + \beta \mathbf{z}) = \alpha I_{\mathcal{T}}(\mathbf{y}) + \beta I_{\mathcal{T}}(\mathbf{z})$, for every $\alpha, \beta \in \mathbb{R}$, $\mathbf{y}, \mathbf{z} \in \mathbb{R}^{n+1}$,
- iii) (positive) monotonicity preserving: if $y_j \leq y_{j+1}$, $\forall j = 0, ..., n-1$, then $I_{\mathcal{T}}(\mathbf{y})$ is a monotonic non-decreasing function.

Then, for any $\mathbf{y} \in \mathbb{R}^{n+1}$, the derivative of $I_{\mathcal{T}}(\mathbf{y})$ in the interior nodes vanishes:

$$(I_{\mathcal{T}}(\mathbf{y}))'(t_j) = 0$$
, $\forall j \in \{1, \dots, n-1\}$, $\forall \mathbf{y} \in \mathbb{R}^{n+1}$. (5.7.2)

(5-7.a) (15 min.) Show that the assertion (5.7.2) holds for the special data vectors

$$\mathbf{y} = \mathbf{y}^{(j)} := [\underbrace{0,\ldots,0}_{j \text{ times}},1,\ldots,1]^{ op}, \quad j=0,\ldots,n \;.$$

HIDDEN HINT 1 for (5-7.a) \rightarrow 5-7-1-0:s1h1.pdf

SOLUTION for (5-7.a) $\rightarrow 5-7-1-1:s1.pdf$

SOLUTION for (5-7.b) \rightarrow 5-7-2-0:s1a.pdf

(5-7.c) (15 min.) [depends on Sub-problem (5-7.a) and Sub-problem (5-7.b)]

Invoke the insight gained in Sub-problem (5-7.a) to complete the proof of the theorem.

SOLUTION for (5-7.c) \rightarrow 5-7-3-0:s2.pdf

End Problem 5-7, 40 min.

Problem 5-8: Approximating π by extrapolation

In [Lecture \to Section 5.2.3.3] we learned about the approximation of a limit $\lim_{h\to 0} \psi(h)$ based on the evaluation of the function ψ "far away" from 0. In this exercise we will practice the underlying technique of "Lagrangian polynomial extrapolation" for a geometric approximation problem.

Study [Lecture \rightarrow Section 5.2.3.3] and, in particular [Lecture \rightarrow Code 5.2.3.19], before your tackle this problem.

In this problem we encounter the situation that a quantity of interest is defined as a limit of a sequence

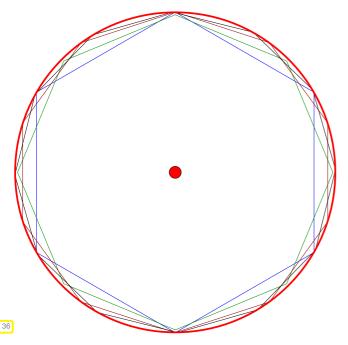
$$x^* = \lim_{n \to \infty} T(n) ,$$

where the function $T:\{n_{\min},n_{\min}+1,\ldots\}\mapsto\mathbb{R}$ may be given in procedural form as **double** \mathbb{T} (**int** n) only. However, invoking \mathbb{T} (\mathbb{I} nf) will usually result in an error and for very large arguments n the implementation of \mathbb{T} () may not yield reliable results, cf. [Lecture \to Ex. 1.5.4.7] and [Lecture \to § 5.2.3.16].

The idea of **extrapolation** is, firstly, to compute a few values $T(n_0), T(n_1), \ldots, T(n_k), k \in \mathbb{N}$, and to consider them as the values $g(1/n_0), g(1/n_1), \ldots, g(1/n_k)$ of a continuous function

$$g:(0,1/n_{\min}]\mapsto \mathbb{R}$$
, for which, obviously $x^*=\lim_{h\to 0}g(h)$.

Secondly, the function g is approximated by an interpolating polynomial $p_k \in \mathcal{P}_k$ with $p_k(n_j^{-1}) = T(n_j)$, $j = 0, \ldots, k$. In many cases we can expect that $p_k(0)$ will provide a good approximation for x^* .



The unit circle can be approximated by inscribed regular polygons with k edges. Thus, the length of half circumference (i.e., π) can be approximated by the half the length of the perimeters s_k of such polygons. These values $s_k/2$ can be calculated by elementary trigonometry and a closed form formula is

$$\frac{1}{2}s_k = k\sin(\pi/k) \ . \tag{5.8.1}$$

(5-8.a) (10 min.)

Show that for any $L \in \mathbb{N}$ and for $k \to \infty$

$$\frac{1}{2}s_k = \pi + \sum_{\ell=1}^{L} c_{\ell} k^{-2\ell} + O(k^{-2(L+1)}) \quad \text{for some} \quad c_{\ell} \in \mathbb{R} \quad \text{independent of} \quad L \; . \tag{5.8.2}$$

We say that s_k has an asymptotic expansion in k^{-2} .

Remark. The heuristics behind polynomial extrapolation in this example is that the existence of an asymptotic expansion (5.8.2) suggests that $\frac{1}{2}s_{\infty}$ can be well approximated by p(0), where p is an interpolating polynomial in k^{-1} or k^{-2} .

HIDDEN HINT 1 for (5-8.a) $\rightarrow 5-8-1-0$:s0h1.pdf

SOLUTION for (5-8.a)
$$\rightarrow 5-8-1-1:s0.pdf$$

Implement a C++ function

double extrapolate_to_pi(const unsigned int k);

that uses the *Aitken-Neville scheme*, see [Lecture \rightarrow Code 5.2.3.10], to approximate π by extrapolation from the data points $(j^{-1}, \frac{1}{2}s_i)$, for $j = 2, \dots, k$. Use the values $\frac{1}{2}s_i$ as given in (5.8.1).

Solution for (5-8.b)
$$\rightarrow 5-8-2-0:s1.pdf$$

(5-8.c) (30 min.) Write a C++ function

void plotExtrapolationError(const unsigned int kmax);

that generates a *linear-logarithmic* plot of the extrapolation errors

$$err(k) := |\pi - p_k(0)|, \quad k = 2, ..., k_{max},$$

versus k and also tabulates the errors. To create the plot use the functions of MATPLOTLIBCPP. Do not forget axis annotations and a meaningful title for the plot.

In main.cpp, plotExtrapolationError (10) is called for $k_{\text{max}} = 10$.

SOLUTION for (5-8.c)
$$\rightarrow 5-8-3-0$$
:s2.pdf

Which kind of convergence of $err(k) \to 0$ for $k \to \infty$ can you conclude from the data generated in Sub-problem (5-8.c)? The main types of asymptotic convergence to 0 are introduced in [Lecture \rightarrow Def. 6.2.2.7].

HIDDEN HINT 1 for (5-8.d) \rightarrow 5-8-4-0:s3h1.pdf

SOLUTION for (5-8.d)
$$\rightarrow 5-8-4-1:s3.pdf$$

(5-8.e) (15 min.) In main.cpp, plotExtrapolationError(30) is called which plots and tabulates err_k for k up to 30. Describe and explain your observations of the plot found in cx_out/pi_error_30.png.

Solution for (5-8.e)
$$\rightarrow 5-8-5-0$$
:s4.pdf

End Problem 5-8, 115 min.

Problem 5-9: Spaces of Piecewise Polynomial Functions

Piecewise polynomial functions, most prominently splines, are widely used for data interpolation and functions approximation. In this problem we examine some aspects of particular specimens.

A purely theoretical problem connected with [Lecture \rightarrow Section 5.4].

For a mesh $\mathcal{M} := \{a = t_0 < t_1 < \ldots < t_{n-1} < t_n = b\}$ we consider the following three spaces of piecewise polynomials:

$$V_{-1} := \left\{ v : [a, b] \to \mathbb{R} : v|_{[t_{i-1}, t_i]} \in \mathcal{P}_2, j = 1, \dots, n \right\}, \tag{5.9.1}$$

$$V_0 := \left\{ v \in \mathbf{C}^0([a,b]) : v|_{[t_{j-1},t_j[} \in \mathcal{P}_2, j = 1,\ldots,n] \right\}, \tag{5.9.2}$$

$$V_1 := \left\{ v \in \mathbf{C}^1([a,b]) : v|_{[t_{j-1},t_j[} \in \mathcal{P}_2, j = 1,\ldots,n] \right\}, \tag{5.9.3}$$

$$V_2 := \left\{ v \in \mathbb{C}^2([a,b]) : v|_{[t_{j-1},t_j[} \in \mathcal{P}_2, j = 1,\dots,n] \right\}.$$
 (5.9.4)

Here, \mathcal{P}_d stands for the space of polynomials $\mathbb{R} \to \mathbb{R}$ of degree $\leq d$.

(5-9.a) (15 min.) Determine the dimensions of the four spaces defined above:

$$\dim V_{-1} = igcup_{,}$$
 $\dim V_0 = igcup_{,}$
 $\dim V_1 = igcup_{,}$
 $\dim V_2 = igcup_{,}$

SOLUTION for (5-9.a)
$$\rightarrow 5-9-1-0:s1.pdf$$

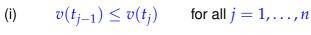
(5-9.b) \odot (10 min.) Which of the spaces V_k , k=-1,0,1, coincides with one of the spline spaces $\mathcal{S}_{d,\mathcal{M}}$ introduced in [Lecture \rightarrow Def. 5.4.1.1] in class?

$$V_{|}=\mathcal{S}_{|}$$
 .

SOLUTION for (5-9.b)
$$\rightarrow 5-9-2-0:s2.pdf$$

(5-9.c) \odot (15 min.) Which of the following conditions for a function $v \in V_1$ are *sufficient* and which are *necessary* for v being non-decreasing, that is, for

$$a \le x \le y \le b$$
 \Rightarrow $v(x) \le v(y)$?

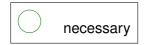




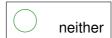


(ii)
$$v'(x) \ge 0$$
 for all $x \in [a, b]$

$$\text{(iii)} \quad v'(\tfrac{1}{2}(t_{j-1}+t_j)) \geq 0 \quad \text{for all } j \in \{1,\dots,n\}$$



(iv)
$$v'(t_j) \geq 0$$
 for all $j \in \{0,\ldots,n\}$



SOLUTION for (5-9.c)
$$\rightarrow 5-9-3-0:s3.pdf$$

In the sequel we consider the mesh (knot set) $\mathcal{M} = \{0, 1, 2\}$.

(5-9.d) \odot (10 min.) Determine the parameters $a, b \in \mathbb{R}$ in the following function definition

$$v(x) := \begin{cases} ax + b & \text{for } 0 \le x < 1 \text{ ,} \\ 1 - x^2 & \text{for } 1 \le x \le 2 \text{ ,} \end{cases}$$

such that $v \in V_1$.

$$a = \boxed{ }$$
 , $b = \boxed{ }$.

SOLUTION for (5-9.d)
$$\rightarrow 5-9-4-0:s4.pdf$$

(5-9.e) \odot (10 min.) Write down a $v \in V_1$ such that

$$v(0) = 0$$
 , $v(1) = 1$, $v(2) = 0$.

$$v(x) = \left\{ egin{array}{c} & ext{for } 0 \leq x < 1 \text{ ,} \\ & ext{for } 1 \leq x \leq 2 \text{ .} \end{array}
ight.$$

SOLUTION for (5-9.e)
$$\rightarrow 5-9-5-0:s5.pdf$$

End Problem 5-9, 60 min.

Problem 5-10: Adaptive Interpolation of Closed Curves

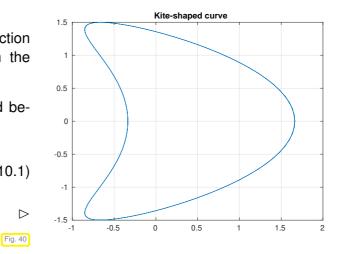
Often, for the sake of efficiency, closed curves (loops) need to be approximated by piecewise polynomials. This problem presents a simple adaptive algorithm for constructing such an approximation.

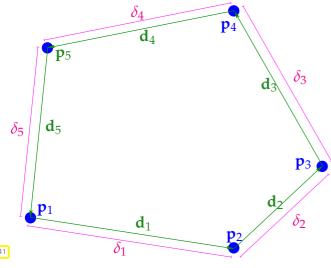
Related to [Lecture → Section 5.3.3]; the idea of adaptive approximation is borrowed from [Lecture \rightarrow Section 7.6].

A 1-periodic continuously differentiable function $\mathbf{c}: \mathbb{R} \to \mathbb{R}^2$ parameterizes a closed curve in the plane.

For instance, the "kite-shaped" curve displayed beside is described by

$$\mathbf{c}(t) = \begin{bmatrix} \cos(2\pi t) + \frac{2}{3}\cos(4\pi t) \\ \frac{3}{2}\sin(2\pi t) \end{bmatrix}. \quad (5.10.1)$$





Notations. Based on a sequence of $n \ge 2$ mutually distinct points $(p_1, p_2, ..., p_n)$ in the plane $(\mathbf{p}_k \in \mathbb{R}^2)$ the following notations are defined:

- $\delta_i := \|\mathbf{p}_{i+1} \mathbf{p}_i\|, i = 1, ..., n-1,$ $\delta_n := \|\mathbf{p}_1 - \mathbf{p}_n\|,$

- $\lambda_0 := 0$, $\lambda_k := \sum_{j=1}^k \delta_j, k = 1, ..., n$, $\mathbf{d}_i := \mathbf{p}_{i+1} \mathbf{p}_i, i = 1, ..., n 1$,
- $\mathbf{d}_{n} := \mathbf{p}_{1} \mathbf{p}_{n},$ $\mathbf{t}_{i} := \frac{\mathbf{d}_{i}}{\|\mathbf{d}_{i}\|}, i = 1, \dots, n.$

We study two ways to connect points in the plane by a piecewise polynomial closed curve.

Definition 5.10.2. Closed polygonal interpolant

Given a sequence $\Sigma:=(p_1,\ldots,p_n),\,n\in\mathbb{N},$ of mutually distinct points $\in\mathbb{R}^2$, their interpolating **closed polygon** Π_{Σ} is a function $\Pi_{\Sigma}:[0,\lambda_n]\to\mathbb{R}^2$ that fulfills:

- (i) $\Pi_{\Sigma}(\lambda_k) = \mathbf{p}_{k+1}, k = 1, \dots, n-1$, and $\Pi_{\Sigma}(0) = \Pi_{\Sigma}(\lambda_n) = \mathbf{p}_1$.
- (ii) $\Pi_{\Sigma}|_{[\lambda_{k-1},\lambda_k]}$, $k=1,\ldots,n$, is componentwise an affine linear function, $\Pi_{\Sigma}|_{[\lambda_{k-1},\lambda_k]} \in (\mathcal{P}_1)^2.$

Definition 5.10.3. Interpolating closed cubic Hermite curve

Given a sequence $\Sigma := (\mathbf{p}_1, \dots, \mathbf{p}_n)$, $n \in \mathbb{N}$, of mutually distinct points $\in \mathbb{R}^2$, their interpolating closed cubic Hermite curve \mathbf{H}_{Σ} is a function $\mathbf{H}_{\Sigma} : [0, \lambda_n] \to \mathbb{R}^2$ that satisfies:

```
(i) \mathbf{H}_{\Sigma}(\lambda_k) = \mathbf{p}_{k+1}, k = 1, \dots, n-1, and \mathbf{H}_{\Sigma}(0) = \mathbf{H}_{\Sigma}(\lambda_n) = \mathbf{p}_1.
```

(ii)
$$\mathbf{H}'(\lambda_k) = \frac{\mathbf{s}_k}{\|\mathbf{s}_k\|}$$
, $\mathbf{s}_k := \begin{cases} \frac{\delta_{k+1}}{\delta_k + \delta_{k+1}} \mathbf{t}_k + \frac{\delta_k}{\delta_k + \delta_{k+1}} \mathbf{t}_{k+1} & \text{, if } k = 1, \dots, n-1 \text{,} \\ \frac{\delta_n}{\delta_1 + \delta_n} \mathbf{t}_1 + \frac{\delta_1}{\delta_1 + \delta_n} \mathbf{t}_n & \text{, if } k = 0, n \text{.} \end{cases}$

 (\mathbf{H}') denotes the derivative of \mathbf{H} .

(iii) $\mathbf{H}_{\Sigma}|_{[\lambda_{k-1},\lambda_k]}$, $k=1,\ldots,n$, is componentwise a cubic polynomial, that is $\mathbf{H}_{\Sigma}|_{[\lambda_{k-1},\lambda_k]} \in (\mathcal{P}_3)^2$.

Code an efficient C++ function

```
std::vector<Eigen::Vector2d>
  closedPolygonalInterpolant(
    std::vector<Eigen::Vector2d> &Sigma,
    const Eigen::VectorXd &x);
```

that returns the sequence of coordinate vectors $\Pi_{\Sigma}(x_k\lambda_n)$, $k=1,\ldots,M,M$ the length of the vector \mathbf{x} with *sorted* entries $0 \leq x_1 < x_2 < \ldots x_M < 1$. Π_{Σ} is the interpolating closed polygon for a set Σ of points passed through the vector argument Sigma.

```
Solution for (5-10.a) \rightarrow 5-10-1-0:s1.pdf
```

(5-10.b) (20 min.) [depends on Sub-problem (5-10.a)]

Implement an efficient C++ function

```
std::vector<Eigen::Vector2d>
  closedHermiteInterpolant(
    std::vector<Eigen::Vector2d> &Sigma,
    const Eigen::VectorXd &x);
```

that returns the M vectors $\mathbf{H}_{\Sigma}(x_k\lambda_n) \in \mathbb{R}^2$, $k=1,\ldots,M$, M the length of the vector \mathbf{x} with *sorted* entries $0 \le x_1 < x_2 < \cdots < x_M \le 1$. In this case \mathbf{H}_{Σ} is the interpolating closed cubic Hermite curve according to Def. 5.10.3 for the points passed in the vector Sigma.

You may use the auxiliary function

```
double hermloceval(double t,
    double t1, double t2,
    double y1, double y2,
    double c1, double c2);
```

provided in the file hermloceval.hpp, which implements the formula [Lecture \rightarrow Eq. (5.3.3.4)] for evaluating a cubic Hermite interpolating polynomial, see also [Lecture \rightarrow Code 5.3.3.6].

```
HIDDEN HINT 1 for (5-10.b) \rightarrow 5-10-2-0:h2.pdf
```

```
SOLUTION for (5-10.b) \rightarrow 5-10-2-1:s1.pdf
```

Now, given a 1-periodic function $\mathbf{c}: \mathbb{R} \to \mathbb{R}^2$ providing the parameterization of a closed curve, we pursue the following policy aiming for its adaptive piecewise polynomial approximation.

Pseudocode 5.10.6: Algorithm for adaptive approximation of closed curve

```
Given: tol > 0, n_{\min} \in \mathbb{N} \setminus \{1\}.
<sup>2</sup> \mathcal{M}:=\{t_0:=0,t_1:=rac{1}{n_{\min}},\ldots,t_{n-1}:=rac{n_{\min}-1}{n_{\min}}\} // An ordered set
         n := \sharp \mathcal{M}; // Notation: \mathcal{M} := \{t_0, \ldots, t_{n-1}\}
         \Sigma := (\mathbf{c}(t))_{t \in \mathcal{M}}; // sequence of points; also defines \lambda_k
         // \Pi_{\Sigma}/H_{\Sigma} according to Def. 5.10.2, Def. 5.10.3
        \sigma_k := \left\| (\mathbf{H}_{\Sigma} - \mathbf{\Pi}_{\Sigma}) (\frac{1}{2}(\lambda_{k-1} + \lambda_k)) \right\|, \quad k = 1, \dots, n;
        \alpha := \frac{1}{n} \sum_{k=1}^{n} \sigma_k;
8
         for j = 1, ..., n do {
               if (\sigma_j > 0.9\alpha) \mathcal{M} \leftarrow \mathcal{M} \cup \{\frac{1}{2}(t_{j-1} + t_j)\}; // Convention t_n := 1
               // \mathcal M is always assumed to be sorted!
11
         }
12
    }
13
    until ( max \sigma_k \leq \text{tol});
```

(5-10.c) (30 min.) depends on Sub-problem (5-10.a) and Sub-problem (5-10.b)

Based on the C++ functions closedPolygonalInterpolant() and closedHermiteInterpolant() implement a C++ function

that performs the interpolation of a closed curve using the above algorithm from Code 5.10.6. The 1-periodic continuous function $\mathbf{c}: \mathbb{R} \to \mathbb{R}$ is passed through a functor object \mathbf{c} , which is to feature an evaluation operator

```
Eigen::Vector2d operator () (double t);
```

The argument tol supplies the tolerance for the termination of the adaptive algorithm. The function is to return

- 1. the vectors $\mathbf{H}_{\Sigma}(x_k L)$, $k=1,\ldots,M$, M the length of the vector \mathbf{x} with *sorted* entries $0 \leq x_1 < x_2 < \cdots < x_M \leq 1$. In this case \mathbf{H}_{Σ} is the interpolating closed cubic Hermite curve according to Def. 5.10.3 for the set Σ at the termination of the algorithm,
- 2. the coordinates of the points in the sequence Σ as a second sequence of 2-vectors.

void plotKite(const char *filename, double tol = 1.0e-3);

```
SOLUTION for (5-10.c) \rightarrow 5-10-3-0:s3.pdf  
(5-10.d) \odot (30 min.) [depends on Sub-problem (5-10.c)]

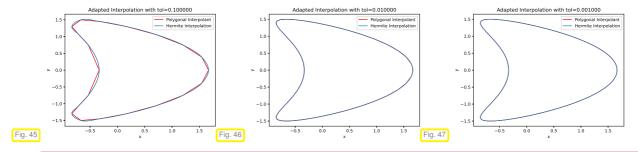
Using MATPLOTLIBCPP and based on your implementation of adaptedHermiteInterpolant(), write a C++ function
```

that, for the kite-shaped curve $\mathbf{c}:[0,1]\to\mathbb{R}^2$ as defined in (5.10.1), and with the choice $n_{\min}=6$, generates a rendering of the closed polygonal interpolant (\to Def. 5.10.2) and the interpolating closed cubic Hermite curve (\to Def. 5.10.3) produced by a call to the function adaptedHermiteInterpolant () with the given tolerance tol. The plots should be created using N:=10000 equidistant sampling points $x_k:=\frac{k}{N}, k=0,\ldots,N-1$.

The graphics is to be saved in PNG format in a file <filename>.png.

SOLUTION for (5-10.d)
$$\rightarrow 5-10-4-0:s4.pdf$$

Figures of kite-shaped curves produced by the code above code:



End Problem 5-10, 100 min.

Problem 5-11: A cubic spline

Since they mimic the behavior of an elastic rod pinned at fixed points, see [Lecture \rightarrow § 5.4.3.1], cubic splines are very popular for creating "aesthetically pleasing" interpolating functions. However, in this problem we look at a cubic spline from the perspective of its defining properties, see [Lecture \rightarrow Def. 5.4.1.1], in order to become more familiar with the concept of spline function and the consequences of the smoothness required by the definition.

Related to [Lecture \rightarrow Section 5.4.1]

(5-11.a) \Box (15 min.) For parameters $\alpha, \beta, \gamma, \delta, \sigma, \rho \in \mathbb{R}$ we define the function

$$s(t) = \begin{cases} -1 & \text{for } -1 \le t < 0 \text{ ,} \\ \alpha t^3 + \beta t^2 + \gamma t + \delta & \text{for } 0 \le t < 1 \text{ ,} \\ -t^3 + \sigma t^2 + \rho t + 2 & \text{for } 1 \le t \le 2 \text{ .} \end{cases}$$

Determine the parameters $\alpha, \beta, \gamma, \delta, \sigma, \rho \in \mathbb{R}$ such that s is a cubic spline with respect to the knot set $\mathcal{M} := \{-1, 0, 1, 2\}$.

SOLUTION for (5-11.a) $\rightarrow 5-11-1-0:sx1.pdf$

(5-11.b) \square (15 min.) For parameters $\alpha, \beta, \gamma, \delta, \sigma, \rho \in \mathbb{R}$ we define the function

$$s(t) = \begin{cases} 2 & \text{for } -1 \le t < 0 \text{ ,} \\ \sigma t^3 + \rho t^2 + \gamma t + \delta & \text{for } 0 \le t < 1 \text{ ,} \\ -t^3 + \alpha t^2 + \beta t - 1 & \text{for } 1 \le t \le 2 \text{ .} \end{cases}$$

Determine the parameters $\alpha, \beta, \gamma, \delta, \sigma, \rho \in \mathbb{R}$ such that s is a cubic spline with respect to the knot set $\mathcal{M} := \{-1, 0, 1, 2\}$.

SOLUTION for (5-11.b) $\rightarrow 5-11-2-0:sx2.pdf$

(5-11.c) \square (15 min.) For parameters $\alpha, \beta, \gamma, \delta, \sigma, \rho \in \mathbb{R}$ we define the function

$$s(t) = \begin{cases} 1 & \text{for } -1 \le t < 0 \text{,} \\ \alpha t^3 + \beta t^2 + \gamma t + \delta & \text{for } 0 \le t < 2 \text{,} \\ -t^3 + \rho t^2 + \sigma t - 7 & \text{for } 2 \le t \le 3 \text{.} \end{cases}$$

Determine the parameters $\alpha, \beta, \gamma, \delta, \sigma, \rho \in \mathbb{R}$ such that s is a cubic spline with respect to the knot set $\mathcal{M} := \{-1, 0, 2, 3\}$.

SOLUTION for (5-11.c) $\rightarrow 5-11-3-0:sx3.pdf$

End Problem 5-11, 45 min.

Problem 5-12: Not-a-knot Cubic Splines

As we have seen in [Lecture \rightarrow Section 5.4.2] cubic splines can be used for one-dimensional data interpolation, which is a popular technique in computer aided geometric design (CAGD). Unfortunately, extra conditions have to be imposed to render the cubic-spline interpolant unique as discussed in [Lecture \rightarrow § 5.4.2.11]. Not-a-knot cubic splines are another way to arrive at a unique cubic-spline interpolant.

This problem builds on [Lecture \rightarrow Section 5.4.2] and contents of [Lecture \rightarrow § 5.4.2.5] should be studied in detail before tackling this problem.

Recall what is a cubic spline interpolant,

Definition [Lecture \rightarrow Def. 5.4.2.3]. Cubic-spline interpolant

Given a node set/knot set $\mathcal{M} := \{t_0 < t_1 < \cdots < t_n\}, n \in \mathbb{N}$, and data values $y_0, \ldots, y_n \in \mathbb{R}$, an associated **cubic spline interpolant** is a function $s \in \mathcal{S}_{3,\mathcal{M}}$ that complies with the interpolation conditions

$$s(t_i) = y_i$$
 , $j = 0, ..., n$. (5.12.1)

and, more fundamentally, the definition of spline functions as special piecewise polynomial functions:

Definition [Lecture → Def. 5.4.1.1]. Splines

Given an interval $I := [a, b] \subset \mathbb{R}$ and a knot sequence $\mathcal{M} := \{a = t_0 < t_1 < \ldots < t_{n-1} < t_n = b\}$, $n \in \mathbb{N}$, the vector space $\mathcal{S}_{d,\mathcal{M}}$ of the spline functions of degree d (or order d+1), $d \in \mathbb{N}_0$, is defined by

$$S_{d,\mathcal{M}} := \{ s \in C^{d-1}(I) : s_j := s_{|[t_{i-1},t_i]} \in \mathcal{P}_d \ \forall j = 1,\ldots,n \} \ .$$

Now we study a specialization of [Lecture \rightarrow Def. 5.4.2.3]:

Definition 5.12.2. Not-a-knot cubic spline

Given a node set/knot set $\mathcal{M} := \{t_0 < t_1 < \dots < t_n\}, n \in \mathbb{N}, n \geq 2$, and data values $y_0, \dots, y_n \in \mathbb{R}$, an associated **not-a-knot cubic spline interpolant** is a cubic spline interpolant s for the data points (t_i, y_i) , $i = 0, \dots, n$, that satisfies the extra conditions that the *third derivative* of s is *continuous at* t_1 *and* t_{n-1} .

The following theorem of De Boor confirms that the not-a-knot constructions resolves the non-uniqueness of cubic spline interpolants.

Theorem 5.12.3. Existence and uniqueness of not-a-knot cubic spline interpolant

For any $n \in \mathbb{N}$, n > 2, and any node set/knot set $\mathcal{M} := \{t_0 < t_1 < \cdots < t_n\}$ and data values $y_i \in \mathbb{R}$, $i = 0, \ldots, n$, there exists a unique not-a-knot cubic spline interpolant according to Def. 5.12.2.

(5-12.a) \odot (10 min.) Why is the assumption n > 2 necessary in the statement of Thm. 5.12.3?

SOLUTION for (5-12.a) $\rightarrow 5-12-1-0$: knks1.pdf

You are supposed to implement a C++ class

```
class NotAKnotCubicSpline {
   public:
   NotAKnotCubicSpline(Eigen::VectorXd t,
    Eigen::VectorXd y);
   ~NotAKnotCubicSpline() = default;
   [[nodiscard]] double eval(double t) const;
   private:
   [[nodiscard]] int getPtIdx(double t) const;
   Eigen::VectorXd t_; // knot sequence
   Eigen::VectorXd y_; // data values
   Eigen::VectorXd c_; // slopes
};
```

The public member function eval () is supposed to evaluate the not-a-knot cubic spline interpolant at a single point $t \in [t_0, t_n]$.

Fortunately, you found a complete implementation of a class NaturalCubicSpline with a 100% analogous definition for natural cubic spline interpolation. You stored it in the files notaknotcubicspline.hpp and ran a text replacement NaturalCubicSpline \rightarrow NotAKnotCubicSpline on it. And this is where you are standing.

(5-12.b) \odot (30 min.) Each of the continuity requirements for the third derivative of the cubic spline interpolant at the knots t_1 and t_{n-1} (assume n > 2, see Def. 5.12.2 for notations) induces one linear relationship between the components of the vectors y_{-} and c_{-} . Derive and state those relationships.

```
SOLUTION for (5-12.b) \rightarrow 5-12-2-0:.pdf 

(5-12.c) (30 min.) [depends on Sub-problem (5-12.b)]
```

In the file notaknotcubicspline.hpp write efficient code for the constructor

that is supposed to initialize the data members y and c in a way that does not break the already implemented eval () member function.

```
HIDDEN HINT 1 for (5-12.c) \rightarrow 5-12-3-0:knk2h1.pdf
SOLUTION for (5-12.c) \rightarrow 5-12-3-1:.pdf
```

```
double evalderivative(double t) const;
```

that returns the value of the first derivative of the not-a-knot cubic spline interpolant at a single point $t \in]t_0, t_n[$.

Add this member function to the class definition and provide an implementation in the file notaknotcubicspline.hpp.

```
SOLUTION for (5-12.d) \rightarrow 5-12-4-0:s3ss.pdf
```

End Problem 5-12, 90 min.

Problem 5-13: Various Aspects of Interpolation by Global Polynomials

[Lecture \rightarrow Chapter 5] teaches the mathematical foundations and algorithmic realizations of interpolation by means of global polynomials. This problems reviews some of these aspects.

This is a purely theoretical problem connected with [Lecture \rightarrow Section 5.2]. It should be solved without resorting to the lecture document or tablet notes.

Throughout this problem we write $\mathcal{P}_d(\mathbb{R})$ for the space of uni-variate polynomials of degree $\leq d$:

$$\mathcal{P}_k := \{t \mapsto \alpha_k t^k + \alpha_{k-1} t^{k-1} + \dots + \alpha_1 t + \alpha_0 \cdot 1, \alpha_i \in \mathbb{R}\}$$
. [Lecture \to Eq. (5.2.1.1)]

(5-13.a) \Box (15 min.) What are the dimensions of the following subspaces of $\mathcal{P}_d(\mathbb{R}), d \in \mathbb{N}$?

(i)
$$V_1 := \{ p \in \mathcal{P}_d(\mathbb{R}) : \int_{-1}^1 p(t) \, \mathrm{d}t = 0 \}$$
: $\dim V_1 = 0$.

(ii)
$$V_2 := \{ p \in \mathcal{P}_d(\mathbb{R}) : \ p(1) = p(-1) = 0 \}$$
: $\dim V_2 = \left| \right|$

Here $p^{(3)}$ stands for the third derivative of p.

SOLUTION for (5-13.a) $\rightarrow 5-13-1-0$:spi1.pdf

(5-13.b) \bigcirc (10 min.) Let $\{t_0, t_1, \dots, t_n\} \subset \mathbb{R}, n \in \mathbb{N}$. The following sets of functions provide bases for \mathcal{P}_d :

$$\mathfrak{B}_1 := \left\{ t \mapsto \prod_{\substack{j=0 \ j \neq i}}^n \frac{t-t_j}{t_i-t_j}, \quad i=0,\ldots,n \right\},$$

$$\mathfrak{B}_2 := \left\{ t \mapsto \prod_{j=0}^{i-1} (t-t_j), \quad i=0,\ldots,n \right\},$$

$$\mathfrak{B}_3:=\left\{t\mapsto t^i\,,\ i=0,\ldots,n\right\}.$$

Note that an empty product is defined to be $\equiv 1$.

Put the right letter in the box indicating the name of the basis:

 $\mathbf{A} \stackrel{.}{=} \text{monomial basis} \quad , \quad \mathbf{B} \stackrel{.}{=} \text{Newton basis} \quad , \quad \mathbf{C} \stackrel{.}{=} \text{Lagrangian basis} \quad .$

SOLUTION for (5-13.b) $\rightarrow 5-13-2-0$: spip2.pdf

(5-13.c) Given data points (t_i, y_i) , i = 0, ..., n, $n \in \mathbb{N}$, we write $p_{k,\ell}$, $0 \le k \le \ell \le n$, for the polynomial $p_{k,\ell} \in \mathcal{P}_{\ell-k}$ interpolating through $(t_i, y_i)_{i=k}^{\ell}$: $p_{k,\ell}(t_i) = y_i$, $i = k, ..., \ell$. Supplement the missing parts of the following recursion:

$$p_{k,\ell}(t) = rac{igg(igg) p_{k+1,\ell}(t) + igg(igg) p_{k,\ell-1}(t)}{igg(igg)} \ , \quad t \in \mathbb{R} \ , \quad 0 \le k < \ell \le n \ .$$

SOLUTION for (5-13.c) \rightarrow 5-13-3-0:pips3.pdf

End Problem 5-13, 35 min.

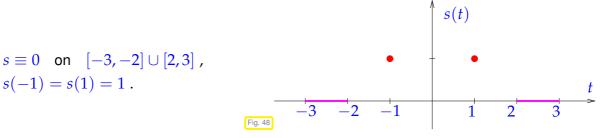
Problem 5-14: Local Representations of Splines

Locally on knot intervals splines coincide with polynomials of a fixed degree. In this problem we compute those local polynomial representations in special cases.

This is a purely theoretical problem based on [Lecture \rightarrow Section 5.4.1] and [Lecture \rightarrow Section 5.4.2].

(5-14.a) □ (6 min.)

Let s denote a *quadratic spline*, $s \in S_{2,\mathcal{M}}$ with respect to the knot set $\mathcal{M} := \{-3, -2, -1, 1, 2, 3\}$ and satisfying



Determine the unknown real coefficients in the local representations

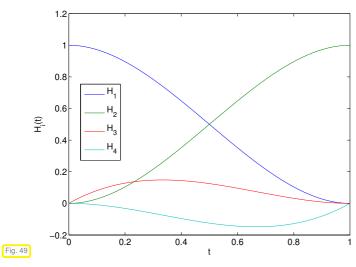
$$s(t)=$$
 t^2+ for $t\in [1,2]$, $s(t)=$ t^2+ $t+$ for $t\in [-1,1]$.

SOLUTION for (5-14.a) $\rightarrow 5-14-1-0$:spls1.pdf

(5-14.b) $\ \ \, \ \ \,$ (12 min.) Let the knot set $\mathcal{M}:=\{t_0 < t_1 < t_2 < \cdots < t_n\} \subset \mathbb{R}$ be given. The local representation of a cubic spline $s \in \mathcal{S}_{3,\mathcal{M}}$ is given by $(h_j:=t_j-t_{j-1},j\in\{1,\ldots,n\})$

$$s|_{[t_{j-1},t_j]}(t) = \alpha_j H_1(\tau) + \beta_j H_2(\tau) + h_j \gamma_j H_3(\tau) + h_j \delta_j H_4(\tau) , \quad \tau := \frac{t - t_{j-1}}{h_j} , \quad \alpha_j, \beta_j, \gamma_j, \delta_j \in \mathbb{R} ,$$

where H_1, H_2, H_3, H_4 are the cardinal basis functions for Hermite interpolation on [0, 1],



$$H_1(\tau) = 1 - 3\tau^2 + 2\tau^3$$
,
 $H_2(\tau) = 3\tau^2 - 2\tau^3$,
 $H_3(\tau) = \tau - 2\tau^2 + \tau^3$,
 $H_4(\tau) = -\tau^2 + \tau^3$.

Inserting the midpoints of knot intervals gives us the extended knot set

$$\widetilde{\mathcal{M}} := \left\{ \widetilde{t}_0 < \widetilde{t}_1 < \widetilde{t}_2 < \dots < \widetilde{t}_{2n} \right\}, \quad \widetilde{t}_{2j} := t_j, \qquad j \in \{0, \dots, n\}, \\ \widetilde{t}_{2j-1} := \frac{1}{2}(t_j + t_{j-1}), \quad j \in \{1, \dots, n\}.$$

On the knot intervals of $\widetilde{\mathcal{M}}$ the spline s has the local representation $(\widetilde{h}_\ell := \widetilde{t}_\ell - \widetilde{t}_{\ell-1}, \, \ell \in \{1, \dots, 2n\})$

$$s|_{[\widetilde{t}_{\ell-1},\widetilde{t}_\ell]}(t) = \widetilde{\alpha}_\ell H_1(\tau) + \widetilde{\beta}_\ell H_2(\tau) + \widetilde{h}_\ell \widetilde{\gamma}_\ell H_3(\tau) + \widetilde{h}_\ell \widetilde{\delta}_\ell H_4(\tau) , \ \tau := \frac{t - \widetilde{t}_{\ell-1}}{\widetilde{h}_\ell} , \ \widetilde{\alpha}_\ell, \widetilde{\beta}_\ell, \widetilde{\gamma}_\ell, \widetilde{\delta}_\ell \in \mathbb{R} .$$

Express the coefficients $\widetilde{\alpha}_{\ell}$, $\widetilde{\beta}_{\ell}$, $\widetilde{\gamma}_{\ell}$, $\widetilde{\delta}_{\ell}$ in terms of α_{j} , β_{j} , γ_{j} , δ_{j} :



SOLUTION for (5-14.b) $\rightarrow 5-14-2-0$:splss2.pdf

End Problem 5-14, 18 min.

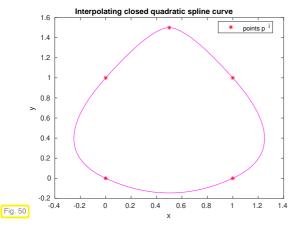
Problem 5-15: Closed Curve Interpolation with Quadratic Splines

The reconstruction of curves and surfaces from given points is a central task in computer-aided design (CAD). Here we focus on the problem of finding a closed and C^1 -smooth curve running through a sequence of given points. That curve should be modeled by means of a periodic quadratic spline function.

This problem relies on [Lecture \rightarrow Section 5.4] and has a focus on implementation in C++ using EIGEN.

We are given $n, n \in \mathbb{N}$, mutually different points $\mathbf{p}^1, \dots, \mathbf{p}^n \in \mathbb{R}^2$ in the plane. Defining a knot sequence $\mathcal{M} := \{0 = t_0 < t_1 < \dots < t_{n-1} < t_n = 1\}$ as ¹

$$t_{j} := \frac{\sum_{k=1}^{j} \|\mathbf{p}^{k} - \mathbf{p}^{k-1}\|_{2}}{\sum_{k=1}^{n} \|\mathbf{p}^{k} - \mathbf{p}^{k-1}\|_{2}} \quad (\mathbf{p}^{0} := \mathbf{p}^{n}), \quad j = 0, \dots, n,$$
 (5.15.1)



we want to find a spline curve $\mathbf{s}:[0,1]\to\mathbb{R}^2$ satisfying

$$\mathbf{s} \in (\mathcal{S}_{2,\mathcal{M}})^2$$
 , c.f. [Lecture \rightarrow Def. 5.4.1.1], (5.15.2)

$$\mathbf{s}(t_i) = \mathbf{p}^j$$
 for all $j = 1, ..., n$, (5.15.3)

$$\mathbf{s}(t_0) = \mathbf{p}^n$$
 and $\mathbf{s}'(t_0) = \mathbf{s}'(t_n)$, (5.15.4)

where ' indicates differentiation with respect to the curve parameter. The requirement (5.15.4) ensures that \mathbf{s} describes a *closed* C^1 -*smooth* curve.

Throughout, for **S** use the following local representation on knot intervals:

$$\mathbf{s}|_{]t_{j-1},t_{j}[}(t) = \mathbf{y}^{j-1}(1-\tau) + \mathbf{x}^{j}h_{j}\tau(1-\tau) + \mathbf{y}^{j}\tau , \quad \tau := \frac{t-t_{j-1}}{h_{i}}, \quad \frac{t_{j-1} < t < t_{j}}{h_{i} := t_{j} - t_{j-1}}, \quad (5.15.5)$$

with unknown vector-valued coefficients $\mathbf{y}^j \in \mathbb{R}^2$, $j = 0, \dots, n$, and $\mathbf{x}^j \in \mathbb{R}^2$, $j = 1, \dots, n$.

(5-15.a) (6 min.) What are the benefits of using the local representation (5.15.5)?

SOLUTION for (5-15.a)
$$\rightarrow 5-15-1-0$$
:pqss1.pdf

(5-15.b) (30 min.)

In explicit form state the linear system of equations

- whose coefficient matrix and right-hand side vector are to be given in terms of the point locations $\mathbf{p}^j \in \mathbb{R}^2$, j = 1, ..., n and the knot positions t_k , k = 0, ..., n, and
- whose solution provides the vector $[(\mathbf{x}^1)_1, (\mathbf{x}^2)_1, \dots, (\mathbf{x}^n)_1]^{\top} \in \mathbb{R}^n$ of first components of the vector-valued coefficients $\mathbf{x}^j \in \mathbb{R}^2$, $j = 1, \dots, n$, in (5.15.5).

¹Sums whose lower summation bound exceeds the upper are supposed to evaluate to zero.

```
\begin{bmatrix} (\mathbf{x}^1)_1 \\ \vdots \end{bmatrix} = \begin{bmatrix} (\mathbf{x}^n)_1 \\ \vdots \\ (\mathbf{x}^n)_1 \end{bmatrix}
```

SOLUTION for (5-15.b) $\rightarrow 5-15-2-0$:pqss2.pdf

In an object-oriented C++ code, for handling interpolating closed C^1 -smooth curves as specified above we rely on the following data type:

```
C++ code 5.15.20: Class definition of ClosedQuadraticSplineCurve
   class ClosedQuadraticSplineCurve {
    public:
3
    // Constructor
     explicit ClosedQuadraticSplineCurve(const Eigen::Matrix2Xd &p);
5
     ~ClosedQuadraticSplineCurve() = default;
6
     ClosedQuadraticSplineCurve() = delete;
     ClosedQuadraticSplineCurve (const ClosedQuadraticSplineCurve &) = delete;
     ClosedQuadraticSplineCurve (const ClosedQuadraticSplineCurve &&) = delete;
     ClosedQuadraticSplineCurve & operator = (const ClosedQuadraticSplineCurve &) =
10
         delete;
11
     // Point evaluation operator for sorted parameter arguments
12
13
     [[nodiscard]] Eigen::Matrix2Xd curve_points(const Eigen::VectorXd &v) const;
     // Curvature evaluation operator for sorted parameter arguments
14
     [[nodiscard]] Eigen::VectorXd local_curvatures(
15
         const Eigen::VectorXd &v) const;
16
     // Approximate computation of the length of the curve
17
     [[nodiscard]] double length(double rtol = 1E-6) const;
18
19
    private:
     [[nodiscard]] bool checkC1(double tol = 1.0E-4) const;
21
     // Number of points to be interpolated
22
23
    unsigned int n_;
    // Coordinates of interpolated points, duplicate: \mathbf{p}^0 = \mathbf{p}^n
24
     Eigen::Matrix2Xd p_;
25
     // Knot sequence
     Eigen::VectorXd t_;
27
     // Coefficients in local representation
28
     Eigen::Matrix2Xd x ;
29
   };
```

Get it on ₩ GitLab (periodicquadraticsplines.hpp).

(5-15.c) **(20 min.)** [depends on Sub-problem (5-15.b)]

In the file periodicquadraticsplines.hpp provide an *efficient* implementation of the constructor

```
ClosedQuadraticSplineCurve::ClosedQuadraticSplineCurve(
const Eigen::Matrix<double, 2, Eigen::Dynamic> &p);
```

which is supposed to initialize the member variable t_{-} with the elements of the knot sequence \mathcal{M} and the columns of x_{-} with the vector coefficients \mathbf{x}^{j} , $j=1,\ldots,n$. The matrix argument p passes the coordinate vectors \mathbf{p}^{j} , $j=1,\ldots,n$, in its columns.

HIDDEN HINT 1 for (5-15.c) \rightarrow 5-15-3-0:pqsh3.pdf

```
SOLUTION for (5-15.c) \rightarrow 5-15-3-1:.pdf
```

(5-15.d) \odot (15 min.) Complete the code in periodicquadraticsplines.hpp to achieve an *efficient* implementation, that is, with cost scaling linearly in the number of data and in the number n of points, of the member function

```
Eigen::Matrix<double, 2, Eigen::Dynamic>
ClosedQuadraticSplineCurve::curve_points(const Eigen::VectorXd &v)
    const;
```

that takes a *sorted* sequence $(0 \le v_1 \le v_2 \le \cdots \le v_{N-1} \le v_N \le 1)$, $N \in \mathbb{N}$, of parameter values as v and returns the sequence $(\mathbf{s}(v_1), \mathbf{s}(v_2), \ldots, \mathbf{s}(v_N))$ as the columns of a $2 \times N$ -matrix.

```
SOLUTION for (5-15.d) \rightarrow 5-15-4-0:pqss3a.pdf
```

(5-15.e) (20 min.) Devise an *efficient* implementation of the member function

```
Eigen::VectorXd ClosedQuadraticSplineCurve::local_curvatures(
const Eigen::VectorXd &v) const;
```

that takes a *sorted* sequence $(0 \le v_1 < v_2 < \cdots < v_{N-1} < v_N \le 1)$, $N \in \mathbb{N}$, of parameter values as \forall and returns a vector of local curvatures of the curve \mathbf{s} .

The **curvature** of $\mathbf{s}:[0,1]\to\mathbb{R}^2$ is a function $\kappa:[0,1]\to\mathbb{R}$ given by

$$\kappa(t) := \frac{\det \left[\mathbf{s}'(t), \mathbf{s}''_j(t) \right]}{\| \mathbf{s}'(t) \|_2}, \quad \mathbf{s}_j := \left. \mathbf{s} \right|_{]t_{j-1}, t_j]}, \quad t_{j-1} < t \le t_j, \quad j = 1, \ldots, n.$$

HIDDEN HINT 1 for (5-15.e) $\rightarrow 5-15-5-0$:pqsh4.pdf

```
Solution for (5-15.e) \rightarrow 5-15-5-1: .pdf
```

(5-15.f) oxdots Based on approximation of the curve $t \mapsto \mathbf{s}(t)$ by means of polygons, implement the following member function of **ClosedQuadraticSplineCurve**:

```
double ClosedQuadraticSplineCurve::length(double rtol) const;
```

It is supposed to return an approximation of the length of the curve with a relative error bounded by rtol.

```
SOLUTION for (5-15.f) \rightarrow 5-15-6-0:pgss5.pdf
```

End Problem 5-15, 91 min.

Problem 5-16: Some Aspects of Interpolation by Global Polynomials

Interpolation by global polynomials is a fundamental numerical technique that forms the basis of many other algorithms. Therefore its various features and related concepts are must-knows for everybody involved in computational science and engineering.

This problem is related to [Lecture \rightarrow Section 5.2].

$$L_i(t_j) = \delta_{ij} := \begin{cases} 1 & \text{if } i = j, \\ 0 & \text{else.} \end{cases}$$
 (5.16.1)

Give an explicit formula for L_i ($\prod = \text{product}$, equal to 1, when empty):

$$L_i(t) = \prod_{i=0,\ldots,n} \frac{t-t_i}{t_i-t_i}$$
 , $i=0,\ldots,n$.

SOLUTION for (5-16.a) $\rightarrow 5-16-1-0$:pintps1.pdf

(5-16.b) • (2 min.) Given a node set $\{t_0, \ldots, t_n\}$ describe the Newton basis $\mathcal{N}_{\mathcal{T}}$ of the space \mathcal{P}_n of uni-variate polynomials of degree $\leq n, n \in \mathbb{N}$ belonging to the sequence $\mathcal{T} := (t_0, t_1, \ldots, t_n)$:

$$\mathcal{N}_{\mathcal{T}} = \left\{ igcap_{i=0}^n .
ight.$$

```
HIDDEN HINT 1 for (5-16.b) \rightarrow 5-16-2-0:pipdeth1.pdf Solution for (5-16.b) \rightarrow 5-16-2-1:.pdf
```

(5-16.c) (3 min.) The following classes are meant to evaluate global polynomial interpolants in different settings:

```
class InterpPolyEval_A {
   public:
    // Constructor taking the evaluation point x as argument
   InterpPolyEval_A(double x);
   // Add another data point (t,v) ∈ ℝ² and update internal
        information
   void addPoint(double t, double y);
   // Value of current interpolating polynomial at x
   double eval(void) const;
   // ... private data and member functions
};
```

Which algorithm would you recommend for the efficient implementation of the class InterpPolyEval_A,
if one has to deal with many calls to the member functions <pre>InterpPolyEval_A::addPoint()</pre>
<pre>and InterpPolyEval_A::eval()? (Pick only one)</pre>

Barycentric interpolation formula	Aitken-Neville scheme		Divided differences
class InterpPolyEval	_ B {		
public:			
// Constructors ta	-		-
InterpPolyEval_B(-		
// Evaluation oper		$y_0,\ldots,y_n);$	computes
$//p(x_k)$ for x_k s pa		77 77	
Eigen::VectorXd ev		: VectorXd	&Υ,
<pre>const Eigen::Vecto // private data</pre>		ations	
<pre>};</pre>	and member run	CLIOIIS	
J /			
_	-		he class InterpPolyEval_B in t
Which algorithm is most suitable ase $n\gg 1$ and many invocation Barycentric interpolation formula	-		
Barycentric interpolation formula class InterpPolyEval public: // Idle Constructor	Aitken-Neville scheme C {		1 ()? (Pick only one)
ase $n \gg 1$ and many invocation Barycentric interpolation formula class InterpPolyEval public:	Aitken-Neville scheme C {	ral_B::eva	1 () ? (Pick only one) Divided differences
Barycentric interpolation formula class InterpPolyEval public: // Idle Constructo InterpPolyEval_C() // Add another dat	Aitken-Neville scheme $ - \mathbf{C} \{$ $\mathbf{c}_{x} = \mathbf{c}_{y} (t,y) \in \mathbb{R} $	al_B::eva	1 () ? (Pick only one) Divided differences
Barycentric interpolation formula class InterpPolyEval public: // Idle Constructo InterpPolyEval_C() // Add another dat information void addPoint(doul // Evaluation of construction)	Aitken-Neville scheme -C { or it a point $(t,y) \in \mathbb{R}$ ole t, double y)	al_B::eva	1 () ? (Pick only one) Divided differences
Barycentric interpolation formula class InterpPolyEval public: // Idle Constructo InterpPolyEval_C() // Add another dat information void addPoint(doul	Aitken-Neville scheme C { or it a point $(t,y) \in \mathbb{R}$ ble t, double y) current interpol	al_B::eva and upda ating poly	Divided differences ate internal ynomial at many points
Barycentric interpolation formula class InterpPolyEval public: // Idle Constructo InterpPolyEval_C() // Add another dat information void addPoint(doul // Evaluation of construct on passed in x	Aitken-Neville scheme C { or ; ca point $(t,y) \in \mathbb{R}$ ble t, double y) current interpol ral(const Eigen:	2 and upda; ating poly :VectorXd	Divided differences ate internal ynomial at many points

in arbitrary order? (Pick only one)

Barycentric Aitken-Neville Divided differences interpolation formula scheme SOLUTION for (5-16.c) $\rightarrow 5-16-3-0$:.pdf

End Problem 5-16, 7 min.

Problem 5-17: Algorithms for Polynomial Interpolation

Global polynomials are the most important set of functions used for interpolation. This problem recalls key algorithms for the evaluating interpolating polynomials.

The problem is connected with [Lecture \rightarrow Section 5.2.3.4] and [Lecture \rightarrow Section 5.2.3.3].

```
(5-17.a) ☑ (15 min.) In Code 5.17.1 complete the implementation of the C++ function double evalNewtonForm(const std::vector<double> &t, const std::vector<double> &c, double x);
```

which evaluates a polynomial given through the coefficients of its Newton basis expansion at a single point $x \in \mathbb{R}$ (argument x). More precisely, setting n := c.size() - 1 the function should return

$$\begin{split} p(x) := \sum_{j=0}^{n} \mathtt{c}[\mathtt{j}] \, N_{j}(x) \,, \quad N_{j}(t) := \prod_{i=0}^{j-1} (t - \mathtt{t}[\mathtt{i}]) \,, \quad j = 0, \dots, n \,, \\ p(x) = \mathtt{c}[\mathtt{0}] + \mathtt{c}[\mathtt{1}](x - \mathtt{t}[\mathtt{0}]) + \mathtt{c}[\mathtt{2}](x - \mathtt{t}[\mathtt{0}])(x - \mathtt{t}[\mathtt{1}]) + \dots \\ + \mathtt{c}[\mathtt{n}](x - \mathtt{t}[\mathtt{0}]) \cdot \dots \cdot (x - \mathtt{t}[\mathtt{n} - \mathtt{1}]) \,, \end{split}$$

where an "empty product" is meant to evaluate to 1.

```
HIDDEN HINT 1 for (5-17.a) → 5-17-1-0:nhehal.pdf

SOLUTION for (5-17.a) → 5-17-1-1:.pdf

(5-17.b) ② (15 min.) After initialization or reset () an object of the type

class ConvergentSequence {
    public:
        ConvergentSequence ();
        double next();
        void reset();
    private:
        // Some private data
}:
```

produces the members of a sequence $(s_n)_{n=1}^{\infty}$ one after another by calls to the next() member function. The first value returned is s_1 . The sequence is known to converge for $n \to \infty$ with limit $s^* \in \mathbb{R}$: $s^* = \lim_{n \to \infty} s_n$.

In Listing 5.17.3 write valid C++ code in the blanks to complete the code of the function

```
double limitByExtrapolation(
   SEQGEN &&seq, double rtol = 1.0E-6,
   double atol = 1.0E-8, unsigned int maxdeg = 20);
```

which is to compute the limit s^* by extrapolation to zero. The type **SEQGEN** is compatible with **ConvergentSequence** and must offer a member function next () fitting the specification given before.

```
C++ code 5.17.3: Function limitByExtrapolation()
  template < class SEQGEN>
  double limitByExtrapolation (SEQGEN &&seq, double rtol = 1.0E-6,
                                double atol = 1.0E-8, unsigned int maxdeg = 20) {
     seq.reset();
4
     std::vector<double> h{
5
                                    };
     std::vector<double> s{seq.next()};
     for (unsigned i = 1; i < maxdeg; ++i) {
8
       h.push_back(
                                );
       s.push_back(seq.next());
9
                                            ; --k)
       for (int k =
                            ; k >=
10
         s[k] =
11
       const double errest = std::abs(s[1] - s[0]);
12
       if (errest < rtol * std::abs(s[0]) || errest < atol)</pre>
13
         return s[0];
14
15
     throw std::runtime_error("limitByExtrapolation failed to converge!");
16
     return std::nan("");
17
18
```

```
HIDDEN HINT 1 for (5-17.b) \rightarrow 5-17-2-0:nhehl.pdf
Solution for (5-17.b) \rightarrow 5-17-2-1:.pdf
```

End Problem 5-17, 30 min.

Chapter 6

Approximation of Functions in 1D

Problem 6-1: Piecewise linear approximation on graded meshes

One of the main point made in [Lecture \rightarrow Section 6.2.3] is that the quality of an interpolant depends heavily on the choice of the interpolation nodes. If the function to be interpolated has a "bad behavior" in a small part of the domain, for instance it has very large derivatives of high order, more interpolation points are required in that area of the domain. Commonly used tools to cope with this task are *graded meshes*, which will be the topic of this problem.

Substantial implementation in C++ is requested.

Given a mesh $\mathcal{M} = \{0 = t_0 < t_1 < \dots < t_{n-1} < t_n = 1\}$ on the unit interval I = [0,1], $n \in \mathbb{N}$, we define the *piecewise linear interpolant*:

$$I_{\mathcal{M}}: C^{0}(I) \to \mathcal{P}_{1,\mathcal{M}} = \{ s \in C^{0}(I), \ s_{|[t_{j-1},t_{j}]} \in \mathcal{P}_{1} \ \forall \ j \}, \quad \text{s.t.} \quad (I_{\mathcal{M}}f)(t_{j}) = f(t_{j}), \quad j = 0, \dots, n$$
(6.1.1)

See also [Lecture \rightarrow Section 5.3.2].

(6-1.a) If we choose the uniform mesh $\mathcal{M}=\{t_j\}_{j=0}^n$ with $t_j=j/n$, given a function $f\in C^2(I)$ what is the asymptotic behavior of the error $\|f-\mathsf{I}_{\mathcal{M}}f\|_{L^\infty(I)}$ when $n\to\infty$?

HIDDEN HINT 1 for (6-1.a) \rightarrow 6-1-1-0:GradedMeshes1h.pdf

SOLUTION for (6-1.a) \rightarrow 6-1-1-1: GradedMeshes1s.pdf

HIDDEN HINT 1 for (6-1.b) \rightarrow 6-1-2-0:GradedMeshes2h.pdf

SOLUTION for (6-1.b) \rightarrow 6-1-2-1:GradedMeshes2s.pdf

(6-1.c) (30 min.) Implement a templated C++ function

```
template <typename FUNCTION>
double pwlintpMaxError(FUNCTION &&f, const Eigen::VectorXd &t);
```

that computes $\|f - \mathsf{I}_{\mathcal{M}} f\|_{L^\infty(I)}$ when f is passed through the functor f (with evaluation operator double operator () (double) const) and the mesh \mathcal{M} through its *sorted* vector $[t_0, t_1, \ldots, t_n]^\top \in \mathbb{R}^{n+1}, n \in \mathbb{N}$, of node positions, which also defines the interval $I := [t_0, t_n]$.

The maximum norm should be approximated by sampling in 10^5 equidistant points in I, see, e.g., [Lecture \rightarrow Ex. 6.6.1.7].

SOLUTION for (6-1.c) \rightarrow 6-1-3-0:sa.pdf

We consider the root functions $f(t) := t^{\alpha}$, $\alpha > 0$. Implement a C++ function

void cvgplotEquidistantMesh(const Eigen::VectorXd &alpha);

that uses MATPLOTLIBCPP to create log-log plots of the the error norms $\|f - I_{\mathcal{M}} f\|_{L^{\infty}(I)}$ as a function of $n := \sharp \mathcal{M} - 1 \to \infty$ for families of *equidistant meshes*. Use meshes with $n = 2^k$, $k = 5, 6, \ldots, 12$, and values for α passed in alpha. Do not forget axis labels, plot title, and a meaningful legend. Save the plot in <code>cvgplotequidistant.png</code>.

Describe your observations qualitatively.

```
SOLUTION for (6-1.d) \rightarrow 6-1-4-0:sal.pdf
```

Create a C++ function

```
Eigen::VectorXd cvgrateEquidistantMesh(const Eigen::VectorXd
    &alpha);
```

that estimates rates of asymptotic algebraic convergence (\rightarrow [Lecture \rightarrow § 6.2.2.9]) of $||f - I_{\mathcal{M}}f||_{L^{\infty}(I)}$ for $f(t) := t^{\alpha}$ in terms of $n := \sharp \mathcal{M} - 1 \rightarrow \infty$ for families of *equidistant meshes* for the passed values of the parameter α and returns the estimates. These regression-based estimates should rely on meshes with $n = 2^k, k = 5, 6, \ldots, 12$.

```
HIDDEN HINT 1 for (6-1.e) \rightarrow 6-1-5-0:gm3h1.pdf
```

```
SOLUTION for (6-1.e) \rightarrow 6-1-5-1:gsmb.pdf
```

(6-1.f) \odot (15 min.) Tabulate and plot the empiric rate of asymptotic n-convergence (in maximum norm) of the piecewise linear interpolant of $f(t)=t^{\alpha}$, $0<\alpha<2$, on *equidistant meshes*. To conduct this experiment implement a C++ function

```
void testcvgEquidistantMesh();
```

that plots (use MatPlotLibCpp) and tabulates the measured convergence rates based on $\alpha=0.05,0.15,0.25,\ldots,2.95$. Save the plot in the file <code>cvgRateEquidistant.png</code>.

```
Solution for (6-1.f) \rightarrow 6-1-6-0: GradedMeshes3s.pdf
```

(6-1.g) \odot (5 min.) In which mesh interval do you expect $|f - I_M f|$ to attain its maximum?

HIDDEN HINT 1 for (6-1.g) \rightarrow 6-1-7-0: GradedMeshes4hb.pdf

```
SOLUTION for (6-1.g) \rightarrow 6-1-7-1:GradedMeshes4s.pdf
```

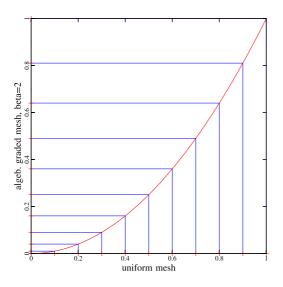
(6-1.h) □ (20 min.) [depends on Sub-problem (6-1.g)]

Compute ("on paper") the exact value of $||f - I_{\mathcal{M}} f||_{L^{\infty}(I)}$ for $f(t) := t^{\alpha}$, $0 < \alpha < 2$. Compare the order of convergence obtained with the one observed numerically in (6-1.b).

```
HIDDEN HINT 1 for (6-1.h) \rightarrow 6-1-8-0:GradedMeshes5h.pdf
```

SOLUTION for (6-1.h) \rightarrow 6-1-8-1:GradedMeshes5s.pdf

(6-1.i) © (20 min.) Since the interpolation error is concentrated in the left part of the domain, it seems reasonable to use a finer mesh only in this part.



A suitable choice is an **algebraically graded mesh**, defined as

$$\mathcal{G} = \left\{ t_j = \left(\frac{j}{n} \right)^{\beta}, j = 0, \dots, n \right\}$$
 (6.1.8)

for the grading parameter $\beta > 1$.

 \triangleleft Algebraically graded mesh for $\beta = 2$.

Fig. 53

Realize a C++ function

that estimates rates of asymptotic algebraic convergence (\rightarrow [Lecture \rightarrow § 6.2.2.9]) of $\|f - \mathsf{I}_{\mathcal{G}} f\|_{L^{\infty}(I)}$ for $f(t) := t^{\alpha}$ in terms of $n := \sharp \mathcal{M} - 1 \rightarrow \infty$ for families of *graded meshes* for the passed values of the parameter α (argument alpha) and of the grading parameter $\beta > 0$ (argument beta). The function should return the measured rates as entries of a matrix, a column for each α -value, a row for each β -value. The regression-based estimates should rely on meshes with $n = 2^k, k = 5, 6, \ldots, 12$.

SOLUTION for (6-1.i)
$$\rightarrow$$
 6-1-9-0:.pdf

(6-1.j) □ (20 min.) [depends on Sub-problem (6-1.i)]

When using algebraically graded meshes for the approximation of $t\mapsto t^{\alpha}$, $0<\alpha<2$, by means of piecewise linear interpolation, how do you have to choose the grading parameter $\beta=\beta(\alpha)$ as a function of α in order to recover an asymptotic algebraic convergence like $O(n^{-2})$ of $\|f-I_{\mathcal{G}}f\|_{L^{\infty}(I)}$ for $\sharp \mathcal{G} \sim n \to \infty$?

To arrive at a conjecture rely on the function cvgrateGradedMesh () from Sub-problem (6-1.i) and write another function

void testcvgGradedMesh();

that calls cvgrateGradedMesh() for

$$\alpha \in \{0.05, 0.15, 0.25, \dots, 2.95\},$$

 $\beta \in \{0.1, 0.2, 0.3, \dots, 2.0\},$

and creates a useful plot of the data (unsing MATPLOTLIBCPP) and stores it in the file alphabeta.png.

Solution for (6-1.j)
$$\rightarrow$$
 6-1-10-0:cls.pdf

End Problem 6-1, 170 min.

Problem 6-2: Piecewise linear interpolation with knots different from nodes

In [Lecture \rightarrow Ex. 5.1.0.15] we examined piecewise linear interpolation in the case where the interpolation nodes coincide with the abscissas of the corners of the interpolating polygon. However, that one is a very special situation. In this problem we consider a more general setting for piecewise linear interpolation.

Involves moderate implementation in C++.

We are given two ordered sets of real numbers to be read as points on the real line:

- (I) the knots $x_0 < x_1 < x_2 < \cdots < x_n$,
- (II) the (interpolation) nodes $t_0 < t_1 < \cdots < t_n$, $n \in \mathbb{N}$, satisfying $x_0 \le t_j \le x_n$, $j = 0, \ldots, n$.

The space of piecewise linear continuous functions with respect to the knot set $\mathcal{N} := \{x_j\}_{j=0}^n$ is defined as:

$$S_{1,\mathcal{N}} := \{ s \in C^0([x_0, x_n]) : s(t) = \gamma_j t + \beta_j \text{ for } t \in]x_{j-1}, x_j], \ \gamma_j, \beta_j \in \mathbb{R} \ i = 1, \dots, n \} \ . \tag{6.2.1}$$

Compare [Lecture \rightarrow Ex. 5.1.0.15] for the special case $x_i := t_i$.

On the function space $S_{1,N} \subset C^0(I)$, $I := [x_0, x_n]$, we consider the **interpolation problem**:

Find $s \in \mathcal{S}_{1,\mathcal{N}}$ such that $s(t_i) = y_i$, $i = 0, \ldots, n$, for given values $y_i \in \mathbb{R}$.

This fits the general framework of [Lecture \rightarrow § 5.1.0.21].

Below we set $I_0 := [x_0, x_1], I_j :=]x_{j-1}, x_{j+1}], j = 1, ..., n-1, I_n := [x_{n-1}, x_n].$

(6-2.a) • Show that the interpolation problem may not have a solution for some values y_j if a single interval $[x_j, x_{j+1}], j = 0, \ldots, n-1$, contains more than 2 nodes t_i .

Solution for (6-2.a)
$$\rightarrow$$
 6-2-1-0:lips1.pdf

(6-2.b) Show that the interpolation problem can have a unique solution for any data values y_j only if each interval I_j , j = 0, ..., n, contains at least one node.

```
HIDDEN HINT 1 for (6-2.b) \rightarrow 6-2-2-0:glgh1.pdf
```

SOLUTION for (6-2.b)
$$\rightarrow$$
 6-2-2-1:lips2.pdf

(6-2.c) Show that the interpolation problem has a unique solution for all data values y_j if each of the closed intervals $]x_0, x_1[,]x_1, x_2[,]x_2, x_3[,...,]x_{n-1}, x_n[$ contains at least one node t_j .

HIDDEN HINT 1 for (6-2.c) \rightarrow 6-2-3-0:glgh1.pdf

SOLUTION for (6-2.c)
$$\rightarrow$$
 6-2-3-1:lips3.pdf

(6-2.d) ☐ Implement a C++ function

that returns the vector of values $s(x_j)$ of the interpolant s of the data points (t_i, y_i) , $i = 0, \ldots, n$ in $S_{1,\mathcal{N}}$ in the knots x_j , $j = 0, \ldots, n$. The argument vectors x, t, and y pass the x_j , t_j , and values y_j , respectively.

The function should first check whether the condition formulated in Sub-problem (6-2.c) is satisified. You must not take for granted that knots or nodes are sorted already.

```
HIDDEN HINT 1 for (6-2.d) → 6-2-4-0:glgh1.pdf

SOLUTION for (6-2.d) → 6-2-4-1:glgs1.pdf

(6-2.e)  Implement a C++ class

    class PwLinIP {
    public:
        PwLinIP (const Eigen::VectorXd &x, const Eigen::VectorXd &t, const Eigen::VectorXd &y);
        double operator() (double arg) const;
    private:
        ...
    };
```

that realizes an interpolator class in the spirit of [Lecture \rightarrow Rem. 5.1.0.11]. The meanings of the arguments of the constructor are the same as for the function tentBasCoeff from Sub-problem (6-2.d). Pay attention to the efficient implementation of the evaluation operator!

```
HIDDEN HINT 1 for (6-2.e) \rightarrow 6-2-5-0:lipx1.pdf

HIDDEN HINT 2 for (6-2.e) \rightarrow 6-2-5-1:lipzz.pdf

Solution for (6-2.e) \rightarrow 6-2-5-2:dsfg1.pdf
```

(6-2.f) \Box Implement a C++ code that creates plots of the cardinal basis functions for interpolation in $\mathcal{S}_{1,\mathcal{N}}$

- with know set $\mathcal{N} := \{0, 1, 2, 3, \dots, 9, 10\}$
- and for interpolation nodes $t_0 := 0$, $t_j = j \frac{1}{2}$, $j = 1, \dots, 10$.

Recall that the cardinal basis function $b_k \in \mathcal{S}_{1,\mathcal{N}}$ is characterized by the conditions:

$$b_k(t_j) = \delta_{kj}, \quad k, j = 0, \dots, 10$$
 (6.2.6)

SOLUTION for (6-2.f) \rightarrow 6-2-6-0:glgs1.pdf

End Problem 6-2

Problem 6-3: Adaptive polynomial interpolation

In [Lecture \rightarrow Section 6.2.3] we have seen that the *a priori* placement of interpolation nodes is key to a good approximation by a polynomial interpolant. This problem deals with an *a posteriori* adaptive strategy that controls the placement of interpolation nodes depending on the interpolant. It employs a **greedy algorithm** to grow the node set based on an intermediate interpolant.

This strategy has recently gained prominence for a wide range of approximation problems, see [TEM08].

Considerable implementation in C++ is requested in this problem.

A description of the greedy algorithm for adaptive polynomial interpolation is as follows:

Given a function $f:[a,b]\mapsto \mathbb{R}$ one starts $\mathcal{T}_0:=\{\frac{1}{2}(b+a)\}$. Based on a fixed finite set $\mathcal{S}\subset [a,b]$ of *sampling points* one augments the set of nodes according to

$$\mathcal{T}_{n+1} = \mathcal{T}_n \cup \left\{ \underset{t \in \mathcal{S}}{\operatorname{argmax}} |f(t) - I_{\mathcal{T}_n}(t)| \right\}, \tag{6.3.1}$$

where $I_{\mathcal{T}_n}$ is the polynomial interpolation operator (\rightarrow [Lecture \rightarrow Cor. 5.2.2.8]) for the node set \mathcal{T}_n , until the relative interpolation error (in maximum norm) drops below a prescribed tolerance tol > 0:

$$\max_{t \in \mathcal{S}} |f(t) - I_{\mathcal{T}_n}(t)| \le \operatorname{tol} \cdot \max_{t \in \mathcal{S}} |f(t)|. \tag{6.3.2}$$

Throughout this exercise you may rely on a C++ function

```
// IN: t: vector of nodes t_0, \ldots, t_n

// y: vector of data y_0, \ldots, y_n

// x: vector of evaluation points x_1, \ldots, x_N

// OUT: p: interpolant evaluated at x

Eigen::VectorXd intpolyval(const Eigen::VectorXd& t, const Eigen::VectorXd& y, const Eigen::VectorXd& x);
```

that computes the values $p(x_i)$, i = 1, ..., N, where $p \in \mathcal{P}_n$ is the global polynomial interpolant through the data points (t_i, y_i) , j = 0, ..., n, see [Lecture \rightarrow Code 5.2.3.7].

```
C++ code 6.3.3: Function intpolyval()
```

```
Eigen::VectorXd intpolyval(const Eigen::VectorXd& t, const Eigen::VectorXd& y,
                            const Eigen::VectorXd& x) {
3
    const unsigned int
4
        n = t.size(), // no. of interpolation nodes = deg. of polynomial
        N = x.size(); // no. of evaluation points
6
    Eigen::VectorXd p = Eigen::VectorXd(N);
8
9
    // Precompute the weights \lambda_i with effort O(n^2)
10
    Eigen::VectorXd lambda(n);
11
12
    for (unsigned int k = 0; k < n; ++k) {
      // little workaround: cannot subtract a vector from a scalar
13
      // -> multiply scalar by vector of ones
14
      lambda(k) =
15
```

```
16
           ((t(k) * Eigen::VectorXd::Ones(k) - t.head(k)).prod() *
17
            (t(k) * Eigen:: VectorXd:: Ones(n - k - 1) - t.tail(n - k - 1)).prod());
19
    // Compute quotient of weighted sums of \frac{\lambda_i}{t-t_i},
20
     // effort O(n)
21
     for (unsigned int i = 0; i < N; ++i) {
22
       Eigen::VectorXd z = (x(i) * Eigen::VectorXd::Ones(n) - t);
23
       // check if we want to evaluate at a node <-> avoid division by zero
25
       double* ptr = std::find(z.data(), z.data() + n, 0.0);
       if (ptr != z.data() + n) \{ // if ptr = z.data + n = z.end no zero was
         p(i) = y(ptr - z.data()); // ptr - z.data gives the position of the
28
       } else {
29
         Eigen::VectorXd mu = lambda.cwiseQuotient(z);
         p(i) = (mu.cwiseProduct(y)).sum() / mu.sum();
31
       }
32
     }
     return p;
```

that implements the algorithm described above.

The function arguments are: the functor object f, the interval bounds a, b, the relative tolerance tol, the number N of *equidistant* sampling points, that is

$$S := \{ a + (b-a)j/N : j = 0, \dots, N \}, \qquad (6.3.4)$$

and t will be used as return parameter for interpolation nodes (i.e. the final set \mathcal{T}_n). The type **Function** defines a function and is supposed to provide an operator **double operator** () (**double**) **const**. A suitable lambda function can satisfy this requirement. The function should return the final set of interpolation nodes as proposed by the greedy algorithm.

```
HIDDEN HINT 1 for (6-3.a) \rightarrow 6-3-1-0:s1h1.pdf

SOLUTION for (6-3.a) \rightarrow 6-3-1-1:solfile.pdf

(6-3.b) (15 min.) [depends on Sub-problem (6-3.a)]

Extend the function adaptive polyintp() from Sub-problem (6-3.a) to
```

```
template <class Function>
  Eigen::VectorXd
  adaptivepolyintp(Function &&f,
        double a, double b,
        double tol, unsigned int N,
        std::vector<double> *errortab = nullptr);
```

so that it also reports (and returns in errortab) the quantity

$$\epsilon_n := \max_{t \in \mathcal{S}} |f(t) - \mathsf{T}_{\mathcal{T}_n}(t)| \tag{6.3.5}$$

for each intermediate node set \mathcal{T}_n .

SOLUTION for (6-3.b) \rightarrow 6-3-2-0:solfile.pdf

For $f_1(t) := \sin(e^{2t})$ and $f_2(t) = \frac{\sqrt{t}}{1+16t^2}$ write a C++ function

void plotInterpolationError();

that creates (employing MATPLOTLIBCPP) and stores (in intperrplot.png) a plot of ϵ_n versus n (the number of interpolation nodes). Choose axis scales that reveal the qualitative decay (types of convergence as given in [Lecture \rightarrow Def. 6.2.2.7]) of this error norm as the number of interpolation nodes is increased. Use interval [a,b] = [0,1], N=1000 sampling points, tolerance tol = 1e-6.

SOLUTION for (6-3.c) \rightarrow 6-3-3-0:solfile.pdf

End Problem 6-3, 105 min.

Problem 6-4: Chebyshev polynomials and their properties

Chebyshev polynomials as defined in [Lecture \rightarrow Def. 6.2.3.3] are a sequence of orthogonal polynomials, which can be defined recursively, see [Lecture \rightarrow Thm. 6.2.3.4]. In this problem we will examine these polynomials and a few of their many properties.

Involves a little implementation in C++

Let $T_n \in \mathcal{P}_n$ be the n-th Chebyshev polynomial, as defined in [Lecture \to Def. 6.2.3.3] and $\xi_0^{(n)}, \ldots, \xi_{n-1}^{(n)}$ be the n zeros of T_n , the Chebychev nodes. According to [Lecture \to Eq. (6.2.3.10)], these are given by:

$$\xi_j^{(n)} = \cos\left(\frac{2j+1}{2n}\pi\right), \quad j = 0, \dots, n-1.$$
 (6.4.1)

Recall the notion of a "almost" inner product.

Definition [Lecture → Def. 6.3.1.1]. (Semi-)inner product [NIS02]

Let V be a vector space over the field \mathbb{K} . A mapping $b: V \times V \to \mathbb{K}$ is called an inner product on V, if it satisfies

- (i) b is linear in the first argument: $b(\alpha v + \beta w, u) = \alpha b(v, u) + \beta b(w, u)$ for all $\alpha, \beta \in \mathbb{K}$, $u, v, w \in V$,
- (ii) b is (anti-)symmetric: $b(v, w) = \overline{b(w, v)}$ (- \hat{b} complex conjugation),
- (iii) b is positive definite: $v \neq 0 \Leftrightarrow b(v, v) > 0$.

b is a semi-inner product, if it still complies with (i) and (ii), but is only positive semi-definite: $b(v,v) \ge 0$ for all $v \in V$.

We define the family of discrete L^2 semi-inner products [Lecture \rightarrow Eq. (6.3.2.4)]:

$$(f,g)_n := \sum_{j=0}^{n-1} f(\xi_j^{(n)}) g(\xi_j^{(n)}), \quad f,g \in C^0([-1,1]) . \tag{6.4.2}$$

We also define the special weighted L^2 inner product:

$$(f,g)_w := \int_{-1}^1 \frac{1}{\sqrt{1-t^2}} f(t)g(t) dt \quad f,g \in C^0([-1,1]) . \tag{6.4.3}$$

(6-4.a) \odot (15 min.) Show that the weight function in (6.4.3) $w(t):=\left(1-t^2\right)^{-\frac{1}{2}}$ satisfies $\int_{-1}^{1}|w(t)|\,\mathrm{d}t<\infty$.

HIDDEN HINT 1 for (6-4.a) \rightarrow 6-4-1-0:s0h1.pdf

Solution for (6-4.a)
$$\rightarrow$$
 6-4-1-1:.pdf

(6-4.b) ☐ (15 min.) [depends on Sub-problem (6-4.a)]

Show that the Chebyshev polynomials are an orthogonal family of polynomials with respect to the inner product defined in (6.4.3) in the sense of [Lecture \rightarrow Def. 6.3.2.7]:

$$(T_k, T_\ell)_w = 0$$
 for all $k, \ell \in \mathbb{N}_0$, $k \neq \ell$. (6.4.4)

HIDDEN HINT 1 for (6-4.b) \rightarrow 6-4-2-0: ChebPolyProperties1h.pdf

SOLUTION for (6-4.b)
$$\rightarrow$$
 6-4-2-1:ChebPolyProperties1s.pdf

Recall the concept of an orthonormal basis:

Definition [Lecture \rightarrow Def. 6.3.1.14]. Orthonormal basis

A subset $\{b_1, \ldots, b_N\}$ of an N-dimensional vector space V with inner product (\rightarrow [Lecture \rightarrow Def. 6.3.1.1]) $(\cdot, \cdot)_V$ is an orthonormal basis (ONB), if $(b_k, b_j)_V = \delta_{kj}$.

A basis of $\{b_1, \ldots, b_N\}$ of V is called orthogonal, if $(b_k, b_j)_V = 0$ for $k \neq j$.

The following is a surprising result.

Theorem 6.4.5. Discrete orthogonality of Chebychev polynomials

The family of polynomials $\{T_0, \ldots, T_n\}$ is an orthogonal basis [Lecture \to Def. 6.3.1.14] of \mathcal{P}_n with respect to the semi-inner product $(\ ,\)_{n+1}$ defined in (6.4.2).

(6-4.c) In (10 min.) Show that for every $n \in \mathbb{N}_0$ $(\cdot, \cdot)_{n+1}$ as defined in (6.4.2) is an inner product on the (real) vector space \mathcal{P}_n of polynomials of degree $\leq n$.

Solution for (6-4.c)
$$\rightarrow$$
 6-4-3-0:sy.pdf

(6-4.d) (45 min.) Write a C++ function

bool checkDiscreteOrthogonality(unsigned int n);

meant to test the assertion of Thm. 6.4.5 numerically.

HIDDEN HINT 1 for (6-4.d) \rightarrow 6-4-4-0: ChebPolyProperties2h.pdf

HIDDEN HINT 2 for (6-4.d) \rightarrow 6-4-4-1:2h2.pdf

SOLUTION for (6-4.d) \rightarrow 6-4-4-2: ChebPolyProperties2s.pdf

(6-4.e) (30 min.) Prove Thm. 6.4.5.

HIDDEN HINT 1 for (6-4.e) \rightarrow 6-4-5-0: ChebPolyProperties3h.pdf

Solution for (6-4.e) \rightarrow 6-4-5-1:ChebPolyProperties3s.pdf

(6-4.f) **(30 min.)** [depends on Sub-problem (6-4.c) and Thm. 6.4.5]

Given $n, m \in \mathbb{N}$, $m \le n$, and a function $f \in C^0([-1,1])$, find an expression for the best approximant $q_m \in \mathcal{P}_m$ of f in the discrete L^2 -norm:

$$q_m = \underset{p \in \mathcal{P}_m}{\operatorname{argmin}} ||f - p||_{n+1},$$
 (6.4.9)

where $\|\cdot\|_{n+1}$ is the norm induced by the scalar product $(\cdot,\cdot)_{n+1}$ from (6.4.2). You should represent q_m through a Chebychev expansion of the form:

$$q_m = \sum_{j=0}^m \alpha_j T_j , \quad \alpha_j \in \mathbb{R} , \qquad (6.4.10)$$

see also [Lecture \rightarrow Section 6.2.3.3].

HIDDEN HINT 1 for (6-4.f) \rightarrow 6-4-6-0: ChebPolyProperties4h.pdf

Solution for (6-4.f) \rightarrow 6-4-6-1: ChebPolyProperties 4s.pdf

Predict q_m as defined through (6.4.9) in the case m = n.

Solution for (6-4.g) \rightarrow 6-4-7-0:sw.pdf

Write a C++ function

template <typename Function>

Eigen::VectorXd

bestpolchebnodes (Function &&f,

unsigned int n, unsigned int m);

that returns the vector of coefficients α_j , $j=0,\ldots,m$, in (6.4.10) given the function f. The degree m of the polynomial and the number n+1 of the Chebychev nodes entering the inner product are passed through the arguments m and n. The input f is a functor of signature std::function<double(double)>, for instance,

auto f = [] (**double** & x) { return
$$1/(pow(5*x,2)+1)$$
; };

SOLUTION for (6-4.h) \rightarrow 6-4-8-0: ChebPolyProperties5s.pdf

We test the implementation of bestpolchebnodes () with the function $f(x) = \frac{1}{(5x)^2+1}$. To that end write a C++ function

std::vector<double> testBestPolyChebNodes(unsigned int n);

that takes the integer n as argument and both tabulates and returns the sequence of the approximation error norms $||f - q_m||_{L^{\infty}(]-1,1[)}$, $m = 0, \ldots, n$, q_m as in (6.4.9).

Approximate the supremum norm of the approximation error by sampling on an equidistant grid with 10^6 points.

HIDDEN HINT 1 for (6-4.i) \rightarrow 6-4-9-0: ChebPolyProperties 6h.pdf

SOLUTION for (6-4.i) \rightarrow 6-4-9-1: ChebPolyProperties6s.pdf

(6-4.j) (30 min.) [depends on Sub-problem (6-4.f)]

Let $L_j \in \mathcal{P}_n$, $j=0,\ldots,n$, be the Lagrange polynomials [Lecture \to § 5.2.2.3] associated with the Chebychev nodes $t_j := \xi_j^{(n+1)}$ in [-1,1]. Show that:

$$L_j(t) = \frac{1}{n+1} + \frac{2}{n+1} \sum_{l=1}^{n} T_l(\xi_j^{(n+1)}) T_l(t) , \quad j = 0, \dots, n.$$
 (6.4.15)

HIDDEN HINT 1 for (6-4.j) \rightarrow 6-4-10-0: ChebPolyProperties7h.pdf

SOLUTION for (6-4.j) \rightarrow 6-4-10-1: ChebPolyProperties7s.pdf

End Problem 6-4, 260 min.

Problem 6-5: Chebychev interpolation of analytic functions

This problem concerns Chebychev interpolation as discussed in [Lecture \rightarrow Section 6.2.3] and takes a close look at the interpolation error estimates derived for analytic interpolands in [Lecture \rightarrow Section 6.2.2.3].

This problem is closely connected to [Lecture \rightarrow Rem. 6.2.3.26] and involves implementation in C++ based on EIGEN. You are expected to be able to compute with complex numbers.

Using techniques from complex analysis, notably the residue theorem [Lecture \rightarrow Thm. 6.2.2.50], in class we derived an expression for the interpolation error [Lecture \rightarrow Eq. (6.2.2.57)] and from it an error bound [Lecture \rightarrow Eq. (6.2.2.58)], as much sharper alternative to [Lecture \rightarrow Thm. 6.2.2.15] and [Lecture \rightarrow Lemma 6.2.2.20] for *analytic* interpolands. The bound tells us that for all $t \in [a, b]$

$$|f(t) - \mathsf{L}_{\mathcal{T}} f(t)| \leq \left| \frac{w(x)}{2\pi i} \int_{\gamma} \frac{f(z)}{(z-t)w(z)} dz \right| \leq \frac{|\gamma|}{2\pi} \frac{\max_{a \leq \tau \leq b} |w(\tau)|}{\min_{z \in \gamma} |w(z)|} \frac{\max_{z \in \gamma} |f(z)|}{d([a,b],\gamma)},$$

where $d([a,b],\gamma)$ is the geometric distance of the integration contour $\gamma\subset\mathbb{C}$ from the interval $[a,b]\subset\mathbb{C}$ in the complex plane. The contour γ must be contractible in the domain D of analyticity of f and must wind around [a,b] exactly once, see [Lecture \to Fig. 213].

Now we consider the interval [-1,1]. Following [Lecture \rightarrow Rem. 6.2.3.26], our task is to find an upper bound for this expression, in the case where f possesses an analytical extension to a complex neighbourhood of [-1,1].

For the analysis of the Chebychev interpolation of analytic functions we used the elliptical contours, see [Lecture \rightarrow Fig. 227],

double lengthEllipse(double rho, unsigned int N);

that approximates the length of the ellipse γ_{ρ} by sampling (6.5.1) for $N \in \mathbb{N}$ (equidistant) arguments θ and returns the result.

SOLUTION for (6-5.a)
$$\rightarrow$$
 6-5-1-0:sa.pdf

From now on we consider the *S*-curve function (the logistic function):

$$f(t) := \frac{1}{1 + e^{-3t}}, \quad t \in \mathbb{R}.$$
 (6.5.3)

(6-5.b) \bigcirc (20 min.) Determine the maximal domain of analyticity of the extension of f to the complex plane \bigcirc .

HIDDEN HINT 1 for (6-5.b) \rightarrow 6-5-2-0:s2h1.pdf

SOLUTION for (6-5.b)
$$\rightarrow$$
 6-5-2-1:solfile.pdf

(6-5.c) □ (15 min.) [depends on Sub-problem (6-5.b)]

Find the least upper bound for the set of ellipse parameters

$$\{\rho > 1: \gamma_{\rho} \subset D\}$$
,

where $D \subset \mathbb{C}$ is the domain of analyticity found in Sub-problem (6-5.b).

HIDDEN HINT 1 for (6-5.c) \rightarrow 6-5-3-0:s2ah1.pdf

Solution for (6-5.c)
$$\rightarrow$$
 6-5-3-1:s2a.pdf

Write a C++ function

that computes an approximation M of

$$\min_{\rho>1} \frac{|\gamma_{\rho}| \max_{z \in \gamma_{\rho}} |f(z)|}{\pi d([-1,1], \gamma_{\rho})}, \tag{6.5.4}$$

by sampling 1000 equidistant ρ -values in the interval $(1, \rho_{\text{max}}]$, where ρ_{max} is the value obtained in Sub-problem (6-5.c). For sampling along the integration path γ_{ρ} use 1000 points corresponding to equidistant θ -values, cf. Sub-problem (6-5.a). The function should return the pair $(M, \rho) \in \mathbb{R}^2$, where ρ stands for the ρ -value for which the minimum is attained.

The distance of [a, b] from γ_{ρ} is formally defined as

$$d([a,b],\gamma) := \inf\{|z-t| \mid z \in \gamma, t \in [a,b]\}. \tag{6.5.5}$$

HIDDEN HINT 1 for (6-5.d) \rightarrow prbfile.pdf

SOLUTION for (6-5.d)
$$\rightarrow$$
 6-5-4-1:solfile.pdf

Based on the result of Sub-problem (6-5.d) and [Lecture \rightarrow Rem. 6.2.3.26], give an "optimal" bound for

$$||f - L_n f||_{L^{\infty}([-1,1])},$$

where $L_n: C^0([-1,1]) \to \mathcal{P}_n$ is the operator of Chebychev interpolation, that is, polynomial Lagrange interpolation based on Chebychev nodes [Lecture \to Eq. (6.2.3.10)], on [-1,1] into the space of polynomials of degree $\le n$.

Solution for (6-5.e)
$$\rightarrow$$
 6-5-5-0:solfile.pdf

(6-5.f) (20 min.) [depends on Sub-problem (6-5.e)]

Implement a C++ function

```
std::pair<std::vector<double>, std::vector<double> >
compareErrorAndBound(unsigned int n_max);
```

that computes and tabulates

- 1. estimates for $||f L_n f||_{L^{\infty}([-1,1])}$ obtained by sampling in 10000 equidistant points $\in [-1,1]$,
- 2. the bounds derived in Sub-problem (6-5.e),

for polynomial degrees $n=4,\ldots,n_{\max}$, where n_{\max} is passed through the argument n_max.

You may rely on the provided function intpolyval() for the evaluation of a polynomial interpolant at many points.

```
HIDDEN HINT 1 for (6-5.f) \rightarrow 6-5-6-0:s5h1.pdf
SOLUTION for (6-5.f) \rightarrow 6-5-6-1:solfile.pdf
```

Solution for (6-5.g) \rightarrow 6-5-7-0:solfile.pdf

End Problem 6-5, 130 min.

Problem 6-6: Polynomial Best Approximation

The famous Jackson Theorem [Lecture \rightarrow Thm. 6.2.1.11] gives detailed information about the best approximation of C^r-functions by means of polynomials. In this exercise we recall that theorem and study best-approximation errors in general.

Purely theoretical problem related to [Lecture \rightarrow Section 6.2.1]

Throughout this problems we write Here

(6-6.a) • (7 min.) The following is the assertion of Jackson's theorem of [Lecture \rightarrow Thm. 6.2.1.11] with some parts missing. Fill either n or r in the green boxes and supplement the right subscript for the norm in the white boxes.

$$\inf_{p \in \mathcal{P}_n} \|f - p\| \leq \left(1 + \frac{\pi^2}{2}\right)^{\frac{1}{2}} \quad \frac{\left(1 - \frac{1}{2}\right)!}{\frac{1}{2}!} \left\|f^{(1)}\right\|_{L^{\infty}}$$

for all $n, r \in \mathbb{N}$, $n \ge r$, and $f \in C^r([-1,1])$. Here \mathcal{P}_n designates the space of univariate polynomials of degree < n and $f^{(k)}$ is a short notation for the k-th derivative of a function.

Solution for (6-6.a)
$$\rightarrow$$
 6-6-1-0:s1.pdf

(6-6.b) □ (8 min.) For $n \in \mathbb{N}$ we abbreviate

$$E_n(f) := \inf_{p \in \mathcal{P}_n} \|f - p\|_{L^{\infty}([-1,1])} \quad \text{for} \quad f \in C^0([-1,1]) \ . \tag{6.6.2}$$

Which of the following assertions on E_n are correct, which are wrong? Throughout, f and g denote arbitrary continuous functions on [-1,1] and $n \in \mathbb{N}$

HIDDEN HINT 1 for (6-6.b) \rightarrow 6-6-2-0:2norm.pdf

SOLUTION for (6-6.b) \rightarrow 6-6-2-1:s2.pdf

End Problem 6-6, 15 min.

(iv)

Problem 6-7: Piecewise cubic Hermite Interpolation of Functions

The variant of piecewise cubic Hermite interpolation studied in [Lecture \rightarrow Section 6.6.2] was a generalized interpolation scheme in the sense that it required point values of the derivative of the interpoland. Conversely, for piecewise cubic Hermite data interpolation we had to use reconstructed slopes. This problem examines the use of reconstructed slopes also for the purpose of function approximation.

Relies on [Lecture \rightarrow Section 5.3.3] and [Lecture \rightarrow Section 6.6.2] and involves C++ coding.

Piecewise cubic Hermite interpolation of a function $f \in C^1([a,b])$ with exact slopes on a mesh

$$\mathcal{M} := \{ a = x_0 < x_1 < \dots < x_m = b \}$$

was defined in [Lecture \to Def. 6.6.2.1]. For $f \in C^4([a,b])$ it enjoys h-convergence with rate 4 as we have seen in [Lecture \to Exp. 6.6.2.2] and which is stated in [Lecture \to Thm. 6.6.2.3].

Theorem [Lecture o Thm. 6.6.2.3]. Convergence of approximation by cubic Hermite interpolation

Let s be the cubic Hermite interpolant of $f \in C^4([a,b])$ on a mesh $\mathcal{M} := \{a = x_0 < x_1 < \ldots < x_{m-1} < x_m = b\}$ according to [Lecture \to Def. 6.6.2.1]. Then

$$||f - s||_{L^{\infty}([a,b])} \le \frac{1}{4!} h_{\mathcal{M}}^{4} ||f^{(4)}||_{L^{\infty}([a,b])},$$

with the meshwidth $h_{\mathcal{M}} := \max_{j} |x_j - x_{j-1}|$.

Now we consider cases, where perturbed or reconstructed slopes are used. For instance, this was done in the context of monotonicity preserving piecewise cubic Hermite interpolation as discussed in [Lecture \rightarrow Section 5.3.3.2].

$$s(x_i) = f(x_i)$$
 , $s'(x_i) = f'(x_i) + \delta_i$, (6.7.1)

where the perturbations $\delta_i \in \mathbb{R}$ may depend on \mathcal{M} , too.

Now we consider a sequence of meshes with meshwidth $h_{\mathcal{M}} \to 0$. Which rate of asymptotic h-convergence of the maximum norm $\|f-s\|_{L^{\infty}(]a,b[)}$ of the approximation error can be expected, if we know that for all j that

$$|\delta_i| = O(h_{\mathcal{M}}^{\beta})$$
 as $h_{\mathcal{M}} \to 0$ for some $\beta \in \mathbb{N}_0$, (6.7.2)

for mesh-width $h_{\mathcal{M}} \to 0$.

HIDDEN HINT 1 for (6-7.a) \rightarrow 6-7-1-0:s1h1.pdf

SOLUTION for (6-7.a)
$$\rightarrow$$
 6-7-1-1:s1.pdf

(6-7.b) We define a "strange" piecewise cubic interpolant:

Definition 6.7.3. Flat-point piecewise cubic interpolant

Given $f \in C^0([a,b])$ and a mesh $\mathcal{M} := \{a = x_0 < x_1 < \ldots < x_{m-1} < x_m = b\}$ the flat-point piecewise cubic interpolant $s : [a,b] \to \mathbb{R}$ is defined as

$$s_{|[x_{j-1},x_j]} \in \mathcal{P}_3$$
, $j=1,\ldots,m$, $s(x_j)=f(x_j)$, $s'(x_j)=0$, $j=0,\ldots,m$.

Implement a C++ function

template <typename Function>

```
Eigen::VectorXd fppchip(Function &&f,
   const Eigen::VectorXd &x, const Eigen::VectorXd &t);
```

that returns the vector $[s(t_k)]_{k=1}^N$, where

- the argument \times passes the m+1 (ordered) nodes of a mesh $\mathcal{M}:=\{a=x_0< x_1<\ldots< x_{m-1}< x_m=b\},$
- s is the flat-point piecewise cubic interpolant of f according to Def. 6.7.3,
- the (sorted!) points $t_k \in [x_0, x_m]$ are supplied as components of the vector t.

Remark. Please refer to [Lecture \rightarrow Thm. 5.3.3.16] to understand the motivation for introducing the "flat-point piecewise cubic interpolant". This is the only linear interpolation operator into C^1 -functions that is locally monotonicity preserving.

SOLUTION for (6-7.b)
$$\rightarrow$$
 6-7-2-0:sx2.pdf

For the function $f(t) = \frac{1}{1+t^2}$ and the interval [-5,5] empirically determine (qualitatively and quantitatively) the asymptotic convergence (in terms of meshwidth $h_{\mathcal{M}} \to 0$) of the flat-point piecewise cubic interpolants according to Def. 6.7.3 on sequences of *equidistant meshes*.

To that end study the maximum norm of the interpolation error in numerical experiments. Approximate those maximum norms by sampling in 10000 equidistant points. Implement a C++ function

```
std::vector<double> fppchipConvergence();
```

that tabulates and returns the error norms for meshes with 4, 8, 16, 32, 64, 128, 256, 512 nodes.

Solution for (6-7.c)
$$\rightarrow$$
 6-7-3-0:sz3.pdf

```
(6-7.d) (10 min.) [depends on Sub-problem (6-7.a) and Sub-problem (6-7.c)]
```

Invoke the result of Sub-problem (6-7.a) to give a theoretical explanation for the observations made in Sub-problem (6-7.c).

```
Solution for (6-7.d) \rightarrow 6-7-4-0:s3a.pdf
```

We now define a piecewise cubic Hermite interpolant with slopes reconstructed from sampled function values.

Definition 6.7.6. Piecewise cubic Hermite interpolant with reconstructed slopes

Given $f \in C^0([a,b])$ and a mesh $\mathcal{M} := \{a = x_0 < x_1 < \ldots < x_{m-1} < x_m = b\}$ the piecewise cubic Hermite interpolant with reconstructed slopes $s : [a,b] \to \mathbb{R}$ is defined as

$$s'(x_j) = \begin{cases} \mathcal{P}_3, & j = 1, \dots, m \\ \frac{f(x_{j-1}, x_j)}{2h} & for j = 0, \dots, m, \end{cases}$$

$$s'(x_j) = \begin{cases} \frac{-f(x_2) + 4f(x_1) - 3f(x_0)}{2h} & for j = 0, \\ \frac{f(x_{j+1}) - f(x_{j-1})}{2h} & for j = 1, \dots, m-1, \\ \frac{3f(x_m) - 4f(x_{m-1}) + f(x_{m-2})}{2h} & for j = m. \end{cases}$$

A class for piecewise cubic Hermite interpolation of continuous functions according to Def. 6.7.6 is defined as follows:

```
class CubicHermiteInterpolant {
public:
    template <typename Function>
    CubicHermiteInterpolant(Function &&f, const Eigen::VectorXd &t);
    virtual ~CubicHermiteInterpolant(void) = default;
    Eigen::VectorXd eval(const Eigen::VectorXd &x) const;
private:
    Eigen::VectorXd t_; // sorted mesh nodes
    Eigen::VectorXd y_; // function values at mesh nodes
    Eigen::VectorXd c_; // slopes at mesh nodes
};
```

The constructor expects a functor f with signature **std**::function<**double** (**double**) and a (sorted!) vector t of nodes. Its evaluation method CubicHermiteInterpolant::eval() returns $s(x_k)$ for the (sorted!) components of the vector passed in x. Here $s \in C^1([a,b])$ is the piecewise cubic Hermite interpolant as defined in Def. 6.7.6.

(6-7.e) (25 min.) Provide an implementation of the class CubicHermiteInterpolant.

```
Solution for (6-7.e) \rightarrow 6-7-5-0:s4.pdf
```

Based on the class **CubicHermiteInterpolant** realize a C++ function

```
std::vector<double> rspchipConververgence();
```

that answers the questions of Sub-problem (6-7.c) for the interpolant according to Def. 6.7.6.

```
Solution for (6-7.f) \rightarrow 6-7-6-0:.pdf
```

Based on the result of Sub-problem (6-7.a) prove that

$$f \in C^3([a,b]) \Rightarrow \exists C > 0: \|f - I_H f\|_{L^{\infty}([a,b[))} \le Ch_{\mathcal{M}}^3$$

with C>0 independent of the mesh and $h_{\mathcal{M}}$ standing for the mesh width of \mathcal{M} .

```
HIDDEN HINT 1 for (6-7.g) \rightarrow 6-7-7-0:s6h1.pdf
```

Solution for (6-7.g) \rightarrow 6-7-7-1:.pdf

End Problem 6-7, 125 min.

Problem 6-8: Monomial representation of Chebychev polynomials

Chebychev polynomials are important both for theoretical analysis and algorithm design. For isntance, as we have seen in [Lecture \rightarrow Section 6.2.3.3] Chebyshev interpolants should internally be represented through their Chebychev expansion. In this problem we will study an algorithm for converting Chebychev expansions into monomial expansions.

This problem relies about elementary facts about Chebychev polynomials as introduced in [Lecture \rightarrow Section 6.2.3.1].

On [-1,1] the n-th Chebychev polynomial T_n , $n \in \mathbb{N}_0$ is defined as [Lecture \to Def. 6.2.3.3]

$$T_n(t) := \cos(n \arccos t) \quad -1 \le t \le 1.$$
 (6.8.1)

That T_n is a polynomial of degree $\leq n$ is clear from the 3-term recursion

$$T_{n+1}(t)=2t\,T_n(t)-T_{n-1}(t)$$
 , $T_0\equiv 1$, $T_1(t)=t$, $n\in\mathbb{N}$. [Lecture o Eq. (6.2.3.5)]

The set $\{T_0, \ldots, T_n\}$ forms a basis of \mathcal{P}_n .

(6-8.a) (8 min.) Below is an incomplete implementation of a C++ function

```
Eigen::VectorXd chebexpToMonom(const Eigen::VectorXd &a);
```

that expects the coefficients $a_i \in \mathbb{R}$, $j = 0, \ldots, n$, of the Chebyshev expansion

$$p(t) = \sum_{j=0}^{n} a_j T_j(t)$$
 , $t \in \mathbb{R}$,

of a polynomial $p \in \mathcal{P}_n$ to be passed in the vector a and returns the coefficients r_j , $j = 0, \ldots, n$ with respect to the monomial basis,

$$p(t) = r_0 + r_1 t + r_2 t^2 + \dots + r_n t^n$$
,

collected in a vector $[r_0, r_1, \ldots, r_n]^{\top}$.

Fill in the missing pieces of the following code, which relies on [Lecture \rightarrow Eq. (6.2.3.5)].

```
Eigen::VectorXd chebexpToMonom(const Eigen::VectorXd &a) {
  const int n = a.size() -1; // degree of polynomial
  assert(n >= 0);
  Eigen::VectorXd r{ Eigen::VectorXd::Zero(n+1) };
  r[0] = A;
  if (n > 0) {
    r[1] = B;
    if (n > 1) {
        // two vectors temporarily storing the monomial coefficients
        // of Chebychev polynomials of two consecutive degrees
        std::array<Eigen::VectorXd, 2> c
        { Eigen::VectorXd::Zero(n+1), Eigen::VectorXd::Zero(n+1) };
        c[0][0] = 1.0; c[1][1] = 1.0;
        for (int j=2; j<= C ; ++j) {
            c[j%2][0] *= -1.0;</pre>
```

SOLUTION for (6-8.a) \rightarrow 6-8-1-0:s1.pdf

End Problem 6-8, 8 min.

Problem 6-9: Finding domains of analyticity

The developments of [Lecture \rightarrow Lemma 6.2.2.53], [Lecture \rightarrow Rem. 6.2.3.26], and [Lecture \rightarrow Section 6.5.3] make clear that knowledge about that maximal subdomain $D \subset \mathbb{C}$ of the complex plane to which an analytic function can be extended from its original real interval of definition is key to predicting the speed of exponential convergence of approximation schemes like Chebychev interpolation and trigonometric interpolation. In this problem we discuss the analytic extension of 1-periodic functions.

Depends on [Lecture \rightarrow Rem. 6.2.2.67] and is related to [Lecture \rightarrow § 6.5.3.16]. This problem does not ask for C++ coding.

In order to analyze the trigonometric functions occurring in this problem, the following identities can be useful:

$$\sin(x + iy) = \sin(x)\cosh(y) + i\cos(x)\sinh(y) \quad \forall x, y \in \mathbb{R},
\cos(x + iy) = \cos(x)\cosh(y) - i\sin(x)\sinh(y) \quad \forall x, y \in \mathbb{R},
\sin^{2}(z) + \cos^{2}(z) = 1 \quad \forall z \in \mathbb{Z}.$$
(6.9.1)

The following result gives domains of analyticity of special functions:

Theorem 6.9.2. domains of analyticity of special functions

- Every polynomial is analytic on ℂ.
- Every rational function, that is, a quotient of two polynomials, is analytic in the complement of the set of zeros of its denominator.
- The functions \exp , \sin , and \cos are analytic on \mathbb{C} .
- The functions $\log: \mathbb{R}^+ \to \mathbb{R}$ and $\sqrt{\mathbb{R}^+_0} \to \mathbb{R}$ can be extended to analytic functions on $\mathbb{C} \setminus \mathbb{R}^-_0$.

(6-9.a) (15 min.) We consider the 1-periodic function

$$f(t) = b \log(a + \sin(2\pi t)), \quad t \in \mathbb{R}, \quad b > 0, \, a > 1.$$
 (6.9.3)

Determine the largest possible subset D of \mathbb{C} to which f can be extended analytically.

$$D=\mathbb{C}\setminus\left\{ egin{array}{c} \operatorname{Re}(z)\in & & & \\ z\in\mathbb{C}: & & \\ \operatorname{Im}(z)\in & & & \end{array}
ight.$$

SOLUTION for (6-9.a) \rightarrow 6-9-1-0:s1.pdf

(6-9.b) □ (15 min.) We consider the 1-periodic function

$$f(t) = a \log(b - \cos(2\pi t)), \quad t \in \mathbb{R}, \quad b > 1, a > 1.$$
 (6.9.4)

Determine the largest possible subset D of \mathbb{C} to which f can be extended analytically.

$$D = \mathbb{C} \setminus \left\{ z \in \mathbb{C} : \\ \operatorname{Im}(z) \in \mathbb{C} \right\}$$

SOLUTION for (6-9.b) \rightarrow 6-9-2-0:s2.pdf

190

(6-9.c) (15 min.) We consider the 1-periodic function

$$f(t) = a\sqrt{b - \sin(2\pi t)}$$
, $t \in \mathbb{R}$, $b > 1$, $a > 1$. (6.9.5)

Determine the largest possible subset D of \mathbb{C} to which f can be extended analytically.

$$D = \mathbb{C} \setminus \left\{ z \in \mathbb{C} : \\ \operatorname{Im}(z) \in \mathbb{C} \right\}$$

SOLUTION for (6-9.c) \rightarrow 6-9-3-0:s2.pdf

End Problem 6-9, 45 min.

Chapter 7

Numerical Quadrature

Problem 7-1: Smooth integrand by transformation

In [Lecture \rightarrow Rem. 7.4.3.10] we saw how knowledge about the structure of a non-smooth integrand can be used to restore its smoothness by transformation. The very same idea can be put to use in this problem.

This problem practises the use of Gauss-Legendre quadrature and of the techniques introduced in [Lecture \rightarrow Rem. 7.4.3.10].

Given a smooth, odd function $f: [-1,1] \to \mathbb{R}$, consider the integral

$$I := \int_{-1}^{1} \arcsin(t) f(t) dt.$$
 (7.1.1)

We want to approximate this integral using global Gauss-Legendre quadrature as introduced in [Lecture \rightarrow Section 7.4.2].

The nodes and the weights of the n-point Gauss-Legendre quadrature rule on [-1,1] can be computed using the provided C++ function

```
QuadRule gaussquad(const unsigned n);
```

which initializes a data structure describing a quadrature rule

```
struct QuadRule {
   Eigen::VectorXd nodes_;
   Eigen::VectorXd weights_;
};
```

The names of the fields tell their function. The function <code>gaussquad()</code> can be found in <code>gaussquad.hpp</code>.

```
(7-1.a) ☑ (30 min.) Based on MATPLOTLIBCPP write a C++ function
```

```
template <class Function>
double gaussConv(const Function& fh, const double I_ex, const
  unsigned N);
```

that produces an *appropriate* plot of the quadrature error of the Gauss-Legendre quadrature rules versus the number n = 1, ..., 50 of quadrature points. The function should output the best approximation of

the integral, i.e. the approximated value at the maximum number of quadrature points. Here f is a functor object [Lecture \rightarrow Section 0.3.3] providing the function f in (7.1.1). I_ex is the exact solution of the integral and N is the maximum number of quadrature points.

Save your convergence plot for $f(t) = \sinh(t)$ as GaussConv.png.

```
HIDDEN HINT 1 for (7-1.a) \rightarrow 7-1-1-0:hcv.pdf
```

SOLUTION for (7-1.a)
$$\rightarrow$$
 7-1-1-1:gaug1.pdf

(7-1.b) \Box (15 min.) Describe qualitatively (type of convergence) and quantitatively (speed of convergence) the asymptotic (for $n \to \infty$) convergence of the quadrature error you observe in (7-1.a).

```
HIDDEN HINT 1 for (7-1.b) \rightarrow 7-1-2-0:h1.pdf
```

```
SOLUTION for (7-1.b) \rightarrow 7-1-2-1:gauq2.pdf
```

```
HIDDEN HINT 1 for (7-1.c) \rightarrow 7-1-3-0:sp3h0.pdf
```

HIDDEN HINT 2 for (7-1.c) \rightarrow 7-1-3-1:sp3h1.pdf

Solution for (7-1.c) \rightarrow 7-1-3-2:gauq3.pdf

Now, write a C++ function

```
template <class Function>
double gaussConvCV(const Function& f, const double I_ex, const
  unsigned N);
```

which plots (employing MATPLOTLIBCPP) the quadrature error versus the number $n=1,\ldots,50$ of quadrature points for the integral obtained in the previous subtask and returns the best approximation of the integral, i.e. the approximated value at the maximum number of quadrature points.

Again, as in Sub-problem (7-1.a), choose $f(t) = \sinh(t)$ and save your convergence plot as GaussConvCV.png.

SOLUTION for (7-1.d) \rightarrow 7-1-4-0:gauq4.pdf

Analogous question as in Sub-problem (7-1.b): describe qualitatively the asymptotic (for $n \to \infty$) convergence of the quadrature error you observe in Sub-problem (7-1.d).

SOLUTION for (7-1.e)
$$\rightarrow$$
 7-1-5-0: gaug2.pdf

Explain the difference between observations made in Sub-problem (7-1.b) and Sub-problem (7-1.d).

SOLUTION for (7-1.f)
$$\rightarrow$$
 7-1-6-0: gaug5.pdf

End Problem 7-1, 100 min.

Problem 7-2: Generalized "Hermite-type" quadrature formula

In this exercise we will meet a new class of quadrature formulas and then use one of them in an application.

This is a purely theoretical exercise. You may want to have a look at the methods used in [Lecture \rightarrow Ex. 7.4.2.2], as they will come handy.

(7-2.a) \square Determine $A, B, C, x_1 \in \mathbb{R}$ such that the quadrature formula

$$\int_0^1 f(x)dx \approx Af(0) + Bf'(0) + Cf(x_1)$$
 (7.2.1)

is exact for polynomials of highest possible degree.

SOLUTION for (7-2.a)
$$\rightarrow$$
 7-2-1-0:herm1s.pdf

(7-2.b) \Box Compute an approximation of z(2) using the quadrature rule of Eq. (7.2.1), where the function z is defined as the solution of the initial value problem

$$z'(t) = \frac{t}{1+t^2}$$
 , $z(1) = 1$. (7.2.2)

HIDDEN HINT 1 for (7-2.b) \rightarrow 7-2-2-0:herm2h.pdf

SOLUTION for (7-2.b)
$$\rightarrow$$
 7-2-2-1:herm2s.pdf

End Problem 7-2

Problem 7-3: Numerical Quadrature of Improper Integrals

In analysis you have seen so-called improper integrals over unbounded intervals, which, due to a fast decay of the integrand have a finite value. Occasionally such integrals have to evaluated numerically and this problem will be concerned with corresponding quadrature rules.

This problem makes extensive use of substitution rules for integrals, implementation in C++ is requested.

We want to devise a numerical method for the computation of improper integrals of the form $\int_{-\infty}^{\infty} f(t)dt$ for continuous functions $f: \mathbb{R} \to \mathbb{R}$ that decay sufficiently fast for $|t| \to \infty$ (such that they are integrable on \mathbb{R}).

A first option is the truncation of the domain to a bounded interval [-b,b], $b \le \infty$, that is, we approximate:

$$\int_{-\infty}^{\infty} f(t)dt \approx \int_{-b}^{b} f(t)dt$$

and then use a standard quadrature rule like Gauss-Legendre quadrature [Lecture \rightarrow Def. 7.4.2.18] on [-b,b].

(7-3.a) \Box (10 min.) For the integrand $g(t) := 1/(1+t^2)$ determine b such that the truncation error E_T satisfies:

$$E_T := \left| \int_{-\infty}^{\infty} g(t)dt - \int_{-b}^{b} g(t)dt \right| \le 10^{-6} . \tag{7.3.1}$$

HIDDEN HINT 1 for (7-3.a) \rightarrow 7-3-1-0:imp1s1.pdf

SOLUTION for (7-3.a)
$$\rightarrow$$
 7-3-1-1:op1sol.pdf

(7-3.b) (5 min.) What is the algorithmic difficulty faced in the implementation of the truncation approach for a generic integrand?

SOLUTION for (7-3.b)
$$\rightarrow$$
 7-3-2-0:op1dif.pdf

A second option is the transformation of the improper integral to a bounded domain by substitution. For instance, we may use the map $t = \cot(s) := \frac{\cos s}{\sin s}$, $\cot :]0, \pi[\to] - \infty, \infty[$.

(7-3.c) \odot (15 min.) Into which integral does the substitution $t = \cot(s)$ convert $\int_{-\infty}^{\infty} f(t) dt$?

SOLUTION for (7-3.c)
$$\rightarrow$$
 7-3-3-0:op2sub.pdf

(7-3.d) \odot (10 min.) Write down the transformed integral explicitly for the integrand $g(t) := \frac{1}{1+t^2}$. Simplify the integrand.

HIDDEN HINT 1 for (7-3.d) \rightarrow 7-3-4-0:h1ti.pdf

SOLUTION for (7-3.d)
$$\rightarrow$$
 7-3-4-1:op2com.pdf

Write a C++ function that uses the previous transformation together with the n-point Gauss-Legendre quadrature rule to evaluate $\int_{-\infty}^{\infty} f(t)dt$:

```
template <typename Function>
double quadinf(int n, Function && f);
```

f passes an object that provides an evaluation operator of the form:

```
double operator() (double x) const;
```

You can use the function golubwelsh () implemented in the file golubwelsh.hpp that computes nodes and weights for Gauss quadrature rules using the Golub-Welsch algorithm discussed in [Lecture \rightarrow Rem. 7.4.2.23], see also [Lecture \rightarrow Code 7.4.2.24].

```
HIDDEN HINT 1 for (7-3.e) \rightarrow 7-3-5-0:hlf1.pdf
```

Solution for (7-3.e)
$$\rightarrow$$
 7-3-5-1:op2imp.pdf

```
void cvgQuadInf(void);
```

in order to study the convergence for $n \to \infty$ of the quadrature method implemented in the previous subproblem for the integrand $h(t) := \exp(-(t-1)^2)$ (shifted Gaussian). To compute the error you can use that $\int_{-\infty}^{\infty} h(t)dt = \sqrt{\pi}$.

Use MATPLOTLIBCPP to create a suitable plot of the quadrature error versus the number of quadrature nodes. What kind of convergence do you observe?

SOLUTION for (7-3.f)
$$\rightarrow$$
 7-3-6-0:testsol.pdf

End Problem 7-3, 75 min.

Problem 7-4: Nested numerical quadrature

This problem addresses numerical quadrature over a two-dimensional domain. It turns out that this problem can be reduced to two one-dimensional integrals, which are amenable to the techniques covered in [Lecture \rightarrow Chapter 7].

Implementation in C++ requested connected with [Lecture \rightarrow Chapter 7].

A laser beam has intensity

$$I(x,y) := \exp(-\alpha((x-p)^2 + (y-q)^2)), \quad x,y \in \mathbb{R}, \quad \alpha > 0,$$
 (7.4.1)

on the plane orthogonal to the direction of the beam. The beam is directed orthogonal to the x-y-plane.

Note that the radiant power absorbed by a surface is the integral of the intensity over the surface.

(7-4.a) (10 min.) Write down a formula involving only 1D integrals and giving the radiant power absorbed by the triangle

$$\Delta := \{(x, y)^T \in \mathbb{R}^2 \mid x \ge 0, y \ge 0, x + y \le 1\}.$$

as a double integral.

SOLUTION for (7-4.a) \rightarrow 7-4-1-0:neq1s.pdf

that evaluates an n-point quadrature formula as introduced in [Lecture \rightarrow Section 7.2] for an integrand passed in f on domain [a,b]. It should rely on the quadrature rule on the reference interval [-1,1] that is supplied through an object of type **QuadRule**:

```
struct QuadRule { Eigen::VectorXd nodes_, weights_; };
```

(The vectors weights and nodes denote the weights and nodes, respectively, of a quadrature rule on the reference interval [-1,1], see [Lecture \rightarrow Def. 7.2.0.1].)

```
Solution for (7-4.b) \rightarrow 7-4-2-0:neq2s.pdf
```

(7-4.c) □ (30 min.) [depends on Sub-problem (7-4.a)]

Write a C++ function

```
template <class Function>
double gaussquadtriangle(const Function &f, const unsigned N)
```

for the computation of the integral

$$\int_{\Lambda} f(x,y) \, dx \, dy \tag{7.4.3}$$

using nested n-point, 1D Gauss quadrature [Lecture \rightarrow Section 7.4] in combination with the function evalquad() of Sub-problem (7-4.b). Use the function gaussquad() from gaussquad.hpp to obtain weights and nodes of Gauss quadrature formulas on [-1,1].

```
HIDDEN HINT 1 for (7-4.c) \rightarrow 7-4-3-0:neq3h1.pdf

HIDDEN HINT 2 for (7-4.c) \rightarrow 7-4-3-1:neq3h2.pdf

SOLUTION for (7-4.c) \rightarrow 7-4-3-2:neq3s.pdf

(7-4.d) ① (30 min.) Write a C++ function

void convtest2DQuad (unsigned int nmax = 20);
```

that applies the function <code>gaussquadtriangle()</code> from Sub-problem (7-4.c) to approximate $\int_{\Delta} I(x,y) \, \mathrm{d}x \mathrm{d}y$ using the parameters $\alpha=1$, p=0, q=0. Tabulate and plot (using MATPLOTLIBCPP, storing the figure in <code>convergence.png</code>) the quadrature error w.r.t. to the number of nodes $n=1,2,3,\ldots,20$ using the "exact" value of the integral $I\approx0.366046550000405$. What kind of convergence do you observe and why?

```
HIDDEN HINT 1 for (7-4.d) \rightarrow 7-4-4-0:neq4h.pdf
SOLUTION for (7-4.d) \rightarrow 7-4-4-1:neq4s.pdf
```

End Problem 7-4, 85 min.

Problem 7-5: Quadrature plots

In this problem, we will try and deduce the type of quadrature and the integrand from an error plot.

This is a purely theoretical problem discussing convergence of quadrature rules, see [Lecture \rightarrow Lemma 7.4.3.6], [Lecture \rightarrow Ex. 7.4.3.9], [Lecture \rightarrow § 7.5.0.9], [Lecture \rightarrow Exp. 7.5.0.12].

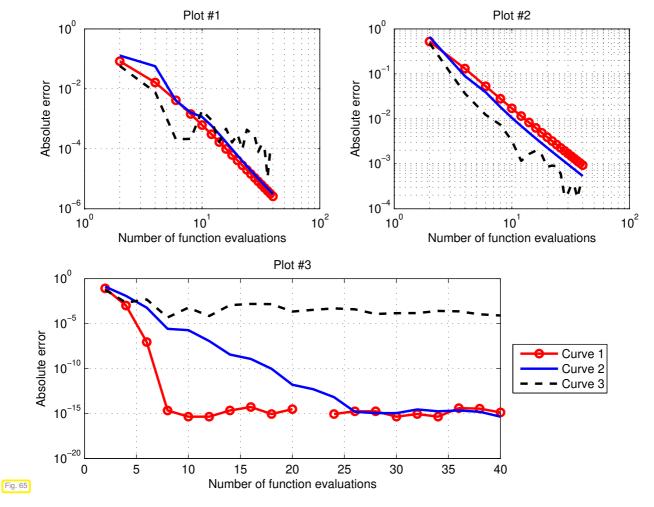
We consider three different functions on the interval I = [0, 1], which have the following properties:

function A: $f_A \in C^{\infty}(I)$, $f_A \notin \mathcal{P}_k \ \forall \ k \in \mathbb{N}$; function B: $f_B \in C^0(I)$, $f_B \notin C^1(I)$; function C: $f_C \in \mathcal{P}_{12}$,

where \mathcal{P}_k is the space of the polynomials of degree at most k defined on I. The following quadrature rules are applied to these functions:

- quadrature rule A is a global Gauss quadrature;
- quadrature rule B is a composite trapezoidal rule;
- quadrature rule C is a composite 2-point Gauss quadrature.

The corresponding absolute values of the quadrature errors are plotted against the number of function evaluations in the figure below. Notice that only the quadrature errors obtained with an even number of function evaluations are shown.



(7-5.a) (15 min.) Match the three plots (plot #1, #2 and #3) with the three quadrature rules (quadrature rule A, B, and C). Justify your answer.

HIDDEN HINT 1 for (7-5.a) \rightarrow 7-5-1-0:plot1h.pdf

SOLUTION for (7-5.a) \rightarrow 7-5-1-1:plot1s.pdf

 \blacktriangle

(7-5.b) \odot (15 min.) The quadrature error curves for a particular function f_A , f_B and f_C are plotted in the same style (curve 1 as red line with small circles, curve 2 means the blue solid line, curve 3 is the black dashed line). Which curve corresponds to which function (f_A , f_B , f_C)? Justify your answer.

SOLUTION for (7-5.b) \rightarrow 7-5-2-0:plot2s.pdf

 \blacktriangle

End Problem 7-5, 30 min.

Problem 7-6: Weighted Gauss quadrature

In [Lecture \rightarrow Rem. 7.4.3.10] we saw, how, in special cases, integrals with singular integrands can be transformed into integrals with smooth integrands, which allow much faster convergence of Gauss quadrature. This problem examines an alternative method: If the type of the singular behavior of the integrand is known, special quadrature rules, called weighted Gauss rules, can be designed that enjoy fast convergence.

This problem relies on the material of [Lecture \rightarrow Section 7.4] and requests a little implementation in C++.

The development of an alternative quadrature formula for a particular singular integral relies on the Chebyshev polynomials of the second kind U_n , defined as

$$U_n(t) = \frac{\sin((n+1)\arccos t)}{\sin(\arccos t)}, \quad n \in \mathbb{N}.$$

The U_n are **orthogonal polynomials** with respect to a weighted L^2 inner product (see [Lecture \to Eq. (6.3.2.3)]) with weight function given by $w(\tau) = \sqrt{1-\tau^2}$.

(7-6.b) \odot (15 min.) Show that the U_n satisfy the 3-term recursion

$$U_{n+1}(t) = 2tU_n(t) - U_{n-1}(t), \qquad U_0(t) = 1, \qquad U_1(t) = 2t,$$

for every $n \ge 1$.

Solution for (7-6.b) \rightarrow 7-6-2-0:wgg1s.pdf

(7-6.c) \odot (5 min.) Show that $U_n \in \mathcal{P}_n$ and has leading coefficient 2^n .

SOLUTION for (7-6.c) \rightarrow 7-6-3-0:wgq2s.pdf

(7-6.d) \odot (15 min.) Show that for every $m, n \in \mathbb{N}_0$, we have

$$\int_{-1}^{1} \sqrt{1-t^2} \, U_m(t) U_n(t) \, dt = \frac{\pi}{2} \delta_{mn}.$$

SOLUTION for (7-6.d) \rightarrow 7-6-4-0:wgq3s.pdf

(7-6.e) (10 min.) What are the zeros ξ_j^n (j = 1, ..., n) of U_n , $n \ge 1$? Give an explicit formula similar to the formula for the Chebyshev nodes in [-1, 1].

SOLUTION for (7-6.e)
$$\rightarrow$$
 7-6-5-0:wgg4s.pdf

Show that the choice of weights

$$w_j = \frac{\pi}{n+1} \sin^2\left(\frac{j}{n+1}\pi\right), \quad j = 1, \dots, n,$$

ensures that the quadrature formula

$$Q_n^U(f) = \sum_{j=1}^n w_j f(\xi_j^n)$$
 (7.6.2)

provides the exact value of the integral

$$I := \int_{-1}^{1} \sqrt{1 - t^2} f(t) \, \mathrm{d}t \tag{7.6.3}$$

for $f \in \mathcal{P}_{n-1}$ (assuming exact arithmetic).

SOLUTION for (7-6.f) \rightarrow 7-6-6-0:wgq5s.pdf

(7-6.g) Show that the quadrature formula (7.6.2) gives the exact value of (7.6.3) even for every $f \in \mathcal{P}_{2n-1}$.

HIDDEN HINT 1 for (7-6.g) \rightarrow 7-6-7-0:wgq6h.pdf

SOLUTION for (7-6.g)
$$\rightarrow$$
 7-6-7-1:wqg6s.pdf

(7-6.h) Show that the quadrature error

$$|Q_n^U(f) - W(f)|$$

decays to 0 exponentially as $n \to \infty$ for every $f \in C^{\infty}([-1,1])$ that admits an analytic extension to an open subset of the complex plane.

HIDDEN HINT 1 for (7-6.h) \rightarrow 7-6-8-0:wgq7h.pdf

SOLUTION for (7-6.h) \rightarrow 7-6-8-1:wqg7s.pdf

(7-6.i) Write a C++ function

```
template < typename Function >
double quadU(const Function & f, const unsigned n);
```

that gives $Q_n^U(f)$ as output, where f is an object with an evaluation operator, like a lambda function, representing f, e.g.

```
auto f = [] (double & t) { return 1/(2 + \exp(3*t)); };
```

SOLUTION for (7-6.i) \rightarrow 7-6-9-0:wqg8s.pdf

(7-6.j) Devise a C++ function

```
void testQuadU(unsigned int nmax = 20);
```

that, for the function $f(t) = 1/(2 + e^{3t})$ and $n = 1, \ldots, n_{\text{max}}$ plots (employing MATPLOTLIBCPP) the quadrature error $E_n(f) = |W(f) - Q_n^U(f)|$ using the "exact" value W(f) = 0.483296828976607. Estimate the parameter $0 \le q < 1$ in the asymptotic decay law $E_n(f) \approx Cq^n$ characterising (sharp) exponential convergence, see [Lecture \rightarrow Def. 6.2.2.7].

SOLUTION for (7-6.j) \rightarrow 7-6-10-0:wgq9s.pdf

End Problem 7-6, 85 min.

Problem 7-7: Quadrature by transformation

In [Lecture \rightarrow Rem. 7.4.3.10] a suitable transformation converted a singular integrand into a smooth one. In this problem we see another example. Also Problem 7-9 and Problem 7-1 cover this topic.

This problem does not contain an implementation part. Study [Lecture \rightarrow Rem. 7.4.3.10] to learn about regularizing transformations.

For $f \in C^0([0,2])$ we consider the definite improper integral with *singular integrand*

$$I(f) := \int_0^2 \sqrt{\frac{2-t}{t}} f(t) \, \mathrm{d}t \,. \tag{7.7.1}$$

HIDDEN HINT 1 for (7-7.a) \rightarrow 7-7-1-0:hsi1.pdf

HIDDEN HINT 2 for (7-7.a) \rightarrow 7-7-1-1:hsi5.pdf

SOLUTION for (7-7.a) \rightarrow 7-7-1-2:s1t.pdf

(7-7.b) \square (20 min.) For $n \in \mathbb{N}^*$, derive a family of 2n-point quadrature formulas

$$Q_n(f) = \sum_{j=1}^{2n} w_j^n f(c_j^n), \qquad f \in C([0,2])$$

for which $Q_n(f)$ converges to I(f) exponentially as $n \to \infty$, if $f \in C^{\infty}([0,2])$.

Remark. A more precise statement of the problem would be ", if f has an analytic extension to a complex neighborhood of [0,2]". This is, how the statement $f \in C^{\infty}([0,2])$ above should be read.

SOLUTION for (7-7.b) \rightarrow 7-7-2-0:qtf1.pdf

(7-7.c) \Box (10 min.) Given n > 0, determine the maximal $d \in \mathbb{N}$ such that $I(p) = Q_n(p)$ for every polynomial $p \in \mathcal{P}_{d-1}$.

Solution for (7-7.c) \rightarrow 7-7-3-0:1tf2ZerosLegendre5s.pdf

End Problem 7-7, 30 min.

Problem 7-8: Discretization of an integral operator

In this problem we consider a completely new concept, an integral operator of the form

$$(\mathsf{T}f)(x) = \int_I k(x,y)f(y)\,\mathrm{d}y\,,\quad x\in I\subset\mathbb{R}\,,\tag{7.8.1}$$

where the kernel k is continuous on $I \times I$, $I \subset \mathbb{R}$ a closed interval. It maps a function on I to another function on I. Such integral operators are very common in models of continuous physics. Of course, their numerical treatment heavily draws on numerical quadrature.

Requires knowledge about Gauss quadrature from [Lecture \rightarrow Section 7.4] and involves substantial implementation in C++.

Given $f \in C^0([0,1])$, the integral

$$g(x) := \int_0^1 e^{|x-y|} f(y) \, dy$$

defines a function $g \in C^0([0,1])$. We approximate the integral by means of n-point Gauss quadrature, which yields a function $g_n(x)$.

(7-8.a) $oxed{:}$ Let $\{\xi_j^n\}_{j=1}^n$ be the nodes for the Gauss quadrature on the interval [0,1]. Assume that the nodes are ordered, namely $\xi_j^n < \xi_{j+1}^n$ for every $j=1,\ldots,n-1$. We can write

$$(g_n(\xi_l^n))_{l=1}^n = M(f(\xi_j^n))_{j=1}^n,$$

for a suitable matrix $M \in \mathbb{R}^{n,n}$. Give a formula for the entries of M.

SOLUTION for (7-8.a)
$$\rightarrow$$
 7-8-1-0:ongp1.pdf

(7-8.b) Implement a function

template < typename Function>

```
Eigen::VectorXd comp_g_gausspts(Function f, unsigned int n)
```

that computes the *n*-vector $(g_n(\xi_l^n))_{l=1}^n$ with optimal complexity O(n) (excluding the computation of the nodes and weights), where f is an object with an evaluation operator, e.g. a lambda function, that represents the function f.

You may use the provided function

```
void gaussrule(int n, Eigen::VectorXd & w, Eigen::VectorXd & xi)
```

that computes the weights w and the ordered nodes xi relative to the n-point Gauss quadrature on the interval [-1,1].

Solution for (7-8.b)
$$\rightarrow$$
 7-8-2-0:ongp2.pdf

(7-8.c) Test your implementation and write a C++ function

double testCompGGaussPts();

that computes $g(\xi_{11}^{21})$ for $f(y) = e^{-|0.5-y|}$. What result do you expect?

SOLUTION for (7-8.c) \rightarrow 7-8-3-0:ongp3.pdf

End Problem 7-8

Problem 7-9: Efficient quadrature of singular integrands

We know that smoothness of the integrand is essential for fast covergence of high-order numerical quadrature like Gauss quadrature rules. In some cases smart transformation can convert the integral into another with a smooth integrand. This problem uses this trick for the sake of efficient numerical quadrature of non-smooth integrands with a special structure.

Before you tackle this problem, read about regularization of integrands by transformation, cf. [Lecture \rightarrow Rem. 7.4.3.10]. For other problems on the same topic see Problem 7-7 and Problem 7-1.

Our task is to develop quadrature formulas for integrals of the form:

$$W(f) := \int_{-1}^{1} \sqrt{1 - t^2} f(t) dt, \tag{7.9.1}$$

where f possesses an analytic extension to a complex neighbourhood of [-1,1].

(7-9.a) \bigcirc (10 min.) Understand how to use the C++ function

```
QuadRule gauleg (unsigned int n);
```

(contained in gauleg.hpp), which returns a structure QuadRule containing nodes (x_j) and weights (w_j) of a Gauss-Legendre quadrature [Lecture \rightarrow Def. 7.4.2.18] on [-1,1] with n nodes.

(7-9.b) Study [Lecture \rightarrow § 7.4.3.2] in order to learn about the convergence of Gauss-Legendre quadrature. What can you say about the least expected convergence for an integrand $f \in C^r([a,b])$? What about convergence in the case $f \in C^\infty([a,b])$?

```
SOLUTION for (7-9.b) \rightarrow 7-9-2-0:effs1.pdf
```

(7-9.c) (40 min.) Based on the function gauleg(), implement a C++ function

```
template <class Function>
double quadsingint(const Function& f, const unsigned n);
```

that approximately evaluates W(f) using 2n evaluations of f. An object of type Function must provide an evaluation operator

```
double operator (double t) const;
```

Ensure that, as $n \to \infty$, the error of your implementation has asymptotic exponential convergence to zero.

```
HIDDEN HINT 1 for (7-9.c) \rightarrow 7-9-3-0:h21.pdf
```

HIDDEN HINT 2 for (7-9.c) \rightarrow 7-9-3-1:h23.pdf

```
Solution for (7-9.c) \rightarrow 7-9-3-2:effs2.pdf
```

(7-9.d) \odot Give formulas for the nodes c_j and weights \tilde{w}_j of a 2n-point quadrature rule on [-1,1], whose application to the integrand f will produce the same results as the function quadsingint that you implemented in the previous subproblem.

```
SOLUTION for (7-9.d) \rightarrow 7-9-4-0:effs3.pdf
```

(7-9.e) 2 Implement a C++ function

```
void tabAndPlotQuadErr()
```

that tabulates and plots the quadrature error:

$$\epsilon_i := |W(f) - ext{quadsingint(f,n)}|$$

for $f(t) := \frac{1}{2 + \exp(3t)}$ and n = 1, 2, ..., 25. Then describe the type of convergence observed, see [Lecture \rightarrow Def. 6.2.2.7].

SOLUTION for (7-9.e) \rightarrow 7-9-5-0:effs5.pdf

End Problem 7-9, 65 min.

Problem 7-10: Zeros of orthogonal polynomials

This problem combines elementary methods for finding zeros from [Lecture \rightarrow Section 8.4] and 3-term recursions satisfied by orthogonal polynomials, *cf.* [Lecture \rightarrow Thm. 6.3.2.14]. This will permit us to find the zeros of Legendre polynomials, the so-called Gauss nodes (see [Lecture \rightarrow Def. 7.4.2.18]).

Connected with [Lecture \rightarrow Section 7.4] and [Lecture \rightarrow Section 8.4]

The zeros of the Legendre polynomial P_n (see [Lecture \to Def. 7.4.2.16]) are the n Gauss points ξ_j^n , $j=1,\ldots,n$. In this problem we compute the Gauss points by zero-finding methods applied to P_n . The 3-term recursion [Lecture \to Eq. (7.4.2.21)] for Legendre polynomials will play an essential role. Moreover, recall that, by definition, the Legendre polynomials are $L^2(]-1,1[)$ -orthogonal.

(7-10.a) $oxed{\square}$ Prove the following interleaving property of the zeros of the Legendre polynomials. For all $n \in \mathbb{N}_0$ we have:

$$-1 < \xi_j^n < \xi_j^{n-1} < \xi_{j+1}^n < 1, \quad j = 1, \dots, n-1.$$

HIDDEN HINT 1 for (7-10.a) \rightarrow 7-10-1-0:ZerosLegendre1h.pdf

SOLUTION for (7-10.a) \rightarrow 7-10-1-1:ZerosLegendre1s.pdf

(7-10.b) • By differentiating [Lecture \rightarrow Eq. (7.4.2.21)] derive a combined 3-term recursion for the sequences $(P_n)_n$ and $(P'_n)_n$.

SOLUTION for (7-10.b) \rightarrow 7-10-2-0: ZerosLegendre2s.pdf

(7-10.c) Use the recursions obtained in (7-10.b) to write a C++ function

```
void legvals(const Eigen::VectorXd& x,
Eigen::MatrixXd& Lx, Eigen::MatrixXd& DLx);
```

that fills the matrices Lx and DLx in $\mathbb{R}^{N\times(n+1)}$ with the values $\{P_k(x_j)\}_{jk}$ and $\{P'_k(x_j)\}_{jk}$, $j=0,\ldots,N-1$ and $k=0,\ldots,n$, for an input vector $x\in\mathbb{R}^N$ (passed in x):

SOLUTION for (7-10.c) \rightarrow 7-10-3-0: ZerosLegendre3s.pdf

(7-10.d) $oldsymbol{\cdot}$ We can compute the zeros of P_k , $k=1,\ldots,n$, by means of the secant rule (see [Lecture \rightarrow § 8.4.2.28]) using the endpoints $\{-1,1\}$ of the interval and the zeros of the previous Legendre polynomial as initial guesses; see (7-10.a). We opt for a correction-based termination criterion (see [Lecture \rightarrow Section 8.2.3]) based on prescribed relative and absolute tolerance (see [Lecture \rightarrow Code 8.4.2.31]).

Write a C++ function that computes the Gauss points $\xi_j^k \in [-1,1]$, $j=1,\ldots,k$ and $k=1,\ldots,n$, using the zero finding approach outlined above:

The Gauss points should be returned in an upper triangular $n \times n$ -matrix.

HIDDEN HINT 1 for (7-10.d) \rightarrow 7-10-4-0:ZerosLegendre4h.pdf

SOLUTION for (7-10.d) \rightarrow 7-10-4-1:ZerosLegendre4s.pdf

(7-10.e) Validate your implementation of the function gaussPts with n = 8 by computing the values of the Legendre polynomials in the zeros obtained (use the function leggrals). Explain the failure of

of the Legendre polynomials in the zeros obtained (use the function legvals). Explain the failure of the method.

```
HIDDEN HINT 1 for (7-10.e) \rightarrow 7-10-5-0:ZerosLegendre5h.pdf SOLUTION for (7-10.e) \rightarrow 7-10-5-1:ZerosLegendre5s.pdf
```

(7-10.f) Fix your function <code>gaussPts</code> taking into account the above considerations. You should use the *regula falsi*, that is a variant of the secant method in which, at each step, we choose the old iterate to keep depending on the signs of the function. More precisely, given two approximations $x^{(k)}$, $x^{(k-1)}$ of a zero in which the function f has different signs, compute another approximation $x^{(k+1)}$ as zero of the secant. Use this as the next iterate, but then chose as $x^{(k)}$ the value $z \in \{x^{(k)}, x^{(k-1)}\}$, for which $sign\left[f(x^{(k+1)})\right] \neq sign[f(z)]$. This ensures that f has always a different sign in the last two iterates.

The regula falsi variation of the secant method can be easily implemented with a little modification of the code.

SOLUTION for (7-10.f) \rightarrow 7-10-6-0:ZerosLegendre6s.pdf

End Problem 7-10

Problem 7-11: Quadrature rules and quadrature errors

When sequences of quadrature rules are applied to evaluate $\int_a^b f(t) dt$, then the asymptotic convergence of quadrature error in terms of $n \to \infty$, n = 1 the number of quadrature nodes, is determined both by

- (i) the smoothness of the integrand f,
- (ii) family of quadrature rules used.

This problems asks you to guess the family of quadrature rules from error curves.

```
To prepare for this problem study [Lecture \rightarrow Section 7.4], in particular [Lecture \rightarrow Ex. 7.4.3.9], and also [Lecture \rightarrow Exp. 7.5.0.12], [Lecture \rightarrow Exp. 7.5.0.16]
```

In this problem we consider three different families of quadrature rules, for which we use the following abbreviations:

- **(TR)** $\hat{=}$ Equidistant trapezoidal rule, see [Lecture \rightarrow Eq. (7.5.0.17)],
- **(SR)** \triangleq Equidistant Simpson rule, see [Lecture \rightarrow Eq. (7.5.0.5)],
- **(GR)** $\hat{=}$ Gauss-Legendre quadrature rules, as defined in [Lecture \rightarrow Def. 7.4.2.18].
- (7-11.a) \odot (15 min.) The following C++ function implements the n-point equidistant trapezoidal rule, $n \in \mathbb{N}$, $n \ge 2$, for the approximate evaluation of $\int_a^b f(t) \, \mathrm{d}t$. The integrand is passed through the functor f

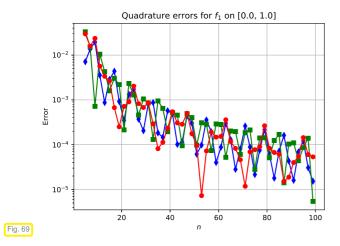
Supplement the missing C++ code in the boxes.

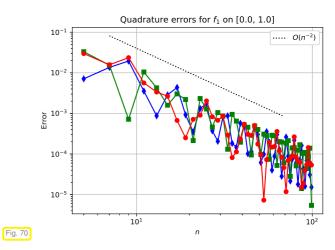
```
SOLUTION for (7-11.a) \rightarrow 7-11-1-0:sa.pdf
```

(7-11.b) (15 min.) The following plots display the quadrature errors for the approximation of the integral

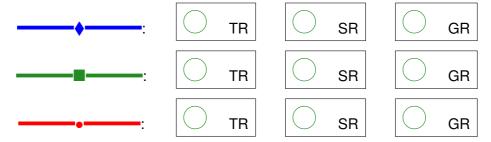
$$\int_0^1 f_1(t) \, \mathrm{d}t \quad , \quad f_1(t) := |\sin(5t)| \ ,$$

by the three families of quadrature rules and numbers $n = 4, \dots, 100$ of quadrature points.





Tick the families of quadrature rules, if the curve correctly represents the expected behavior of its quadrature error as a function of n.

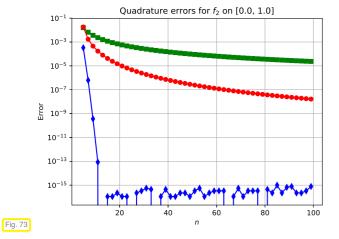


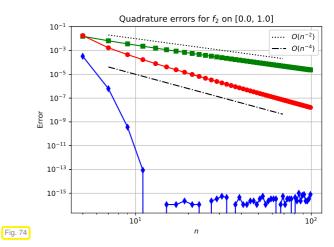
SOLUTION for (7-11.b) \rightarrow 7-11-2-0:s1.pdf

Below we have plotted the quadrature errors for the approximation of the

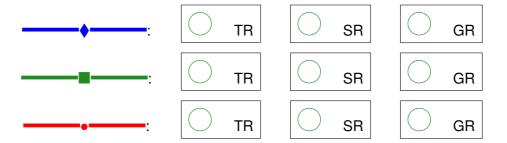
$$\int_0^1 f_2(t) dt , f_2(t) := |\sin(5t)|^2,$$

for the three families of quadrature rules and numbers n = 4, ..., 100 of quadrature points.





Select the quadrature rule, if the curve correctly represents the expected behavior of its quadrature error as a function of n.

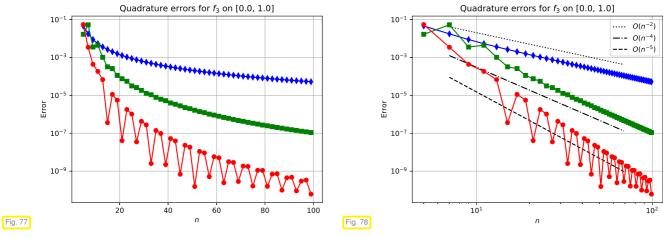


SOLUTION for (7-11.c) \rightarrow 7-11-3-0:s2.pdf

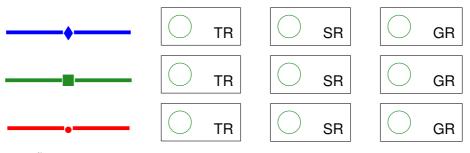
(7-11.d) (15 min.) The plots below display the quadrature errors for

$$\int_0^1 f_3(t) \, \mathrm{d}t \, , \quad f_3(t) := |\sin(5t)|^5 \, ,$$

and for equidistant trapezoidal rules (TR), equidistant Simpson rules (SR), and Gauss-Legendre rules (GR) as a function of the number n of quadrature points, $n = 4, \ldots, 100$.



Select those quadrature rules, for which the curve correctly represents the expected behavior of its quadrature error as a function of n.



SOLUTION for (7-11.d) \rightarrow 7-11-4-0:s3.pdf

(7-11.e) (20 min.) Decide whether the following assertions about quadrature rules are true or false

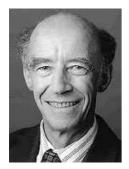
- (i) For any given set of n quadrature nodes $\in [0,1]$ there is a quadrature rule of order n with those nodes.
 - true false
- (ii) Let Q_n designate an n-point quadrature rule of order n on [0,1]. Then

$$f \geq 0 \quad \Rightarrow \quad Q_n(f) \geq 0 \; .$$

true false

End Pro	blom 7-11 90 n	nin					
SOLUTION	N for (7-11.e) \rightarrow	7-11-5-0	:sx.pdf				A
		true				false	
	family of equidista uadrature nodes fo			enjoys exponentia	ıl conv	ergence in t	the number <i>n</i>
		true				false	
(iii) An n	An n -point quadrature rule of order $2n-1$ has all positive weights						

Problem 7-12: Clenshaw-Curtis-Fejer Quadrature Formula



In [J. Waldvogel, Fast construction of the Fejér and Clenshaw-Curtis quadrature rules, BIT, 46 (2006), pp. 195–202.] Jörg Waldvogel, former professor at the Seminar of Applied Mathematics of ETH Zürich, elaborates formulas for the weights and nodes of so-called Clenshaw-Curtis-Fejer quadrature rules, $\it cf.$ [Lecture \rightarrow Rem. 7.3.0.8]. This problem is concerned with their efficient implementation.

This problem requires familiarity with the discrete Fourier transform (DFT) from [Lecture \rightarrow Section 4.2] and the FFT algorithm [Lecture \rightarrow Section 4.3] and its use in Eigen. Knowledge about the concept and the order of quadrature formulas is taken for granted.

The following class provides the weights and nodes of a quadrature formula of the Clenshaw-Curtis-Fejer family on the interval [-1,1].

C++ code 7.12.1: Definition of quadrature class

```
class CCFQuadRule {
    public:
3
     explicit CCFQuadRule(unsigned int n);
     ~CCFQuadRule() = default;
     const Eigen::VectorXd &nodes() const { return nodes_; }
     const Eigen::VectorXd &weights() const { return weights_; }
8
9
    private:
10
     Eigen::VectorXd nodes_;
11
     Eigen::VectorXd weights_;
  };
13
```

Get it on ₩ GitLab (clenshawcurtisfejer.hpp).



Note that passing an integer n to the constructor builds an n+1-point quadrature formula!

The methods nodes () and weights () return vectors whose entries supply the nodes c_j and weights w_i , i = 0, ..., n, of the quadrature formula (C++ indexing used!).

This is the listing of the implementation of the constructor:

C++ code 7.12.2: Constructor of CCFQuadRule

```
CCFQuadRule::CCFQuadRule(unsigned int n) : weights_(n + 1) {
2
     assert(n > 0);
3
     nodes = (Eigen::ArrayXd::LinSpaced(n + 1, 0, n) * M PI / n).cos().matrix();
     const unsigned int m = n / 2; // Integer division!
5
     const Eigen::ArrayXd idx = Eigen::ArrayXd::LinSpaced(m, 1.0, m);
6
     const Eigen::ArrayXd cos_arg = 2.0 * M_PI * idx / n;
     Eigen::ArrayXd fac = 2.0 / (4 * idx.pow(2) - 1.0);
8
     if (n \% 2 == 0) fac [m - 1] /= 2.0;
9
     weights_{0} = (1.0 - fac.sum()) / n;
10
11
     for (unsigned int j = 1; j < n; ++j) {
       weights_[j] = (1.0 - (fac * (j * cos_arg).cos()).sum()) * 2.0 / n;
12
     }
13
```

```
weights_[n] = weights_[0];

Get it on  GitLab (clenshawcurtisfejer.hpp).
```

Recall that

- Eigen::ArrayXd::LinSpaced(n,a,b) builds a sequence of n values $\left(a + \frac{b-a}{n-1}j\right)_{j=0}^{n-1}$,
- the cos () method of Eigen::ArrayXd applies the cosine function to all entries of the sequence,
 and
- the call of pow (2) for Eigen::ArrayXd squares all members of the sequence,
- that the sum () method computes the sum of all entries of a vector/matrix in EIGEN.

```
SOLUTION for (7-12.a) \rightarrow 7-12-1-0:ccfs1.pdf
```

```
SOLUTION for (7-12.b) \rightarrow 7-12-2-0:.pdf (7-12.c) \odot (15 min.) The C++ function unsigned int determineOrderCCF (unsigned int n);
```

is supposed to return the **order** [Lecture \rightarrow Def. 7.4.1.1] of the quadrature formulas provided by an object of the class **CCFQuadRule** when constructed with a parameter n.

Fill in the missing parts of the code in the boxes (valid C++ syntax!):

```
unsigned int determineOrderCCF(unsigned int n) {
  const CCFQuadRule ccfqr(n);
  unsigned int d = 0;
                            // Degree of monomial
  double quaderr; // Quadrature error
  double I_ex;
                 // Exact integral value
                  // Value computed with CCF quadrature rule
  double I_qr;
  const double abstol = std::numeric_limits<double>::epsilon() * n
  const double reltol = std::numeric_limits<double>::epsilon() * n
    * 20;
  do {
    const Eigen::ArrayXd pdvals = ccfqr.
       .array().pow(d);
    I_qr =
                                                                 ;
    I ex =
    quaderr = std::abs(I_qr - I_ex);
    d = d + 1;
  } while ((quaderr <= abstol) || (quaderr <= reltol * I_ex));</pre>
  return
```

}

SOLUTION for (7-12.c)
$$\rightarrow$$
 7-12-3-0:cfqrs3.pdf (7-12.d) \odot (90 min.) [depends on Sub-problem (7-12.b)]

Your task is to devise a more efficient implementation of (the constructor of) the class **CCFQuadRule**. To that end you have created a copy, renamed it to **CCFQuadRule_Fast** and deleted parts of the constructor. Now you should re-implement the constructor in the file clenshawcurtisfejer.hpp so that the asymptotic effort for computing the nodes and weights of the Clenshaw-Curtis-Fejer quadrature rule becomes $O(n \log n)$ for $n \to \infty$.

Solution for (7-12.d)
$$\rightarrow$$
 7-12-4-0:spx1.pdf

End Problem 7-12, 130 min.

Problem 7-13: Smooth integrand by transformation (ArcCos)

Integrals with non-smooth integrands can *sometimes* be treated by a suitable transformation, which converts the integrand into a smooth function, see [Lecture \rightarrow Rem. 7.4.3.10]. In this problem we study an example.

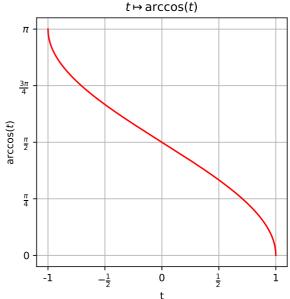
You should be familiar with the transformation of quadrature rules [Lecture \rightarrow Rem. 7.2.0.4] and Gaussian-Legendre quadrature formulas [Lecture \rightarrow Section 7.4], [Lecture \rightarrow Section 7.4.2]. You must be able to perform integration by substitution.

Given a smooth function $f:[-1,1]\to\mathbb{R}$ our task is to come up with a numerical approximation of

$$I(f) := \int_{-1}^{1} \arccos(t) f(t) dt$$
. (7.13.1)

We want to approximate this integral using global Gauss-Legendre quadrature as introduced in [Lecture \rightarrow Section 7.4.2].

The function $t\mapsto \arccos t \rhd$ ($\hat{=}$ the inverse of \cos on $[0,\pi]$)



The nodes and the weights of the n-point Gauss-Legendre quadrature rule on [-1,1] can be computed for $n \le 256$ using the provided C++ function

Fig. 81

```
QuadRule gaussquad(const unsigned int n);
```

which initializes a data structure describing a quadrature rule

```
struct QuadRule {
   QuadRule() = default;
   explicit QuadRule(unsigned int n) : nodes_(n), weights_(n) {}
   Eigen::VectorXd nodes_;
   Eigen::VectorXd weights_;
};
```

The names of the data members tell their function. The function <code>gaussquad()</code> is declared in <code>gaussquad.hpp</code>.

```
(7-13.a) ☑ (15 min.) In the file arccosquad.hpp implement a C++ function void testConvGaussQuad();
```

that prints a table that allows you to predict *qualitatively and quantitatively* the *asymptotic behavior* (w.r.t. $n \to \infty$) of the quadrature error of n-point Gauss-Legendre numerical quadrature, when directly applied to (7.13.1) for $f(t) = \frac{1}{1+e^t}$.

Based on that table describe the observed empiric convergence of the quadrature error as a function of the number of quadrature nodes qualitatively and quantitatively. Explain, how you arrive at your

conclusions.

```
HIDDEN HINT 1 for (7-13.a) \rightarrow 7-13-1-0:slhacos.pdf

HIDDEN HINT 2 for (7-13.a) \rightarrow 7-13-1-1:slref.pdf

Solution for (7-13.a) \rightarrow 7-13-1-2:.pdf
```

```
SOLUTION for (7-13.b) \rightarrow 7-13-2-0:s2trf.pdf
```

In the file arccosquad.hpp complete the code of the C++ function

that

- expects the functor argument f to pass a function $f:[-1,1] \to \mathbb{R}$ and to provide an evaluation operator **operator double** () (**double**) **const**,
- uses *n* evaluations of the function *f* to compute an approximation $I_n(f)$ of I(f) from (7.13.1),
- that achieves exponential asymptotic convergence $I_n(f) \to I(f)$ for $n \to \infty$, if the function f possesses an analytic extension beyond [-1,1].

```
SOLUTION for (7-13.c) \rightarrow 7-13-3-0:s3sol.pdf
```

Analoguous to Sub-problem (7-13.a) in the file arccosquad.hpp write a C++ function

```
void testConvTrfGaussQuad();
```

that outputs (in a suitable table) information that allows you to predict *qualitatively and quantitatively* the *asymptotic behavior* (w.r.t. $n \to \infty$) of the approximation error of your implementation of arccosWeightedQuad() for $f(t) = \frac{1}{1+e^t}$.

Characterize the observed empiric convergence of the quadrature error as a function of n qualitatively and quantitatively. Justify your conclusions.

```
Solution for (7-13.d) \rightarrow 7-13-4-0:.pdf
```

End Problem 7-13, 50 min.

Problem 7-14: Aspects of Numerical Quadrature

A quadrature rule approximates the integral of a (continuous) function by means of a weighted sum of function values at so-called quadrature nodes/points. Numerical analysis studies the order of quadrature rules and the asymptotic convergence of the quadrature error as a function of the number of quadrature points.

This problem is linked with [Lecture \rightarrow Section 7.2], [Lecture \rightarrow Section 7.4.1], and [Lecture \rightarrow Section 7.5] and also requires familiarity with [Lecture \rightarrow § 7.5.0.18].

(7-14.a) (4 min.) Complete the definition of the order of a quadrature rule:

```
Definition cf. [Lecture \rightarrow Def. 7.4.1.1]. Order of a quadrature rule
```

The **order** of a quadrature rule $Q_n : C^0([a,b]) \to \mathbb{R}$ is defined as

$$\operatorname{order}(Q_n) := \left\{ m \in \mathbb{N}_0 \colon \ Q_n(\right] \right\}$$

```
SOLUTION for (7-14.a) \rightarrow 7-14-1-0:bqns1.pdf
```

expects that the input vectors c and w have the same length and contain the weights and nodes of a quadrature formula on [-1,1]. It returns the order of that quadrature formula. Supplement the missing parts of the following listing by writing valid C++ code in the boxes.

```
unsigned int checkQuadOrder(
const Eigen::VectorXd &c, const Eigen::VectorXd &w,
const double tol = 1.0E-10) {
  const unsigned long N = c.size();
  assert(N == w.size());
  for (unsigned int d = 0; d <
                                          ; ++d) {
    double s = 0.0;
    for (int j = 0; j <
                                 ; ++j) {
                    * std::pow(
      s +=
                                        , d);
    double val = (d \% 2 == 0) ?
                                                            ;
    if (std::abs(s - val) > tol) {
      return
                     ;
    }
  return
```

```
}
```

HIDDEN HINT 1 for $(7-14.b) \rightarrow 7-14-2-0: nqms2h1.pdf$ SOLUTION for $(7-14.b) \rightarrow 7-14-2-1: nqm2s.pdf$ (7-14.c) \bigcirc (3 min.) The C++ function

template <typename FUNCTOR>

double trapezoidalRule(
FUNCTOR &&f, double a, double b, unsigned int N);

evaluates the **equidistant trapezoidal quadrature rule** with \mathbb{N} nodes on the interval [a,b] for a (continuous) function passed through the functor object f.

Write valid C++ code in the boxes of the following listing so that the function meets this specification.

```
template <typename FUNCTOR>
double trapezoidalRule(
FUNCTOR &&f, double a, double b, unsigned int N) {
  assert(a < b);</pre>
  assert (N >= 2);
  const double h =
  double t = a+h;
  double s =
  for (int j = 1; j <
                                ; ++ †) {
                           ;
    t += h;
  s +=
  return
                    ;
}
```

HIDDEN HINT 1 for (7-14.c) \rightarrow 7-14-3-0:nqm3h1.pdf SOLUTION for (7-14.c) \rightarrow 7-14-3-1:nqms3.pdf

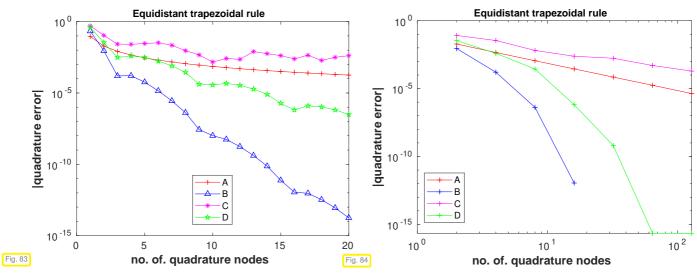
(7-14.d) (4 min.) The equidistant trapezoidal quadrature rule as implemented in Subproblem (7-14.c) is used for the approximate evaluation of the integrals

$$\int_{a}^{b} f_{\alpha,\beta}(t) dt$$
 , $a = 0.33$, $b = 1.33$,

where the family of functions $f_{\alpha,\beta}: \mathbb{R} \to \mathbb{R}$ is defined as

$$f_{\alpha,\beta}(t) := \sqrt{|1 + \alpha \sin(\beta t)|}, \quad \alpha, \beta, t \in \mathbb{R}.$$

The following plots display the modulus of the quadrature error for certain members of this family of functions as a function of the number of quadrature nodes.



State, which graph (A,B,C, or D) belongs to which pair of parameters α , β (defining the integrand $f_{\alpha,\beta}$).

Parameters	\leftrightarrow	Graph
(1) $\alpha=0.9,~~\beta=2\pi$	\leftrightarrow	
(2) $\alpha = 0.5$, $\beta = \pi$	\leftrightarrow	
(3) $\alpha=1.1$, $\beta=2\pi$	\leftrightarrow	
(4) $\alpha=0.5$, $\beta=2\pi$	\leftrightarrow	

SOLUTION for (7-14.d) \rightarrow 7-14-4-0:nqms4.pdf

End Problem 7-14, 15 min.

Problem 7-15: Convolution Quadrature

Convolution quadrature is a special numerical technique for the evaluation of the convolution of two causal functions, when one is given in the form of its Laplace transform. This problem studies formulas for the weights of convolution quadrature and their efficient numerical computation.

This problem requires knowledge of the discrete Fourier transform (DFT) [Lecture \rightarrow Def. 4.2.1.18], its implementation in EIGEN, some complex analysis, and numerical quadrature for periodic analytic functions [Lecture \rightarrow § 7.5.0.18].

The **convolution** of two causal integrable functions $f,g:\mathbb{R}^+:=[0,\infty[\to\mathbb{R}],$ written as f*g, is a function $\mathbb{R}^+\to\mathbb{R}$ defined as

$$(f * g)(t) := \int_0^t f(t - \tau)g(\tau) d\tau, \quad t \ge 0.$$
 (7.15.1)

Given a fixed "timestep size" $\tau > 0$, we approximate f * g at equidistant times $t_n := \tau n$, $n \in \mathbb{N}_0$, by

$$(f * g)(t_n) \approx \sum_{\ell=0}^{n} w_{n-\ell}^{\tau} g_{\ell} , \quad g_{\ell} := g(t_{\ell}) , \quad n \in \mathbb{N}_0 .$$
 (7.15.2)

with suitable convolution quadrature weights $w_j^{\tau} \in \mathbb{C}, j \in \mathbb{N}_0$.

In the sequel let the function f be fixed and denote by $F:D\subset\mathbb{C}\to\mathbb{C}$ its **Laplace transform**, which is supposed to be known.

Assumption 7.15.3. Holomorphy of the Laplace transform

The Laplace transform $F:D\subset\mathbb{C}\to\mathbb{C}$ of f is *analytic* (holomorphic) on an open set $D\subset\mathbb{C}$ with

$$\mathbb{C}^+ := \{ z \in \mathbb{C} : \operatorname{Re} z > 0 \} \subset D ,$$

that is, F is analytic in the right half-plane of \mathbb{C} .

As explained in [HAS16] the convolution quadrature weights w_j^{τ} in (7.15.2) are given by the following complex contour integrals: for any 0 < r < 1

$$w_j^{\tau} = \frac{1}{2\pi \imath} (-1)^j \int_{|z|=r} \frac{1}{z^{j+1}} F\left(\frac{1-z}{\tau}\right) dz, \quad j \in \mathbb{Z}.$$
 (7.15.4)

With the parameterization $z(\varphi):=re^{2\pi\imath\varphi},\,\varphi\in[0,1[$, we convert (7.15.4) into

$$w_j^{\tau} = r^{-j} \int_0^1 \psi_j(\varphi) \, \mathrm{d}\varphi \,, \quad \psi_j(\varphi) := e^{-2\pi \imath j \varphi} \cdot F\left(\frac{1 - r e^{2\pi \imath \varphi}}{\tau}\right) \,, \quad j \in \mathbb{Z} \,. \tag{7.15.5}$$

(7-15.a) (10 min.)

We identify $\psi_j : \mathbb{R} \to \mathbb{C}$ from (7.15.5) with its extension to a neighborhood of the real axis in the complex plane. Merely taking for granted Ass. 7.15.3, what is the *maximal* $\rho > 0$ such that ψ_j is analytic in the strip

$$S_{\rho}:=\{z\in\mathbb{C}:\ |\operatorname{Im}z|<\rho\}$$
 ?

Give that maximal ρ as a function of r:

$$\rho(r) =$$

SOLUTION for (7-15.a)
$$\rightarrow$$
 7-15-1-0:cqs1.pdf

$$w_j^{ au} pprox w_j^{ au,N} =$$

SOLUTION for (7-15.b) \rightarrow 7-15-2-0:cqs2.pdf

In the file convolution quadrature.hpp implement an efficient C++ function

```
template <typename FFUNCTOR>
```

```
Eigen::VectorXcd compute_cq_weights(
   FFUNCTOR&& F, unsigned int N, double tau, double r = 0.0);
```

that computes the approximate convolution quadrature weights $w_j^{\tau,N}$ by means of the (N+1)-point equidistant trapezoidal rule and for $j=0,\ldots,N$, and returns them as the components of a complex column vector. The functor object $\mathbb F$ passes the Laplace transform F. The parameter $\mathbb F$ specifies the radius F of the integration contour. If $\mathbb F=0.0$, then a default value is chosen, see the code template.

"Efficient" means that the function should enjoy an asymptotic computational cost of $O(N \log N)$ for $N \to \infty$.

SOLUTION for (7-15.c)
$$\rightarrow$$
 7-15-3-0:cqs3.pdf

(7-15.d) □ (5 min.) [depends on Sub-problem (7-15.a), Sub-problem (7-15.b)]

Assume computations in exact arithmetic. How will the error introduced into the computation of the convolution quadrature weights by numerical quadrature by means of the (N+1)-point equidistant trapezoidal rule, that is $|w_j^{\tau} - w_j^{\tau,N}|$, behave as a function of N asymptotically for $N \to \infty$? Explain your answer.

Solution for (7-15.d)
$$\rightarrow$$
 7-15-4-0:qc4s4.pdf

End Problem 7-15, 44 min.

Problem 7-16: On Quadrature Formulas

This problem addresses some key ideas of numerical quadrature, offering a repetition of considerations already covered in class.

```
This problem is linked to [Lecture \rightarrow Section 7.2] and [Lecture \rightarrow Section 7.4.1].
```

Descendants of the following abstract base class are supposed to provide weights and nodes for a quadrature formula on [-1,1].

```
struct QuadFormulaRef {
    virtual const std::vector<double>& nodes() const = 0;
    virtual const std::vector<double>& weights() const = 0;
  };
The following class is a type for a quadrature formula on a general interval
 [a,b]\subset \mathbb{R}.
    class QuadFormula {
     public:
      QuadFormula(const QuadFormulaRef& qf, double a, double b);
      std::pair<double, double> interval() const { return {a_, b_}; }
      const std::vector<double>& nodes() const { return nodes_; }
      const std::vector<double>& weights() const { return weights_; }
     private:
      double a_;
      double b_;
      std::vector<double> nodes ;
      std::vector<double> weights_;
    };
```

The constructor takes a quadrature formula object qf defined for the reference interval [-1,1] plus the interval bounds a and b and initializes the class member variables nodes_ and weights_ with the nodes and weights of the quadrature formula transformed to [a,b]. Fill the blanks in Code 7.16.1 to complete the implementation of the constructor.

```
C++ code 7.16.1: Constructor of QuadFormula
  QuadFormula::QuadFormula(const QuadFormulaRef& qf, double a, double b)
       : a_(a), b_(b), nodes_(qf.nodes()), weights_(qf.weights()) {
2
     const double len = b - a;
3
     const unsigned int n_pts = nodes_.size();
     for (unsigned int j = 0; j < n_pts; ++j) {
6
7
       nodes_[j] =
8
       weights_[j] *=
9
    }
  }
10
```

```
SOLUTION for (7-16.a) \rightarrow 7-16-1-0:.pdf 

(7-16.b) \square (20 min.) The C++ function checkQuadOrder () whose incomplete listing is given
```

as Code 7.16.3 determines and returns the order of a quadrature formula on [-1,1] given to it via the qf argument.

Complete the listing so that the function can perform its function.

C++ code 7.16.3: Code for function checkQuadOrder() unsigned int checkQuadOrder(const QuadFormulaRef& qf, const double to I = 1.0E-10) { 2 const std::vector<double>& c = qf.nodes(); 3 const std::vector<double>& w = qf.weights(); const unsigned int n_pts = c.size(); 6 std::vector<double> monom_vals(n_pts, 1.0); 7 for (unsigned int d = 0; $d < n_pts$; ++d) { double s = 0.0; for (unsigned int j = 0; $j < n_pts$; ++j) $s += w[j] * monom_vals[j]$; 10 **double** val = 2.0 / (2*d + 1); 11 > tol) return 2 * d; 12 13 for (unsigned int j = 0; j < n_pts; ++j) monom_vals[j] *=</pre> 14 for (unsigned int j = 0; $j < n_pts$; ++j) $s += w[j] * monom_vals[j]$; 15 if > tol) return 2 * d + 1; 16 for (unsigned int j = 0; j < n_pts; ++j) monom_vals[j] *=</pre> 17 18 return 2 * n_pts; 19 } 20

```
HIDDEN HINT 1 for (7-16.b) \rightarrow 7-16-2-0:qfsbh1.pdf
```

SOLUTION for (7-16.b) \rightarrow 7-16-2-1:.pdf

End Problem 7-16, 30 min.

Chapter 8

Iterative Methods for Non-Linear Systems of Equations

Problem 8-1: Convergent Newton iteration

As explained in [Lecture \rightarrow Section 8.4.2.1], the convergence of Newton's method in 1D may only be local. This problem investigates a particular setting, in which global convergence can be expected.

This is a purely theoretical problem practising "intuitive mathematical reasoning".

We recall the notion of a *convex function* [Lecture \to Def. 5.3.1.4] and its geometric meaning [Lecture \to Fig. 158]: A differentiable function $f:[a,b] \mapsto \mathbb{R}$ is convex if and only if its graph lies on or above its tangent at any point. Equivalently, differentiable function $f:[a,b] \mapsto \mathbb{R}$ is convex, if and only if its derivative is non-decreasing.

(8-1.a) (20 min.) Give a "graphical proof" of the following statement:

Theorem 8.1.1. Global convergence of Newton's method in 1D for convex monotone functions

If F(x) belongs to $C^2(\mathbb{R})$, is strictly increasing, is convex, and has a unique zero, then the Newton iteration [Lecture $\to Eq. (8.4.2.1)$] for F(x) = 0 is well defined and will converge to the zero of F(x) for any initial guess $x^{(0)} \in \mathbb{R}$.

SOLUTION for (8-1.a) \rightarrow 8-1-1-0:Conv1s.pdf

End Problem 8-1, (25 min.)

Problem 8-2: Code quiz

A frequently encountered drudgery in scientific computing is the use and modification of poorly documented code. This makes it necessary to understand the ideas behind the code first. Now we practice this in the case of a simple iterative method.

Related to [Lecture \rightarrow Section 8.4.2.1], involves implementation in C++.

C++ code 8.2.1: Undocumented function

```
double myfunction (double x) {
     constexpr double dy = 0.693147180559945; // = std::log(2)
3
     double y = 0.;
4
     while (x > 2. * std :: sqrt(2.)) {
5
       x /= 2.;
6
       y += dy;
     while (x < std :: sqrt(2.)) {
       x *= 2.;
10
       y -= dy;
11
12
     double z = x - 1.; //
13
     double dz = x * std :: exp(-z) - 1.;
14
     while (std::abs(dz / z) > std::numeric_limits < double > ::epsilon()) {
15
       z += dz;
       dz = x * std :: exp(-z) - 1.;
17
18
     return y + z + dz; //
19
20
  }
```

```
HIDDEN HINT 1 for (8-2.a) \rightarrow 8-2-1-0:Code1ha.pdf
```

HIDDEN HINT 2 for (8-2.a) \rightarrow 8-2-1-1:Code1hb.pdf

SOLUTION for (8-2.a) \rightarrow 8-2-1-2:Code1s.pdf

```
Solution for (8-2.b) \rightarrow 8-2-2-0:Code2s.pdf
```

(8-2.c) (10 min.) Explain what the while loop in lines 13–18 of Code 8.2.1 does.

```
Solution for (8-2.c) \rightarrow 8-2-3-0:Code3s.pdf
```

(8-2.d) (10 min.) Explain the conditional expression in the last while loop.

```
SOLUTION for (8-2.d) \rightarrow 8-2-4-0:Code4s.pdf
```

Write a C++ function

```
double myfunction_modified(double x);
```

where you replace the last while-loop with a fixed number of iterations that, nevertheless, guarantee that the result has a relative accuracy EPS. Derive that required minimal number of iterations!

```
HIDDEN HINT 1 for (8-2.e) \rightarrow 8-2-5-0:s5h1.pdf
Solution for (8-2.e) \rightarrow 8-2-5-1:Code5s.pdf
```

 \blacktriangle

End Problem 8-2, 110 min.

Problem 8-3: Newton's method for $F(x) := \arctan x = 0$

The merely local convergence of Newton's method is notorious, see [Lecture \rightarrow Section 8.5.2] and [Lecture \rightarrow Ex. 8.5.4.1]. The failure of the convergence is often caused by the overshooting of Newton correction. In this problem we try to understand the observations made in [Lecture \rightarrow Ex. 8.5.4.1].

Moderate implementation in C++ is requested.

```
HIDDEN HINT 1 for (8-3.a) \rightarrow 8-3-1-0:NewtonArctan1ha.pdf
```

HIDDEN HINT 2 for (8-3.a) \rightarrow 8-3-1-1:NewtonArctan1hb.pdf

SOLUTION for (8-3.a) \rightarrow 8-3-1-2:NewtonArctan1s.pdf

```
double newton_arctan(double x0_ = 2.0);
```

that uses Newton's method to find an approximation of such $x^{(0)}$. The argument x_0 can be used to supply an initial guess for the iteration.

The considerations from Sub-problem (8-3.a) suggest that we use an initial guess in [1,2], we opt for $x^{(0)} = 2$.

SOLUTION for (8-3.b) \rightarrow 8-3-2-0:NewtonArctan2s.pdf

End Problem 8-3, 40 min.

Problem 8-4: A derivative-free iterative scheme for finding zeros

[Lecture \rightarrow Rem. 8.2.2.12] shows how to detect the order of convergence of an iterative method from a numerical experiment. In this problem we study the so-called Steffensen's method, which is a derivative-free iterative method for finding zeros of functions in 1D.

Related to [Lecture \rightarrow Section 8.2.2] and requests simple C++ coding.

Let $f:[a,b]\mapsto \mathbb{R}$ be twice continuously differentiable with $f(x^*)=0$ and $f'(x^*)\neq 0$. Consider the iteration defined by

$$x^{(n+1)} := x^{(n)} - \frac{f(x^{(n)})}{g(x^{(n)})}$$
, where $g(x) = \frac{f(x+f(x)) - f(x)}{f(x)}$. (8.4.1)

(8-4.a) • (15 min.) Write a C++ function for the Steffensen's method:

```
template <class Function>
double steffensen(Function &&f, double x0)
```

f is a functor object or lambda function providing f, x0 is the initial guess $x^{(0)}$. Terminate the iteration when the computed sequence of approximations becomes stationary, see [Lecture \rightarrow Code 8.2.3.6] for an example.

```
SOLUTION for (8-4.a) \rightarrow 8-4-1-0:Quad1s.pdf
```

Write a C++ function

```
void testSteffensen();
```

that applies your implementation of steffensen() to find the zero of the function $f(x) = xe^x - 1$ (see [Lecture \rightarrow Exp. 8.3.1.3]). Use $x^{(0)} = 1$ as initial guess.

```
Solution for (8-4.b) \rightarrow 8-4-2-0:sq2.pdf
```

We now want to investigate the order of Steffensen's method. To do so extend your implementation of steffensen() to

```
template <typename Function>
Eigen::VectorXd steffensen_log(Function &&f, double x0, LOGGER
    &&log = [](double) -> void {});
```

such that it logs each iterate using a callable \log () as described in [Lecture \rightarrow § 0.3.3.4]. Then use this new facility to implement a function

```
void orderSteffensen();
```

that tabulates values from which you can read of the order of Steffensen's method when applied to the test case of Sub-problem (8-4.b).

```
HIDDEN HINT 1 for (8-4.c) \rightarrow 8-4-3-0:steffh1.pdf
```

Solution for (8-4.c)
$$\rightarrow$$
 8-4-3-1:sq3.pdf

(8-4.d) \odot (10 min.) For the test case of Sub-problem (8-4.b) the function g(x) from (8.4.1) contains a term like e^{xe^x} . Therefore it grows very fast in x and the method cannot start for large $x^{(0)}$. How can you modify the function f (keeping the same zero) in order to allow the choice of a larger initial guess?

HIDDEN HINT 1 for (8-4.d) \rightarrow 8-4-4-0:Quad2h.pdf

SOLUTION for (8-4.d) \rightarrow 8-4-4-1:Quad2s.pdf

End Problem 8-4, 50 min.

Problem 8-5: Order-p convergent iterations

In [Lecture \to Section 8.2.2] we investigated the speed of convergence of iterative methods for the solution of a general non-linear problem $F(\mathbf{x})=0$ and introduced the notion of **convergence** of order $p\geq 1$, see [Lecture \to Def. 8.2.2.10]. This problem highlights the fact that for p>1 convergence may not be guaranteed, even if the error norm estimate of [Lecture \to Def. 8.2.2.10] may hold for some $\mathbf{x}^*\in\mathbb{R}^n$ and all iterates $\mathbf{x}^{(k)}\in\mathbb{R}^n$.

Problem with a theoretical focus and a little C++ coding. Relies on the concepts introduced in [Lecture \rightarrow Section 8.2.2]

Given $\mathbf{x}^* \in \mathbb{R}^n$, suppose that a sequence $\mathbf{x}^{(k)}$ satisfies [Lecture \rightarrow Def. 8.2.2.10]:

$$\exists C > 0: \|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| \le C\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^p \quad \forall k \text{ and } p > 1.$$
 (8.5.1)

(8-5.a) \odot (20 min.) Determine $\epsilon_0 > 0$ as large as possible such that we have

$$\|\mathbf{x}^{(0)} - \mathbf{x}^*\| \le \epsilon_0 \quad \Longrightarrow \quad \lim_{k \to \infty} \mathbf{x}^{(k)} = \mathbf{x}^*.$$
 (8.5.2)

In other words, ϵ_0 tells us which distance of the initial guess from \mathbf{x}^* still guarantees local convergence.

HIDDEN HINT 1 for (8-5.a) \rightarrow 8-5-1-0:ocvh1.pdf

SOLUTION for (8-5.a)
$$\rightarrow$$
 8-5-1-1:OrdC1s.pdf

SOLUTION for (8-5.b)
$$\rightarrow$$
 8-5-2-0:OrdC2s.pdf

void kminplot();

that uses MATPLOTLIBCPP's plt::imshow() to create a image of $k_{\min}(\epsilon_0,\tau)$ for the values $p=1.5,\ C=2$ and

$$\epsilon_0 \in \left[0, C^{\frac{1}{1-p}}\right]^2$$
 , $\tau \in [10^{-5}, 10^{-1}]$.

Use a logarithmic axis scale for τ .

SOLUTION for (8-5.c)
$$\rightarrow$$
 8-5-3-0:OrdC3s.pdf

End Problem 8-5, 65 min.

Problem 8-6: Order of convergence from error recursion

In [Lecture \rightarrow Exp. 8.4.2.32] we have observed *fractional* orders of convergence ([Lecture \rightarrow Def. 8.2.2.10]) for both the secant method and the quadratic inverse interpolation method. This is fairly typical for 2-point methods in 1D and arises from the underlying recursions for error bounds. This problem addresses how to determine the order of convergence from an abstract error recursion.

A similar analysis is elaborated for the secant method in [Lecture \rightarrow Rem. 8.4.2.33], where a linearised error recursion is given in [Lecture \rightarrow Eq. (8.4.2.37)].

For an iterative scheme producing the sequence $(x^{(n)})_{n\in\mathbb{N}}$, $x^{(n)}\in\mathbb{R}$ we suppose a recursive bound for the norms of the iteration errors of the form

$$||e^{(n+1)}|| \le ||e^{(n)}|| \sqrt{||e^{(n-1)}||}$$
, (8.6.1)

where $e^{(n)} = x^{(n)} - x^*$ is the error of *n*-th iterate.

double testOrder(const unsigned int n = 20);

that guesses the maximal order of convergence of the method from a numerical experiment with n iterations.

```
HIDDEN HINT 1 for (8-6.a) \rightarrow 8-6-1-0:RecursionOrder1h.pdf
SOLUTION for (8-6.a) \rightarrow 8-6-1-1:RecursionOrder1s.pdf
```

```
HIDDEN HINT 1 for (8-6.b) \rightarrow 8-6-2-0:RecursionOrder2ha.pdf
```

HIDDEN HINT 2 for (8-6.b) \rightarrow 8-6-2-1:RecursionOrder2hb.pdf

HIDDEN HINT 3 for (8-6.b) \rightarrow 8-6-2-2: RecursionOrder2hc.pdf

SOLUTION for (8-6.b) \rightarrow 8-6-2-3:RecursionOrder2s.pdf

End Problem 8-6, 50 min.

Problem 8-7: Nonlinear electric circuit

[Lecture \rightarrow Ex. 2.1.0.3] discusses electric circuits with elements that give rise to linear voltage—current dependence, see [Lecture \rightarrow Ex. 2.1.0.3] and [Lecture \rightarrow Ex. 2.8.0.1]. The principles of nodal analysis were explained in these cases.

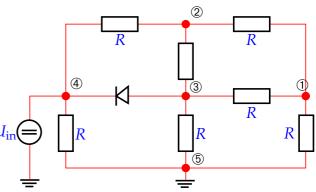
However, the electrical circuits encountered in practise usually feature elements with a *non-linear* current-voltage characteristic. Then nodal analysis leads to non-linear systems of equations as was elaborated in [Lecture \rightarrow Ex. 8.1.0.1]. Please note that transformation to frequency domain is not possible for non-linear circuits so that we will always study the direct current (DC) situation.

In this problem we deal with a very simple non-linear circuit element, a diode. The current I through a diode to which a voltage U is applied can be modelled by the relationship

$$I = \alpha \left(e^{\beta \frac{U}{U_T}} - 1 \right)$$

with suitable parameters $\alpha, \beta > 0$ and the thermal $U_{\rm in}$ voltage U_T .

Now we consider the circuit depicted in Fig. 93 and assume that all resistors have non-dimensional resistance R=1.



The circuit is driven by imposing a voltage U_{in} through an ideal voltage source.

(8-7.a) \odot (20 min.) Study again [Lecture \rightarrow Ex. 2.1.0.3] and [Lecture \rightarrow Ex. 8.1.0.1] to refresh your knowledge about the nodal analysis of electric circuits.

SOLUTION for (8-7.b)
$$\rightarrow$$
 8-7-2-0:NonL1s.pdf

In the file nonlinear_circuit.hpp write an EIGEN-based C++ function

```
void circuit (const alpha, double beta,
const Eigen::VectorXd &Uin, Eigen::VectorXd &Uout)
```

that computes the output voltages Uout, which is defined to be the voltage at node 1 in Fig. 93 for a sorted vector of input voltages Uin (at node 4) and for a thermal voltage $U_T = 0.5$. The parameters alpha, beta pass the (non-dimensional) diode parameters α , β .

Use Newton's method to solve $\mathbf{F}(\mathbf{u}) = \mathbf{0}$ with a relative tolerance of $\tau = 10^{-6}$.

SOLUTION for (8-7.c)
$$\rightarrow$$
 8-7-3-0:NonL2s.pdf

We are interested in the nonlinear effects introduced by the diode. Using MATPLOTLIBCPP implement a C++ function

```
void plotU();
```

that creates a plot of the output voltage $U_{\text{out}} = U_{\text{out}}(U_{\text{in}})$ as a function of the variable input voltage $U_{\text{in}} \in [0, 20]$ (for non-dimensional parameters $\alpha = 8$, $\beta = 1$ and for a thermal voltage $U_T = 0.5$). How can you discern the non-linearity of the circuit in the plot?

Solution for (8-7.d) \rightarrow 8-7-4-0:NonL3s.pdf

End Problem 8-7, 75 min.

Problem 8-8: Julia Set

Julia sets are famous fractal shapes in the complex plane. They are constructed from the basins of attraction of zeros of complex functions when the Newton method is applied to find them.

This problem treats Newton's method in 2D and is related to [Lecture \rightarrow Section 8.5]. You may watch this video by 3Blue1Brown to learn some background, see also here.

In the space C of complex numbers the equation

$$z^3 = 1 (8.8.1)$$

has three solutions:

$$z_1 = 1$$
 , $z_2 = -\frac{1}{2} + \frac{1}{2}\sqrt{3}\imath$, $z_3 = -\frac{1}{2} - \frac{1}{2}\sqrt{3}\imath$, (8.8.2)

the cubic roots of unity.

(8-8.a) • (15 min.) As you know from the analysis course, the complex plane $\mathbb C$ can be identified with $\mathbb R^2$ via $(x,y)\mapsto z=x+\imath y$. Using this identification, convert Eq. (8.8.1) into a system of equations $\mathbf F(\mathbf x)=\mathbf 0$ for a suitable function $\mathbf F:\mathbb R^2\mapsto\mathbb R^2$.

SOLUTION for (8-8.a)
$$\rightarrow$$
 8-8-1-0:JuliaSet1s.pdf

Formulate the Newton iteration [Lecture \rightarrow Eq. (8.5.1.6)] for the non-linear equation $\mathbf{F}(\mathbf{x}) = \mathbf{0}$ with $\mathbf{x} = [x, y]^T$ and \mathbf{F} from Sub-problem (8-8.a).

SOLUTION for (8-8.b)
$$\rightarrow$$
 8-8-2-0:JuliaSet2s.pdf

In the file julia.hpp implement two C++ functions

```
Eigen::Vector2d juliaF(const Eigen::Vector2d& x);
Eigen::Matrix2d juliaDF(const Eigen::Vector2d& x);
```

that return $F(\mathbf{x})$ and the Jacobian $DF(\mathbf{x})$ for for the function $F: \mathbb{R}^2 \to \mathbb{R}^2$ found in Subproblem (8-8.a).

SOLUTION for (8-8.c)
$$\rightarrow$$
 8-8-3-0:sal.pdf

(8-8.d) □ (60 min.) [depends on Sub-problem (8-8.b)]

Denote by $\mathbf{x}^{(k)}$ the iterates produced by Newton method from the previous subproblem with some initial vector $\mathbf{x}^{(0)} \in \mathbb{R}^2$. Depending on $\mathbf{x}^{(0)}$, the sequence $\mathbf{x}^{(k)}$ will either diverge or converge to one of the three cubic roots of unity.

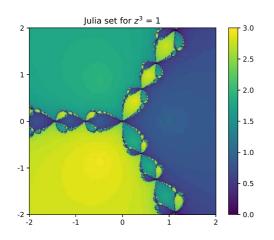
Analyze the behavior of Newton iterations by means of a C++ code using the following approach:

- Use equally spaced points on the domain $[-2,2]^2 \subset \mathbb{R}^2$ as starting points of Newton iterations.
- Color the starting points differently depending on which of the three roots is the limit of the sequence $\mathbf{x}^{(k)}$.

Concretely complete the implementation of the C++ function

```
void julia();
```

supplied in julia.hpp. It is creates and saves a plot in julia.png.



This is how your final plot should look like.

Remark. The boundary of the three so-called domain of attraction of the three different zeros is a fractal set: At whatever scale you look at it, it displays the same complex shape.

Fig. 95

HIDDEN HINT 1 for (8-8.d) \rightarrow 8-8-4-0:JuliaSet3h.pdf Solution for (8-8.d) \rightarrow 8-8-4-1:JuliaSet3s.pdf

End Problem 8-8, 105 min.

Problem 8-9: Modified Newton method

The following problem consists in EIGEN implementation of a modified version of the Newton method (in one dimension [Lecture \rightarrow Section 8.4.2.1] and many dimensions [Lecture \rightarrow Section 8.5]) for the solution of a nonlinear system.

Involves coding in C++. Refresh your knowledge of stopping criteria for iterative methods [Lecture \rightarrow Section 8.2.3].

For the solution of the non-linear system of equations $\mathbf{F}(\mathbf{x}) = \mathbf{0}$ (with $\mathbf{F} : \mathbb{R}^n \to \mathbb{R}^n$), the following iterative method has been proposed:

$$\mathbf{y}^{(k)} = \mathbf{x}^{(k)} + D\mathbf{F}(\mathbf{x}^{(k)})^{-1}\mathbf{F}(\mathbf{x}^{(k)}),$$

$$\mathbf{x}^{(k+1)} = \mathbf{y}^{(k)} - D\mathbf{F}(\mathbf{x}^{(k)})^{-1}\mathbf{F}(\mathbf{y}^{(k)}),$$
(8.9.1)

where $D \mathbf{F}(\mathbf{x}) \in \mathbb{R}^{n,n}$ is the Jacobian matrix of \mathbf{F} evaluated in the point \mathbf{x} .

$$\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$$
, $\mathbf{F}(\mathbf{x}) = \mathbf{F}(\mathbf{y}) \Rightarrow \mathbf{x} = \mathbf{y}$. (8.9.2)

Show that under this assumption the iteration (8.9.1) is *consistent* with $\mathbf{F}(\mathbf{x}) = \mathbf{0}$ in the sense of [Lecture \rightarrow Def. 8.3.1.1], that is, show that, under the assumption that $\mathbf{D} \mathbf{F}(\mathbf{x}^*)$ is regular for every fixed point \mathbf{x}^* of the iteration, $\mathbf{x}^{(k)} = \mathbf{x}^{(0)}$ for every $k \in \mathbb{N}$ if and only if $\mathbf{F}(\mathbf{x}^{(0)}) = \mathbf{0}$.

SOLUTION for (8-9.a)
$$\rightarrow$$
 8-9-1-0:Modils.pdf

that computes a step of the modified Newton method for a *scalar* function \mathbf{F} , that is, for the case n=1.

Here, f is a functor object of type Function passing the function $F: \mathbb{R} \to \mathbb{R}$ and df a functor object of type Jacobian passing the derivative $F': \mathbb{R} \to \mathbb{R}$. Both require an appropriate evaluation operator operator () and have to conform with std::function<double(double)>.

SOLUTION for (8-9.b)
$$\rightarrow$$
 8-9-2-0:Modi2s.pdf

(8-9.c) (30 min.) What is the order of convergence of the method described in (8.9.1)?

To investigate it, write a C++ function

```
void mod_newt_ord();
```

that

- uses the function mod_newt_step to the following scalar equation: $\arctan(x) 0.123 = 0$,
- generates a suitable terminal output that makes it possible to empirically determine the order of convergence, in the sense of [Lecture → Rem. 8.2.2.12],
- implements meaningful stopping criteria, see [Lecture → Section 8.2.3].

Use $x_0 = 5$ as initial guess and tabulate data from which you can read off the order.

```
HIDDEN HINT 1 for (8-9.c) \rightarrow 8-9-3-0:Modi3ha.pdf
```

HIDDEN HINT 2 for (8-9.c) \rightarrow 8-9-3-1:Modi3hb.pdf

SOLUTION for (8-9.c)
$$\rightarrow$$
 8-9-3-2:Modi3s.pdf

(8-9.d) □ (20 min.) Implement a templated C++ function

that *efficiently* realizes a step of the iteration (8.9.1) for a function $\mathbf{F}:\mathbb{R}^n\to\mathbb{R}^n$ passed through the functor \mathbf{f} , which complies with $\mathbf{std}: \mathtt{function} < \mathtt{Vector}(\mathtt{const}\ \mathtt{Vector}\ \mathtt{\&}) >$; The \mathtt{Vector} type can be assumed to have the capabilities of $\mathtt{Eigen}: \mathtt{VectorXd}$, whereas the type $\mathtt{Jacobian}$ must have an evaluation operator that returns an object compatible with $\mathtt{Eigen}: \mathtt{MatrixXd}$.

HIDDEN HINT 1 for (8-9.d) \rightarrow 8-9-4-0:3xh1.pdf

SOLUTION for (8-9.d)
$$\rightarrow$$
 8-9-4-1:s3x.pdf

In the sequel we consider the non-linear system of equations

$$\mathbf{F}(\mathbf{x}) := \mathbf{A}\mathbf{x} + \begin{bmatrix} c_1 e^{x_1} \\ \vdots \\ c_n e^{x_n} \end{bmatrix} = \mathbf{0} , \qquad (8.9.8)$$

where $\mathbf{A} \in \mathbb{R}^{n,n}$ is symmetric positive definite, $\mathbf{x} = [x_j]_{j=1}^n \in \mathbb{R}^n$ is the solution vector, and $c_i \geq 0$, $i = 1, \ldots, n$.

(8-9.e) Compute the Jacobian matrix $D \mathbf{F}(\mathbf{x})$ for \mathbf{F} as given in (8.9.8).

SOLUTION for (8-9.e)
$$\rightarrow 8-9-5-0$$
:s4x.pdf

Based on your implementation of mod_newt_step_system(), create a C++ function

that uses the modified Newton iteration (8.9.1) to solve (8.9.8). The argument A passes the matrix $\bf A$ and $\bf c$ supplies the values c_i .

Stop the iteration and return the approximate solution when the Euclidean norm of the increment $\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}$ relative to the norm of $\mathbf{x}^{(k+1)}$ is smaller than the tolerance passed in tol or if maxit steps of the iterations have been carried out. Use the zero vector as initial guess.

SOLUTION for (8-9.f)
$$\rightarrow$$
 8-9-6-0:s4m.pdf

(8-9.g) (30 min.) [depends on Sub-problem (8-9.f)]

Write a C++ function

```
void mod_newt_sys_test();
```

that tests your implementation of mod_newt_sys() for

$$\mathbf{A} = \begin{bmatrix} 2 & -1 & 0 & 0 \\ -1 & 2 & -1 & 0 \\ 0 & -1 & 2 & -1 \\ 0 & 0 & -1 & 2 \end{bmatrix} , \mathbf{c} = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix},$$

and determines the order of convergence from suitable tabulated values.

HIDDEN HINT 1 for (8-9.g) \rightarrow 8-9-7-0:s4h1.pdf

SOLUTION for (8-9.g) $\rightarrow 8-9-7-1$: Modi4s.pdf

lack

End Problem 8-9, 145 min.

Problem 8-10: Solving a quasi-linear system

In [Lecture \rightarrow § 8.5.1.25] we studied Newton's method for a so-called quasi-linear system of equations, see [Lecture \rightarrow Eq. (8.5.1.26)]. In [Lecture \rightarrow Ex. 8.5.1.30] we then dealt with concrete quasi-linear system of equations and in this problem we will supplement the theoretical considerations from class by implementation in Eigen. We will also learn about a simple fixed point iteration for that system, see [Lecture \rightarrow Section 8.3].

Refresh yourself about the relevant parts of the lecture. You should also try to recall the Sherman-Morrison-Woodbury formula [Lecture \rightarrow Lemma 2.6.0.21].

Consider the *nonlinear* (quasi-linear) system:

$$\mathbf{A}(\mathbf{x})\mathbf{x} = \mathbf{b} , \qquad (8.10.1)$$

as in [Lecture \to Ex. 8.5.1.30]. Here, $\mathbf{A}: \mathbb{R}^n \to \mathbb{R}^{n,n}$ is a matrix-valued function:

where $\|\cdot\|_2$ is the Euclidean norm.

```
SOLUTION for (8-10.a) \rightarrow 8-10-1-0:s1.pdf
```

We consider the fixed point iteration derived in Sub-problem (8-10.a). Implement an *efficient* EIGEN-based C++ function

that computes the iterate $\mathbf{x}^{(k+1)}$ from $\mathbf{x}^{(k)}$.

HIDDEN HINT 1 for (8-10.b) \rightarrow 8-10-2-0:q2h1.pdf

SOLUTION for (8-10.b)
$$\rightarrow$$
 8-10-2-1:QuasiLinear2s.pdf

that finds the solution \mathbf{x}^* of (8.10.1), (8.10.2) with the fixed point method applied to the previous quasi-linear system. Use $\mathbf{x}^{(0)} = \mathbf{b}$ as initial guess. Supply it with a suitable *correction based stopping criterion* as discussed in [Lecture \rightarrow Section 8.2.3] and pass absolute and relative tolerance as arguments atol and rtol.

```
SOLUTION for (8-10.c) \rightarrow 8-10-3-0:QuasiLinear3s.pdf
```

(8-10.d) \Box (10 min.) Let $\mathbf{b} \in \mathbb{R}^n$ be given. Write the recursion formula for the solution of

$$\mathbf{A}(\mathbf{x})\mathbf{x} = \mathbf{b} \tag{8.10.6}$$

with the Newton method.

```
SOLUTION for (8-10.d) \rightarrow 8-10-4-0:QuasiLinear4s.pdf
```

(8-10.e) \odot (15 min.) The matrix $\mathbf{A}(\mathbf{x})$, being symmetric and tri-diagonal, is cheap to invert, see, e.g., [Lecture \rightarrow Rem. 3.3.4.4]. Rewrite the previous iteration from Sub-problem (8-10.d) exploiting the Sherman-Morrison-Woodbury inversion formula for rank-one modifications from [Lecture \rightarrow Lemma 2.6.0.21].

```
SOLUTION for (8-10.e) \rightarrow 8-10-5-0:QuasiLinear5s.pdf
```

(8-10.f) (30 min.) [depends on Sub-problem (8-10.e)]

Write an EIGEN-based C++ function

the realizes a single step of the Newton method as derived in Sub-problem (8-10.e).

```
HIDDEN HINT 1 for (8-10.f) \rightarrow 8-10-6-0:QuasiLinear6ha.pdf
```

```
SOLUTION for (8-10.f) \rightarrow 8-10-6-1:QuasiLinear6s.pdf
```

Repeat Sub-problem (8-10.c) for the Newton method. To that end implement a C++ function

that solves (8.10.1), (8.10.2) with Newton's method. As initial guess use $\mathbf{x}^{(0)} = \mathbf{b}$. The parameters atol and rtol enter the stopping criterion as discussed in [Lecture \rightarrow Section 8.5.3].

```
SOLUTION for (8-10.g) \rightarrow 8-10-7-0:QuasiLinear7s.pdf
```

End Problem 8-10, 95 min.

Problem 8-11: Radioactive Decay

Least-squares estimation is still the main mathematical tool for parameter estimation. This exercise studies a particular (non-linear) least-squares parameter estimation problem.

Requires familiarity with [Lecture → Section 8.7] and involves moderate coding in C++.

For radioactive substances A and B consider the decay chain

$$A \xrightarrow{\lambda_1} B \xrightarrow{\lambda_2} C$$
.

That is A decays into B with a decay rate λ_1 . B itself decays into a third substance C with decay rate λ_2 . The evolution of the amounts $\Phi_A(t)$, $\Phi_B(t)$ of the substances A and B can be modelled by the following coupled system of ordinary differential equations, see [STRLN09]:

$$\frac{d\Phi_A}{dt}(t) = -\lambda_1 \Phi_A(t)
\frac{d\Phi_B}{dt}(t) = -\lambda_2 \Phi_B(t) + \lambda_1 \Phi_A(t) .$$
(8.11.1)

(8-11.a) (30 min.) Calculate the analytical solution of the system (8.11.1) using the method of variation of constants, see [STRNL09].

HIDDEN HINT 1 for (8-11.a) \rightarrow 8-11-1-0:ra1h1.pdf

HIDDEN HINT 2 for (8-11.a) \rightarrow 8-11-1-1:ra1h2.pdf

SOLUTION for (8-11.a) \rightarrow 8-11-1-2:rd1.pdf

(8-11.b) (15 min.) [depends on Sub-problem (8-11.a)]

Unfortunately, only the quantity $\Phi_B(t)$ can be measured and the available values are affected by significant measurement errors. Nevertheless, we would like to estimate the initial quantities A_0 , B_0 and the decay rates λ_1 and λ_2 . This leads to an overdetermined *non-linear* system of equations $\mathbf{F}(\mathbf{x}) = \mathbf{0}$, which has to be solved in least squares sense, see [Lecture \rightarrow Section 8.7]. Here, the vector \mathbf{x} of unknown parameters is $[A_0, B_0, \lambda_1, \lambda_2]^T \in \mathbb{R}^4$.

The measurements (t_i, m_i) , $i \in \{1, ..., N\}$, are stored in the lines of the file decay.txt: t_i are the times of measurements and m_i are the measured values of Φ_B .

Write down the function **F** for the current problem.

Solution for (8-11.b)
$$\rightarrow$$
 8-11-2-0:rd2.pdf

We want to use the Gauss-Newton method (see [Lecture \rightarrow Section 8.7.2]) for solving the non-linear problem described above.

Write down the concrete iteration formulas [Lecture \rightarrow Eq. (8.7.2.1)] for our special problem. In order to do this, compute the Jacobian DF.

Which linear least squares problem is solved in each Gauss-Newton step?

Solution for (8-11.c)
$$\rightarrow$$
 8-11-3-0:rd3.pdf

(8-11.d) (30 min.) [depends on Sub-problem (8-11.c)]

Write a C++ function

```
template <typename F, typename DF>
std::vector<double> GaussNewton(Eigen::Vector4d& x, F&& f, DF&& df,
const double tol = 1e-14)
```

that implements the Gauss-Newton method for estimating the parameters A_0 , B_0 , λ_1 , and λ_2 . The function outputs a vector containing L_{∞} norm of the updates at each step of the iteration. The estimated parameters are to be stored in the object x.

Describe qualitatively and quantitatively the observed convergence, when choosing as initial guess $\mathbf{x}^{(0)} = [1, 1, 1, 0.1]^T$.

```
HIDDEN HINT 1 for (8-11.d) \rightarrow 8-11-4-0:radsh2.pdf

HIDDEN HINT 2 for (8-11.d) \rightarrow 8-11-4-1:Radioactive4h.pdf

SOLUTION for (8-11.d) \rightarrow 8-11-4-2:rd4.pdf
```

End Problem 8-11, 95 min.

Problem 8-12: Approximation of a circle

In this problem we study how a fitting problem arising in computational geometry can be solved using a least squares approach in several ways, leading to different results.

Relies on [Lecture \rightarrow Section 8.7] and involves substantial implementation in C++ based on EIGEN.

Let us consider a sequence of N of points $\in \mathbb{R}^2$ approximately located on a circle (N=8):

A circle with center $[m_1, m_2]^{\top} \in \mathbb{R}^2$ and radius r > 0 is described by

$$(x - m_1)^2 + (y - m_2)^2 = r^2. (8.12.1)$$

8-12.1: linear algebraic fit

(8-12.a) • (15 min.) Plugging the point coordinates (x_i, y_i) into (8.12.1) we obtain an overdetermined linear system with three unknowns

$$m_1, m_2, c := r^2 - m_1^2 - m_2^2$$

Specify the system matrix A and the right-hand side vector b of the corresponding *linear* least squares problem [Lecture \rightarrow Eq. (3.1.3.7)].

SOLUTION for (8-12.a)
$$\rightarrow$$
 8-12-1-0:cala1.pdf

Write a C++ function

```
Eigen::Vector3d circl_alg_fit(const Eigen::VectorXd& x,
const Eigen::VectorXd& y);
```

that receives point coordinates in the vectors \mathbf{x} and \mathbf{y} and solves the overdetermined system in least squares sense and returns a vector $[m_1, m_2, r]^{\top} \in \mathbb{R}^3$ of estimated circle parameters.

SOLUTION for (8-12.b)
$$\rightarrow$$
 8-12-2-0:cala2.pdf

8-12.II: geometric fit

The algebraic approach lacks an intuitive geometrical meaning: minimising the equation residual of (8.12.1) in least squares sense does not necessarily yield the best circle fit in aesthetic sense.

A more appealing approach consist in the minimization of the distances between the data points and the (unknown) circle

$$d_i = \left| \sqrt{(m_1 - x_i)^2 + (m_2 - y_i)^2} - r \right| \qquad i = 1, \dots, N,$$

in the sense of least squares, i.e., determine m_1, m_2 and r such that the sum $\sum_{i=1}^n d_i^2$ is minimal. This is a *non-linear* least squares problem of the form

$$\mathbf{z}^* = \underset{\mathbf{z}}{\operatorname{argmin}} \|\mathbf{F}(\mathbf{z})\|^2$$

see [Lecture \rightarrow Section 8.7].

(8-12.c) \Box (15 min.) Write down the concrete function $\mathbf{F}: \mathbb{R}^{n_1} \to \mathbb{R}^{n_2}$ for this non-linear least squares problem.

You can use the auxiliary values

$$R_i = \sqrt{(x_i - m_1)^2 + (y_i - m_2)^2}$$
, $i = 1, ..., N$.

SOLUTION for (8-12.c) \rightarrow 8-12-3-0:cage1.pdf

Define the functional

$$\Phi(\mathbf{z}) := \frac{1}{2} \|\mathbf{F}(\mathbf{z})\|^2.$$

Compute the Jacobian D **F** of **F**, the gradient $\operatorname{grad} \Phi(\mathbf{z})$ of Φ and the Hessian $\operatorname{H} \Phi(\mathbf{z})$.

SOLUTION for (8-12.d)
$$\rightarrow$$
 8-12-4-0:cage2.pdf

Use a C++ code to find the circle that fit best the data given in the table above, according to the distances d_i . In order to do this, implement a C++ function

that uses the Gauss-Newton method introduced in [Lecture \rightarrow Section 8.7.2] to minimize the functional Φ . The parameters x, y and the return values in z are the same as for circl_alg_fit from Subproblem (8-12.a). The recorder object linked to the pointer err should be passed the current estimate in each step of the Gauss-Newton iteration.

```
SOLUTION for (8-12.e) \rightarrow 8-12-5-0:cage3.pdf
```

(8-12.f) **(30 min.)** [depends on Sub-problem (8-12.d)]

Now solve the same problem using the Newton method for least squares equations as described in [Lecture \rightarrow Section 8.7.1]. To that end, in the file <code>circle_approx.hpp</code>, implement a C++ function

The parameters and return values obey the same specification as those for circl_geo_fit() from Sub-problem (8-12.e).

```
SOLUTION for (8-12.f) \rightarrow 8-12-6-0:cage4.pdf
```

(8-12.g) Compare the convergence of Gauss-Newton and Newton methods implemented in the previous sub-problems. You can use the parameters determined by the algebraic fit as initial guess.

Rely on the error in the parameters measured in the maximum norm.

To that end, in the file circle_approx.hpp, implement a C++ function

SOLUTION for (8-12.g) \rightarrow 8-12-7-0:cage5.pdf

•

8-12.III: Constrained non-linear fitting

As you will know, for $a \neq 0$ the solution set (for the unknown vector x) of the quadratic equation

$$a \mathbf{x}^T \mathbf{x} + \mathbf{b}^T \mathbf{x} + c = 0, \quad \mathbf{b} \in \mathbb{R}^2, \ a, c \in \mathbb{R}.$$
 (8.12.9)

(8-12.h) \bigcirc (20 min.) Derive an expression in terms of a, b, c for the center m and the radius r of the circle defined by (8.12.9).

```
SOLUTION for (8-12.h) \rightarrow 8-12-8-0:casvdl.pdf
```

(8-12.i) \boxdot (60 min.) According to equation (8.12.9), the same circle can be defined by different parameter vectors $\mathbf{v} = (a, b_1, b_2, c)^T$ and $\mathbf{v}' = (a', b'_1, b'_2, c')^T$, when $\mathbf{v}' = \lambda \mathbf{v}$, for every $\lambda \in \mathbb{R}$, $\lambda \neq 0$. Thus, this equation, for different data values (x_i, y_i) , can be solved in a least squares sense if supplemented by a **non-linear constraint**:

$$a\mathbf{x}_{i}^{T}\mathbf{x}_{i} + \mathbf{b}^{T}\mathbf{x}_{i} + c = 0 \quad i = 1,...,N, \qquad \|(a,b_{1},b_{2},c)^{T}\|_{2}^{2} = 1.$$

Write a C++ function

that solves this constrained overdetermined linear system of equations in least squares sense.

To learn how to use SVD to solve this problem, study carefully the hyperplane fitting problem [Lecture \rightarrow Ex. 3.4.4.5] and [Lecture \rightarrow Eq. (3.4.4.1)].

```
Solution for (8-12.i) \rightarrow 8-12-9-0:casvd2.pdf
```

8-12.IV: comparison of the results

(8-12.j) • Draw the data points and the fitted circles computed in the previous subtasks. Compare the centers and the radii.

To that end, in the file circle_approx.hpp, implement a C++ function

```
void plot(const Eigen::VectorXd& x, const Eigen::VectorXd& y,
const Eigen::Vector3d& z_alg, const Eigen::Vector3d& z_geo_GN,
const Eigen::Vector3d& z_geo_N, const Eigen::Vector3d& z_svd)
```

SOLUTION for (8-12.j) \rightarrow 8-12-10-0:caco.pdf

End Problem 8-12, 230 min.

Problem 8-13: Computing a Level Set

For a real-valued function $f: D \subset \mathbb{R}^d \to \mathbb{R}$ the level sets $\mathcal{L}_c := \{x \in D: f(x) \leq c\}, c \in \mathbb{R}$, are often used to characterize subsets of \mathbb{R}^d . This problem examines a numerical method for approximately computing the boundary of level sets of convex functions for d = 2.

This problem relies on methods for finding zeros of functions in 1D, see [Lecture \rightarrow Section 8.4].

Throughout this problem let $f: \mathbb{R}^2 \to \mathbb{R}$ be a continuously differentiable *strictly convex* function which has a *global minimum* in x = 0 and satisfies the growth condition $|f(x)| \to \infty$ uniformly for $||x|| \to \infty$

Then the level sets

$$\mathcal{L}_c := \{ x \in \mathbb{R}^2 : f(x) \le c \}, \quad c > 0,$$
 (8.13.1)

will be closed and convex subsets of \mathbb{R}^2 , and every ray

$$R_{\mathbf{d}} := \{ \xi \mathbf{d} : \xi \ge 0 \}, \quad \mathbf{d} \in \mathbb{R}^2 \setminus \{ 0 \}, \tag{8.13.2}$$

will intersect the boundary of \mathcal{L}_c ,

$$\partial \mathcal{L}_c := \{ x \in \mathbb{R}^2 : f(x) = c \} , \qquad (8.13.3)$$

in exactly one point.

Solution for (8-13.a)
$$\rightarrow$$
 8-13-1-0:s1.pdf

(8-13.b) Now, write down the recursion for the secant method [Lecture \rightarrow § 8.4.2.28] for solving the *scalar* non-linear equation tackled in Sub-problem (8-13.a). Your formula should involve only f-evaluations and $\mathbf{d} \in \mathbb{R}^2 \setminus \{\mathbf{0}\}$ and c > 0 as parameters.

Solution for (8-13.b)
$$\rightarrow$$
 8-13-2-0:s3.pdf

(8-13.c) □ (15 min.) [depends on Sub-problem (8-13.b)]

Implement an efficient C++ function

```
template <typename Functor>
  Eigen::VectorXd pointLevelSet(
   Functor &&f, const Eigen::Vector2d &d, double c,
   const Eigen::Vector2d &x0,
   double rtol = 1E-10, double atol = 1E-16);
```

that uses the secant method to find the unique point in the intersection of R_d and $\partial \mathcal{L}_c$, c > 0. The function f is given in procedural form by a functor object f. The coordinates of the intersection point are returned.

The vector argument $\mathbf{x}0$ passes (a guesses for) the coordinates of a point $\mathbf{x}_0 \in \partial \mathcal{L}_c$. This information may be used to obtain initial guesses for the secant iteration, which should start from the points $\frac{\|\mathbf{x}_0\|_2}{\|\mathbf{d}\|_2}\mathbf{d}$ and $1.1 \cdot \frac{\|\mathbf{x}_0\|_2}{\|\mathbf{d}\|_2}\mathbf{d}$.

Employ a correction-based termination criterion with relative tolerance and absolute tolerance supplied by the arguments rtol and atol.

```
SOLUTION for (8-13.c) \rightarrow 8-13-3-0:s4.pdf
```

(8-13.d) (20 min.) [depends on Sub-problem (8-13.c)]

Based on pointLevelSet () write an efficient C++ function

template <typename Functor>
 double areaLevelSet (Functor &&f, unsigned int n, double c);

that returns the "Archimedean approximation" of the area of \mathcal{L}_c , c>0, for the function f passed in f.

That Archimedean approximation replaces the set \mathcal{L}_c with a *convex polygonal domain* through the n points (n > 2)

$$\mathbf{p}_j \in \mathbb{R}^2$$
: $\{\mathbf{p}_j\} = R_{\mathbf{d}_j} \cap \partial \mathcal{L}_c$, $\mathbf{d}_j = \begin{bmatrix} \cos(2\pi i/n) \\ \sin(2\pi i/n) \end{bmatrix}$, $j = 0, \ldots, n-1$.

HIDDEN HINT 1 for (8-13.d) \rightarrow 8-13-4-0:hcp.pdf

HIDDEN HINT 2 for (8-13.d) \rightarrow 8-13-4-1:hcp2.pdf

SOLUTION for (8-13.d) \rightarrow 8-13-4-2:s4a.pdf

(8-13.e) :: (10 min.) We consider

$$f(\mathbf{x}) = x_1^2 + 2x_2^4$$
, $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \in \mathbb{R}^2$.

n	error in area
5	0.7322814883
10	0.2020820848
15	0.0943700007
20	0.0539834484
25	0.0346797600
30	0.0241767180
40	0.0136526713
50	0.0087539878
80	0.0034265044

The table lists the error in the approximate area of \mathcal{L}_1 returned by areaLevelSet () as $n \to \infty$, using relative tolerance rtol = 10^{-10} and absolute tolerance atol = 10^{-16} .

Describe qualitatively and quantitatively the empiric convergence of the approximation of the area in terms of $\frac{1}{n} \to 0$.

Solution for (8-13.e) \rightarrow 8-13-5-0:s5.pdf

End Problem 8-13, 60 min.

Problem 8-14: Symmetric Rank-1 Best Approximation of a Matrix

The approximation of matrices by data-sparse matrices of low rank has become a fundamental technique in modern computational mathematics and data science, leading to very efficient methods for the compression of non-local operators and multi-dimensional data.

This problems is based on [Lecture \rightarrow Section 8.7], in particular [Lecture \rightarrow Section 8.7.2], and also draws on [Lecture \rightarrow Section 3.4.4.2].

Given a square matrix $\mathbf{M} \in \mathbb{R}^{n,n}$, $\mathbf{M} \neq \mathbf{O}$, we want to find a symmetric rank-1 matrix closest in Frobenius norm $\|\cdot\|_F$ [Lecture \rightarrow Def. 3.4.4.17]:

$$\mathbf{x} \in \underset{\mathbf{z} \in \mathbb{R}^n}{\operatorname{argmin}} \left\| \mathbf{M} - \mathbf{z} \mathbf{z}^{\top} \right\|_{F} = \underset{\mathbf{z} \in \mathbb{R}^n}{\operatorname{argmin}} \left\| \mathbf{\Phi}(\mathbf{z}) \right\|_{2},$$
 (8.14.1)

with
$$\Phi:\mathbb{R}^n o\mathbb{R}^{n^2}$$
 , $\Phi(x):=\mathrm{vec}(\mathbf{M})-x\otimes x$, $x\in\mathbb{R}^n$, (8.14.2)

where $\|\cdot\|_2$ denotes the Euclidean vector norm, vec the vectorization of matrices

$$\operatorname{vec}: \mathbb{R}^{n,m} \to \mathbb{R}^{n \cdot m}$$
 , $\operatorname{vec}(\mathbf{A}) := \begin{bmatrix} (\mathbf{A})_{:,1} \\ (\mathbf{A})_{:,2} \\ \vdots \\ (\mathbf{A})_{:,m} \end{bmatrix} \in \mathbb{R}^{n \cdot m}$, [Lecture \to Eq. (1.2.3.5)]

as introduced in [Lecture \to Rem. 1.2.3.4], and \otimes the Kronecker product of two matrices/vectors from [Lecture \to Def. 1.4.3.7].

(8-14.a) : (30 min.)

In the file symrank1.hpp implement an EIGEN-based C++ function

that solves (8.14.1) for a *symmetric positive semi-definite* [Lecture \to Def. 1.1.2.6] matrix \mathbf{M} , that is, provided that $\mathbf{M}^{\top} = \mathbf{M}$ and $\mathbf{x}^{\top} \mathbf{M} \mathbf{x} > 0$ for all $\mathbf{x} \in \mathbb{R}^n$.

HIDDEN HINT 1 for (8-14.a) \rightarrow 8-14-1-0:s1h1.pdf

Solution for (8-14.a)
$$\rightarrow$$
 8-14-1-1:s1.pdf

For general M we intend to solve (8.14.1) by means of the **Gauss-Newton method**, that is, by means of the iteration [Lecture \rightarrow Eq. (8.7.2.1)]

$$\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} + \underset{\mathbf{h} \in \mathbb{R}^n}{\operatorname{argmin}} \left\| \mathbf{\Phi}(\mathbf{x}^{(k)}) + \mathsf{D} \, \mathbf{\Phi}(\mathbf{x}^{(k)}) \mathbf{h} \right\|_{2}, \tag{8.14.4}$$

with **1** from (8.14.2).

(8-14.b) (30 min.)

Give an expression for the linear mapping $\mathbf{h} \in \mathbb{R}^n \mapsto \mathsf{D}\, \mathbf{\Phi}(\mathbf{x})\mathbf{h}$, $\mathbf{x} \in \mathbb{R}^n$ fixed.

HIDDEN HINT 1 for (8-14.b) \rightarrow 8-14-2-0:kpbl.pdf

SOLUTION for (8-14.b)
$$\rightarrow$$
 8-14-2-1:s2a.pdf

(8-14.c) (30 min.) [depends on Sub-problem (8-14.b)]

Derive a formula for the Jacobian $D \Phi(\mathbf{x}) \in \mathbb{R}^{n^2,n}$ of Φ as defined in (8.14.2).

Write down D $\Phi(\mathbf{x})$ explicitly for n=3, abbreviating $x_i=(\mathbf{x})_i, i=1,2,3$.

HIDDEN HINT 1 for (8-14.c) \rightarrow 8-14-3-0:handy2.pdf

SOLUTION for (8-14.c)
$$\rightarrow$$
 8-14-3-1:s2b.pdf

(8-14.d) (30 min.) [depends on Sub-problem (8-14.c)]

The minimization problem in (8.14.4) involves a linear least-squares problem of the form

$$\Delta \mathbf{x} := \underset{\mathbf{b} \in \mathbb{R}^n}{\operatorname{argmin}} \|\mathbf{A}\mathbf{h} - \mathbf{b}\|_2 , \qquad (8.14.11)$$

for some matrix $\mathbf{A} \in \mathbb{R}^{n^2,n}$ and vector $\mathbf{b} \in \mathbb{R}^{n^2}$.

For general n compute an explicit formula for the system matrix $\mathbf{T} \in \mathbb{R}^{n,n}$ of the **normal equations** for the linear least squares problem (8.14.11).

HIDDEN HINT 1 for (8-14.d) \rightarrow 8-14-4-0:3hI.pdf

SOLUTION for (8-14.d)
$$\rightarrow$$
 8-14-4-1:s4.pdf

Give a simple explicit formula for the inverse T^{-1} of the system matrix of the normal equations.

HIDDEN HINT 1 for (8-14.e) \rightarrow 8-14-5-0:smw.pdf

HIDDEN HINT 2 for (8-14.e) \rightarrow 8-14-5-1:fails.pdf

SOLUTION for (8-14.e)
$$\rightarrow 8-14-5-2:s5.pdf$$

(8-14.f) (30 min.)

Code an efficient C++ function (in the file symrank1.hpp)

```
Eigen::VectorXd computeKronProdVecMult(const Eigen::VectorXd &v,
    const Eigen::VectorXd &b);
```

that evaluates the expression

$$(\mathbf{v}^{\top} \otimes \mathbf{I}_n + \mathbf{I}_n \otimes \mathbf{v}^{\top})\mathbf{b}$$
 for $\mathbf{v} \in \mathbb{R}^n$, $\mathbf{b} \in \mathbb{R}^{n^2}$, (8.14.14)

for any $n \in \mathbb{N}$.

HIDDEN HINT 1 for (8-14.f) \rightarrow 8-14-6-0:hrs.pdf

```
SOLUTION for (8-14.f) \rightarrow 8-14-6-1:sx.pdf
```

(8-14.g) ☐ (45 min.) [depends on Sub-problem (8-14.f) and Sub-problem (8-14.e)]

Based on the normal equation method for the occurring linear least squares problems, in the file symrank1.hpp implement an efficient C++ function

```
Eigen::VectorXd symmRankOneApprox(
   const Eigen::MatrixXd &M
   double rtol = 1E-6, atol = 1.0E-8);
```

that implements the **Gauss-Newton iteration** (8.14.4) for solving (8.14.1) using a correction-based termination criterion with relative tolerance rtol and absolute tolerance atol. The parameter M passes the matrix $\mathbf{M} \in \mathbb{R}^{n,n}$.

As initial guess for z we use that column of the symmetric part $\frac{1}{2}(M^{\top} + M)$ of M with the largest Euclidean norm.

```
HIDDEN HINT 1 for (8-14.g) \rightarrow 8-14-7-0:mcgn.pdf

HIDDEN HINT 2 for (8-14.g) \rightarrow 8-14-7-1:mcgn2.pdf

HIDDEN HINT 3 for (8-14.g) \rightarrow 8-14-7-2:fail2.pdf

SOLUTION for (8-14.g) \rightarrow 8-14-7-3:s6.pdf
```

End Problem 8-14, 215 min.

Problem 8-15: Non-linear Least-Squares Location Estimation

Least-squares techniques are the main tool for estimating unknown parameters from measurements. In this problem we apply them to an overdetermined non-linear system of equations arising from the problem of estimating the location of an acoustic source. As solution method we use the iterative **Gauss-Newton method**.

This problem requires knowledge about the Gauss-Newton method from [Lecture \rightarrow Section 8.7.2] and skills in multi-dimensional differentiation. Implementation in C++ based on EIGEN is part of the problem.

A source of acoustic waves at an unknown location $\mathbf{p}^* \in \mathbb{R}^3$ in a room starts emitting a loud and short noise at an unknown time t^* . At times $t_j > t^*$, $j = 1, \ldots, n$, this signal is detected by each of n > 4 microphones placed in known locations $\mathbf{q}^j \in \mathbb{R}^3$, $j = 1, \ldots, n$, $\mathbf{q}^i \neq \mathbf{q}^j$ for $i \neq j$. Writing c (physical units $[c] = 1 \frac{m}{s}$) for the speed of sound in air, it is straightforward that the arrival times t_j of the signal satisfy

$$\|\mathbf{q}^{j} - \mathbf{p}^{*}\|_{2} = c(t_{j} - t^{*}), \quad j = 1, ..., n.$$
 (8.15.1)

From the measured arrival times t_j we want to reconstruct the location \mathbf{p}^* of the noise source and the time t^* , when it went off.

(8-15.a) • (5 min.)

Rewrite (8.15.1) as an overdetermined non-linear system of equations in the standard form $F(\mathbf{x}) = \mathbf{0}$. Specify \mathbf{x} and F:

$$F: \left\{ \begin{array}{c} \mathbb{R} \\ \mathbf{x} := \\ \end{array} \right. \rightarrow \mathbb{R}$$

$$\mapsto F(\mathbf{x}) := \left[\begin{array}{c} \mathbb{R} \\ \mathbb{R} \\ \end{array} \right]$$

SOLUTION for (8-15.a)
$$\rightarrow$$
 8-15-1-0:locs1.pdf

Compute the Jacobian DF(x) for the function F found in Sub-problem (8-15.a). Determine its maximal domain of definition, that is, the set of x's for which F is differentiable.

$$\mathsf{D}\,F(\mathbf{x}) = \mathsf{D}\,F\left(\left[\right] \right) =$$

for
$$\mathbf{x} \in$$

SOLUTION for (8-15.b) \rightarrow 8-15-2-0:locs2.pdf

A

Assume a choice of physical units such that c=1. In the file locationestimation.hpp write a C++ function

```
std::pair<Eigen::Vector3d, double>
source_estimation(
  const Eigen::Matrix<double, 3, Eigen::Dynamic> &Q,
  const Eigen::VectorXd &ta);
```

that solves the non-linear system of equations (8.15.1) in least-squares sense using the Gauss-Newton method introduced in [Lecture \rightarrow Section 8.7.2]. The $3 \times n$ -matrix $\mathbb Q$ passes the vectors $\mathbf q^j$, $j=1,\ldots,n$, in its columns, whereas the vector ta contains the arrival times t_j . The function should return an estimate for the pair $(\mathbf p^*,t^*)$.

As initial guess use

$$\mathbf{p}^{(0)} := \frac{1}{n} \sum_{j=1}^{n} \mathbf{q}^{j}$$
 , $t^{(0)} := t_k - \left\| \mathbf{q}^k - \mathbf{p}^{(0)} \right\|_2$, k such that $t_k = \min\{t_j, j = 1, \dots, n\}$.

As in [Lecture \to Code 8.7.2.2] rely on a correction-based stopping rule in the spirit of [Lecture \to Eq. (8.5.3.2)] with relative tolerance $\tau_{\rm rel}=10^{-6}$ and absolute tolerance $\tau_{\rm abs}=10^{-8}$.

SOLUTION for (8-15.c)
$$\rightarrow$$
 8-15-3-0:locsol3.pdf

•

(8-15.d) \odot (5 min.) The following table displays the history of the Gauss-Newton iteration as used in source_location () for $t^* = 0$,

$$\begin{bmatrix} \mathbf{q}^1, \dots, \mathbf{q}^{11} \end{bmatrix} = \begin{bmatrix} 0 & 0 & 0 & 1 & 1 & 1 & 1 & 0.5 & 0.0 & 1.0 & 0.5 \\ 1 & 0 & 1 & 1 & 0 & 0 & 1 & 0.0 & 0.5 & 0.5 & 1.0 \\ 0 & 1 & 1 & 0 & 1 & 0 & 1 & 0.5 & 0.5 & 0.5 & 0.5 \end{bmatrix}$$

and the arrival times t_i fitting the source position $\mathbf{p}^* = [0.1, 0.1, 0.9]$.

k	$\left\ \mathbf{x}^{(k)}-\mathbf{x}^* ight\ _2$	$\left\ \mathbf{s}\right\ _2$	$\left\ F(\mathbf{x}^{(k)}) ight\ $
1	1.021791e+00	1.135326e+00	3.442935e+00
2	2.588111e-01	2.527947e-01	6.041095e-01
3	7.550346e-03	7.504712e-03	1.559289e-02
4	6.192137e-05	6.192055e-05	9.969019e-05
5	9.310988e-10	9.310988e-10	1.862087e-09
6	1.240973e-16	7.295652e-17	3.159274e-16

Characterize the observed convergence and explain your answer.

SOLUTION for (8-15.d) \rightarrow 8-15-4-0:locsol4.pdf

End Problem 8-15, 45 min.

Problem 8-16: Periodic Collocation for a Non-Linear Differential Equation

This problem examines the so-called collocation method for the approximate solution of differential equations. Since we look for periodic solutions, it is natural to use linear combinations of trigonometric polynomials as trial space. Equations are obtained by requiring that the equations remains satisfied in finitely many points.

This problem requires the techniques from [Lecture \rightarrow Section 4.2] and [Lecture \rightarrow Section 8.5] and involves substantial implementation in C++ using EIGEN.

We look for 1-periodic solutions $u:[0,1]\to\mathbb{C}$, u=u(t), of the differential equation

$$-\frac{d^2u}{dt^2}(t) + u(t)^3 = \sin(\pi t) , \quad 0 \le t < 1 .$$

To approximate them we replace u with

$$u_N(t):=\sum_{j=0}^N x_j\cos(2\pi jt) \quad ext{for some} \quad N\in\mathbb{N} \; ,$$
 (8.16.1)

and try to determine the unknown coefficients $x_j \in \mathbb{R}$, j = 0, ..., N, by solving the collocation equations

$$-\frac{d^2u_N}{dt^2}(t_k) + u_N(t_k)^3 = \sin(\pi t_k) , \quad k = 0, \dots, N , \quad t_k := \frac{k}{N+1} . \tag{8.16.2}$$

Eigen::VectorXd eval_uN(const Eigen::VectorXd &x, unsigned int M);

which, given the coefficients $x_j \in \mathbb{R}$, j = 0, ..., N, in (8.16.1) through the vector x and a number M > N in M, returns the vector

$$\left[u_N(\frac{k}{M})\right]_{k=0}^{M-1} \in \mathbb{R}^M.$$

"Efficient" stipulates that the computational cost of eval_uN() must scale like $O(M \log M)$ for $M \to \infty$.

HIDDEN HINT 1 for (8-16.a) \rightarrow 8-16-1-0:pcollh1.pdf

Solution for (8-16.a)
$$\rightarrow$$
 8-16-1-1:pces1.pdf

(8-16.b) \odot (12 min.) Write the collocation equations (8.16.2) in the standard form $F(\mathbf{x}) = \mathbf{0}$ of a non-linear system of equations for an *explicitly given* $F : \mathbb{R}^n \to \mathbb{R}^n$ and suitable \mathbf{x} and $n \in \mathbb{N}$.

$$n = \boxed{ }$$
 , $\mathbf{x} = \boxed{ }$,

```
F(\mathbf{x}) =
```

```
HIDDEN HINT 1 for (8-16.b) \rightarrow 8-16-2-0:ceh2.pdf
Solution for (8-16.b) \rightarrow 8-16-2-1:pcolls2.pdf
```

Relying on the function $eval_uN$ () from Sub-problem (8-16.a) implement another C++ function in periodiccollocation.hpp,

```
Eigen::VectorXd eval_F(const Eigen::VectorXd &x);
```

that returns the vector $F(\mathbf{x})$ when given the argument \mathbf{x} in \mathbf{x} .

```
SOLUTION for (8-16.c) \rightarrow 8-16-3-0:pcolls3.pdf
```

We want to try Newton's method to solve $F(\mathbf{x}) = \mathbf{0}$ with F from Sub-problem (8-16.b). To that end we need a function

```
Eigen::MatrixXd eval_DF(const Eigen::VectorXd &x);
```

that returns the Jacobian DF(x), when x contains the argument x. Write such a function in the file periodiccollocation.hpp

```
HIDDEN HINT 1 for (8-16.d) \rightarrow 8-16-4-0:pch4.pdf
```

```
SOLUTION for (8-16.d) \rightarrow 8-16-4-1:pces4.pdf
```

End Problem 8-16, 57 min.

Problem 8-17: Aspects of Iterative Methods for Non-Linear Equations

This problem revisits central concepts and algorithms connected with iterative methods for non-linear (systems of) equations.

This problem assumes familiarity with [Lecture \rightarrow Section 8.2.2], [Lecture \rightarrow Section 8.4.2.1], and [Lecture \rightarrow Section 8.5.3].

(8-17.a) (5 min.) Fill in the blanks in the following definition:

SOLUTION for (8-17.a) \rightarrow 8-17-1-0:sa2.pdf

(8-17.b) \odot (10 min.) The templated C++ function

implements a generic Newton iteration for solving $F(\mathbf{x}) = \mathbf{0}$ with correction-based termination relying on the simplified Newton correction:

```
\begin{split} \textbf{STOP}, \text{ as soon as } & \left\| \Delta \bar{\mathbf{x}}^{(k)} \right\| \leq \mathtt{rtol} \left\| \mathbf{x}^{(k)} \right\| \quad \textbf{or} \quad \left\| \Delta \bar{\mathbf{x}}^{(k)} \right\| \leq \mathtt{atol} \;, \\ \text{with } & \mathbf{simplified Newton correction} \quad \Delta \bar{\mathbf{x}}^{(k)} := \mathsf{D} \, F(\mathbf{x}^{(k-1)})^{-1} F(\mathbf{x}^{(k)}). \\ & \left[ \mathsf{Lecture} \to \mathsf{Eq.} \; (8.5.3.5) \right] \end{split}
```

The functor arguments F and DF pass the function $F: \mathbb{R}^n \to \mathbb{R}^n$ and its derivative (Jacobian) $DF: \mathbb{R}^n \to \mathbb{R}^{n,n}$, whereas x supplies the initial guess and also returns the final results. The user has to specify the absolute tolerance atol and relative tolerance rtol to control termination.

The type **JacType** must be a *functor type* whose return type has capabilities like a matrix type of EIGEN, whereas **VecType** is a matching vector type, for instance a vector type of EIGEN.

Fill in the blanks in Code 8.17.1 so that the resulting code is a valid C++ implementation of newton ()

C++ code 8.17.1: Generic driver function for Newton's method template <typename FuncType, typename JacType, typename VecType> void newton(const FuncType &F, const JacType &DF, VecType &x, 2 double rtol, double atol) { 3 double sn; **do** { auto jacfac = . **lu**(); x -= jacfac.solve(); sn = jacfac.solve().norm();) && (sn > while ((sn >)); 10 11 }

SOLUTION for (8-17.b) \rightarrow 8-17-2-0:imsb.pdf

(8-17.c) :: (10 min.)

Let f denote the function

$$f: \left\{ \begin{array}{ccc} [0,1] & \to & [0,1] \\ x & \mapsto & xe^{x-1} \end{array} \right.$$

Explicitly state the Newton iteration converging to g(y) for a given $y \in [0, 1]$ and sufficiently good initial guess $x^{(0)}$, where $g := f^{-1}$ is the *inverse function* of f.

$$x^{(k+1)} = x^{(k)} -$$
 , $k \in \mathbb{N}_0$.

SOLUTION for (8-17.c) \rightarrow 8-17-3-0:imsc.pdf

End Problem 8-17, 25 min.

Chapter 9

Computation of Eigenvalues and Eigenvectors

Chapter 10

Krylov Methods for Linear Systems of Equations

Chapter 11

Numerical Integration – Single Step Methods

Problem 11-1: Linear ODE in spaces of matrices

In this problem we consider initial value problems (IVPs) for linear ordinary differential equations (ODEs) whose state space is a vector space of $n \times n$ matrices. Such ODEs arise when modelling the dynamics of rigid bodies in classical mechanics. A related problem is Problem 11-5.

Mainly relies on the information contained in [Lecture \rightarrow Section 11.2].

We consider the *linear* matrix differential equation

$$\dot{\mathbf{Y}} = \mathbf{A}\mathbf{Y} =: \mathbf{f}(\mathbf{Y}) \quad \text{with} \quad \mathbf{A} \in \mathbb{R}^{n,n}.$$
 (11.1.1)

whose solutions are *matrix-valued functions* $\mathbf{Y}: \mathbb{R} \to \mathbb{R}^{n,n}$ and whose right-hand-side function maps $\mathbb{R}^{n,n} \to \mathbb{R}^{n,n}$.

(11-1.a) \odot (15 min.) Show that for *skew-symmetric* \mathbf{A} , i.e. $\mathbf{A} = -\mathbf{A}^{\top}$ we have:

```
\mathbf{Y}(0) orthogonal \Longrightarrow \mathbf{Y}(t) orthogonal \forall t.
```

```
\label{eq:hint1} \mbox{Hidden Hint 1 for (11-1.a)} \ \ \rightarrow \mbox{11-1-1-0:MatO1h1.pdf}
```

HIDDEN HINT 2 for (11-1.a) \rightarrow 11-1-1-1:MatO1h2.pdf

SOLUTION for (11-1.a) \rightarrow 11-1-1-2:Mat01s.pdf

1. a single step of the explicit Euler method, see [Lecture \rightarrow Section 11.2.1]:

2. a single step of the implicit Euler method, see [Lecture \rightarrow Section 11.2.2],

3. a single step of the implicit mid-point method, see [Lecture \rightarrow Section 11.2.3].

which compute, for a given initial value $\mathbf{Y}(t_0) = \mathbf{Y}_0$ and for given step size h, approximations for $\mathbf{Y}(t_0 + h)$ using one step of the corresponding method for the approximation of the ODE (11.1.1)

HIDDEN HINT 1 for (11-1.b) \rightarrow 11-1-2-0:MatO2h.pdf

SOLUTION for (11-1.b)
$$\rightarrow$$
 11-1-2-1:MatO2s.pdf

Investigate numerically, which one of the implemented methods preserves orthogonality for the ODE (11.1.1) and which one doesn't. To that end, consider the matrix

$$\mathbf{M} := \begin{bmatrix} 8 & 1 & 6 \\ 3 & 5 & 7 \\ 9 & 9 & 2 \end{bmatrix}$$

and use the matrix Q arising from the QR-decomposition (\rightarrow [Lecture \rightarrow Section 3.3.3.4]) of M as initial data Y_0 . As matrix A, use the skew-symmetric matrix

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

Perform N=20 time steps of size h=0.01 with each method and compute and tabulate the Frobenius norm of $\mathbf{Y}_k^{\top}\mathbf{Y}_k - \mathbf{I}$.

Do all this with a C++ function

std::tuple<double, double, double> checkOrthogonality();

that also returns the Frobenius norms $\mathbf{Y}_{N}^{\mathsf{T}}\mathbf{Y}_{N} - \mathbf{I}$ for each of the three methods.

Solution for (11-1.c)
$$\rightarrow$$
 11-1-3-0:Mat03s.pdf

End Problem 11-1, 75 min.

Problem 11-2: Explicit Runge-Kutta methods

The most widely used class of numerical integrators for IVPs is that of *explicit* Runge-Kutta (RK) methods as defined in [Lecture \rightarrow Def. 11.4.0.11]. They are usually described by giving their coefficients in the form of a Butcher scheme [Lecture \rightarrow Eq. (11.4.0.13)].

Related to [Lecture \rightarrow Section 11.4], [Lecture \rightarrow Def. 11.4.0.11], and Butcher schemes as defined in [Lecture \rightarrow Eq. (11.4.0.13)]

which provides a generic **explicit Runge-Kutta single-step method** (\rightarrow [Lecture \rightarrow Def. 11.4.0.11]) given by a Butcher scheme [Lecture \rightarrow Eq. (11.4.0.13)] to solve the autonomous initial-value problem $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y}), \ \mathbf{y}(t_0) = \mathbf{y}_0$. The matrix $\mathfrak{A} \in \mathbb{R}^{s,s}$ and the vector $\mathbf{b} \in \mathbb{R}^s$ from the Butcher scheme are passed to the constructor, whereas the right-hand side vector field \mathbf{f} , the final time T > 0, the number N of *equidistant timesteps*, and the intial state \mathbf{y}_0 are arguments of the $\mathtt{solve}()$ method, which is supposed to return the sequence $(\mathbf{y}_k)_{k=0}^N$ of states generated by the explicit Runge-Kutta method.

```
SOLUTION for (11-2.a) \rightarrow 11-2-1-0:RK3P1s.pdf

(11-2.b) (30 min.) [depends on Sub-problem (11-2.a)]

In the template file rk3prey.hpp implement a C++ function double RK3prey();
```

in order to test your implementation of the RK methods with the following test case:

As autonomous initial value problem, consider the predator/prey model (cf. [Lecture \rightarrow Ex. 11.1.2.5]):

$$\dot{y}_1(t) = (\alpha_1 - \beta_1 y_2(t)) y_1(t) , \qquad (11.2.2)$$

$$\dot{y}_2(t) = (\beta_2 y_1(t) - \alpha_2) y_2(t)$$
, (11.2.3)

$$\mathbf{y}(0) = [100, 5]^{\top},$$
 (11.2.4)

with coefficients $\alpha_1 = 3$, $\alpha_2 = 2$, $\beta_1 = \beta_2 = 0.1$.

Use the explicit 3-stage Runge-Kutta single step method described by the following **Butcher scheme** (*cf.* [Lecture \rightarrow Def. 11.4.0.11]):

Compute an approximated solution up to time T=10 for the number of equidistant time steps $N=2^j$, $j=7,\ldots,14$.

As reference solution, use $\mathbf{y}(10) = [0.319465882659820, 9.730809352326228]^{\top}.$

Tabulate the error at final time and estimate the empiric rate of algebraic convergence of the method by means of linear regression using the supplied function polyfit (file polyfit.hpp). Your function should return the estimated convergence rate.

SOLUTION for (11-2.b) \rightarrow 11-2-2-0:RK3P2s.pdf

End Problem 11-2, 60 min.

•

Problem 11-3: Extrapolation of evolution operators

In [Lecture \rightarrow § 11.3.1.1] we have seen how discrete evolution operators can describe single-step methods for the numerical integration of ODEs. This task will study a way to combine discrete evolution operators of known order in order to build a single-step method with increased order. The method can be generalized to an extrapolation construction of higher-order single-step methods. These can be used for time-local stepsize control following the policy of [Lecture \rightarrow § 11.5.2.11].

This problem assumes familiarity with the [Lecture \rightarrow Section 11.3] and, in particular, the concept of discrete evolution, *cf.* [Lecture \rightarrow Def. 11.3.1.5]. It also addresses the adaptive timestepping strategy presented in [Lecture \rightarrow Section 11.5].

Let Ψ^h define the discrete evolution of an order p Runge-Kutta single step method for the autonomous ODE $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$, $\mathbf{f} : D \subseteq \mathbb{R}^d \to \mathbb{R}^d$. We define a new discrete evolution operator:

$$\widetilde{\mathbf{Y}}^h := \frac{1}{1 - 2^p} \left(\mathbf{Y}^h - 2^p \cdot (\mathbf{Y}^{h/2} \circ \mathbf{Y}^{h/2}) \right)$$
(11.3.1)

where o denotes the composition of mappings.

```
SOLUTION for (11-3.a) \rightarrow 11-3-1-0:ODESolve1s.pdf
```

(11-3.b) ☐ (15 min.) Implement an EIGEN-based C++ function

that returns $\widetilde{\Psi}^h y_0$ when given the underlying Ψ through the functor Psi. Objects of type **DiscEvIOp** must provide an evaluation operator:

```
Eigen::VectorXd operator()(double h, const Eigen::VectorXd &y);
```

providing the result of $\Psi^h(y)$. For instance, suitable C++ lambda functions satisfy this requirement, see [Lecture \rightarrow Section 0.3.3].

```
Solution for (11-3.b) \rightarrow 11-3-2-0:ODESolve2s.pdf
```

that carries out N equidistant steps of size $\tau := T/N$ of a single step method described by the discrete evolution Ψ passed through the functor object Psi. The function should return the sequence $(\mathbf{y}_k)_{k=0}^N$ produced by the single-step method when started with initial value \mathbf{y}_0 (given as argument y_0).

```
HIDDEN HINT 1 for (11-3.c) \rightarrow 11-3-3-0:ODESolve3h.pdf
Solution for (11-3.c) \rightarrow 11-3-3-1:ODESolve3s.pdf
```

For the particular scalar initial-value problem (IVP)

$$\dot{y} = 1 + y^2$$
 , $y(0) = 0$, (11.3.4)

with exact solution $y(t) = \tan(t)$, determine empirically the order of the single step method induced by $\widetilde{\Psi}^h$ from (11.3.1), when Ψ arises from the explicit Euler method, recall Sub-problem (11-3.a). For that purpose, write a C++ function

```
double testcvpExtrapolatedEuler();
```

that monitors and tabulates the error $|y_N(1) - y_{ex}(1)|$ at final time T = 1 for N uniforms steps with $N = 2^q, q = 2, \ldots, 12$, also returns an estimate of the rate of algebraic convergence based on linear regression.

```
HIDDEN HINT 1 for (11-3.d) \rightarrow 11-3-4-0:ODESolve4h.pdf
Solution for (11-3.d) \rightarrow 11-3-4-1:ODESolve4s.pdf
```

Complete the implementation of a function

for the approximation of the solution of an IVP by means of adaptive timestepping based on Ψ and $\widetilde{\Psi}$, where Ψ is passed through the argument Psi.

Step rejection and stepsize correction and prediction as explained in [Lecture \rightarrow Section 11.5] is to be employed as in [Lecture \rightarrow Code 11.5.2.18]. The argument T supplies the final time, y0 the initial state, h0 an initial stepsize, p the order of the discrete evolution Ψ , reltol and abstol the respective tolerances, and hmin a minimal stepsize that will trigger premature termination.

The function should return a vector of tuples (t_k, \mathbf{y}_k) , k = 0, ..., M, where $M \in \mathbb{N}$ is the total number of timesteps, t_k are the knots of the temporal mesh created by the adaptive integrator, and $(\mathbf{y}_k)_k$ the sequence of states generated by the single-step method.

```
HIDDEN HINT 1 for (11-3.e) \rightarrow 11-3-5-0:EPh.pdf

SOLUTION for (11-3.e) \rightarrow 11-3-5-1:ODESolve5s.pdf

(11-3.f) ① (20 min.) [ depends on Sub-problem (11-3.e) ]
```

Implement a C++ function

```
void solveTangentIVP()
```

that relies on odeintssctrl from Sub-problem (11-3.e) to compute the solution of the IVP (11.3.4) up to time T=1.5 with the following data: h0=1/100, reltol=10e-4, abstol=10e-6, hmin=10e-5. Finally, the function should use MATPLOTLIBCPP to plot the approximated solution and store the plot in the file tangent.png. Do not forget to add axis labels.

```
Solution for (11-3.f) \rightarrow 11-3-6-0:sx.pdf
```

End Problem 11-3, 100 min.

Problem 11-4: System of second-order ODEs

In this problem we practise the conversion of a second-order ODE into a first-order system in the case of a large linear system of ODEs.

This problem assumes familiarity with Runge-Kutta single-step method as introduced in [Lecture \rightarrow Section 11.4]

Consider the following initial value problem for an (implicit) second-order system of ordinary differential equations in the time interval [0, T]:

$$2\ddot{u}_{1} - \ddot{u}_{2} = u_{1}(u_{2} + u_{1}),$$

$$-\ddot{u}_{i-1} + 2\ddot{u}_{i} - \ddot{u}_{i+1} = u_{i}(u_{i-1} + u_{i+1}), \qquad i = 2, \dots, n-1,$$

$$2\ddot{u}_{n} - \ddot{u}_{n-1} = u_{n}(u_{n} + u_{n-1}),$$

$$u_{i}(0) = u_{0,i} \qquad i = 1, \dots, n,$$

$$\dot{u}_{i}(0) = v_{0,i} \qquad i = 1, \dots, n.$$

$$(11.4.1)$$

Here the notation \ddot{w} designates the second derivative of a time-dependent function $t \mapsto w(t)$.

(11-4.a) \odot (15 min.) Write (11.4.1) as a first-order IVP of the form $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$, $\mathbf{y}(0) = \mathbf{y}_0$.

```
HIDDEN HINT 1 for (11-4.a) \rightarrow 11-4-1-0:Syst1h1.pdf
```

HIDDEN HINT 2 for (11-4.a) \rightarrow 11-4-1-1:Syst1h2.pdf

HIDDEN HINT 3 for (11-4.a) \rightarrow 11-4-1-2:Syst1h3.pdf

```
SOLUTION for (11-4.a) \rightarrow 11-4-1-3:Syst1s.pdf
```

```
HIDDEN HINT 1 for (11-4.b) \rightarrow 11-4-2-0:h2.pdf
```

```
SOLUTION for (11-4.b) \rightarrow 11-4-2-1:Syst2s.pdf
```

Implement a C++ function

that performs a single step of stepsize h of the classical Runge-Kutta method of order 4 for the first-order ODE obtained in Sub-problem (11-4.a) (associated right-hand side function passed via the functor odefun). y0 contains the state before the step and the function has to return the state after the step in y1.

```
Solution for (11-4.c) \rightarrow 11-4-3-0:Syst3s.pdf
```

Implement a function

```
double testcvgRK4()
```

that applies the classical Runge-Kutta method of order 4 to solve a particular initial-value problem for the ODE derived in Sub-problem (11-4.a).

Use the parameters and initial values

$$n = 5$$
, $u_i(0) = i/n$, $v_i(0) = -1$, $T = 1$,

and tabulate the Euclidean norm of the error at final time T=1 for numbers $N=2,4,8,\ldots,1024$ of uniform timesteps. As reference solution use the value obtained with $N=2^{12}$ equidistant timesteps. The function should also return an estimate of the rate of algebraic convergence obtained by means of linear regression for which you can use the supplied function polyfit () from polyfit.hpp.

Note that you should use EIGEN's sparse matrix data type for any sparse matrix encountered. Refer to [Lecture \rightarrow Section 2.7.2] for details.

SOLUTION for (11-4.d) \rightarrow 11-4-4-0:Syst4s.pdf

•

End Problem 11-4, 60 min.

Problem 11-5: Non-linear Evolutions in Spaces of Matrices

In this problem we consider initial value problems (IVPs) for ordinary differential equations whose state space is a vector space of $n \times n$ matrices. Such IVPs occur in mathematical models of discrete mechanical systems.

Related to this problem is Problem 11-1.

We consider a *non-linear* ODE in the state space of $n \times n$ matrices and study the associated initial value problem

$$\dot{\mathbf{Y}} = -(\mathbf{Y} - \mathbf{Y}^{\top})\mathbf{Y} =: \mathbf{f}(\mathbf{Y}) , \quad \mathbf{Y}(0) = \mathbf{Y}_0 \in \mathbb{R}^{n,n},$$
 (11.5.1)

whose solution is given by a *matrix-valued function* $t \mapsto \mathbf{Y}(t) \in \mathbb{R}^{n,n}$.

```
(11-5.a) ② (15 min.) In the file NLmatode.hpp, implement the C++ function

Eigen::MatrixXd matode(const Eigen::MatrixXd & Y0, double T)
```

which solves (11.5.1) on [0, T] using the C++ header-only class ode 45 (in the file ode 45 hpp). The initial value should be given by a $n \times n$ EIGEN matrix Y0. Set the absolute tolerance to 10^{-10} and the relative tolerance to 10^{-8} . The output should be an approximation of $\mathbf{Y}(T) \in \mathbb{R}^{n,n}$.

Instructions on the use of ode45, see [Lecture \rightarrow ??].

The class ode45 is header-only, meaning you just include the file and use it right away (no linking required). The file ode45.hpp defines the class ode45 implementing an embedded Runge-Kutta-Fehlberg method of order 4(5), see [Lecture \rightarrow Section 11.5.3] and [Lecture \rightarrow Ex. 11.5.3.2], with an adaptive stepsize control as presented in [Lecture \rightarrow § 11.5.2.11].

1. Construct an object of **ode45** type: create an instance of the class, passing the right-hand-side function **f** as a functor object to the constructor

```
template <class StateType,
class RhsType = std::function<StateType(const StateType &)>>
class ode45 {
   public:
   ode45(const RhsType &rhs);
   // ......
}
```

Template parameters are

- **StateType**: type of initial data and solution (state space), the only requirement is that the type possesses a normed vector-space structure, that is, it must implement the operations +, *, *=, += and assignment/copy operators. Moreover a norm() method must be available. EIGEN's vector and matrix types, as well as fundamental types are eligible as **StateType**.
- RhsType: type of rhs function (automatically deduced).

The argument rhs must be of a functor type that provides an evaluation operator

```
StateType operator()(const StateType & vec);
```

It can also be a lambda function.

2. (optional) Set the integration options: set data members of the data structure **ode45.options** to configure the solver:

```
O.options.<option_you_want_to_set> = <value>;
```

Examples:

- rtol: relative tolerance for error control (default is 1e-6)
- atol: absolute tolerance for error control (default is 1e-8)

e.g.:

```
O.options.rtol = 1e-5;
```

3. Solve stage: invoke the single-step method through calling the method

```
template <class NormFunc = decltype(_norm<StateType>)>
std::vector<std::pair<StateType, double>>
solve(const StateType &y0, double T,
const NormFunc &norm = _norm<StateType>);
```

The type **NormType** should provide a norm for vectors of type **StateType**. However, this type can be deduced automatically and the argument norm is optional. The other arguments are

- y0: initial value of type StateType ($\mathbf{y}_0 = y0$)
- T: final time of integration
- norm: (optional) norm function to call for objects of StateType, for the computation of the error

Return value The function returns the solution of the IVP, as a std::vector of std::pair (y(t), t) for every snapshot.

For more explanations and details, please consult the in-class documentation provided in the comments.

```
SOLUTION for (11-5.a) \rightarrow 11-5-1-0:NMatO1s.pdf
```

(11-5.b) \Box (10 min.) Show that the function $t \mapsto \mathbf{Y}^{\top}(t)\mathbf{Y}(t)$ is constant for the exact solution $\mathbf{Y}(t)$ of (11.5.1).

```
HIDDEN HINT 1 for (11-5.b) \rightarrow 11-5-2-0:NMat02h1.pdf
```

HIDDEN HINT 2 for (11-5.b) \rightarrow 11-5-2-1:NMat02h2.pdf

SOLUTION for (11-5.b)
$$\rightarrow$$
 11-5-2-2:NMat02s.pdf

In the file NLmatode.hpp, implement the C++ function

```
bool checkinvariant(const Eigen::MatrixXd & M, double T);
```

which (numerically) determines if the invariant $t \mapsto \mathbf{Y}(t)^{\top}\mathbf{Y}(t)$ is preserved for approximate solutions of (11.5.1) as output by $\mathtt{matode}()$. You must take into account round-off errors. The function's arguments should be the same as that of $\mathtt{matode}()$.

SOLUTION for (11-5.c)
$$\rightarrow$$
 11-5-3-0:NMat03s.pdf

The so-called discrete gradient method for (11.5.1) reads

$$\mathbf{Y}_* = \mathbf{Y}_0 + \frac{1}{2}hf(\mathbf{Y}_0)$$
 , $\mathbf{Y}_1 = (\mathbf{I} + \frac{1}{2}h(\mathbf{Y}_* - \mathbf{Y}_*^{\top}))^{-1}(\mathbf{I} - \frac{1}{2}h(\mathbf{Y}_* - \mathbf{Y}_*^{\top}))\mathbf{Y}_0$. (11.5.5)

Using the solution produced by matode() from Sub-problem (11-5.a) as a substitute for an exact solution, determine the order of convergence of the discrete gradient rule in a numerical experiment that uses $M \in \{10, 20, 40, 80, 160, 320, 640, 1280\}$ equidistant integration steps for solving (11.5.1) approximately. As initial value \mathbf{Y}_0 use the *orthogonal* "left-shift" matrix

$$\mathbf{Y}_0 := egin{bmatrix} 0 & 1 & 0 & 0 & 0 \ 0 & 0 & 1 & 0 & 0 \ 0 & 0 & 0 & 1 & 0 \ 0 & 0 & 0 & 0 & 1 \ 1 & 0 & 0 & 0 & 0 \end{bmatrix}.$$

Concretely, implement a C++ function

double cvgDiscreteGradientMethod();

that tabulates the Frobenius norm of the discretization error at final time T=1 and returns an estimate of the rate of algebraic convergence obtained by linear regression using the supplied function polyfit (file polyfit.hpp).

SOLUTION for (11-5.d) \rightarrow 11-5-4-0:s5.pdf

End Problem 11-5, 60 min.

Problem 11-6: Order is not everything

In [Lecture \rightarrow Section 11.3.2] we have seen that Runge-Kutta single step methods when applied to initial value problems with sufficiently smooth solutions will converge algebraically (with respect to the maximum error in the mesh points) with a rate given by their intrinsic order, see [Lecture \rightarrow Def. 11.3.2.8].

This problem studies the convergence of some Runge-Kutta single-step methods as discussed in [Lecture \rightarrow Section 11.3.2]. It relies on a class implemented in Problem 11-2.

In this problem we perform empiric investigations of orders of convergence of several explicit Runge-Kutta single step methods. We rely on two IVPs, one of which has a perfectly smooth solution, whereas the second has a solution that is merely piecewise smooth. Thus in the second case the smoothness assumptions of the convergence theory for RK-SSMs might be violated and it is interesting to study the consequences.

Remark. For this problem you can reuse the class **RKIntegrator** implemented in Problem 11-2. You first need to construct an object of this class, passing as arguments the Butcher tableau matrices A and b, for instance:

```
RKIntegrator<double> rk(A,b);
```

After that, call the methods solve, with parameters: r.h.s. function, final time, initial value and number of steps. For instance:

```
rk.solve(f,T,y0,n);
```

The output of this function will be **std**: **:vector**<**double**> containing the solution at each equidistant time step.

(11-6.a) (30 min.) Consider the autonomous ODE

$$\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y}), \quad \mathbf{y}(0) = \mathbf{y}_0,$$
 (11.6.1)

where $f: \mathbb{R} \to \mathbb{R}$ and $\mathbf{y}_0 \in \mathbb{R}$. Using the class **RKIntegrate** write a C++ function

```
template <class Function>
double testCvgRKSSM(const Function &f, double T, double y0,
const Eigen::MatrixXd &A, const Eigen::VectorXd &b)
```

in order_not_all.hpp that computes an approximated solution \mathbf{y}_M of (11.6.1) up to time T by means of an explicit Runge-Kutta method with $M=2^k$, $k=1,\ldots,12$, uniform timesteps. The method is defined by the Butcher scheme described by the inputs A (Butcher matrix) and b (weight vector), see [Lecture \rightarrow Eq. (11.4.0.13)]. The input f is an object with an evaluation operator (e.g. a lambda function) for arguments of type **double** representing \mathbf{f} . The input y0 passes the initial value \mathbf{y}_0 .

For each k, the function should print the error at the final point $E_M = |\mathbf{y}_M(T) - \mathbf{y}_{2^{15}}(T)|$, $M = 2^k$, $k = 1, \ldots, 12$, accepting $\mathbf{y}_{2^{15}}(T)$ as exact value. Assuming algebraic convergence for $E_M \approx CM^{-r}$, at each $k = 2, \ldots, 12$ also print an approximation of the order of convergence r_k (recall that $M = 2^k$). This will be an expression involving E_M and $E_{M/2}$.

Finally, compute and return an approximate order of convergence by linear regression. Only take into account results for which the error is larger than 10^{-14} in order not to be misled by the impact of round-off errors.

```
SOLUTION for (11-6.a) \rightarrow 11-6-1-0:OrdN1h.pdf
```

$$\dot{y} = (1 - y)y, \quad y(0) = 1/2,$$
 (11.6.4)

and of the initial value problem

$$\dot{y} = |1.1 - y| + 1, \quad y(0) = 1.$$
 (11.6.5)

HIDDEN HINT 1 for (11-6.b) \rightarrow 11-6-2-0:ONA2s.pdf

HIDDEN HINT 2 for (11-6.b) \rightarrow 11-6-2-1:OMA6s.pdf

HIDDEN HINT 3 for (11-6.b) \rightarrow 11-6-2-2:OMA8s.pdf

SOLUTION for (11-6.b)
$$\rightarrow$$
 11-6-2-3:OrdN2h.pdf

Using testCvgRKSSM() write in order_not_all.hpp a C++ function

void cmpCvgRKSSM();

that empirically determines and prints the rates of convergence of

- the explicit Euler method [Lecture → Eq. (11.2.1.5)], a RK single step method of order 1,
- the explicit trapezoidal rule [Lecture \rightarrow Eq. (11.4.0.8)], a RK single step method of order 2,
- an RK method of order 3 given by the Butcher tableau [Lecture \rightarrow Eq. (11.4.0.13)]

$$\begin{array}{c|ccccc}
0 & & & \\
1/2 & 1/2 & & \\
\hline
1 & -1 & 2 & \\
\hline
1/6 & 2/3 & 1/6 & \\
\end{array}$$

• the classical RK method of order 4, see [Lecture \rightarrow Ex. 11.4.0.17] for details.

when applied for the numerical integration of the initial-value problems (11.6.4) and (11.6.5). Use final time T=1 in each case.

Comment on the calculated order of convergence for the different methods and the two different initial value problems.

SOLUTION for (11-6.c)
$$\rightarrow$$
 11-6-3-0:x1s.pdf

End Problem 11-6, 80 min.

Problem 11-7: Initial Condition for Lotka-Volterra ODE

In this problem we will face a situation, where we need to compute the derivative of the solution of an initial value problem with respect to the initial state in order to apply Newton's method. So thus exercise covers both numerical integration and the solution of non-linear systems of equations.

You should grasp the abstract view of ODEs from [Lecture \rightarrow Section 11.1.4] and still remember Newton's method from [Lecture \rightarrow Section 8.5].

A differential equation for the derivative w.r.t. initial state. We consider IVPs for the autonomous ODE

$$\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y}) \tag{11.7.1}$$

with smooth right hand side $\mathbf{f} \colon D \to \mathbb{R}^d$, where $D \subseteq \mathbb{R}^d$ is the state space. We take for granted that for all initial states, solutions exist for all times (global solutions, see [Lecture \to Ass. 11.1.4.1]).

By its very definition given in [Lecture \rightarrow Def. 11.1.4.3], the evolution operator

$$\Phi \colon \mathbb{R} \times D \to D$$
, $(t, \mathbf{y}) \mapsto \Phi(t, \mathbf{y})$

satisfies

$$\frac{\partial \mathbf{\Phi}}{\partial t}(t, \mathbf{y}) = \mathbf{f}(\mathbf{\Phi}(t, \mathbf{y})).$$

Next, we can differentiate this identity with respect to the state variable y. We assume that all derivatives can be interchanged, which can be justified by rigorous arguments (which we won't do here). Thus, by the chain rule, we obtain, after swapping partial derivatives $\frac{\partial}{\partial t}$ and D_{y} ,

$$\frac{\partial \mathsf{D}_{\mathbf{y}} \, \mathbf{\Phi}}{\partial t}(t, \mathbf{y}) = \mathsf{D}_{\mathbf{y}} \, \frac{\partial \mathbf{\Phi}}{\partial t}(t, \mathbf{y}) = \mathsf{D}_{\mathbf{y}}(\mathbf{f}(\mathbf{\Phi}(t, \mathbf{y}))) = \mathsf{D} \, \mathbf{f}(\mathbf{\Phi}(t, \mathbf{y})) \, \mathsf{D}_{\mathbf{y}} \, \mathbf{\Phi}(t, \mathbf{y}).$$

Abbreviating $W(t, y) := D_y \Phi(t, y)$ we can rewrite this as the non-autonomous ODE

$$\dot{\mathbf{W}} = \mathsf{D}\,\mathbf{f}(\mathbf{\Phi}(t,\mathbf{y}))\mathbf{W} \quad .. \tag{11.7.2}$$

Here, the state y can be regarded as a parameter. Since $\Phi(0, y) = y$, we also know W(0, y) = I (identity matrix), which supplies an initial condition for (11.7.2). In fact, we can even merge (11.7.1) and (11.7.2) into the ODE

$$\frac{d}{dt}[\mathbf{y}(\cdot), \mathbf{W}(\cdot, \mathbf{y}_0)] = [\mathbf{f}(\mathbf{y}(t)), \mathbf{D}\mathbf{f}(\mathbf{y}(t))\mathbf{W}(t, \mathbf{y}_0)], \qquad (11.7.3)$$

which is autonomous again.

Now let us apply (11.7.2)/(11.7.3). As in [Lecture \rightarrow Ex. 11.1.2.5], we consider the following autonomous Lotka-Volterra differential equation of a predator-prey model

$$\dot{u} = (2 - v)u,
\dot{v} = (u - 1)v,$$
(11.7.4)

on the state space $D = \mathbb{R}^2_+$, $\mathbb{R}_+ = \{\xi \in \mathbb{R} : \xi > 0\}$. All the solutions of (11.7.4) are periodic and their period depends on the initial state $[u(0), v(0)]^T$. In this exercise we want to develop a numerical method which computes a suitable initial condition for a given period.

(11-7.a) \odot (5 min.) For fixed state $\mathbf{v} \in D$, (11.7.2) represents an ODE. What is its state space?

SOLUTION for (11-7.a)
$$\rightarrow$$
 11-7-1-0:Init1h.pdf

(11-7.b) (10 min.) What is the right hand side function for the ODE (11.7.2), when the underlying autonomous ODE $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$ is the Lotka-Volterra ODE (11.7.4)? You may write u(t), v(t) for solutions of (11.7.4).

SOLUTION for (11-7.b)
$$\rightarrow$$
 11-7-2-0:Init2h.pdf

$$\mathbf{y}_0 = \left[\begin{array}{c} u(0) \\ v(0) \end{array} \right]$$

is equal to a given value $T_P > 0$.

SOLUTION for (11-7.c)
$$\rightarrow$$
 11-7-3-0:Init4h.pdf

We write $W(T, y_0)$, $T \ge 0$, $y_0 \in \mathbb{R}^2_+$ for the solution of (11.7.2) for the underlying ODE (11.7.4). Express the Jacobian of F by means of W.

SOLUTION for (11-7.d)
$$\rightarrow$$
 11-7-4-0:Init5h.pdf

(11-7.e) ☐ (10 min.) [depends on Sub-problem (11-7.c)]

Argue, why the solution of $\mathbf{F}(\mathbf{y}) = 0$ will, in gneneral, not be unique. When will it be unique?

```
HIDDEN HINT 1 for (11-7.e) \rightarrow 11-7-5-0:Init1s.pdf
```

```
SOLUTION for (11-7.e) \rightarrow 11-7-5-1:Init6h.pdf
```

A C++ implementation **ode45** of an adaptive embedded Runge-Kutta method has been presented in [Lecture \rightarrow ??].

Instructions on the use of ode45, see [Lecture \rightarrow ??].

The class ode45 is header-only, meaning you just include the file and use it right away (no linking required). The file ode45.hpp defines the class ode45 implementing an embedded Runge-Kutta-Fehlberg method of order 4(5), see [Lecture \rightarrow Section 11.5.3] and [Lecture \rightarrow Ex. 11.5.3.2], with an adaptive stepsize control as presented in [Lecture \rightarrow § 11.5.2.11].

 Construct an object of ode45 type: create an instance of the class, passing the right-hand-side function f as a functor object to the constructor

```
template <class StateType,
class RhsType = std::function<StateType(const StateType &)>>
class ode45 {
  public:
    ode45(const RhsType &rhs);
    // ......
}
```

Template parameters are

- **StateType**: type of initial data and solution (state space), the only requirement is that the type possesses a normed vector-space structure, that is, it must implement the operations +, *, *=, += and assignment/copy operators. Moreover a norm() method must be available. EIGEN's vector and matrix types, as well as fundamental types are eligible as **StateType**.
- RhsType: type of rhs function (automatically deduced).

The argument rhs must be of a functor type that provides an evaluation operator

```
StateType operator()(const StateType & vec);
```

It can also be a lambda function.

2. (optional) Set the integration options: set data members of the data structure **ode45.options** to configure the solver:

```
0.options.<option_you_want_to_set> = <value>;
```

Examples:

- rtol: relative tolerance for error control (default is 1e-6)
- atol: absolute tolerance for error control (default is 1e-8)

e.g.:

```
O.options.rtol = 1e-5;
```

3. Solve stage: invoke the single-step method through calling the method

```
template <class NormFunc = decltype(_norm<StateType>)>
std::vector<std::pair<StateType, double>>
solve(const StateType &y0, double T,
const NormFunc &norm = _norm<StateType>);
```

The type **NormType** should provide a norm for vectors of type **StateType**. However, this type can be deduced automatically and the argument norm is optional. The other arguments are

- y0: initial value of type **StateType** ($\mathbf{v}_0 = y0$)
- T: final time of integration
- norm: (optional) norm function to call for objects of StateType, for the computation of the error

Return value The function returns the solution of the IVP, as a std::vector of std::pair (y(t),t) for every snapshot.

For more explanations and details, please consult the in-class documentation provided in the comments.

that computes $\Phi(T, [u_0, v_0]^T)$ and $\mathbf{W}(T, [u_0, v_0]^T)$. The first component of the output pair should contain $\Phi(T, [u_0, v_0]^T)$ and the second component the matrix $\mathbf{W}(T, [u_0, v_0]^T)$.

Set relative and absolute tolerances of **ode45** to 10^{-14} and 10^{-12} , respectively.

```
HIDDEN HINT 1 for (11-7.f) \rightarrow 11-7-6-0:Init2s.pdf

SOLUTION for (11-7.f) \rightarrow 11-7-6-1:Init7h.pdf

(11-7.g) ① (25 min.) [ depends on Sub-problem (11-7.f) and Sub-problem (11-7.d) ]

Relyiong on PhiAndW(), write a C++ routine

std::pair<double, double> findInitCond();
```

that determines and returns initial conditions u(0) and v(0) such that the solution of the system (11.7.4) has period $T_P = 5$. Use the multi-dimensional **Newton method** for $\mathbf{F}(\mathbf{y}) = 0$ with \mathbf{F} . As your initial approximation, use $[3,2]^T$. Terminate the Newton iteration as soon as $|\mathbf{F}(\mathbf{y})| \leq 10^{-5}$. Validate your implementation by comparing the obtained initial data \mathbf{y} with $\Phi(100,\mathbf{y})$.

Remark. The residual based termination criterion recommended above [Lecture \rightarrow § 8.2.3.4] is appropriate for this particular application and, in general, should not be used for Newton's method. Better termination criteria are proposed in [Lecture \rightarrow Section 8.5.3].

```
HIDDEN HINT 1 for (11-7.g) \rightarrow 11-7-7-0:lvxh.pdf
Solution for (11-7.g) \rightarrow 11-7-7-1:Init8h.pdf
```

End Problem 11-7, 90 min.

Problem 11-8: Integrating ODEs using the Taylor expansion method

In [Lecture \to Chapter 11] of the course we studied single step methods for the integration of initial value problems for ordinary differential equations $\dot{y}=f(y)$, [Lecture \to Def. 11.3.1.5]. Explicit single step methods have the advantage that they only rely on *point evaluations* of the right hand side f.

However, if derivatives of \mathbf{f} are also available, one have more options and this problem examines a class of single-step methods that also rely on $D \mathbf{f}$.

The new class of methods is obtained by the following reasoning: if the right hand side $\mathbf{f}: \mathbb{R}^n \to \mathbb{R}^n$ of an autonomous initial value problem

$$\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y}) , \qquad \mathbf{y}(0) = \mathbf{y}_0 , \qquad (11.8.1)$$

with solution $\mathbf{y}: \mathbb{R} \to \mathbb{R}^n$ is smooth, the solution $\mathbf{y}(t)$ will also be regular and it is possible to expand it into a Taylor sum at t = 0, see [Lecture \to Thm. 8.3.2.12],

$$\mathbf{y}(t) = \sum_{k=0}^{m} \frac{\mathbf{y}^{(k)}(0)}{k!} t^k + R_m(t) , \qquad (11.8.2)$$

with remainder term $R_m(t) = O(t^{m+1})$ for $t \to 0$.

A single step method for the numerical integration of (11.8.1) can be obtained by choosing m=3 in (11.8.2), neglecting the remainder term and taking the remaining sum as an approximation of $\mathbf{y}(h)$, that is

$$\mathbf{y}(h) \approx \mathbf{y}_1 := \mathbf{y}(0) + \frac{d\mathbf{y}(0)}{dt}h + \frac{1}{2}\frac{d^2\mathbf{y}(0)}{dt^2}h^2 + \frac{1}{6}\frac{d^3\mathbf{y}(0)}{dt^3}h^3$$
.

Subsequently, one uses the ODE and the initial condition to replace the temporal derivatives $\frac{d^{t}\mathbf{y}}{dt^{i}}$ with expressions in terms of (derivatives of) \mathbf{f} . This yields a single step integration method called *Taylor (expansion) method*.

(11-8.a) (15 min.) Express $\frac{d\mathbf{y}}{dt}(t)$ and $\frac{d^2\mathbf{y}}{dt^2}(t)$ in terms of \mathbf{f} and its Jacobian \mathbf{Df} .

HIDDEN HINT 1 for (11-8.a) \rightarrow 11-8-1-0: Tayl1h.pdf

SOLUTION for (11-8.a)
$$\rightarrow$$
 11-8-1-1: Tayl1s.pdf

(11-8.b) **③** (30 min.) [depends on Sub-problem (11-8.a)]

Verify the formula

$$\frac{d^3\mathbf{y}}{dt^3}(0) = \mathbf{D}^2\mathbf{f}(\mathbf{y}_0)(\mathbf{f}(\mathbf{y}_0), \mathbf{f}(\mathbf{y}_0)) + \mathbf{D}\mathbf{f}(\mathbf{y}_0)^2\mathbf{f}(\mathbf{y}_0).$$
(11.8.3)

HIDDEN HINT 1 for (11-8.b) \rightarrow 11-8-2-0: Tayl0h.pdf

HIDDEN HINT 2 for (11-8.b) \rightarrow 11-8-2-1: Tayl2h.pdf

HIDDEN HINT 3 for (11-8.b) \rightarrow 11-8-2-2: Tayl24h.pdf

SOLUTION for (11-8.b) \rightarrow 11-8-2-3: Tayl2s.pdf

(11-8.c) : Now we apply the Taylor expansion method introduced above to the following ODE for the predator-prey model, as introduced in [Lecture \rightarrow Ex. 11.1.2.5]:

$$\dot{y}_1(t) = (\alpha_1 - \beta_1 y_2(t)) y_1(t) , \qquad (11.8.4a)$$

$$\dot{y}_2(t) = (\beta_2 y_1(t) - \alpha_2) y_2(t)$$
, (11.8.4b)

$$\mathbf{y}(0) = [100, 5]^{\top}$$
 (11.8.4c)

To this end, in the template file taylorintegrator.hpp implement a C++ function

```
std::vector<Eigen::Vector2d> solvePredPreyTaylor(
  double alpha1, double beta1, double alpha2, double beta2,
  double T, const Eigen::Vector2d& y0, unsigned int M);
```

for the numerical integration of initial-value problems for (11.8.4) using the Taylor expansion method with $M, M \in \mathbb{N}$ uniform timesteps on the temporal interval [0, T]. The arguments alpha1, beta1, alpha2, beta2 provide the parameters of the ODE-based model (11.8.4), while y0 passes the initial value, and T the final time. M specifies the number of equidistant timesteps to be carried out. The function should return a sequence of approximate states \mathbf{y}_k , $k = 0, \dots, M$.

SOLUTION for (11-8.c)
$$\rightarrow$$
 11-8-3-0: Tayl3s.pdf

(11-8.d) (20 min.) [depends on Sub-problem (11-8.c)]

Based on the function solvePredPreyTaylor() implemented in Sub-problem (11-8.c) experimentally determine the order of convergence of the considered Taylor expansion method when it is applied to solve (11.8.4). Study the behaviour of the error at final time t=10 for the initial data $\mathbf{y}(0) = [100, 5]^{\top}$ and model coefficients $\alpha_1 = 3, \alpha_2 = 2, \beta_1 = \beta_2 = 0.1$. As reference solution, use $\mathbf{y}(10) = [0.319465882659820, 9.730809352326228]^{\top}$. Implement a function

```
double testCvgTaylorMethod();
```

that generates a suitable error table and returns the empiric rate of algebraic convergence in terms of the (uniform) timestep $h \to 0$. This rate should be determined by linear regression.

SOLUTION for (11-8.d)
$$\rightarrow$$
 11-8-4-0: Tayl4s.pdf

(11-8.e) (5 min.) What is the disadvantage of the Taylor's method compared with a Runge-Kutta method?

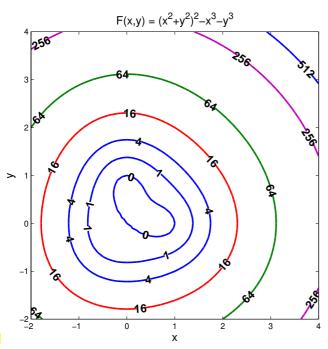
Solution for (11-8.e)
$$\rightarrow$$
 11-8-5-0: Tayl5s.pdf

End Problem 11-8, 70 min.

Problem 11-9: Drawing Isolines of a Function

Isoline/countour plots of functions are a popular way to visualize real-valued functions or data defined on a subset of the plane. In this exercise we will use numerical integration for finding isolines.

This problem assumes familiarity with explicit adaptive Runge-Kutta single-step methods as introduced in [Lecture \rightarrow Section 11.4] and [Lecture \rightarrow Section 11.5]. It relies on the **ode45** integrator class from [Lecture \rightarrow ??].



An isoline is a curve in the plane \mathbb{R}^2 that can be defined as the set of zeros of a continuous function $F: \mathbb{R}^2 \mapsto \mathbb{R}$:

$$\mathcal{C}:=\{x\in\mathbb{R}^2:\,F(x)=0\}\;.$$

In PYTHON (matplotlib) or MATLAB they can be plotted with the function contour

$$\triangleleft$$
 Contour plot of the function $F(x,y) = (x^2 + y^2)^2 - x^3 - y^3$

Fig. 104

If F is continuously differentiable, we can obtain connected components of C as solution curves of the initial value problems

$$\dot{\mathbf{y}}(t) = \frac{\mathbf{grad} \, F(\mathbf{y}(t))^{\perp}}{\|\mathbf{grad} \, F(\mathbf{y}(t))\|} \,, \quad F(\mathbf{y}(0)) = 0 \,. \tag{11.9.1}$$

Here $^{\perp}$ stands for the counterclockwise rotation of a vector $\in \mathbb{R}^2$ by 90° , that is,

$$\begin{bmatrix} x \\ y \end{bmatrix}^{\perp} = \begin{bmatrix} -y \\ x \end{bmatrix}, \quad x, y \in \mathbb{R} . \tag{11.9.2}$$

Then (11.9.1) is a consequence of the fact that an isoline always runs perpendicular to the direction of steepest ascent of F, which is provided by $\operatorname{grad} F$.

(11-9.a) (5 min.)

State the ordinary differential equation (11.9.1) for the concrete case of the function

$$F: \mathbb{R}^2 \to \mathbb{R}$$
, $F(x,y) := (x^2 + y^2)^2 - x^3 - y^3$. (11.9.3)

SOLUTION for (11-9.a) $\rightarrow 11-9-1-0:s1.pdf$

(11-9.b) (30 min.)

Use (with rtol = 0.00001, atol = 1e-9) the provided EIGEN-based integrator class **ode45** that employs an embedded explicit Runge-Kutta pair of methods to implement (in the file contour.hpp) a C++ function

that solves an initial value problem for (11.9.1) with initial state $\mathbf{y}_0 \in \mathbb{R}^2$ passed in y_0 from initial time 0 up to final time T (given in T). The argument gradF supplies a functor object equipped with an evaluation operator

```
Eigen::Vector2d operator(Eigen::Vector2d x) const;
```

that returns $\operatorname{grad} F(x)$. The function $\operatorname{computeIsolinePoints}()$ should return the sequence of computed states $y_0, y_1, \ldots, y_N, N \in \mathbb{N}$, in the columns of an $2 \times (N+1)$ -matrix, where N is the number of steps taken by the adaptive integrator of ode45.

Instructions on the use of ode45, see [Lecture \rightarrow ??].

The class ode45 is header-only, meaning you just include the file and use it right away (no linking required). The file ode45.hpp defines the class ode45 implementing an embedded Runge-Kutta-Fehlberg method of order 4(5), see [Lecture \rightarrow Section 11.5.3] and [Lecture \rightarrow Ex. 11.5.3.2], with an adaptive stepsize control as presented in [Lecture \rightarrow § 11.5.2.11].

1. Construct an object of **ode45** type: create an instance of the class, passing the right-hand-side function **f** as a functor object to the constructor

```
template <class StateType,
class RhsType = std::function<StateType(const StateType &)>>
class ode45 {
   public:
    ode45(const RhsType &rhs);
   // ......
}
```

Template parameters are

- **StateType**: type of initial data and solution (state space), the only requirement is that the type possesses a normed vector-space structure, that is, it must implement the operations +, *, *=, += and assignment/copy operators. Moreover a norm() method must be available. EIGEN's vector and matrix types, as well as fundamental types are eligible as **StateType**.
- RhsType: type of rhs function (automatically deduced).

The argument rhs must be of a functor type that provides an evaluation operator

```
StateType operator()(const StateType & vec);
```

It can also be a lambda function.

2. (optional) Set the integration options: set data members of the data structure **ode45.options** to configure the solver:

```
O.options.<option_you_want_to_set> = <value>;
```

Examples:

- rtol: relative tolerance for error control (default is 1e-6)
- atol: absolute tolerance for error control (default is 1e-8)

```
e.g.:
    O.options.rtol = 1e-5;
```

3. Solve stage: invoke the single-step method through calling the method

```
template <class NormFunc = decltype(_norm<StateType>)>
std::vector<std::pair<StateType, double>>
solve(const StateType &y0, double T,
const NormFunc &norm = _norm<StateType>);
```

The type **NormType** should provide a norm for vectors of type **StateType**. However, this type can be deduced automatically and the argument norm is optional. The other arguments are

- y0: initial value of type StateType ($y_0 = y0$)
- T: final time of integration
- norm: (optional) norm function to call for objects of StateType, for the computation of the error

Return value The function returns the solution of the IVP, as a std::vector of std::pair (y(t),t) for every snapshot.

For more explanations and details, please consult the in-class documentation provided in the comments.

```
Solution for (11-9.b) \rightarrow 11-9-2-0:s2.pdf
```

The 0-isoline of F from (11.9.3)

$$C_{\text{egg}} := \{ x = (x, y) \in \mathbb{R}^2 \setminus \{ \mathbf{0} \} : F(x) := (x^2 + y^2)^2 - x^3 - y^3 = 0 \}.$$

is known as "crooked egg curve". Based on <code>computeIsolinePoints()</code> and with $\mathbf{y}_0 = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \in \mathcal{C}_{egg}$ and T = 4 implement a C++ function (in the file <code>contour.hpp</code>)

```
Eigen::Matrix<double,2,Eigen::Dynamic> crookedEgg();
```

that returns coordinates of densely spaced points $\in \mathcal{C}_{egg}$ in the columns of the produced matrix.

```
HIDDEN HINT 1 for (11-9.c) \rightarrow 11-9-3-0:s3h1.pdf
```

```
Solution for (11-9.c) \rightarrow 11-9-3-1:s3.pdf
```

```
(11-9.d) (11-9.d) (20 min.) [ depends on Sub-problem (11-9.b) ]
```

Again rely on the integrator class ode45 to implement (in the file contour.hpp) a function

that performs the same computation as computeIsolinePoints() from Sub-problem (11-9.b). However, this time the argument F merely contains a functor object providing the function F itself.

To obtain an approximation of $\operatorname{grad} F$ rely on an approximation by symmetric difference quotients, e.g.,

$$\frac{\partial F}{\partial x_1}(x) \approx \frac{F(x+\left[\begin{smallmatrix} h \\ 0 \end{smallmatrix} \right]) - F(x-\left[\begin{smallmatrix} h \\ 0 \end{smallmatrix} \right])}{2h} \quad \text{for some} \quad h>0 \; .$$

Choose a suitable fixed h > 0 and explain your choice in a comment in your code.

HIDDEN HINT 1 for (11-9.d) \rightarrow 11-9-4-0:h4.pdf

Solution for (11-9.d) \rightarrow 11-9-4-1:s4.pdf

End Problem 11-9, 65 min.

Problem 11-10: Symplectic Timestepping for Equations of Motion

This problem examines a special class of timestepping schemes suitable for solving numerically equations of motions from classical mechanics. These timestepping schemes preserve profound stuctural properties of those equations and have been dubbed **symplectic**.

This problem considers a particular instance of a single-step method as introduced in [Lecture \rightarrow Section 11.3].

The authors of [RWR04] consider ordinary differential equations on the state space \mathbb{R}^{2n} , $n \in \mathbb{N}$, of the form

$$\frac{d}{dt} \begin{bmatrix} \mathbf{p} \\ \mathbf{q} \end{bmatrix} = \begin{bmatrix} \mathbf{f}(\mathbf{q}, t) \\ \mathbf{g}(\mathbf{p}) \end{bmatrix} \quad \text{with continuous functions} \quad \begin{aligned} \mathbf{f} : \mathbb{R}^n \times \mathbb{R} \to \mathbb{R}^n \\ \mathbf{g} : \mathbb{R}^n \to \mathbb{R}^n \end{aligned} . \tag{11.10.1}$$

For the temporal discretization of (11.10.1), [RWR04] proposes the following class of s-stage methods, $s \in \mathbb{N}$, with uniform timestep h > 0 for generating sequences $\left(\begin{bmatrix} \mathbf{p}_j \\ \mathbf{q}_j \end{bmatrix}\right)_{j=0,\dots,M}$ of approximations of $\begin{bmatrix} \mathbf{p}(jh) \\ \mathbf{q}(jh) \end{bmatrix}$, starting from initial values $\mathbf{p}_0 := \mathbf{p}(0)$ and $\mathbf{q}_0 := \mathbf{q}(0)$:

Here a_1, \ldots, a_s and b_1, \ldots, b_s are given coefficients, carefully chosen such that the single-step method achieves a certain order. A possible choice discussed in [**RWR04**] is, for s = 3,

$$a_1 := \frac{2}{3}$$
, $a_2 := -\frac{2}{3}$, $a_3 := 1$, $b_1 := \frac{7}{24}$ $b_2 := \frac{3}{4}$, $b_3 := -\frac{1}{24}$. (11.10.3)

(11-10.a) :: (10 min.)

Write a C++ function

void sympTimestep(double tau, Eigen::Vector2d& pq_j);

that implements a single timestep $\begin{bmatrix} \mathbf{p}_j \\ \mathbf{q}_j \end{bmatrix} \mapsto \begin{bmatrix} \mathbf{p}_{j+1} \\ \mathbf{q}_{j+1} \end{bmatrix}$ of the method (11.10.2) with s=3 and coefficients according to (11.10.3) for the "harmonic oscillator" ODE (n=1)

$$\dot{p} = q$$
 , $\dot{q} = -p$. (11.10.4)

The components of the argument vector supply p_i and q_i , and will be replaced with p_{i+1} and q_{i+1} .

SOLUTION for (11-10.a)
$$\rightarrow$$
 11-10-1-0:symp4.pdf

Experimentally determine the order of convergence of the single-step method implemented in Subproblem (11-10.a) in terms of the number of timesteps by

- integrating (11.10.4) over the period $[0, 2\pi]$,
- with initial values p(0) = 0, q(0) = 1,
- and using M = 10, 20, 40, 80, 160, 320, 640 timesteps.

Tabulate the following errors at the final time:

$$\operatorname{err}(M) := |p(2\pi) - p^{(M)}| + |q(2\pi) - q^{(M)}|, \quad M = 10, 20, 40, 80, 160, 320, 640,...$$

where $p^{(M)}$ and $q^{(M)}$ are the approximations of $p(2\pi)$ and $q(2\pi)$ produced by the single-step method. To that end complete the function <code>sympTimesteppingODETest()</code> in the file <code>symplectic.hpp</code>, which is called from <code>main()</code>.

HIDDEN HINT 1 for (11-10.b) \rightarrow 11-10-2-0:smp5h1.pdf

SOLUTION for (11-10.b)
$$\rightarrow$$
 11-10-2-1:symp5.pdf

The symplectic timestepping method has especially been designed for Hamiltonian ODEs

$$\dot{\mathbf{p}}(t) = -\frac{\partial H}{\partial \mathbf{q}}(\mathbf{p}(t), \mathbf{q}(t)),$$

$$\dot{\mathbf{q}}(t) = \frac{\partial H}{\partial \mathbf{p}}(\mathbf{p}(t), \mathbf{q}(t))$$
with $H : \mathbb{R}^n \times \mathbb{R}^n \to \mathbb{R}$. (11.10.7)

The smooth function H is called a **Hamiltonian**.

Notation: The partial derivative vector $\frac{\partial H}{\partial \mathbf{q}}(\mathbf{p}_0,\mathbf{q}_0)\in\mathbb{R}^n$ stands for the gradient of the function $\mathbf{q}\mapsto H(\mathbf{p}_0,\mathbf{q})$ in $\mathbf{q}=\mathbf{q}_0$. $\frac{\partial H}{\partial \mathbf{p}}(\mathbf{p}_0,\mathbf{q}_0)$ is defined analogously.

(11-10.c) $\ \ \Box$ (15 min.) Show that for every solution $t\mapsto (\mathbf{p}(t),\mathbf{q}(t)),\,t\in I\subset\mathbb{R},$ of (11.10.7) the function

$$t \mapsto \underline{E}(t) := H(\mathbf{p}(t), \mathbf{q}(t)) , \quad t \in I , \tag{11.10.8}$$

is constant.

Solution for (11-10.c)
$$\rightarrow 11-10-3-0:s6.pdf$$

(11-10.d) (10 min.) Write down explicitly in the form (11.10.1) the Hamiltonian ODE (11.10.7) for the particular Hamiltonian

$$H(\mathbf{p}, \mathbf{q}) = \frac{1}{2} \|\mathbf{p}\|_{2}^{2} + \|\mathbf{q}\|_{2}^{4}, \quad \mathbf{p}, \mathbf{q} \in \mathbb{R}^{n}, \quad n \in \mathbb{N}.$$
 (11.10.9)

SOLUTION for (11-10.d)
$$\rightarrow$$
 11-10-4-0:s7.pdf

```
Eigen::MatrixXd simulateHamiltonianDynamics(
  const Eigen::VectorXd &p0, const Eigen::VectorXd &q0,
  double T, unsigned int M);
```

that uses M equidistant steps of the sympletic timestepping scheme (11.10.2) + (11.10.3) in order solve the Hamiltonian ODE (11.10.7) with Hamiltonian H as in (11.10.9) up to final time T > 0 and with initial values \mathbf{p}_0 and \mathbf{q}_0 passed in the argument vectors \mathbf{p}_0 and \mathbf{q}_0 of equal length $n \in \mathbb{N}$.

The function should return the numerical solution

$$\begin{bmatrix} \mathbf{p}_0 \\ \mathbf{q}_0 \end{bmatrix}, \begin{bmatrix} \mathbf{p}_1 \\ \mathbf{q}_1 \end{bmatrix}, \dots, \begin{bmatrix} \mathbf{p}_M \\ \mathbf{q}_M \end{bmatrix},$$

in the columns of a $2n \times (M+1)$ -matrix.

Solution for (11-10.e)
$$\rightarrow$$
 11-10-5-0:s8.pdf

(11-10.f) (10 min.) The following table gives the quantities

$$\mathcal{E}(M) := \max\{E(\mathbf{p}_k, \mathbf{q}_k), k = 0, \dots, M\}, \qquad (11.10.12)$$

where E is the "energy" as defined in (11.10.8) and $\left(\begin{bmatrix}\mathbf{p}_k\\\mathbf{q}_k\end{bmatrix}\right)_{k=0}^M$ is the sequence produced by simulateHamiltonianDynamics () for n=2, $\mathbf{p}_0=\mathbf{0}$, $\mathbf{q}_0=\begin{bmatrix}1\\0\end{bmatrix}$, T=5 and M timesteps.

M	$\mathcal{E}(M)$	
10	1.136158638333363	
20	1.009854310088357	
40	1.001187287267320	Based on the data in table describe qualitatively and
80	1.000149549345033	quantitatively the expected asymptotic behavior of
160	1.000018794627444	$\mathcal{E}(M)$ for $M o \infty$.
320	1.000002358738605	
640	1.000000295580299	
1280	1.000000036994599	

SOLUTION for (11-10.f) $\rightarrow 11-10-6-0:s9.pdf$

End Problem 11-10, 80 min.

Problem 11-11: RK-SSMs and Discrete Evolutions

This problem takes a look at the discrete evolution operators spawned by explicit Runge-Kutta single-step methods (RK-SSMs).

This problem relies on the material covered in [Lecture \rightarrow Section 11.3] and [Lecture \rightarrow Section 11.4.

The formulas for an explicit s-stage, $s \in \mathbb{N}$, Runge-Kutta single-step method described by the Butcher scheme

with $c_i, b_i, a_{i,j} \in \mathbb{R}$ are given in the following definition:

Definition [Lecture → Def. 11.4.0.11]. Explicit Runge-Kutta single-step method

For $b_i, a_{i,j} \in \mathbb{R}$, $c_i := \sum_{j=1}^{i-1} a_{ij}, i,j=1,\ldots,s, s \in \mathbb{N}$, an s-stage explicit Runge-Kutta single **step method** (RK-SSM) for the ODE $\dot{\mathbf{y}} = \mathbf{f}(t, \mathbf{y}), \mathbf{f} : \Omega \to \mathbb{R}^N$, is defined by $(\mathbf{y}_0 \in D)$

$$\mathbf{k}_i := \mathbf{f}(t_0 + c_i h, \mathbf{y}_0 + h \sum_{j=1}^{i-1} a_{i,j} \mathbf{k}_j) , \quad i = 1, \dots, s , \quad \mathbf{y}_1 := \mathbf{y}_0 + h \sum_{i=1}^{s} b_i \mathbf{k}_i .$$

The vectors $\mathbf{k}_i \in \mathbb{R}^N$, $i=1,\ldots,s$, are called **increments**, h>0 is the size of the timestep.

(11-11.a) (10 min.) What conditions on the coefficients c_i , b_i , $a_{i,j} \in \mathbb{R}$ guarantee that the explicit Runge-Kutta single-step method according to the above definition is consistent with the autonomous ODE $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$ with Lipschitz continuous right-hand side function \mathbf{f} .

```
c_i, b_i, a_{i,i} such that
                                                                                 RK-SSM consistent.
```

```
HIDDEN HINT 1 for (11-11.a) \rightarrow 11-11-1-0:rkssmhal.pdf
```

SOLUTION for (11-11.a) \rightarrow 11-11-1-1:erkssmsa.pdf The templated class **ExpIRKSSMEvolOp** implements a functor type realizing the discrete evolution operator $(h, \mathbf{y}) \mapsto \Psi(h, \mathbf{y})$ for an explicit s-stage Runge-Kutta single-

step method described by the Butcher scheme $\frac{\mathbf{c} \mid \mathfrak{A}}{\mathbf{b}^T}$, $\mathbf{c}, \mathbf{b} \in \mathbb{R}^s$, $\mathfrak{A} \in \mathbb{R}^{s,s}$ strictly lower triangular,

```
applied to an autonomous ODE \dot{\mathbf{y}} = \mathbf{f}(\mathbf{y}).
```

```
template <typename RHSFunctor>
class ExplRKSSMEvolOp {
 public:
```

ExplRKSSMEvolOp (RHSFunctor& f, const Eigen::MatrixXd& A,

The constructor takes a functor object encoding the right-hand-side function f and $\mathfrak A$ and b as input arguments and stores them in class-local member variables.

As stated above, the evaluation operator **operator** () is supposed to realize the discrete evolution operator $\Psi(h,y)$ for the RK-SSM applied to $\dot{y}=f(y)$. Complete the implementation of Code 11.11.4 so that it meets this specification.

C++ code 11.11.4: Implementation of evaluation operator for ExpIRKSSMEvolOp 1 template <typename RHSFunctor> 2 template <typename State> State ExpIRKSSMEvolOp<RHSFunctor>::operator()(double h, const State& y) const { unsigned int $s = A_.cols();$ State y1{y}; 5 std::vector<State> k{s, y}; 6 for (int i = 0; i <; ++i) { 7 for (int j = 0; j <; ++j) { k[1; 9] += h * A (} 10); 11 12 $y1 += h * b_{[}$] * **k**[]; 13 14 return y1; } 15

```
HIDDEN HINT 1 for (11-11.b) \rightarrow 11-11-2-0:rkssmhb1.pdf
Solution for (11-11.b) \rightarrow 11-11-2-1:.pdf
```

End Problem 11-11, 30 min.

Chapter 12

Single Step Methods for Stiff Initial Value Problems

Problem 12-1: Implicit Runge-Kutta method

This problem addresses the implementation of general implicit Runge-Kutta methods [Lecture \rightarrow Def. 12.3.3.1]. We will adapt an integrator class so far available for explicit Runge-Kutta single-step methods to the implicit case.

Related to Problem 11-2. This problem assumes familiarity with [Lecture \rightarrow Section 12.3], and, especially, [Lecture \rightarrow Section 12.3.3] and [Lecture \rightarrow Rem. 12.3.3.9]. It practices the implementation of a full implicit RK-SSM in C++ based on Eigen.

Problem 11-2 introduced the class **RKIntegrator** that implemented the timestepping for a general explicit Runge-Kutta method according to [Lecture \rightarrow Def. 11.4.0.11]. Keeping the interface we now extend this class so that it realizes a general Runge-Kutta timestepping method as given in [Lecture \rightarrow Def. 12.3.3.1] for solving an autonomous initial value problem $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y}), \mathbf{y}(0) = \mathbf{y}_0$.

(12-1.a) (15 min.) Rederive the stage form

$$\mathbf{g}_i = h \sum_{j=1}^s a_{ij} \mathbf{f}(t_0 + c_j h, \mathbf{y}_0 + \mathbf{g}_j)$$
,
$$\mathbf{y}_1 = \mathbf{y}_0 + h \sum_{i=1}^s b_i \mathbf{f}(t_0 + c_i h, \mathbf{y}_0 + \mathbf{g}_i)$$
, [Lecture \rightarrow Eq. (12.3.3.8)]

of a general (implicit) Runge-Kutta method from the increment equations as given in [Lecture \rightarrow Def. 12.3.3.1].

SOLUTION for (12-1.a)
$$\rightarrow$$
 12-1-1-0:0s.pdf

SOLUTION for (12-1.b)
$$\rightarrow$$
 12-1-2-0:0as.pdf (12-1.c) \odot (45 min.) [depends on Sub-problem (12-1.b)]

By modifying the class **RKIntegrator** for the implementation of explicit Runge-Kutta methods, design in implicit_rkintegrator.hpp a similar header-only C++ class **implicitRKIntegrator** which implements a general implicit RK method given through a Butcher scheme [Lecture \rightarrow Eq. (12.3.3.3)] to solve the autonomous initial value problem $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y}), \mathbf{y}(0) = \mathbf{y}_0$.

This class should implement a generic implicit Runge-Kutta single-step method for an autonomous ODE according to [Lecture \rightarrow Def. 12.3.3.1]. The method is specified through its Butcher matrix $\mathfrak{A} \in \mathbb{R}^{s,s}$ and the weight vector $\mathbf{b} \in \mathbb{R}^{s}$ that are passed to the constructor as arguments \mathbb{A} and \mathbb{b} .

The solve () method carries out N equidistant timesteps of the method for the ODE $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$ on the time interval [0,T], T>0, and with initial value \mathbf{y}_0 . The right-hand side function \mathbf{f} and its Jacobian D \mathbf{f} are passed through the functor objects \mathbf{f} and $\mathbf{J}\mathbf{f}$ equipped with suitable evaluation operators.

In the solve () method the stages \mathbf{g}_i as introduced in ((12-1.a)) are to be computed with the **damped** Newton method (see [Lecture \rightarrow Section 8.5.4]) applied to the nonlinear system of equations satisfied by the stages (see [Lecture \rightarrow Rem. 12.3.3.6] and [Lecture \rightarrow Rem. 12.3.3.9]). Use the provided function

from the file dampnewton.hpp, that is a simplified version of [Lecture \rightarrow Code 8.5.4.5]. Note that we do not use the simplified Newton method as discussed in [Lecture \rightarrow Rem. 12.3.3.9].

In the code template you will find large parts of **implicitRKIntegrator** already implemented. In fact, you only have to write the method step () for the actual implicit RK-SSM timestepping.

```
SOLUTION for (12-1.c) \rightarrow 12-1-3-0:Impl1h.pdf (12-1.d) \odot (15 min.) Examine the following code
```

```
C++11-code 12.1.5: Initialization of implicitRKIntegrator object

// Definition of coefficients in Butcher scheme

constexpr unsigned int s = 2;

Eigen:: MatrixXd A(s, s);

Eigen:: VectorXd b(s);

// What method is this?

A << 5. / 12., -1. / 12., 3. / 4., 1. / 4.;

b << 3. / 4., 1. / 4.;

// Initialize implicit RK with Butcher scheme
```

```
implicitRKIntegrator RK(A, b);
```

Get it on ₩ GitLab (main.cpp).

Write down the complete Butcher scheme according to [Lecture \rightarrow Eq. (12.3.3.3)] for the implicit Runge-Kutta method now available through RK. Which method is it? Is it A-stable [Lecture \rightarrow Def. 12.3.4.9], L-stable [Lecture \rightarrow Def. 12.3.4.15]?

HIDDEN HINT 1 for (12-1.d)
$$\rightarrow$$
 12-1-4-0:Impl1s.pdf

SOLUTION for (12-1.d) \rightarrow 12-1-4-1:Impl2h.pdf

The file main.cpp uses your implementation in implicitRKIntegrator of general implicit RK-SSMs to solve an initial value problem for the Lotka-Volterra ordinary differential equation

$$\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y}) := \begin{bmatrix} y_1(\alpha_1 - \beta y_2) \\ y_2(-\alpha_2 + \beta y_2) \end{bmatrix}. \tag{12.1.6}$$

with the particular paramters $\alpha_1 = 3$, $\alpha_2 = 2$, $\beta = 0.1$. The right-hand-side function and its Jacobian are implemented as suitable lambda functions. Initial state is $\mathbf{y}_0 = \begin{bmatrix} 100 \ 5 \end{bmatrix}^\top$ and final time T = 10. The norm of the error at final time is tabulated.

Run the code. What information can be gleaned from the table?

SOLUTION for (12-1.e)
$$\rightarrow$$
 12-1-5-0: Impl3h.pdf

End Problem 12-1, 105 min.

Problem 12-2: Damped precession of a magnetic needle

This problem deals with a dynamical system from mechanics describing the movement of a rod-like magnet in a strong magnetic field. This can be modelled by an ODE with a particular invariant that can become stiff in the case of large friction.

Assumes an idea about the notion of stiffness [Lecture \rightarrow Notion 12.2.0.7] and knowledge of heuristic criteria for predicting stiffness of an IVP [Lecture \rightarrow § 12.2.0.18]. Also looks at simple implicit and semi-implicit [Lecture \rightarrow Section 12.4] single-step methods for a concrete ODE.

We consider the initial value problem

$$\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y}) := \mathbf{a} \times \mathbf{y} + c\mathbf{y} \times (\mathbf{a} \times \mathbf{y}), \quad \mathbf{y}(0) = \mathbf{y}_0 = [1, 1, 1]^{\top},$$
 (12.2.1)

where c > 0 and $\mathbf{a} \in \mathbb{R}^3$, $\|\mathbf{a}\|_2 = 1$.

Note: $\mathbf{x} \times \mathbf{y}$ denotes the **cross product** of the two vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^3$. It is defined by

$$\mathbf{x} \times \mathbf{y} = \begin{bmatrix} x_2 y_3 - x_3 y_2 \\ x_3 y_1 - x_1 y_3 \\ x_1 y_2 - x_2 y_1 \end{bmatrix}.$$

It satisfies $\mathbf{x} \times \mathbf{y} \perp \mathbf{x}$. In Eigen, it is available as $\mathbf{x}.\mathtt{cross}(\mathbf{y}) \rightarrow \mathbf{x}$ EIGEN documentation, #include <Eigen/Geometry> is required.

(12-2.a) \Box (10 min.) Show that $\|\mathbf{y}(t)\|_2 = \|\mathbf{y}_0\|_2$ for every solution \mathbf{y} of the ODE.

HIDDEN HINT 1 for (12-2.a) \rightarrow 12-2-1-0:Cros1s.pdf

SOLUTION for (12-2.a)
$$\rightarrow$$
 12-2-1-1:Cros1h.pdf

SOLUTION for (12-2.b)
$$\rightarrow$$
 12-2-2-0:Cros2h.pdf

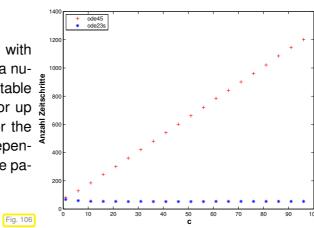
Discuss whether and when IVP for (12.2.1) will be stiff in the case $\mathbf{y}_0 \approx \mathbf{a}$.

HIDDEN HINT 1 for (12-2.c) \rightarrow 12-2-3-0:stiffcrit.pdf

SOLUTION for (12-2.c)
$$\rightarrow$$
 12-2-3-1:s22.pdf

(12-2.d) (10 min.) [depends on Sub-problem (12-2.b)]

For $\mathbf{a}=[1,0,0]^{\top}$, (12.2.1) was solved with the standard explicit adaptive Runge-Kutta numerical integrator **Ode45** and another A-stable semi-implicit adaptive numerical integrator up to the point T=10 (same tolerances for the stepsize control). Explain the different dependence of the total number of steps from the parameter c observed in the figure.



```
SOLUTION for (12-2.d) \rightarrow 12-2-4-0:Cros3h.pdf
```

```
SOLUTION for (12-2.e) \rightarrow 12-2-5-0:Cros4h.pdf (12-2.f) (20 min.) In cross.hpp, implement the C++ function
```

```
void tab_crossprod();
```

which solves (12.2.1) with $\mathbf{a} = [1, 0, 0]^{\top}$, c = 1, up to T = 10, using M = 128 uniform time steps of the implicit mid-point method, and tabulates $\|\mathbf{y}_k\|_2$ for the generated sequence of approximate states. What do you observe?

For this task rely on the helper class **implicitRKIntegrator**:

```
class implicitRKIntegrator {
public:
   implicitRKIntegrator(
      const Eigen::MatrixXd &A,
      const Eigen::VectorXd &b);
   template <class Function, class Jacobian>
   std::vector<Eigen::VectorXd> solve(
      Function &&f, Jacobian &&Jf, double T,
      const Eigen::VectorXd &y0,
      unsigned int M) const;
private:
   ......
};
```

It implements a generic implicit Runge-Kutta single-step method for an autonomous ODE according to [Lecture \rightarrow Def. 12.3.3.1]. The method is specified through its Butcher matrix $\mathfrak{A} \in \mathbb{R}^{s,s}$ and the weight vector $\mathbf{b} \in \mathbb{R}^s$ that are passed to the constructor as arguments \mathbf{A} and \mathbf{b} .

The solve () method carries out M equidistant timesteps of the method for the ODE $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$ on the time interval [0,T], T>0, and with initial value \mathbf{y}_0 . The right-hand side function \mathbf{f} and its Jacobian D \mathbf{f} are passed through the functor objects \mathbf{f} and \mathbf{J} equipped with suitable evaluation operators. The solve () method relies on the damped Newton iteration to solve the stage equations, see [Lecture \rightarrow Rem. 12.3.3.6] and [Lecture \rightarrow Rem. 12.3.3.9].

```
SOLUTION for (12-2.f) \rightarrow 12-2-6-0:Cros5h.pdf
```

(12-2.g) (20 min.) The so-called **linear-implicit mid-point method** can be obtained by a simple *linearization* of the (single) increment equation of the implicit mid-point method around the current state.

Equivalently, we can obtain it from the standard implicit midpoint method by carrying out a single step of the Newton iteration for the increment equations with initial guess **0**.

Give the defining equation of the linear-implicit mid-point method for the general autonomous differential equation

$$\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$$

with smooth right-hand-side function $\mathbf{f}: \mathbb{R}^N \to \mathbb{R}^N$, $N \in \mathbb{N}$.

```
HIDDEN HINT 1 for (12-2.g) \rightarrow 12-2-7-0:linhi.pdf
```

```
SOLUTION for (12-2.g) \rightarrow 12-2-7-1:Cros6h.pdf
```

Implement the linear-implicit midpoint method in the C++ function

```
template <class Function, class Jacobian>
std::vector<Eigen::VectorXd> solve_lin_mid(
   Function &&f, Jacobian &&Jf,
   double T, const Eigen::VectorXd &y0,
   unsigned int M);
```

The signature of this function is the same as that of the solve () method of implicitRKIntegrator.

Extend your implementation of tab_crossprod() so that it uses solve_lin_mid() to solve (12.2.1) with $\mathbf{a} = [1,0,0]^{\mathsf{T}}$, c=1 up to T=10 and M=128. Tabulate $\|\mathbf{y}_k\|_2$ for the sequence of approximate states generated by the linear implicit midpoint method. What do you observe?

SOLUTION for (12-2.h) \rightarrow 12-2-8-0:Cros7h.pdf

End Problem 12-2, 120 min.

Problem 12-3: Singly Diagonally Implicit Runge-Kutta Method

SDIRK methods (Singly Diagonally Implicit Runge-Kutta methods) are distinguished by Butcher schemes with a lower triangular coefficient matrix $\mathfrak A$ whose diagonal entries are all the same:

$$\frac{\mathbf{c} \quad \mathfrak{A}}{\mathbf{b}^{T}} = \begin{array}{c}
c_{1} \quad \gamma & \cdots & 0 \\
c_{2} \quad a_{21} \quad \ddots & \vdots \\
\vdots \quad \vdots & & \vdots \\
\vdots \quad \vdots & & \vdots \\
\vdots \quad \vdots & & \ddots & \vdots \\
c_{s} \quad a_{s1} \quad \cdots \quad a_{s,s-1} \quad \gamma \\
\hline
b_{1} \quad \cdots \quad b_{s-1} \quad b_{s}
\end{array}$$

$$(12.3.1)$$

with $\gamma \neq 0$. Their advantage is that systems of equations of the same structure have to be solved in order to compute the s increments sequentially.

Familiarity with [Lecture \rightarrow Section 12.3.3] is required. This problem has a focus on stability theory.

In this problem, the scalar linear initial value problem of second order

$$\ddot{y} + \dot{y} + y = 0, \quad y(0) = 1, \quad \dot{y}(0) = 0$$
 (12.3.2)

should be solved numerically using a family of 2-stage SDIRK method described by the Butcher scheme

$$\begin{array}{c|cccc}
\gamma & \gamma & 0 \\
1 - \gamma & 1 - 2\gamma & \gamma \\
\hline
& 1/2 & 1/2
\end{array}, (12.3.3)$$

where $\gamma > 0$ is a parameter.

(12-3.a) • (5 min.) Assume you want to solve an initial-value problem for an ordinary differential equation $\dot{\mathbf{y}} = \mathbf{f}(t,\mathbf{y}), \ \mathbf{f}: I \times D \subset \mathbb{R} \times \mathbb{R}^N \to \mathbb{R} \times \mathbb{R}^N$, numerically. Explain the benefit of using SDIRK-SSMs compared to using Gauss-Radau RK-SSMs as introduced in [Lecture \rightarrow Ex. 12.3.4.21]. In what situations will this benefit matter much?

HIDDEN HINT 1 for (12-3.a) \rightarrow 12-3-1-0:SDIR0h.pdf

SOLUTION for (12-3.a)
$$\rightarrow$$
 12-3-1-1:SDIR0s.pdf

State the equations for the increments k_1 and k_2 of the Runge-Kutta method (12.3.3) applied to the initial value problem corresponding to the differential equation $\dot{y} = \mathbf{f}(t, y)$.

Solution for (12-3.b)
$$\rightarrow$$
 12-3-2-0:SDIR1s.pdf

(12-3.c) \odot (20 min.) Show that the stability function S(z) of the SDIRK method (12.3.3) is given by

$$S(z) = \frac{1 + z(1 - 2\gamma) + z^2(1/2 - 2\gamma + \gamma^2)}{(1 - \gamma z)^2}$$

and plot the stability domain S_{Ψ} for the parameter value $\gamma = 1$ using the supplied PYTHON script stabdomSDIRK.py

Solution for (12-3.c)
$$\rightarrow$$
 12-3-3-0:SDIR2s.pdf

(12-3.d) \Box (25 min.) Find out whether for $\gamma = 1$ the SDIRK RK-SSM (12.3.3) is

A-stable.

· L-stable.

Argue, based on the following version of the maximum principle for analytic functions:

Lemma 12.3.4. Maximum principle for rational functions

If a rational function $z \in \mathbb{C} \mapsto R(z)$ has no pole in the left half plane $\mathbb{C}^- := \{z \in \mathbb{C} : \operatorname{Re} z \leq 0\}$, then either

$$\max\{|R(z)|: z \in \mathbb{C}^-\} = \max\{|R(z)|: \operatorname{Re} z = 0\}$$

or

$$\min\{|R(z)|: z \in \mathbb{C}^-\} = \min\{|R(z)|: \operatorname{Re} z = 0\}$$
,

that is, the function attains on the imaginary axis either its maximum or its minimum over the whole left half plane.

```
Solution for (12-3.d) \rightarrow 12-3-4-0:SDIRxs.pdf
```

(12-3.e) (10 min.) Formulate (12.3.2) as an initial value problem for a linear first order system for the function $z(t) = (y(t), \dot{y}(t))^{\top}$.

```
SOLUTION for (12-3.e) \rightarrow 12-3-5-0:SDIR3s.pdf
```

(12-3.f) (20 min.) [depends on Sub-problem (12-3.e)]

Implement a C++ function

```
Eigen:: Vector2d SdirkStep(const Eigen:: Vector2d &z0, double h,
  double gamma);
```

that realizes one step of the method (12.3.3) for the linear differential equation (12.3.2), starting from the initial state z0 and returning the state after a single step of size h.

```
Solution for (12-3.f) \rightarrow 12-3-6-0:SDIR4s.pdf
```

Realize a C++ function

```
double cvgSDIRK();
```

to conduct a numerical experiment, which gives an indication of the order of the method (with $\gamma = \frac{3+\sqrt{3}}{6}$) for the initial value problem from (12.3.3). Choose $\mathbf{z}_0 = [1, 0]^{\mathsf{T}}$ as initial value, $\mathsf{T} = 10$ as end time and M=20, 40, 80, ..., 10240 as numbers of equidistant timesteps. Tabulate the error $|y(T)-y_M|$ $(y := (z)_1)$ at final time and use it as a basis for the estimation of the rate of algebraic convergence by means of the supplied function polyfit ().

```
HIDDEN HINT 1 for (12-3.g) \rightarrow 12-3-7-0:sd5h1.pdf
```

```
SOLUTION for (12-3.g) \rightarrow 12-3-7-1:SDIR5s.pdf
```

End Problem 12-3, 115 min.

Problem 12-4: Semi-implicit Runge-Kutta SSM

General implicit Runge-Kutta methods as introduced in [Lecture \rightarrow Section 12.3.3] entail solving systems of non-linear equations for the increments, see [Lecture \rightarrow Rem. 12.3.3.9]. Semi-implicit Runge-Kutta single step methods, also known as Rosenbrock-Wanner (ROW) methods [Lecture \rightarrow Eq. (12.4.0.9)] just require the solution of linear systems of equations. This problem deals with a concrete ROW method, its stability and aspects of its implementation.

Relies on the contents of [Lecture \rightarrow Section 12.1] and also depends on [Lecture \rightarrow Section 12.3.4].

We consider the autonomous ODE

$$\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y}),\tag{12.4.1}$$

and discretize it with a *semi-implicit* Runge-Kutta SSM: a so-called **Rosenbrock method** given by

$$\begin{cases}
\mathbf{W}\mathbf{k}_{1} = \mathbf{f}(\mathbf{y}_{0}), \\
\mathbf{W}\mathbf{k}_{2} = \mathbf{f}(\mathbf{y}_{0} + \frac{1}{2}h\,\mathbf{k}_{1}) - ah\,\mathbf{J}\,\mathbf{k}_{1}, \\
\mathbf{y}_{1} = \mathbf{y}_{0} + h\mathbf{k}_{2},
\end{cases} (12.4.2)$$

where

$$J = Df(y_0)$$
 , $W = I - ah J$, $a := \frac{1}{2 + \sqrt{2}}$.

$$\begin{split} (\mathbf{I} - ha_{ii}\mathbf{J})\mathbf{k}_i &= \mathbf{f}(\mathbf{y}_0 + h\sum_{j=1}^{i-1}(a_{ij} + d_{ij})\mathbf{k}_j) - h\mathbf{J}\sum_{j=1}^{i-1}d_{ij}\mathbf{k}_j \;, \quad \mathbf{J} = \mathsf{D}\,\mathbf{f}(\mathbf{y}_0) \;, \\ \mathbf{y}_1 &:= \mathbf{y}_0 + h\sum_{j=1}^{s}b_j\mathbf{k}_j \;. \end{split}$$
 [Lecture \to Eq. (12.4.0.9)]

What is s and what are the coefficients a_{ij} , d_{ij} , and b_i for the method (12.4.2)?

SOLUTION for (12-4.a)
$$\rightarrow$$
 12-4-1-0:simplA.pdf

(12-4.b) (15 min.) Refresh your knowledge about how you can compute the stability function of a Runge-Kutta single-step function.

SOLUTION for (12-4.b)
$$\rightarrow$$
 12-4-2-0:.pdf

Compute the stability function S of the Rosenbrock method (12.4.2) for the linear scalar model ODE $\dot{y} = \lambda y, \lambda \in \mathbb{C}$. In other words, compute the (rational) function S(z), such that

$$y_1 = S(z)y_0, \quad z := h\lambda,$$

when the method is applied to perform a single setp of size h starting from y_0 .

Solution for (12-4.c)
$$\rightarrow$$
 12-4-3-0:SemI1s.pdf

Compute the first 4 terms of the Taylor expansion of S(z) around z=0. What is the maximal $q\in\mathbb{N}$ such that

$$|S(z) - \exp(z)| = O(|z|^q)$$

for $|z| \to 0$? Deduce the maximal possible order of the method Eq. (12.4.2).

HIDDEN HINT 1 for (12-4.d) \rightarrow 12-4-4-0:SemI2h.pdf

SOLUTION for (12-4.d) \rightarrow 12-4-4-1: SemI2s.pdf

(12-4.e) ☐ (30 min.) In rosenbrock.hpp, implement a C++ function:

that applies the Rosenbrock method (12.4.2) for solving an initial-value problem for an autonomous ODE $\dot{y} = f(y)$.

It takes as input functor objects for \mathbf{f} and $D\mathbf{f}$ (e.g., as lambda functions), an initial data (vector or scalar) $y0 = \mathbf{y}(0)$, a number of steps M and a final time T. The function returns the sequence of states generated by the single step method up to t = T, using M equidistant steps of the Rosenbrock method.

SOLUTION for (12-4.e)
$$\rightarrow$$
 12-4-5-0:SemI3s.pdf

Explore the order of the method (12.4.2) empirically by applying it to the IVP for the limit cycle [Lecture \rightarrow Ex. 12.2.0.4]:

$$\mathbf{f}(\mathbf{y}) := \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix} \mathbf{y} + \lambda (1 - \|\mathbf{y}\|^2) \mathbf{y} , \qquad (12.4.5)$$

with $\lambda=1$ and initial state $\mathbf{y}_0=[1,1]^{\top}$ on [0,10]. Use uniform timesteps of size $h=2^{-k}, k=4,\ldots,10$ and compute a reference solution $\mathbf{y}_{\mathrm{ref}}$ with timestep size $h=2^{-12}$. Monitor the maximal error on the temporal mesh

$$\max_{i} \|\mathbf{y}_{j} - \mathbf{y}_{\text{ref}}(t_{j})\|_{2}$$
,

and use it to estimate a rate of algebraic convergence by means of linear regression.

To do this write in rosenbrock.hpp a C++ function

double cvgRosenbrock();

that tabulates the errors and returns the estimated rate of convergence.

SOLUTION for (12-4.f)
$$\rightarrow$$
 12-4-6-0:SemI4s.pdf

In complex analysis you might have hear about the maximum principle for holomorphic/analytic functions. The following is a special version of it.

Theorem 12.4.7. Maximum principle for holomorphic functions

Let

$$\mathbb{C}^- := \{ z \in \mathbb{C} \mid Re(z) < 0 \} .$$

Let $f:D\subset\mathbb{C}\to\mathbb{C}$ be non-constant, defined on $\overline{\mathbb{C}^-}$, and analytic in \mathbb{C}^- . Furthermore, assume that $w:=\lim_{|z|\to\infty}f(z)$ exists and $w\in\mathbb{C}$, then:

$$\forall z \in \mathbb{C}^-$$
: $|f(z)| < \sup_{\tau \in \mathbb{R}} |f(i\tau)|$.

HIDDEN HINT 1 for (12-4.g) \rightarrow 12-4-7-0:SemI5h.pdf

Solution for (12-4.g) \rightarrow 12-4-7-1:SemI5s.pdf

End Problem 12-4, 120 min.

Problem 12-5: Exponential integrator

The exponential integrators are a modern class of single step methods developed for special initial value problems (problems that can be regarded as perturbed linear ODEs), see

M. HOCHBRUCK AND A. OSTERMANN, *Exponential integrators*, Acta Numerica, 19 (2010), pp. 209–286.

These methods fit the concept of single step methods as introduced in [Lecture \rightarrow Def. 11.3.1.5] and, usually, converge algebraically according to [Lecture \rightarrow Eq. (11.3.2.7)].

You have to be familiar with fundamental concepts for single-step methods as introduced in [Lecture \rightarrow Section 11.3.1], their asymptic convergence properties from [Lecture \rightarrow Section 11.3.2] and the notion of a stability domain [Lecture \rightarrow Def. 12.1.0.51].

A step with size h of the so-called **exponential Euler** single step method for the autonomous ODE $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$ with continuously differentiable $\mathbf{f} : \mathbb{R}^N \to \mathbb{R}^N$ reads:

$$\mathbf{y}_1 = \mathbf{y}_0 + h \ \varphi(h \ \mathsf{D} \ \mathbf{f}(\mathbf{y}_0)) \ \mathbf{f}(\mathbf{y}_0), \tag{12.5.1}$$

where $D \mathbf{f}(\mathbf{y}) \in \mathbb{R}^{N,N}$ is the Jacobian of \mathbf{f} at $\mathbf{y} \in \mathbb{R}^N$, and the matrix function $\varphi : \mathbb{R}^{N,N} \to \mathbb{R}^{N,N}$ is defined as

$$\varphi(\mathbf{Z}) = (\exp(\mathbf{Z}) - \operatorname{Id}) \mathbf{Z}^{-1}, \quad \mathbf{Z} \in \mathbb{R}^{N,N}.$$
 (12.5.2)

Here, $\exp(\mathbf{Z})$ is the matrix exponential of \mathbf{Z} , a special function $\exp: \mathbb{R}^{N,N} \to \mathbb{R}^{N,N}$, see [Lecture \to Eq. (12.1.0.34)].

The function φ is implemented in the file exponential integrator. hpp as the function

When plugging in the exponential series, it is clear that the function $z \mapsto \varphi(z) := \frac{\exp(z) - 1}{z}$ is analytic on \mathbb{C} with $\varphi(0) = 1$. Thus, $\varphi(\mathbf{Z})$ is well defined for all matrices $\mathbf{Z} \in \mathbb{R}^{N,N}$.

SOLUTION for (12-5.a)
$$\rightarrow$$
 12-5-1-0:Expo1h.pdf

(12-5.b) (15 min.) Show that the exponential Euler single step method defined in (12.5.1) solves the linear initial value problem

$$\dot{\mathbf{y}} = \mathbf{A} \ \mathbf{y} \ , \quad \mathbf{y}(0) = \mathbf{y}_0 \in \mathbb{R}^d \ , \qquad \mathbf{A} \in \mathbb{R}^{d,d} \ ,$$
 (12.5.3)

exactly.

HIDDEN HINT 1 for (12-5.b) \rightarrow 12-5-2-0:Expo1s.pdf

SOLUTION for (12-5.b)
$$\rightarrow$$
 12-5-2-1:Expo2h.pdf

Determine the region of stability of the exponential Euler single step method defined in (12.5.1) (see [Lecture \rightarrow Def. 12.1.0.51]).

SOLUTION for (12-5.c)
$$\rightarrow$$
 12-5-3-0: Expo3h.pdf

(12-5.d) ☐ (20 min.) In exponentialintegrator.hpp, write a C++ function

```
template <class Function, class Jacobian>
Eigen::VectorXd exponentialEulerStep(const Eigen::VectorXd &y0,
    Function &&f, Jacobian &&df, double h)
```

that implements (12.5.1). Here f and df are objects with evaluation operators representing the ODE right-hand side function $\mathbf{f}: \mathbb{R}^N \to \mathbb{R}^N$ and its Jacobian, respectively.

SOLUTION for (12-5.d)
$$\rightarrow$$
 12-5-4-0:Expo4h.pdf

(12-5.e) ☑ (30 min.) What is the order of the single step method (12.5.1)?

To investigate it empirically, in the file exponentialintegrator.hpp, implement a C++ function

```
void testExpEulerLogODE();
```

that applies the method to the scalar logistic ODE

$$\dot{y} = y (1 - y)$$
, $y(0) = 0.1$,

in the time interval [0,1]. Tabulate the error at the final time against the stepsize h=T/M, $M=2^k$ for $k=1,\ldots,15$. Discuss qualitatively and quantitative the empiric asymptotic convergence of the single-step method.

HIDDEN HINT 1 for (12-5.e) \rightarrow 12-5-5-0:Expo3s.pdf

SOLUTION for (12-5.e) \rightarrow 12-5-5-1:Expo5h.pdf

End Problem 12-5, 100 min.

Problem 12-6: Mono-implicit Runge-Kutta single step method

As explained in [Lecture \rightarrow Section 12.3.4] we must resort to implicit single-step methods as solvers for stiff initial-value problems for ODEs. The big challenge about *implicit* Runge-Kutta single-step methods (RK-SSMs) is the need to solve possibly big non-linear systems of equations. Many attempts have been made to alleviate this drawback and this problem presents one of them.

This problem needs the linear model problem analysis as introduced in [Lecture \rightarrow Section 12.1] and [Lecture \rightarrow Section 12.3.4]. It practices the implementation of implicit single step methods with a focus on solving the non-linear systems of equations, *cf.* [Lecture \rightarrow Rem. 12.3.3.9]. A little C++ implementation based on Eigen is requested.

A so-called *s*-stage mono-implicit Runge-Kutta single step method (MIRK) for the autonomous ODE $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y}), \mathbf{f} : D \subset \mathbb{R}^N \to \mathbb{R}^N$, is defined as:

$$\mathbf{g}_{i} := (1 - v_{i})\mathbf{y}_{0} + v_{i}\mathbf{y}_{1} + h \sum_{j=1}^{i-1} d_{i,j}\mathbf{f}(\mathbf{g}_{j}), \quad i = 1, \dots s,$$

$$\mathbf{y}_{1} := \mathbf{y}_{0} + h \sum_{j=1}^{s} b_{j}\mathbf{f}(\mathbf{g}_{j})$$
(12.6.1)

for suitable coefficients $v_i, b_i, d_{i,j} \in \mathbb{R}$, fixed to achieve a desired order.

(12-6.a) \odot (20 min.) Single step methods defined as in Eq. (12.6.1) belong to the class of implicit Runge-Kutta methods [Lecture \rightarrow Def. 12.3.3.1]. Write down the corresponding Butcher scheme [Lecture \rightarrow Eq. (12.3.3.3)] in terms of the coefficients v_i , b_i , $d_{i,j}$.

HIDDEN HINT 1 for (12-6.a) \rightarrow 12-6-1-0:mirk0.pdf

SOLUTION for (12-6.a) \rightarrow 12-6-1-1:MIRK1h.pdf

(12-6.b) \Box (20 min.) Compute the stability function of a MIRK scheme defined as in Eq. (12.6.1) for s=2.

HIDDEN HINT 1 for (12-6.b) \rightarrow 12-6-2-0:MIRK1.pdf

SOLUTION for (12-6.b) \rightarrow 12-6-2-1:MIRK2h.pdf

(12-6.c) $\[\]$ (10 min.) Now, we consider the special case of a scalar ODE (N=1) and s=2. Abbreviating $\mathbf{z} := [g_1, g_2, y_1]^\top$, rewrite Eq. (12.6.1) as a non-linear system of equations in the form $\mathbf{F}(\mathbf{z}) = \mathbf{0}$ for an explicitly specified suitable function $\mathbf{F} : \mathbb{R}^3 \to \mathbb{R}^3$.

SOLUTION for (12-6.c) \rightarrow 12-6-3-0:MIRK3h.pdf

Find the Jacobian $D\mathbf{F}(\mathbf{z})$ of the function \mathbf{F} (in terms of $d_{i,j}, v_i, b_i$ and the derivative of f) from the previous sub-problem.

SOLUTION for (12-6.d) \rightarrow 12-6-4-0:MIRK4h.pdf

In the following sub-problems, we will implement the 2-stage MIRK scheme for the scalar case: N=1, s=2.

In mirk.hpp, implement a function

that approximates the solution of F(z)=0 by performing two steps of the Newton method applied to F(z)=0. Use z to pass the initial guess. Objects of type Func and Jac have to supply suitable evaluation operators operator (), which, for the type Jac must return an EIGEN-compatible matrix.

SOLUTION for (12-6.e) \rightarrow 12-6-5-0:MIRK5h.pdf

(12-6.f) (20 min.) [depends on Sub-problem (12-6.e)]

Consider the particular MIRK scheme given by the coefficients:

$$v_1 = 1$$
, $v_2 = \frac{344}{2025}$, $d_{21} = -\frac{164}{2025}$, $b_1 = \frac{37}{82}$, $b_2 = \frac{45}{82}$. (12.6.6)

Using the function Newton2Steps, implement in mirk.hpp a function

```
template <class Func, class Jac>
double MIRKStep(Func &&f, Jac &&df, double y0, double h);
```

that realizes one step of the MIRK scheme defined by Eq. (12.6.1) for s=2 and a scalar ODE, that is, d=1. The right hand side function f and its Jacobian are passed through f and f. The solution of the nonlinear system arising from Eq. (12.6.1) is approximated using two Newton steps, namely by using the function f Newton2Steps (). The initial guess has to be chosen appropriately!

SOLUTION for (12-6.f)
$$\rightarrow$$
 12-6-6-0:MIRK6h.pdf

Implement in mirk.hpp a function

for the solution of a scalar ODE up to time T, using M equidistant steps of the mono-implicit Runge-Kutta single step method defined by (12.6.6). The initial value is passed in y0.

SOLUTION for (12-6.g)
$$\rightarrow$$
 12-6-7-0:MIRK7h.pdf

(12-6.h) ☑ (30 min.) In the file mirk.hpp implement a C++ function

void cvgMIRK();

that applies your implementation of the MIRK method to the IVP

$$\dot{y} = 1 + y^2$$
 , $y(0) = 0$, (12.6.9)

on [0,1]. The exact solution of (12.6.9) is $y_{ex}(t) := \tan t$. Compute the solution y_M at T=1 with a sequence of uniform temporal meshes with $M=4,\ldots,512$ intervals. Compute the error $|y_M(1)-y_{ex}(1)|$ and output an error table for $M=4,\ldots,512$.

Determine the empiric rate of convergence of the MIRK method for (12.6.9).

SOLUTION for (12-6.h) \rightarrow 12-6-8-0:MIRK8h.pdf

End Problem 12-6, 160 min.

Problem 12-7: Stability of a Runge-Kutta method

This problem is devoted to an empirical study of stability problems haunting explicit Runge-Kutta single-step methods, whose stability functions invariably are polynomials.

This problem is meant to supplement the discussion in [Lecture \rightarrow Section 12.1].

We focus on a 3-stage Runge-Kutta single step method described by the following Butcher-Tableau [Lecture \rightarrow Eq. (11.4.0.13)].

We also consider the following concrete case of the prey/predator model as introduced in [Lecture \rightarrow Ex. 11.1.2.5]:

$$\dot{y}_1(t) = (1 - y_2(t))y_1(t) ,
\dot{y}_2(t) = (y_1(t) - 1)y_2(t) .$$
(12.7.2)

```
Eigen::Vector2d PredPrey(Eigen::Vector2d y0, double T, unsigned
   int M);
```

that uses the RK-SSM (12.7.1) to solve an initial value problem for (12.7.2) with initial value \mathbf{y}_0 and M equidistant timestep up to final time T > 0. It should return $\mathbf{y}_N \approx \mathbf{y}(T)$.

SOLUTION for (12-7.a)
$$\rightarrow$$
 12-7-1-0:Stab0s.pdf

(12-7.b) (20 min.) [depends on Sub-problem (12-7.a)]

Write a C++ function

```
void SimulatePredPrey();
```

to approximate the solution at time T=1 of the IVP for (12.7.2) with initial value $\mathbf{y}_0=\mathbf{y}(0)=[100,1]^{\top}$. Use the RK-SSM (12.7.1) with uniform timestep h>0 and the function from Sub-problem (12-7.a).

Qualitatively and quantitatively determine the type of asymptotic convergence of the method for uniform steps of size 2^{-j} , $j=2,\ldots,13$. As a reference solution, use an approximation obtained with 2^{14} steps. Tabulate the norm of the error at final time.

What do you notice for big step sizes? Try to find the maximum step size for which blow-up of the numerical solution can still be avoided.

HIDDEN HINT 1 for (12-7.b) \rightarrow 12-7-2-0:Stab1h2.pdf

SOLUTION for (12-7.b)
$$\rightarrow$$
 12-7-2-1:Stab1s.pdf

(12-7.c) lacksquare (15 min.) Calculate the stability function $S=S(z), z\in\mathbb{C}$, of the method given by the Butcher scheme Eq. (12.7.1).

Solution for (12-7.c)
$$\rightarrow$$
 12-7-3-0:Stab2s.pdf

End Problem 12-7, 65 min.

Problem 12-8: Implicit Runge-Kutta Method for Gradient-Flow ODEs

In this problem we apply a special class of *implicit* Runge-Kutta single-step methods to a special class of ODEs, many of which give rise to *stiff* initial-value problems.

You need to master [Lecture \rightarrow Section 12.3.3] and [Lecture \rightarrow Section 12.3.4] and should also be familiar with the notion of a "stiff" IVP, see [Lecture \rightarrow Section 12.2].

We consider the autonomous gradient-flow ODE

$$\dot{\mathbf{y}}(t) = \mathbf{f}(\mathbf{y}) := -\operatorname{grad} V(\mathbf{y}(t)), \qquad (12.8.1)$$

where $V:D\subset\mathbb{R}^N\to\mathbb{R}$ is a continuously differentiable function.

Show that for a solution $t \mapsto \mathbf{y}(t)$ of (12.8.1) defined on $I \subset \mathbb{R}$ the function $t \in I \mapsto V(\mathbf{y}(t))$ is non-increasing.

SOLUTION for (12-8.a)
$$\rightarrow$$
 12-8-1-0:s1.pdf

(12-8.b) (10 min.)

The scalar (N=1) linear ODE $\dot{y}=-\lambda y,\ \lambda\in\mathbb{R}$, belongs to the class of gradient-flow ODEs of the form (12.8.1). What is V in this case?

Solution for (12-8.b)
$$\rightarrow$$
 12-8-2-0:s2.pdf

For the numerical integration of gradient-flow ODEs we opt for a so-called **SDIRK** (singly diagonally implicit Runge-Kutta) method, an implicit 5-stage Runge-Kutta method described by the Butcher scheme

In the file gradientflow.hpp you find a function

Eigen::MatrixXd ButcherMatrix();

that gives you the 6×5 -matrix $\begin{bmatrix} \mathfrak{A} \\ \mathbf{b}^{\top} \end{bmatrix}$ of the Butcher scheme (12.8.2). Its bottom row is the vector \mathbf{b}^{\top} .

(12-8.c) (20 min.)

State the non-linear systems of equations of size $N \times N$ that have to be solved when applying the implicit Runge-Kutta single-step method (with timestep h > 0) described by the Butcher scheme (12.8.2) to a generic autonomous ODE $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$ with $\mathbf{f} : D \subset \mathbb{R}^N \mapsto \mathbb{R}^N$, $N \in \mathbb{N}$.

Solution for (12-8.c)
$$\rightarrow$$
 12-8-3-0:s3.pdf

Give the **stage form** [Lecture \rightarrow Eq. (12.3.3.8)] of the increment equations for the implicit Runge-Kutta single-step method (timestep h > 0) described by the Butcher scheme (12.8.2) and applied to a generic autonomous ODE $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$ with $\mathbf{f} : D \subset \mathbb{R}^N \mapsto \mathbb{R}^N$, $N \in \mathbb{N}$.

```
SOLUTION for (12-8.d) \rightarrow 12-8-4-0:s4.pdf
```

In the file gradientflow.hpp implement a C++ function

```
template <typename Functor, typename Jacobian>
std::array<Eigen::VectorXd, 5> computeStages(
Func &&f, Jac &&df,
const Eigen::VectorXd &y, double h,
double rtol = 1E-6, double atol = 1E-8);
```

that uses a standard (undamped) *Newton iteration* (correction-based termination with relative tolerance rtol and absolute tolerance atol) to solve for the *stages* \mathbf{g}_i (\rightarrow [Lecture \rightarrow Eq. (12.3.3.8)]) required for one step of size h>0 of the implicit Runge-Kutta single-step method described by the Butcher scheme (12.8.2) and applied to a generic autonomous ODE $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y})$ with $\mathbf{f}: D \subset \mathbb{R}^N \mapsto \mathbb{R}^N$, $N \in \mathbb{N}$.

The right-hand side function f and its Jacobian are passed through suitable functor objects f and J respectively, whereas g supplies the initial state for the step.

```
HIDDEN HINT 1 for (12-8.e) \rightarrow 12-8-5-0:newt.pdf

SOLUTION for (12-8.e) \rightarrow 12-8-5-1:s5.pdf

(12-8.f) \bigcirc (30 min.) [depends on Sub-problem (12-8.e)]
```

Based on computeStages () in gradientflow.hpp realize a a C++ function

```
template <typename Func, typename Jac>
Eigen::VectorXd discEvolSDIRK(
Func &&f, Jac &&df, const Eigen::VectorXd &y,
double h, double rtol = 1E-6, double atol = 1E-8);
```

that evaluates the discrete evolution operator $\Psi^h y$ for the implicit Runge-Kutta single-step method described by the Butcher scheme (12.8.2) and applied to a generic autonomous ODE $\dot{y} = f(y)$ with $f: D \subset \mathbb{R}^N \mapsto \mathbb{R}^N$, $N \in \mathbb{N}$. Of course, the parameters y and h supply the arguments y and h for Ψ . The meaning of f, f, rtol and atol has been explained before

```
Solution for (12-8.f) \rightarrow 12-8-6-0:s6.pdf
```

Finally we consider the concrete gradient-flow ODE (12.8.1) spawned by the potential function

$$V(\mathbf{y}) = \sin(\|\mathbf{y}\|^2) + \lambda(\mathbf{d}^{\top}\mathbf{y})^2 \;, \quad \lambda \in \mathbb{R} \;, \quad \mathbf{y}, \mathbf{d} \in \mathbb{R}^N \;, \quad \|\mathbf{d}\| = 1 \;. \tag{12.8.8}$$

We tackle associated initial value-problems with initial state vector \mathbf{y}_0 , $\|\mathbf{y}_0\| = 1$, and want to solve them on the time interval [0, 0.1].

```
(12-8.g)  (15 min.)
```

Write down the autonomous gradient-flow ODE induced by $\mathbf{y} \mapsto V(\mathbf{y})$ from (12.8.8).

```
SOLUTION for (12-8.g) \rightarrow 12-8-7-0:s6a.pdf
```

For what values of the parameter $\lambda \in \mathbb{R}$ do we face a *stiff* ODE for states y close to zero: $\|\mathbf{y}\| \ll 1$?

HIDDEN HINT 1 for (12-8.h) \rightarrow 12-8-8-0:h7s.pdf

SOLUTION for (12-8.h)
$$\rightarrow$$
 12-8-8-1:s7.pdf

(12-8.i) (30 min.) [depends on Sub-problem (12-8.f) and Sub-problem (12-8.g)]

Based on your implementation of ${\tt discEvolSDIRK}$ () in ${\tt gradientflow.hpp}$ write a C++ function

```
std::vector<Eigen::VectorXd>
solveGradientFlow(const Eigen::VectorXd &d,
double lambda,
const Eigen::VectorXd &y0,
double T, unsigned int M);
```

that uses M equidistant steps of the SDIRK RK-SSM with the Butcher scheme (12.8.2) to solve the gradient-flow ODE (12.8.1) spawned by V from (12.8.8) up to final time T>0. The arguments ${\tt d}$ and ${\tt lambda}$ supply the parameters, and ${\tt y0}$ the initial state ${\tt y}^{(0)}$. Use the default tolerance parameters for Newton's method as specified in the signature of computeStages ().

HIDDEN HINT 1 for (12-8.i) \rightarrow 12-8-9-0:h8lf.pdf

SOLUTION for (12-8.i)
$$\rightarrow$$
 12-8-9-1:s8.pdf

(12-8.j) (10 min.)

We have solved the gradient-flow ODE (12.8.1) with V from (12.8.8) and N=2, $\mathbf{d}=\begin{bmatrix}1\\0\end{bmatrix}$, $\lambda=10$, $\mathbf{y}_0=\begin{bmatrix}1\\0\end{bmatrix}$, T=0.1.

N	error norm
10	5.99348e-07
20	3.65190e-08
40	2.25261e-09
80	1.39860e-10
160	8.71007e-12
320	5.41081e-13

The table lists the norms of the error $\mathbf{y}_N - \mathbf{y}(T)$ at final time for different numbers N of equidistant timesteps.

Which empiric order of convergence of the SDIRK single-step method do the data suggest?

SOLUTION for (12-8.j) $\rightarrow 12-8-10-0:s9.pdf$

_

End Problem 12-8, 225 min.