

1.1 Fundamentals

1.1.1 Notations

The notations in this course try to adhere to established conventions. Since these may not be universal, idiosyncrasies cannot be avoided completely. Notations in textbooks may be different, beware!

- 💡 notation for generic field of numbers: \mathbb{K}

In this course, \mathbb{K} will designate either \mathbb{R} (real numbers) or \mathbb{C} (complex numbers); complex arithmetic [14, Sect. 2.5] plays a crucial role in many applications, for instance in signal processing.

(1.1.1) Notations for vectors

- ◆ **Vectors** = are n -tuples ($n \in \mathbb{N}$) with components $\in \mathbb{K}$.

vector = one-dimensional array (of real/complex numbers)

- vectors will usually be denoted by small **bold** symbols: $\mathbf{a}, \mathbf{b}, \dots, \mathbf{x}, \mathbf{y}, \mathbf{z}$

- ◆ Default in this lecture: vectors = column vectors

$$\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{K}^n \quad | \quad [x_1 \cdots x_n] \in \mathbb{K}^{1,n}$$

column vector row vector

\mathbb{K}^n \triangleq vector space of *column vectors* with n components in \mathbb{K} .

- notation for column vectors: **bold** small roman letters, e.g. **x, y, z**

- ◆ Transposing: { column vector \mapsto row vector
row vector \mapsto column vector }

$$\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}^T = [x_1 \cdots x_n] \quad , \quad [x_1 \cdots x_n]^T = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$$

- Notation for row vectors: \mathbf{x}^T , \mathbf{y}^T , \mathbf{z}^T

- #### ◆ Addressing vector components:

-  two notations: $\mathbf{x} = [x_1 \dots x_n]^\top \rightarrow x_i, i = 1, \dots, n$
 $\mathbf{x} \in \mathbb{K}^n \rightarrow (\mathbf{x})_i, i = 1, \dots, n$

- ◆ Selecting sub-vectors:

☞ notation: $\mathbf{x} = [x_1 \dots x_n]^\top \succ (\mathbf{x})_{k:l} = (x_k, \dots, x_l)^\top, 1 \leq k \leq l \leq n$

- ◆ j -th unit vector: $\mathbf{e}_j = [0, \dots, 1, \dots, 0]^\top, (\mathbf{e}_j)_i = \delta_{ij}, i, j = 1, \dots, n.$

☞ notation: Kronecker symbol $\delta_{ij} := 1$, if $i = j$, $\delta_{ij} := 0$, if $i \neq j$.

(1.1.2) Notations and notions for matrices

- ◆ Matrices = two-dimensional arrays of real/complex numbers

$$\mathbf{A} := \begin{bmatrix} a_{11} & \dots & a_{1m} \\ \vdots & & \vdots \\ a_{n1} & \dots & a_{nm} \end{bmatrix} \in \mathbb{K}^{n,m}, \quad n, m \in \mathbb{N}.$$

vector space of $n \times m$ -matrices: ($n \hat{=} \text{number of rows}$, $m \hat{=} \text{number of columns}$)

☞ notation: bold CAPITAL roman letters, e.g., $\mathbf{A}, \mathbf{S}, \mathbf{Y}$

$$\mathbb{K}^{n,1} \leftrightarrow \text{column vectors}, \quad \mathbb{K}^{1,n} \leftrightarrow \text{row vectors}$$

- ◆ Writing a matrix as a tuple of its columns or rows

$$\mathbf{c}_i \in \mathbb{K}^n, \quad i = 1, \dots, m \quad \triangleright \quad \mathbf{A} = [\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_m] \in \mathbb{K}^{n,m},$$

$$\mathbf{r}_i \in \mathbb{K}^m, \quad i = 1, \dots, n \quad \triangleright \quad \mathbf{A} = \begin{bmatrix} \mathbf{r}_1^\top \\ \vdots \\ \mathbf{r}_n^\top \end{bmatrix} \in \mathbb{K}^{n,m}.$$

- ◆ Addressing matrix entries & sub-matrices (☞ notations):

$$\mathbf{A} := \begin{bmatrix} a_{11} & \dots & a_{1m} \\ \vdots & & \vdots \\ a_{n1} & \dots & a_{nm} \end{bmatrix} \quad \begin{aligned} &\rightarrow \text{entry } (\mathbf{A})_{i,j} = a_{ij}, \quad 1 \leq i \leq n, 1 \leq j \leq m, \\ &\rightarrow i\text{-th row, } 1 \leq i \leq n: \quad \mathbf{a}_{i,:} = (\mathbf{A})_{i,:}, \\ &\rightarrow j\text{-th column, } 1 \leq j \leq m: \quad \mathbf{a}_{:,j} = (\mathbf{A})_{:,j}, \\ &\rightarrow \text{matrix block} \quad (a_{ij})_{i=k,\dots,l}^{j=r,\dots,s} = (\mathbf{A})_{k:l,r:s}, \quad 1 \leq k \leq l \leq n, \\ &\quad \quad \quad \text{(sub-matrix)} \quad 1 \leq r \leq s \leq m. \end{aligned}$$

The colon (:) range notation is inspired by MATLAB's matrix addressing conventions, see Section 1.2.1. $(\mathbf{A})_{k:l,r:s}$ is a matrix of size $(l - k + 1) \times (s - r + 1)$.

◆ Transposed matrix:

$$\mathbf{A}^\top = \begin{bmatrix} a_{11} & \dots & a_{1m} \\ \vdots & & \vdots \\ a_{n1} & \dots & a_{nm} \end{bmatrix}^\top := \begin{bmatrix} a_{11} & \dots & a_{n1} \\ \vdots & & \vdots \\ a_{1m} & \dots & a_{mn} \end{bmatrix} \in \mathbb{K}^{m,n}.$$

◆ Adjoint matrix (Hermitian transposed):

$$\mathbf{A}^H := \begin{bmatrix} a_{11} & \dots & a_{1m} \\ \vdots & & \vdots \\ a_{n1} & \dots & a_{nm} \end{bmatrix}^H := \begin{bmatrix} \bar{a}_{11} & \dots & \bar{a}_{n1} \\ \vdots & & \vdots \\ \bar{a}_{1m} & \dots & \bar{a}_{mn} \end{bmatrix} \in \mathbb{K}^{m,n}.$$

☞ notation: $\bar{a}_{ij} = \Re(a_{ij}) - i\Im(a_{ij})$ complex conjugate of a_{ij} .

1.1.2 Classes of matrices

Most matrices occurring in mathematical modelling have a special structure. This section presents a few of these. More will come up throughout the remainder of this chapter; see also [2, Sect. 4.3].

(1.1.3) Special matrices

Terminology and notations for a few very special matrices:

Identity matrix: $\mathbf{I} = \begin{bmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & 1 \end{bmatrix} \in \mathbb{K}^{n,n},$

Zero matrix: $\mathbf{O} = \begin{bmatrix} 0 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & 0 \end{bmatrix} \in \mathbb{K}^{n,m},$

Diagonal matrix: $\mathbf{D} = \begin{bmatrix} d_1 & & 0 \\ & \ddots & \\ 0 & & d_n \end{bmatrix} \in \mathbb{K}^{n,n}, \quad d_j \in \mathbb{K}, \quad j = 1, \dots, n.$

The creation of special matrices can usually be done by special commands or functions in the various languages or libraries dedicated to numerical linear algebra, see § 1.2.5, § 1.2.13.

(1.1.4) Diagonal and triangular matrices

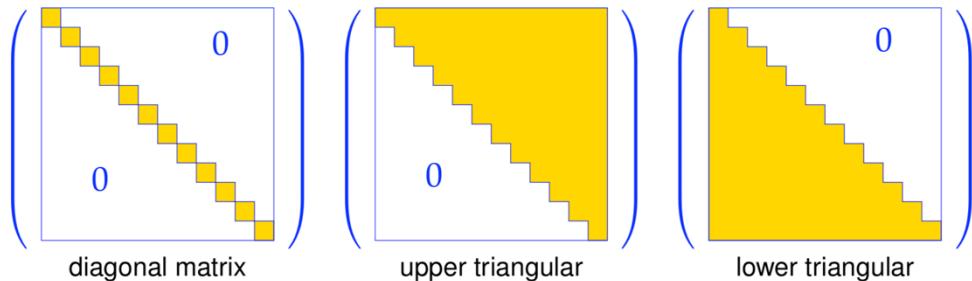
A little terminology to quickly refer to matrices whose non-zero entries occupy special locations:

Definition 1.1.5. Types of matrices

A matrix $\mathbf{A} = (a_{ij}) \in \mathbb{R}^{m,n}$ is

- **diagonal matrix**, if $a_{ij} = 0$ for $i \neq j$,
- **upper triangular matrix** if $a_{ij} = 0$ for $i > j$,
- **lower triangular matrix** if $a_{ij} = 0$ for $i < j$.

A triangular matrix is **normalized**, if $a_{ii} = 1$, $i = 1, \dots, \min\{m, n\}$.



(1.1.6) Symmetric matrices

Definition 1.1.7. Hermitian/symmetric matrices

A matrix $\mathbf{M} \in \mathbb{K}^{n,n}$, $n \in \mathbb{N}$, is **Hermitian**, if $\mathbf{M}^H = \mathbf{M}$. If $\mathbb{K} = \mathbb{R}$, the matrix is called **symmetric**.

Definition 1.1.8. Symmetric positive definite (s.p.d.) matrices → [4, Def. 3.31], [12, Def. 1.22]

$\mathbf{M} \in \mathbb{K}^{n,n}$, $n \in \mathbb{N}$, is **symmetric (Hermitian) positive definite (s.p.d.)**, if

$$\mathbf{M} = \mathbf{M}^H \quad \text{and} \quad \forall \mathbf{x} \in \mathbb{K}^n: \mathbf{x}^H \mathbf{M} \mathbf{x} > 0 \Leftrightarrow \mathbf{x} \neq 0.$$

If $\mathbf{x}^H \mathbf{M} \mathbf{x} \geq 0$ for all $\mathbf{x} \in \mathbb{K}^n$ ▷ \mathbf{M} **positive semi-definite**.

Lemma 1.1.9. Necessary conditions for s.p.d. → [4, Satz 3.33], [12, Prop. 1.18]

For a symmetric/Hermitian positive definite matrix $\mathbf{M} = \mathbf{M}^H \in \mathbb{K}^{n,n}$ holds true:

1. $m_{ii} > 0, i = 1, \dots, n,$
2. $m_{ii}m_{jj} - |m_{ij}|^2 > 0 \quad \forall 1 \leq i < j \leq n,$
3. all eigenvalues of \mathbf{M} are positive. (\leftarrow also sufficient for symmetric/Hermitian \mathbf{M})

Remark 1.1.10 (S.p.d. Hessians)

Recall from analysis: in an isolated local minimum x^* of a C^2 -function $f : \mathbb{R}^n \mapsto \mathbb{R}$ ► Hessian $D^2 f(x^*)$ s.p.d. (see Def. 8.4.11 for the definition of the Hessian)

To compute the minimum of a C^2 -function iteratively by means of Newton's method (→ Sect. 8.4) a linear system of equations with the s.p.d. Hessian as system matrix has to be solved in each step.

The solutions of many equations in science and engineering boils down to finding the minimum of some (energy, entropy, etc.) function, which accounts for the prominent role of s.p.d. linear systems in applications.

?! Review question(s) 1.1.11. (Special matrices)

- We consider two matrices $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n,m}$, both with at most $N \in \mathbb{N}$ non-zero entries. What is the maximal number of non-zero entries of $\mathbf{A} + \mathbf{B}$?
- A matrix $\mathbf{A} \in \mathbb{R}^{n,m}$ enjoys the following property (banded matrix):

$$i \in \{1, \dots, n\}, j \in \{1, \dots, m\}, i - j \notin \{-B_-, \dots, B_+\} \Rightarrow (\mathbf{A})_{ij} = 0,$$

for given $B_-, B_+ \in \mathbb{N}_0$. What is the maximal number of non-zero entries of \mathbf{A} .

1.2 Software and Libraries

Whenever algorithms involve matrices and vectors (in the sense of linear algebra) it is advisable to rely on suitable code libraries or numerical programming environments.

1.2.1 MATLAB

MATLAB (“matrix laboratory”) is a commercial (sold by MathWorks Corp.)

- full fledged high level **programming language** designed for numerical algorithms (domain specific language: **DSL**)
- integrated development environment (**IDE**) offering editor, debugger, profiler, tracing facilities,
- rather comprehensive collection of **numerical libraries**.

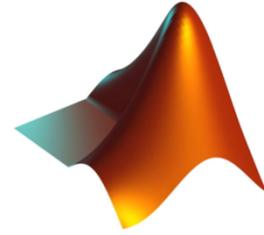


Fig. 25

Many textbooks, for instance [2] and [8] rely on MATLAB to demonstrate the actual implementation of numerical algorithms. So did earlier versions of this course. The current version has dumped MATLAB and, hence, this section **can be skipped** safely.

In its basic form MATLAB is an interpreted scripting language without strict type-binding. This, together with its uniform IDE across many platforms, makes it a very popular tool for *rapid prototyping and testing* in CSE.

Plenty of resources are available for MATLAB’s users, among them

- MATLAB documentation accessible through the Help menu or through this [link](#),
- MATLAB’s help facility through the commands **help <function>** or **doc <function>**,
- A concise [MATLAB primer](#), one of many available online, see also [here](#)

(1.2.1) Fetching the dimensions of a matrix

- ☞ **v = size(A)** yields a row vector **v** of length 2 with **v(1)** containing the number of rows and **v(2)** containing the number of columns of the matrix **A**.
- ☞ **numel(A)** returns the total number of entries of **A**; if **A** is a (row or column) vector, we get the length of **A**.
- [MATLAB-Einführung](#) (in German) by P. Arbenz.

“In MATLAB everything is a matrix”

(Fundamental “data type” in MATLAB = **matrix** of complex numbers)

- In MATLAB vectors are represented as $n \times 1$ -matrices (column vectors) or $1 \times n$ -matrices (row vectors).

Note: The treatment of vectors as special matrices is consistent with the basic operations from matrix calculus.

(1.2.2) Access to matrix and vector components in MATLAB

Access (rvalue & lvalue) to components of a vector and entries of a matrix in MATLAB is possible through the `()`-operator:

- ☞ `r = v(i)`: retrieve i -th entry of vector `v`. i must be an integer and smaller or equal `numel(v)`.
- ☞ `r = A(i, j)`: get matrix entry $(A)_{i,j}$ for two (valid) integer indices i and j .
- ☞ `r = A(end-1, end-2)`: get matrix entry $(A)_{n-1,m-2}$ of an $n \times m$ -matrix `A`.

! In case the matrix `A` is too small to contain an entry $(A)_{i,j}$, write access to `A(i, j)` will automatically trigger a dynamic adjustment of the matrix size to hold the accessed entry. The other new entries are filled with zeros.

Output: `M=`

```
% Caution: matrices are dynamically expanded when
% out of range entries are accessed
M = [1, 2, 3; 4, 5, 6]; M(4, 6) = 1.0; M,
```

1	2	3	0	0	0
4	5	6	0	0	0
0	0	0	0	0	0
0	0	0	0	0	1

➤ Danger of accidental exhaustion of memory!

(1.2.3) Access to submatrices in MATLAB

For any two (row or column) vectors `I, J` of positive integers `A(I, J)` selects the submatrix

$$[(A)_{i,j}]_{\substack{i \in I \\ j \in J}} \in \mathbb{K}^{\#I, \#J}.$$

`A(I, J)` can be used as both r-value and l-value; in the former case the maximal components of `I` and `J` have to be smaller or equal the corresponding matrix dimensions, lest MATLAB issue the error message

Index exceeds matrix dimensions. In the latter case, the size of the matrix is grown, if needed, see § 1.2.2.

(1.2.4) Initialization of matrices in MATLAB by concatenation

Inside square brackets `[]` the following two **matrix construction operators** can be used:

- `,`-operator $\hat{=}$ adding another matrix to the right (horizontal concatenation)
 - `;`-operator $\hat{=}$ adding another matrix at the bottom (vertical concatenation)
- (The `,`-operator binds more strongly than the `;`-operator!)

! Matrices joined by the `,`-operator must have the same number of rows.
Matrices concatenated vertically must have the same number of columns

▷ Filling a small matrix: $A = [1, 2; 3, 4; 5, 6];$ $A = [[1; 3; 5], [2; 4; 6]];$ $\rightarrow 3 \times 2$ matrix $\begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix}.$

▷ Initialization of vectors in MATLAB:

column vectors $x = [1; 2; 3];$
row vectors $y = [1, 2, 3];$

▷ Building a matrix from blocks:

```
% MATLAB script demonstrating the construction of a matrix from blocks
A = [1,2;3,4]; B = [5,6;7,8];
C = [A,B;-B,A], % use concatenation
```

Output: C=

1	2	5	6
3	4	7	8
-5	-6	1	2
-7	-8	3	4

(1.2.5) Special matrices in MATLAB

[Special matrices in MATLAB \rightarrow § 1.1.3]

- ☒ $n \times n$ identity matrix: $I = \text{eye}(n);$
- ☒ $m \times n$ zero matrix: $O = \text{zeros}(m, n);$
- ☒ $m \times n$ random matrix with entries equidistributed in $[0, 1]: R = \text{rand}(m, n);$
- ☒ $n \times n$ diagonal matrix with components of n -vector d (both row or column vectors are possible) on its diagonal: $D = \text{diag}(d);$

(1.2.6) Initialization of equispaced vectors (“loop index vectors”)

In MATLAB `v = (a:s:b)`, where `a, b, s` are real numbers, creates a *row vector* as initialised by the following code

```
if ((b >= a) && (s > 0))
    v = [a];
    while (v(end)+s <= b), v = [v,v(end)+s];
end

elseif ((b <= a) && (s < 0))
    v = [a];
    while (v(end)+s >= b), v = [v,v(end)+s];
end
else v = [];
end
```

Examples:

```
>> v = (3:-0.5:-0.3)
v = 3.0000 2.5000 2.0000 1.5000 1.0000 0.5000 0
>> v = (1:2.5:-13)
v = Empty matrix: 1-by-0
```

These vectors can be used to program loops in MATLAB

```
for i = (a:s:b)
    % Do something with the loop variable i
end
```

In general we could also pass a matrix as “loop index vector”. In this case the loop variable will run through the *columns* of the matrix

<pre>% MATLAB loop over columns of a matrix M = [1,2,3;4,5,6]; for i = M; i, end</pre>	Output: <pre>i = 1 i = 2 i = 3 4 5 6</pre>
--	---

(1.2.7) Special structural operations on matrices in MATLAB

- ◆ $\mathbf{A}' \triangleq$ Hermitian transpose of a matrix \mathbf{A} , transposing without complex conjugation done by `transpose(A)`.
- ◆ `triu(A)` and `tril(A)` return the upper and lower **triangular parts** of a matrix \mathbf{A} as r-value (copy): If $\mathbf{A} \in \mathbb{K}^{m,n}$

$$(\text{triu}(\mathbf{A}))_{i,j} = \begin{cases} (\mathbf{A})_{i,j} & , \text{if } i \leq j \\ 0 & \text{else.} \end{cases} \quad (\text{tril}(\mathbf{A}))_{i,j} = \begin{cases} (\mathbf{A})_{i,j} & , \text{if } i \geq j \\ 0 & \text{else.} \end{cases}$$

- ◆ `diag(A)` for a matrix $\mathbf{A} \in \mathbb{K}^{m,n}$, $\min\{m, n\} \geq 2$ returns the column vector $[(\mathbf{A})_{i,i}]_{i=1,\min\{m,n\}} \in \mathbb{K}^{\min\{m,n\}}$.

1.2.2 PYTHON

PYTHON is a widely used general-purpose and open source programming language. Together with the packages like NUMPY and MATPLOTLIB it delivers similar functionality like MATLAB for free. For interactive computing IPYTHON can be used. All those packages belong to the SciPY ecosystem.

PYTHON features a good documentation and several scientific distributions are available (e.g. Anaconda, Enthought) which contain the most important packages. On most Linux-distributions the SciPY ecosystem is also available in the software repository, as well as many other packages including for example the Spyder IDE delivered with Anaconda.

A good introduction tutorial to numerical PYTHON are the [SciPy-lectures](#). The full documentation of NUMPY and SciPY can be found [here](#). For former MATLAB-users there's also a [guide](#). The scripts in this lecture notes follow the official PYTHON style guide.

Note that in PYTHON we have to import the numerical packages explicitly before use. This is normally done at the beginning of the file with lines like `import numpy as np` and `from matplotlib import pyplot as plt`. Those import statements are often skipped in this lecture notes to focus on the actual computations. But you can always assume the import statements as given here, e.g. `np.ravel(A)` is a call to a NUMPY function and `plt.loglog(x, y)` is a call to a MATPLOTLIB pyplot function.

PYTHON is not used in the current version of the lecture. Nevertheless a few PYTHON codes are supplied in order to convey similarities and differences to implementations in MATLAB and C++.

(1.2.8) Matrices and Vectors in PYTHON

The basic numeric data type in PYTHON are NUMPY's n-dimensional arrays. Vectors are normally implemented as 1D arrays and no distinction is made between row and column vectors. Matrices are represented as 2D arrays.

- ☞ `v = np.array([1, 2, 3])` creates a 1D array with the three elements 1, 2 and 3.
- ☞ `A = np.array([[1, 2], [3, 4]])` creates a 2D array.
- ☞ `A.shape` gives the n-dimensional size of an array.
- ☞ `A.size` gives the total number of entries in an array.

Note: There's also a matrix class in NUMPY with different semantics but its use is officially discouraged and it might even be removed in future release.

(1.2.9) Manipulating arrays in PYTHON

There are many possibilities listed in the documentation how to [create](#), [index](#) and [manipulate](#) arrays.

An important difference to MATLAB is, that all arithmetic operations are normally performed element-wise, e.g. `A * B` is not the matrix-matrix product but element-wise multiplication (in MATLAB: `A.*A`). Also `A * v` does a [broadcasted](#) element-wise product. For the matrix product one has to use `np.dot(A, B)` or `A.dot(B)` explicitly.

1.2.3 EIGEN

Currently, the most widely used programming language for the development of new simulation software in scientific and industrial high-performance computing is C++. In this course we are going to use and discuss [EIGEN](#) as an example for a C++ library for numerical linear algebra (“embedded” domain specific language: DSL).

EIGEN is a [header-only](#) C++ template library designed to enable easy, natural and efficient numerical linear algebra: it provides data structures and a wide range of operations for matrices and vectors, see

below. EIGEN also implements many more fundamental algorithms (see the [documentation page](#) or the discussion below).

EIGEN relies on [expression templates](#) to allow the efficient evaluation of complex expressions involving matrices and vectors. Refer to the [example](#) given in the EIGEN documentation for details.

- ➡ [Link](#) to an “EIGEN Cheat Sheet” (quick reference relating to MATLAB commands)

(1.2.10) Compilation of codes using EIGEN

Compiling and linking on Mac OS X 10.10:

```
clang -D_HAS_CPP0X -std=c++11 -Wall -g \
-Wno-deprecated-register -DEIGEN3_ACTIVATED \
-I/opt/local/include -I/usr/local/include/eigen3 \
-o main.cpp.o -c main.cpp
/usr/bin/c++ -std=c++11 -Wall -g -Wno-deprecated-register \
-DEIGEN3_ACTIVATED -Wl,-search_paths_first \
-Wl,-headerpad_max_install_names main.cpp.o \
-o executable /opt/local/lib/libboost_program_options-mt.dylib
```

Of course, different compilers may be used on different platforms. In all cases, basic EIGEN functionality can be used without linking with a special library. Usually the generation of such elaborate calls of the compiler is left to a build system like CMAKE.

(1.2.11) Matrix and vector data types in EIGEN

A generic matrix data type is given by the templated class

```
Matrix<typename Scalar,  
        int RowsAtCompileTime, int ColsAtCompileTime>
```

Here **Scalar** is the underlying scalar type of the matrix entries, which must support the usual operations '+', '-', '*', '/', and '+=' , '*'= , etc. Usually the scalar type will be either **double**, **float**, or **complex<>**. The cardinal template arguments **RowsAtCompileTime** and **ColsAtCompileTime** can pass a **fixed** size of the matrix, if it is known at compile time. There is a specialization selected by the template argument **Eigen::Dynamic** supporting variable size "dynamic" matrices.

C++11-code 1.2.12: Vector type and their use in EIGEN

```
1 #include <Eigen/Dense >  
2  
3 template<typename Scalar>  
4 void eigenTypeDemo(unsigned int dim)  
5 {  
6     // General dynamic (variable size) matrices  
7     using dynMat_t =  
8         Eigen::Matrix<Scalar, Eigen::Dynamic, Eigen::Dynamic>;  
9     // Dynamic (variable size) column vectors  
10    using dynColVec_t = Eigen::Matrix<Scalar, Eigen::Dynamic, 1>;  
11  
12    // Dynamic (variable size) row vectors  
13    using dynRowVec_t = Eigen::Matrix<Scalar, 1, Eigen::Dynamic>;  
14    using index_t = typename dynMat_t::Index;  
15    using entry_t = typename dynMat_t::Scalar;  
16  
17    // Declare vectors of size 'dim'; not yet initialized  
18    dynColVec_t colvec(dim);  
19    dynRowVec_t rowvec(dim);  
20    // Initialisation through component access  
21    for(index_t i=0; i< colvec.size(); ++i) colvec(i) = (Scalar)i;  
22    for(index_t i=0; i< rowvec.size(); ++i) rowvec(i) = (Scalar)1/(i+1);  
23    colvec[0] = (Scalar)3.14; rowvec[dim-1] = (Scalar)2.718;  
24    // Form tensor product, a matrix, see Section 1.3.1  
25    dynMat_t vecprod = colvec*rowvec;  
26    const int nrows = vecprod.rows();  
27    const int ncols = vecprod.cols();  
28}
```

Note that in Line 23 we could have relied on automatic type deduction via `auto vectprod = ...`. However, often it is safer to forgo this option and specify the type directly.

The following convenience data types are provided by EIGEN, see [documentation](#):

- **MatrixXd** $\hat{=}$ generic variable size matrix with **double** precision entries
- **VectorXd, RowVectorXd** $\hat{=}$ dynamic column and row vectors
($=$ dynamic matrices with one dimension equal to 1)
- **MatrixNd** with $N = 2, 3, 4$ for small fixed size square $N \times N$ -matrices (type **double**)
- **VectorNd** with $N = 2, 3, 4$ for small column vectors with fixed length N .

The **d** in the type name may be replaced with **i** (for **int**), **f** (for **float**), and **cd** (for **complex<double>**) to select another basic scalar type.

All matrix type feature the methods `cols()`, `rows()`, and `size()` telling the number of columns, rows, and total number of entries.

Access to individual matrix entries and vector components, both as Rvalue and Lvalue, is possible through the `()`-operator taking two arguments of type `index_t`. If only one argument is supplied, the matrix is accessed as a linear array according to its memory layout. For vectors, that is, matrices where one dimension is fixed to 1, the `[]`-operator can replace `()` with one argument, see Line 21 of Code 1.2.12.

(1.2.13) Initialization of *dense* matrices in EIGEN

The entry access operator `(int i, int j)` allows the most direct setting of matrix entries; there is hardly any runtime penalty.

Of course, in EIGEN dedicated functions take care of the initialization of the special matrices introduced in `??`:

```

Eigen::MatrixXd I = Eigen::MatrixXd::Identity(n,n);
Eigen::MatrixXd O = Eigen::MatrixXd::Zero(n,m);
Eigen::MatrixXd D = d_vector.asDiagonal();

```

C++11-code 1.2.14: Initializing special matrices in EIGEN

```

1 #include <Eigen/Dense>
2 // Just allocate space for matrix, no initialisation
3 Eigen::MatrixXd A(rows,cols);
4 // Zero matrix. Similar to matlab command zeros(rows,cols);
5 Eigen::MatrixXd B = MatrixXd::Zero(rows, cols);
6 // Ones matrix. Similar to matlab command ones(rows,cols);
7 Eigen::MatrixXd C = MatrixXd::Ones(rows, cols);
8 // Matrix with all entries same as value.
9 Eigen::MatrixXd D = MatrixXd::Constant(rows, cols, value);
10 // Random matrix, entries uniformly distributed in [0,1]
11 Eigen::MatrixXd E = MatrixXd::Random(rows, cols);
12 // (Generalized) identity matrix, 1 on main diagonal
13 Eigen::MatrixXd I = MatrixXd::Identity(rows,cols);
14 std::cout << "size of A = (" << A.rows() << ',' << A.cols() << ')' <<
    std::endl;

```

A versatile way to initialize a matrix relies on a combination of the operators `<<` and `,`, which allows the construction of a matrix from blocks:

```

MatrixXd mat3(6,6);
mat3 <<
    MatrixXd::Constant(4,2,1.5), // top row, first block
    MatrixXd::Constant(4,3,3.5), // top row, second block
    MatrixXd::Constant(4,1,7.5), // top row, third block
    MatrixXd::Constant(2,4,2.5), // bottom row, left block
    MatrixXd::Constant(2,2,4.5); // bottom row, right block

```

The matrix is filled top to bottom left to right, block dimensions have to match (like in MATLAB).

(1.2.15) Access to submatrices in EIGEN (\rightarrow [documentation](#))

The method `block(int i, int j, int p, int q)` returns a reference to the submatrix with upper left corner at position (i, j) and size $p \times q$.

The methods `row(int i)` and `col(int j)` provide a reference to the corresponding row and column of the matrix. Even more specialised access methods are

```
topLeftCorner(p, q), bottomLeftCorner(p, q),  
topRightCorner(p, q), bottomRightCorner(p, q),  
topRows(q), bottomRows(q),  
leftCols(p), and rightCols(q),
```

with obvious purposes.

C++11 code 1.2.16: Demonstration code for access to matrix blocks in EIGEN → GITLAB

```
2 template<typename MatType>
3 void blockAccess(Eigen::MatrixBase<MatType> &M)
4 {
5     using index_t = typename Eigen::MatrixBase<MatType>::Index;
6     using entry_t = typename Eigen::MatrixBase<MatType>::Scalar;
7     const index_t nrows(M.rows()); // No. of rows
8     const index_t ncols(M.cols()); // No. of columns
9
10    cout << "Matrix M = " << endl << M << endl; // Print matrix
11    // Block size half the size of the matrix
12    index_t p = nrows/2, q = ncols/2;
13    // Output submatrix with left upper entry at position (i,i)
14    for(index_t i=0; i < min(p,q); i++)
15        cout << "Block (" << i << ',' << i << ',' << p << ',' << q
16        << ") = " << M.block(i,i,p,q) << endl;
17    // l-value access: modify sub-matrix by adding a constant
18    M.block(1,1,p,q) += Eigen::MatrixBase<MatType>::Constant(p,q,1.0);
19    cout << "M = " << endl << M << endl;
20    // r-value access: extract sub-matrix
21    MatrixXd B = M.block(1,1,p,q);
22    cout << "Isolated modified block = " << endl << B << endl;
23    // Special sub-matrices
24    cout << p << " top rows of m = " << M.topRows(p) << endl;
25    cout << p << " bottom rows of m = " << M.bottomRows(p) << endl;
26    cout << q << " left cols of m = " << M.leftCols(q) << endl;
27    cout << q << " right cols of m = " << M.rightCols(p) << endl;
28    // r-value access to upper triangular part
29    const MatrixXd T = M.template triangularView<Upper>(); //
30    cout << "Upper triangular part = " << endl << T << endl;
31    // l-value access to upper triangular part
32    M.template triangularView<Lower>() *= -1.5; //
33    cout << "Matrix M = " << endl << M << endl;
34 }
```

EIGEN offers [views](#) for access to triangular parts of a matrix, see Line 29 and Line 32, according to

```
M.triangularView<XX>()
```

where XX can stand for one of the following: [Upper](#), [Lower](#), [StrictlyUpper](#), [StrictlyLower](#), [UnitUpper](#), [UnitLower](#), see [documentation](#).

For column and row vectors references to sub-vectors can be obtained by the methods [head](#)([int](#) [length](#)), [tail](#)([int](#) [length](#)), and [segment](#)([int](#) [pos](#), [int](#) [length](#)).

Note: Unless the preprocessor switch [NDEBUG](#) is set, EIGEN performs range checks on all indices.

Since operators like MATLAB's `.*` are not available, EIGEN uses the [Array concept](#) to furnish entry-wise operations on matrices. An EIGEN-Array contains the same data as a matrix, supports the same methods for initialisation and access, but replaces the operators of matrix arithmetic with entry-wise actions. Matrices and arrays can be converted into each other by the [array\(\)](#) and [matrix\(\)](#) methods, see [documentation](#) for details.

C++11 code 1.2.18: Using Array in EIGEN → GITLAB

```
2 void matArray(int nrows, int ncols){  
3     Eigen::MatrixXd m1(nrows, ncols),m2(nrows, ncols);  
4     for(int i = 0; i < m1.rows(); i++)  
5         for(int j=0; j < m1.cols(); j++) {  
6             m1(i,j) = (double)(i+1)/(j+1);  
7             m2(i,j) = (double)(j+1)/(i+1);  
8         }  
9         // Entry-wise product, not a matrix product  
10        Eigen::MatrixXd m3 = (m1.array() * m2.array()).matrix();  
11        // Explicit entry-wise operations on matrices are possible  
12        Eigen::MatrixXd m4(m1.cwiseProduct(m2));  
13        // Entry-wise logarithm  
14        cout << "Log(m1) = " << endl << log(m1.array()) << endl;  
15        // Entry-wise boolean expression, true cases counted  
16        cout << (m1.array() > 3).count() << " entries of m1 > 3" << endl;  
17 }
```

The application of a [functor](#) (→ Section 0.2.3) to all entries of a matrix can also been done via the [unaryExpr\(\)](#) method of a matrix:

```
// Apply a lambda function to all entries of a matrix  
auto fnct = [](double x) { return (x+1.0/x); };  
cout << "f(m1) = " << endl << m1.unaryExpr(fnct) << endl;
```

Remark 1.2.19 (EIGEN in use)

- ☞ EIGEN is used as one of the base libraries for the [Robot Operating System](#) (ROS), an open source project with strong ETH participation.
 - ☞ The geometry processing library [libigl](#) uses EIGEN as its basic linear algebra engine. At ETH it is being used and developed at ETH Zurich, at the [Interactive Geometry Lab](#) and [Advanced Technologies Lab](#).
-

1.2.4 (Dense) Matrix storage formats

All numerical libraries store the entries of a (generic = *dense*) matrix $\mathbf{A} \in \mathbb{K}^{m,n}$ in a *linear* array of length mn (or longer). Accessing entries entails suitable index computations.

Two natural options for “vectorisation” of a matrix: *row major*, *column major*

$$\mathbf{A} = \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{pmatrix}$$

Row major (C-arrays, bitmaps, Python):

A_arr	1	2	3	4	5	6	7	8	9
-------	---	---	---	---	---	---	---	---	---

Column major (Fortran, MATLAB, EIGEN):

A_arr	1	4	7	2	5	8	3	6	9
-------	---	---	---	---	---	---	---	---	---

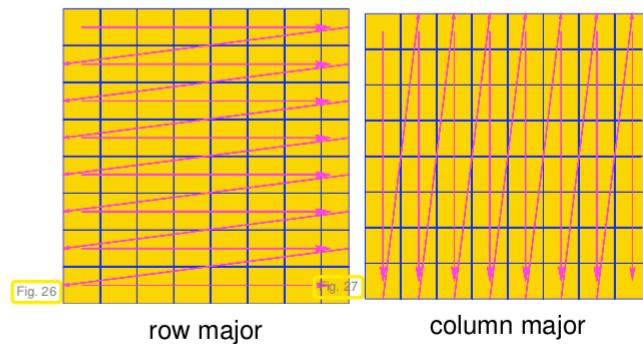
Access to entry $(\mathbf{A})_{ij}$ of $\mathbf{A} \in \mathbb{K}^{n,m}$,
 $i = 1, \dots, n$, $j = 1, \dots, m$:

row major:

$$(\mathbf{A})_{ij} \leftrightarrow A_arr(m*(i-1)+(j-1))$$

column major:

$$(\mathbf{A})_{ij} \leftrightarrow A_arr(n*(j-1)+(i-1))$$



Example 1.2.20 (Accessing matrix data as a vector)

Both in MATLAB and EIGEN the single index access operator relies on the linear data layout: In MATLAB

```
1 A = [1 2 3; 4 5 6; 7 8 9]; A(:)',
```

produces the terminal output

```
1 4 7 2 5 8 3 6 9
```

which clearly reveals the *column major* storage format.

In PYTHON the default data layout is row major, but it can be explicitly set. Further, array transposition does not change any data, but only the memory order and array shape.

In PYTHON the default data layout is row major, but it can be explicitly set. Further, array transposition does not change any data, but only the memory order and array shape.

PYTHON-code 1.2.21: Storage order in PYTHON

```
1 # array creation
2 A = np.array([[1, 2], [3, 4]]) # default (row major) storage
3 B = np.array([[1, 2], [3, 4]], order='F') # column major storage
4
5 # show internal storage
6 np.ravel(A, 'K') # array elements as stored in memory: [1, 2, 3, 4]
7 np.ravel(B, 'K') # array elements as stored in memory: [1, 3, 2, 4]
8
9 # nothing happens to the data on transpose, just the storage order
10 # changes
11 np.ravel(A.T, 'K') # array elements as stored in memory: [1, 2, 3,
12 # 4]
13 np.ravel(B.T, 'K') # array elements as stored in memory: [1, 3, 2,
14 # 4]
15
16 # storage order can be accessed by checking the array's flags
17 A.flags['C_CONTIGUOUS'] # True
18 B.flags['F_CONTIGUOUS'] # True
19 A.T.flags['F_CONTIGUOUS'] # True
20
21 B.T.flags['C_CONTIGUOUS'] # True
```

In EIGEN the data layout can be controlled by a template argument; default is column major.

C++11 code 1.2.22: Single index access of matrix entries in EIGEN → GITLAB

```

2 void storageOrder(int nrows=6,int ncols=7)
3 {
4     cout << "Different matrix storage layouts in Eigen" << endl;
5     // Template parameter ColMajor selects column major data layout
6     Matrix<double, Dynamic, Dynamic, ColMajor> mcm(nrows, ncols);
7     // Template parameter RowMajor selects row major data layout
8     Matrix<double, Dynamic, Dynamic, RowMajor> mrm(nrows, ncols);
9     // Direct initialization; lazy option: use int as index type
10    for (int l=1,i= 0; i< nrows; i++)
11        for (int j= 0; j< ncols; j++,l++)
12            mcm(i,j) = mrm(i,j) = l;
13
14    cout << "Matrix mrm = " << endl << mrm << endl;
15    cout << "mcm linear = ";
16    for (int l=0;l < mcm.size(); l++) cout << mcm(l) << ',';
17    cout << endl;
18
19    cout << "mrm linear = ";
20    for (int l=0;l < mrm.size(); l++) cout << mrm(l) << ',';
21    cout << endl;
22 }
```

The function call `storageOrder(3, 3)`, cf. Code 1.2.22 yields the output

```

1 Different matrix storage layouts in Eigen
2 Matrix mrm =
3 1 2 3
4 4 5 6
5 7 8 9
6 mcm linear = 1,4,7,2,5,8,3,6,9,
7 mrm linear = 1,2,3,4,5,6,7,8,9,
```

Remark 1.2.23 (Vectorisation of a matrix)

Mapping a column-major matrix to a column vector with the same number of entries is called **vectorization** or linearization in numerical linear algebra, in symbols

$$\text{vec} : \mathbb{K}^{n,m} \rightarrow \mathbb{K}^{n \cdot m} , \quad \text{vec}(\mathbf{A}) = \begin{bmatrix} (\mathbf{A})_{:,1} \\ (\mathbf{A})_{:,2} \\ \vdots \\ (\mathbf{A})_{:,m} \end{bmatrix}. \quad (1.2.24)$$

Remark 1.2.25 (MATLAB command `reshape`)

matlab offers the built-in command `reshape` for changing the dimensions of a matrix $\mathbf{A} \in \mathbb{K}^{m,n}$:

```
B = reshape(k, l, A); % error, in case kl ≠ mn
```

This command will create an $k \times l$ -matrix by just reinterpreting the linear array of entries of \mathbf{A} as data for a matrix with k rows and l columns. Regardless of the size and entries of the matrices the following test will always produce a `equal = true` result

```
if ((prod(size(A)) ~= (k*l)), error('Size mismatch'); end  
B = reshape(A, k, l);  
equal = (B(:) == A(:));
```

Remark 1.2.26 (NUMPY function `reshape`)

NUMPY offers the function `np.reshape` for changing the dimensions of a matrix $\mathbf{A} \in \mathbb{K}^{m,n}$:

```
# read elements of A in row major order (default)  
B = np.reshape(A, (k, l)) # error, in case kl ≠ mn  
B = np.reshape(A, (k, l), order='C') # same as above  
# read elements of A in column major order  
B = np.reshape(A, (k, l), order='F')  
# read elements of A as stored in memory  
B = np.reshape(A, (k, l), order='A')
```

This command will create an $k \times l$ -array by reinterpreting the array of entries of \mathbf{A} as data for an array with k rows and l columns. The order in which the elements of \mathbf{A} are read can be set by the `order` argument to row major (default, '`C`'), column major (`'F'`) or \mathbf{A} 's internal storage order, i.e. row major if \mathbf{A} is row major or column major if \mathbf{A} is column major (`'A'`).

Remark 1.2.27 (Reshaping matrices in EIGEN)

If you need a reshaped view of a matrix' data in EIGEN you can obtain it via the raw data vector belonging to the matrix. Then use this information to create a matrix view by means of [Map → documentation](#).

C++11 code 1.2.28: Demonstration on how reshape a matrix in EIGEN → [GITLAB](#)

```
2 template<typename MatType>  
3 void reshapeltest(MatType &M)  
4 {  
5     using index_t = typename MatType::Index;  
6     using entry_t = typename MatType::Scalar;  
7     const index_t nsize(M.size());
```

```

8
9 // reshaping possible only for matrices with non-prime dimensions
10 if ((nsize %2) == 0) {
11     entry_t *Mdat = M.data(); // raw data array for M
12     // Reinterpretation of data of M
13     Map<Eigen::Matrix<entry_t ,Dynamic ,Dynamic>> R(Mdat,2 ,nsize/2 );
14     // (Deep) copy data of M into matrix of different size
15     Eigen::Matrix<entry_t ,Dynamic ,Dynamic> S =
16         Map<Eigen::Matrix<entry_t ,Dynamic ,Dynamic>>(Mdat,2 ,nsize/2 );
17
18     cout << " Matrix M = " << endl << M << endl;
19     cout << "reshaped to " << R.rows() << 'x' << R.cols()
20         << " = " << endl << R << endl;
21     // Modifying R affects M, because they share the data space !
22     R *= -1.5;
23     cout << "Scaled (!) matrix M = " << endl << M << endl;
24     // Matrix S is not affected, because of deep copy
25     cout << "Matrix S = " << endl << S << endl;
26 }
27 }
```

This function has to be called with a mutable (l-value) matrix type object. A sample output is printed next:

```

1 Matrix M =
2 0 -1 -2 -3 -4 -5 -6
3 1 0 -1 -2 -3 -4 -5
4 2 1 0 -1 -2 -3 -4
5 3 2 1 0 -1 -2 -3
6 4 3 2 1 0 -1 -2
7 5 4 3 2 1 0 -1
8 reshaped to 2x21 =
9 0 2 4 -1 1 3 -2 0 2 -3 -1 1 -4 -2 0 -5 -3 -1 -6 -4 -2
10 1 3 5 0 2 4 -1 1 3 -2 0 2 -3 -1 1 -4 -2 0 -5 -3 -1
11 Scaled (!) matrix M =
12 -0 1.5 3 4.5 6 7.5 9
13 -1.5 -0 1.5 3 4.5 6 7.5
14 -3 -1.5 -0 1.5 3 4.5 6
15 -4.5 -3 -1.5 -0 1.5 3 4.5
16 -6 -4.5 -3 -1.5 -0 1.5 3
17 -7.5 -6 -4.5 -3 -1.5 -0 1.5
18 Matrix S =
19 0 2 4 -1 1 3 -2 0 2 -3 -1 1 -4 -2 0 -5 -3 -1 -6 -4 -2
20 1 3 5 0 2 4 -1 1 3 -2 0 2 -3 -1 1 -4 -2 0 -5 -3 -1
```

C++11 code 1.2.31: Timing for row and column oriented matrix access for EIGEN ➔ GITLAB

```
2 void rowcolaccesstiming (void)
3 {
4     const int K = 3; // Number of repetitions
5     const int N_min = 5; // Smalles matrix size 32
6     const int N_max = 13; // Scan until matrix size of 8192
7     unsigned long n = (1L << N_min);
8     Eigen::MatrixXd times (N_max-N_min+1,3);
9
10    for(int l=N_min; l<= N_max; l++, n*=2) {
11        Eigen::MatrixXd A = Eigen::MatrixXd::Random(n,n);
12        double t1 = 1000.0;
13        for(int k=0;k<K;k++) {
14            auto tic = high_resolution_clock::now();
15            for(int j=0; j < n-1; j++) A.row(j+1) -= A.row(j); // row access
16            auto toc = high_resolution_clock::now();
17            double t =
18                (double)duration_cast<microseconds>(toc-tic).count() / 1E6;
19            t1 = std::min(t1,t);
20        }
21        double t2 = 1000.0;
22        for(int k=0;k<K;k++) {
23            auto tic = high_resolution_clock::now();
24            for(int j=0; j < n-1; j++) A.col(j+1) -= A.col(j); //column
25            auto toc = high_resolution_clock::now();
26            double t =
27                (double)duration_cast<microseconds>(toc-tic).count() / 1E6;
28            t2 = std::min(t2,t);
29        }
30        times(l-N_min,0) = n;      times(l-N_min,1) = t1;
31        times(l-N_min,2) = t2;
32    }
33    std::cout << times << std::endl;
34 }
```

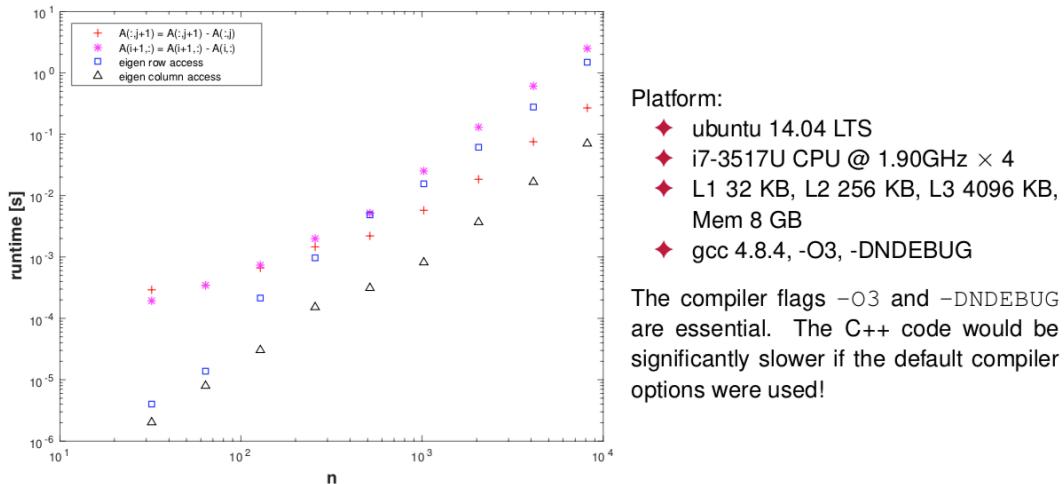
PYTHON-code 1.2.32: Timing for row and column oriented matrix access in PYTHON

```
1 import numpy as np
2 import timeit
3 from matplotlib import pyplot as plt
4
5 def col_wise(A):
6     for j in range(A.shape[1] - 1):
7         A[:, j + 1] -= A[:, j]
8
9 def row_wise(A):
10    for i in range(A.shape[0] - 1):
11        A[i + 1, :] -= A[i, :]
```

```

13 | # Timing for row/column-wise operations on matrix, we conduct k runs in |
|   order |
14 | # to reduce risk of skewed measurements due to OS activity during run. |
15 |
16 k = 3
17 res = []
18 for n in 2**np.mgrid[4:14]:
19     A = np.random.normal(size=(n, n))
20
21     t1 = min(timeit.repeat(lambda: col_wise(A), repeat=k,
22                           number=1))
22     t2 = min(timeit.repeat(lambda: row_wise(A), repeat=k,
23                           number=1))
24
24     res.append((n, t1, t2))
25
26 # plot runtime versus matrix sizes
27 ns, t1s, t2s = np.transpose(res)
28
29 plt.figure()
30 plt.plot(ns, t1s, '+', label='A[:, j + 1] -= A[:, j]')
31 plt.plot(ns, t2s, 'o', label='A[i + 1, :] -= A[i, :]')
32 plt.xlabel(r'n')
33 plt.ylabel(r'runtime [s]')
34 plt.legend(loc='upper left')
35 plt.savefig('../PYTHON_PICTURES/accessrtlin.eps')
36
37 plt.figure()
38 plt.loglog(ns, t1s, '+', label='A[:, j + 1] -= A[:, j]')
39 plt.loglog(ns, t2s, 'o', label='A[i + 1, :] -= A[i, :]')
40 plt.xlabel(r'n')
41 plt.ylabel(r'runtime [s]')
42 plt.legend(loc='upper left')
43 plt.savefig('../PYTHON_PICTURES/accessrtlog.eps')
44
45 plt.show()

```



For both MATLAB and EIGEN codes we observe a glaring discrepancy of CPU time required for accessing entries of a matrix in rowwise or columnwise fashion. This reflects the impact of features of the underlying hardware architecture, like cache size and memory bandwidth:

Interpretation of timings: Since matrices in MATLAB are stored column major all the matrix elements in a column occupy contiguous memory locations, which will all reside in the cache together. Hence, column oriented access will mainly operate on data in the cache even for large matrices. Conversely, row oriented access addresses matrix entries that are stored in distant memory locations, which incurs frequent cache misses ([cache thrashing](#)).

The impact of hardware architecture on the performance of algorithms will **not** be taken into account in this course, because hardware features tend to be both intricate and ephemeral. However, for modern high performance computing it is essential to adapt implementations to the hardware on which the code is supposed to run.

1.3 Basic linear algebra operations

First we refresh the basic rules of vector and matrix calculus. Then we will learn about a very old programming interface for simple dense linear algebra operations.

1.3.1 Elementary matrix-vector calculus

What you should know from linear algebra [10, Sect. 2.2]:

- ◆ vector space operations in matrix space $\mathbb{K}^{m,n}$ (addition, multiplication with scalars)



dot product: $\mathbf{x}, \mathbf{y} \in \mathbb{K}^n, n \in \mathbb{N}$: $\mathbf{x} \cdot \mathbf{y} := \mathbf{x}^H \mathbf{y} = \sum_{i=1}^n \bar{x}_i y_i \in \mathbb{K}$

(in EIGEN: $\mathbf{x} \cdot \text{dot}(\mathbf{y})$ or $\mathbf{x} \cdot \text{adjoint}() * \mathbf{y}$, $\mathbf{x}, \mathbf{y} \hat{=} \text{column vectors}$)



tensor product: $\mathbf{x} \in \mathbb{K}^m, \mathbf{y} \in \mathbb{K}^n, n \in \mathbb{N}$: $\mathbf{x}\mathbf{y}^H = (x_i \bar{y}_j)_{\substack{i=1, \dots, m \\ j=1, \dots, n}} \in \mathbb{K}^{m,n}$

(in EIGEN: $\mathbf{x} * \mathbf{y} \cdot \text{adjoint}()$, $\mathbf{x}, \mathbf{y} \hat{=} \text{column vectors}$)

- ◆ All are special cases of the **matrix product**:

$$\mathbf{A} \in \mathbb{K}^{m,n}, \quad \mathbf{B} \in \mathbb{K}^{n,k}: \quad \mathbf{AB} = \left[\sum_{j=1}^n a_{ij} b_{jl} \right]_{\substack{i=1, \dots, m \\ l=1, \dots, k}} \in \mathbb{R}^{m,k}. \quad (1.3.1)$$

Recall from linear algebra basic properties of the matrix product: for all \mathbb{K} -matrices $\mathbf{A}, \mathbf{B}, \mathbf{C}$ (of suitable sizes), $\alpha, \beta \in \mathbb{K}$

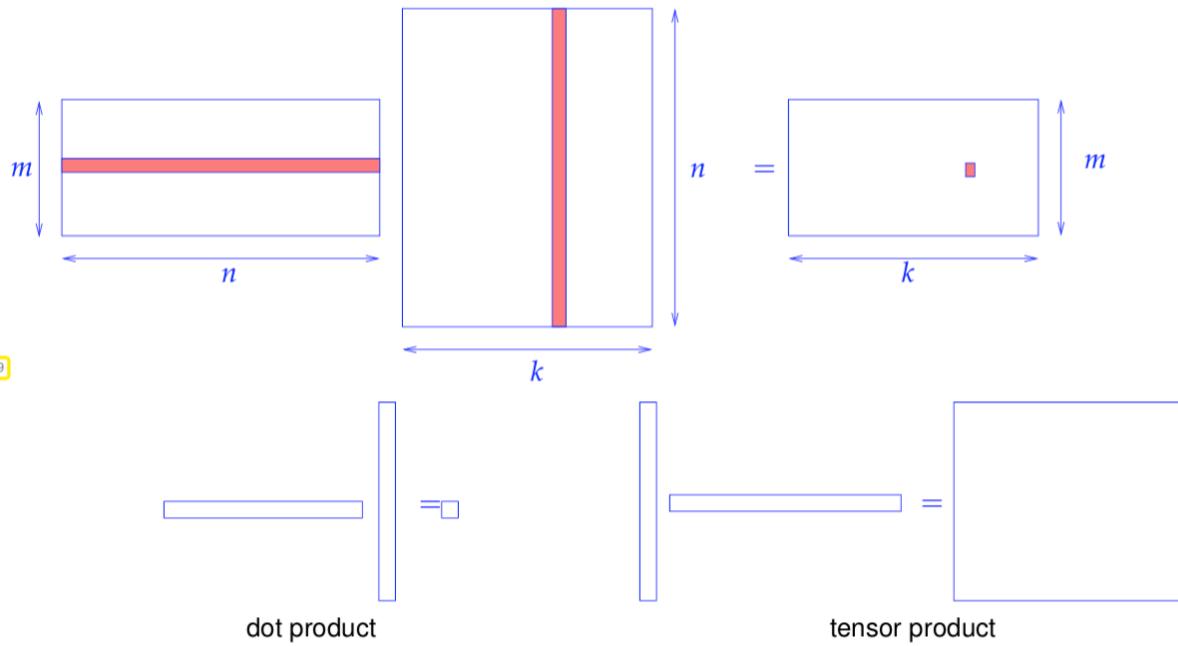
associative: $(\mathbf{AB})\mathbf{C} = \mathbf{A}(\mathbf{BC})$,

bi-linear: $(\alpha\mathbf{A} + \beta\mathbf{B})\mathbf{C} = \alpha(\mathbf{AC}) + \beta(\mathbf{BC})$, $\mathbf{C}(\alpha\mathbf{A} + \beta\mathbf{B}) = \alpha(\mathbf{CA}) + \beta(\mathbf{CB})$,

non-commutative: $\mathbf{AB} \neq \mathbf{BA}$ in general .

(1.3.2) Visualisation of (special) matrix products

Dependency of an entry of a product matrix:



Remark 1.3.3 (Row-wise & column-wise view of matrix product)

To understand what is going on when forming a matrix product, it is often useful to decompose it into matrix \times vector operations in one of the following two ways:

$\mathbf{A} \in \mathbb{K}^{m,n}$, $\mathbf{B} \in \mathbb{K}^{n,k}$:

$$\mathbf{AB} = \begin{bmatrix} \mathbf{A}(\mathbf{B})_{:,1} & \dots & \mathbf{A}(\mathbf{B})_{:,k} \end{bmatrix}, \quad \mathbf{AB} = \begin{bmatrix} (\mathbf{A})_{1,:}\mathbf{B} \\ \vdots \\ (\mathbf{A})_{m,:}\mathbf{B} \end{bmatrix}. \quad (1.3.4)$$

↓ ↓
matrix assembled from columns matrix assembled from rows

For notations refer to Sect. 1.1.1.

Remark 1.3.5 (Understanding the structure of product matrices)

A “mental image” of matrix multiplication is useful for telling special properties of product matrices.

For instance, zero blocks of the product matrix can be predicted easily in the following situations using the idea explained in Rem. 1.3.3 (try to understand how):

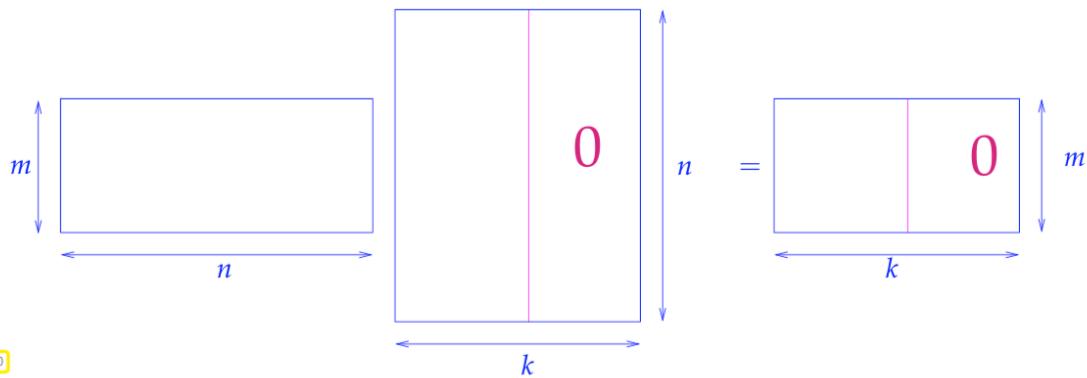


Fig. 30

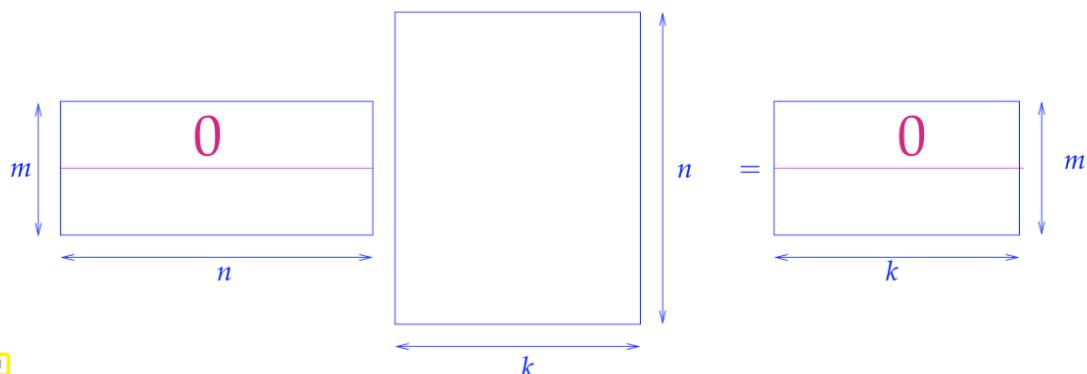


Fig. 31

A clear understanding of matrix multiplication enables you to “see”, which parts of a matrix factor matter in a product:

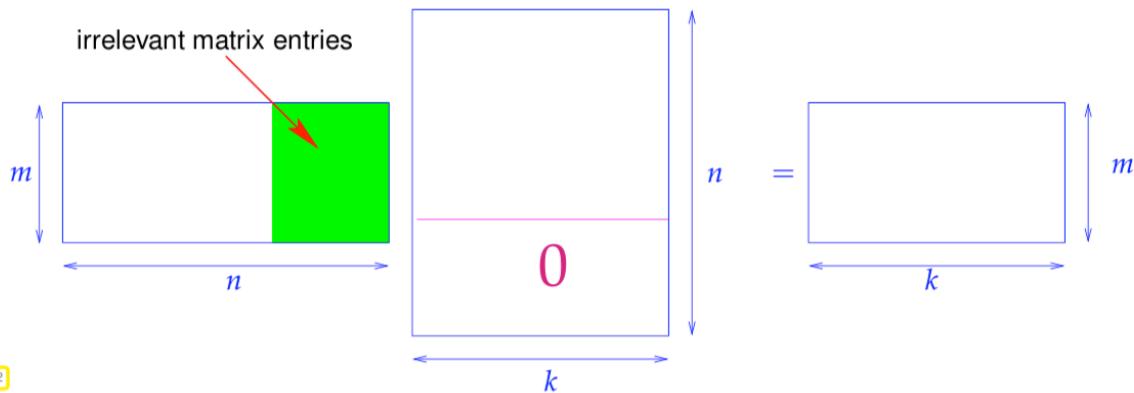
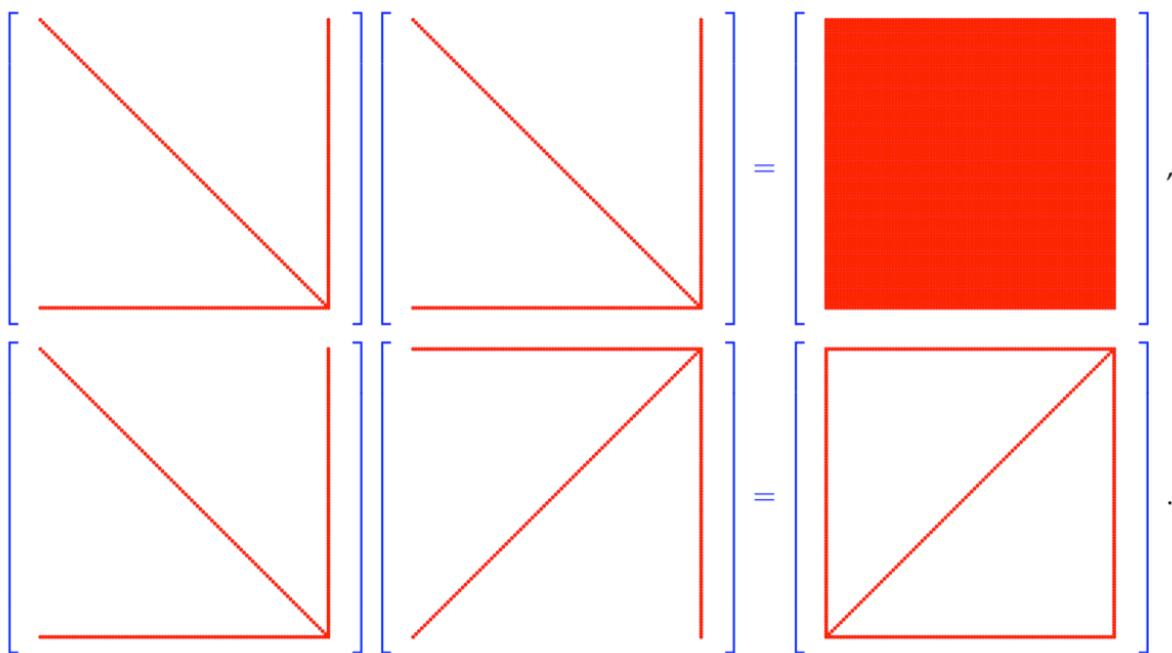


Fig. 32

"Seeing" the structure/pattern of a matrix product:



These nice renderings of the so-called patterns of matrices, that is, the distribution of their non-zero entries have been created by a special EIGEN/**Figure**-command for visualizing the structure of a matrix:
`fig.spy(M)`

C++11 code 1.3.6: Visualizing the structure of matrices in EIGEN → GITLAB

```

2 #include <Eigen/Dense>
3
4 #include <figure/figure.hpp>
5
6 using namespace Eigen;
7
8 int main () {
9     int n = 100;
10    MatrixXd A(n,n), B(n,n); A.setZero(); B.setZero();
11    A.diagonal() = VectorXd::LinSpaced(n,1,n);
12    A.col(n-1) = VectorXd::LinSpaced(n,1,n);

```

```

13 A.row(n-1) = RowVectorXd::LinSpaced(n,1,n);
14 B = A.colwise().reverse();
15 MatrixXd C = A*A, D = A*B;
16 mgl::Figure fig1, fig2, fig3, fig4;
17 fig1.spy(A); fig1.save("Aspy_cpp");
18 fig2.spy(B); fig2.save("Bspy_cpp");
19 fig3.spy(C); fig3.save("Cspy_cpp");
20 fig4.spy(D); fig4.save("Dspy_cpp");
21 return 0;
22 }
```

This code also demonstrates the use of `diagonal()`, `col()`, `row()` for L-value access to parts of a matrix.

PYTHON/MATPLOTLIB-command for visualizing the structure of a matrix: `plt.spy(M)`

PYTHON-code 1.3.7: Visualizing the structure of matrices in PYTHON

```

1 n = 100
2 A = np.diag(np.mgrid[:n])
3 A[:, -1] = A[-1, :] = np.mgrid[:n]
4 plt.spy(A)
5 plt.spy(A[::-1, :])
6 plt.spy(np.dot(A, A))
7 plt.spy(np.dot(A, B))
```

Remark 1.3.8 (Multiplying triangular matrices)

The following result is useful when dealing with matrix decompositions that often involve triangular matrices.

Lemma 1.3.9. Group of regular diagonal/triangular matrices

$$\mathbf{A}, \mathbf{B} \quad \left\{ \begin{array}{l} \text{diagonal} \\ \text{upper triangular} \\ \text{lower triangular} \end{array} \right. \Rightarrow \mathbf{AB} \text{ and } \mathbf{A}^{-1} \quad \left\{ \begin{array}{l} \text{diagonal} \\ \text{upper triangular} \\ \text{lower triangular} \end{array} \right. .$$

(assumes that \mathbf{A} is regular)

"Proof by visualization" → Rem. 1.3.5

Experiment 1.3.10 (Scaling a matrix)

Scaling = multiplication with diagonal matrices (with non-zero diagonal entries):

It is important to know the different effect of multiplying with a diagonal matrix from left or right:

- ◆ multiplication with diagonal matrix *from left* ► **row scaling**

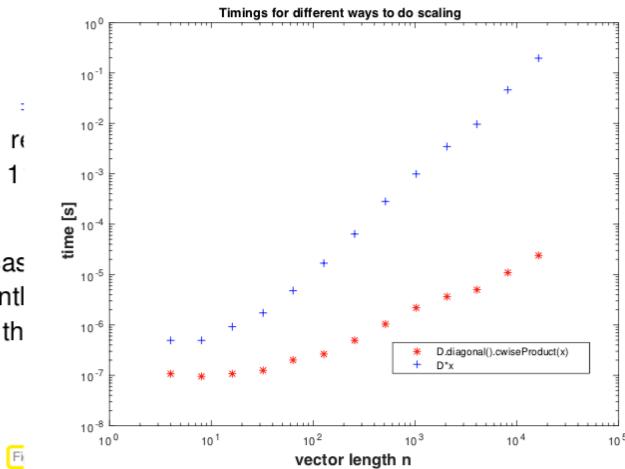
$$\begin{bmatrix} d_1 & 0 & 0 \\ 0 & d_2 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & d_n \end{bmatrix} \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1m} \\ a_{21} & a_{22} & & a_{2m} \\ \vdots & & & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nm} \end{bmatrix} = \begin{bmatrix} d_1 a_{11} & d_1 a_{12} & \dots & d_1 a_{1m} \\ d_2 a_{21} & d_2 a_{22} & \dots & d_2 a_{2m} \\ \vdots & & & \vdots \\ d_n a_{n1} & d_n a_{n2} & \dots & d_n a_{nm} \end{bmatrix} = \begin{bmatrix} d_1 (\mathbf{A})_{1,:} \\ \vdots \\ d_n (\mathbf{A})_{n,:} \end{bmatrix} .$$

- ◆ multiplication with diagonal matrix *from right* ► **column scaling**

$$\begin{bmatrix} a_{11} & a_{12} & \dots & a_{1m} \\ a_{21} & a_{22} & & a_{2m} \\ \vdots & & & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nm} \end{bmatrix} \begin{bmatrix} d_1 & 0 & 0 \\ 0 & d_2 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & d_m \end{bmatrix} = \begin{bmatrix} d_1 a_{11} & d_2 a_{12} & \dots & d_m a_{1m} \\ d_1 a_{21} & d_2 a_{22} & \dots & d_m a_{2m} \\ \vdots & & & \vdots \\ d_1 a_{n1} & d_2 a_{n2} & \dots & d_m a_{nm} \end{bmatrix} = \begin{bmatrix} d_1 (\mathbf{A})_{:,1} & \dots & d_m (\mathbf{A})_{:,m} \end{bmatrix} .$$

Multiplication with a scaling matrix $\mathbf{D} := \text{diag}(d_1, \dots, d_n) \in \mathbb{R}^{n,n}$ in EIGEN can be realised in three ways, see Code 1.3.11, ?? 9-?? 1 ?? 13, and ?? 15.

The code will be slowed down massively in case a temporary dense matrix is created inadvertently. Notice that EIGEN's expression templates avoid the pointless effort, see ?? 15.



C++11 code 1.3.11: Timing multiplication with scaling matrix in EIGEN → GITLAB

```

2   int nruns = 3, minExp = 2, maxExp = 14;
3   MatrixXd tms(maxExp-minExp+1,4);
4   for(int i = 0; i <= maxExp-minExp; ++i){
5     Timer tbad, tgood, topt; // timer class
6     int n = std::pow(2, minExp + i);
7     VectorXd d = VectorXd::Random(n,1), x = VectorXd::Random(n,1),
y(n);
```

```

8   for(int j = 0; j < nruns; ++j) {
9     MatrixXd D = d.asDiagonal(); // 
// matrix vector multiplication
10    tbad.start(); y = D*x; tbad.stop(); //
// componentwise multiplication
11    tgood.start(); y= d.cwiseProduct(x); tgood.stop(); //
// matrix multiplication optimized by Eigen
12    topt.start(); y = d.asDiagonal()*x; topt.stop(); //
13  }
14  tms(i,0)=n;
15  tms(i,1)=tgood.min(); tms(i,2)=tbad.min(); tms(i,3)= topt.min();
16}
17
18}
19}
```

PYTHON-code 1.3.12: Timing multiplication with scaling matrix in PYTHON

```
1 import numpy as np
2 from matplotlib import pyplot as plt
3 import timeit
4
5 # script for timing a smart and foolish way to carry out
6 # multiplication with a scaling matrix
7
8 nruns = 3
9 res = []
10 for n in 2**np.mgrid[2:15]:
11     d = np.random.uniform(size=n)
12     x = np.random.uniform(size=n)
13
14     tbad = min(timeit.repeat(lambda: np.dot(np.diag(d), x),
15                               repeat=nruns, number=1))
16     tgood = min(timeit.repeat(lambda: d * x, repeat=nruns,
17                               number=1))
18
19     res.append((n, tbad, tgood))
20
21 ns, tbads, tgoods = np.transpose(res)
22 plt.figure()
23 plt.loglog(ns, tbads, '+', label='using np.diag')
24 plt.loglog(ns, tgoods, 'o', label='using *')
25 plt.legend(loc='best')
26 plt.title('Timing for different ways to do scaling')
27 plt.savefig('../PYTHON_PICTURES/scaletiming.eps')
28 plt.show()
```

Hardly surprising, the component-wise multiplication of the two vectors is way faster than the intermittent initialisation of a diagonal matrix (main populated by zeros) and the computation of a matrix \times vector product. Nevertheless, such blunders keep on haunting numerical codes. Do not rely solely on EIGEN optimizations!

Remark 1.3.13 (Row and column transformations)

Simple operations on rows/columns of matrices, *cf.* what was done in Exp. 1.2.29, can often be expressed as multiplication with special matrices: For instance, given $\mathbf{A} \in \mathbb{K}^{n,m}$ we obtain \mathbf{B} by adding row $(\mathbf{A})_{j,:}$ to row $(\mathbf{A})_{j+1,:}$, $1 \leq j < n$.

Realisation through matrix product \blacktriangleright

$$\mathbf{B} = \begin{bmatrix} 1 & & & \\ & \ddots & & \\ & & 1 & \\ & & 1 & 1 \\ & & & \ddots \\ & & & & 1 \end{bmatrix} \mathbf{A}.$$

The matrix multiplying \mathbf{A} from the left is a specimen of a **transformation matrix**, a matrix that coincides with the identity matrix \mathbf{I} except for a single off-diagonal entry.

left-multiplication	with transformation matrices	\rightarrow	row transformations
right-multiplication			column transformations

row/column transformations will play a central role in Sect. 2.3

Remark 1.3.14 (Matrix algebra)

A vector space $(V, \mathbb{K}, +, \cdot)$, where V is additionally equipped with a **bi-linear** and associative “multiplication” is called an **algebra**. Hence, the vector space of square matrices $\mathbb{K}^{n,n}$ with matrix multiplication is an algebra with **unit element** \mathbf{I} .

(1.3.15) Block matrix product

Given matrix dimensions $M, N, K \in \mathbb{N}$ block sizes $1 \leq n < N$ ($n' := N - n$), $1 \leq m < M$ ($m' := M - m$), $1 \leq k < K$ ($k' := K - k$) we start from the following matrices:

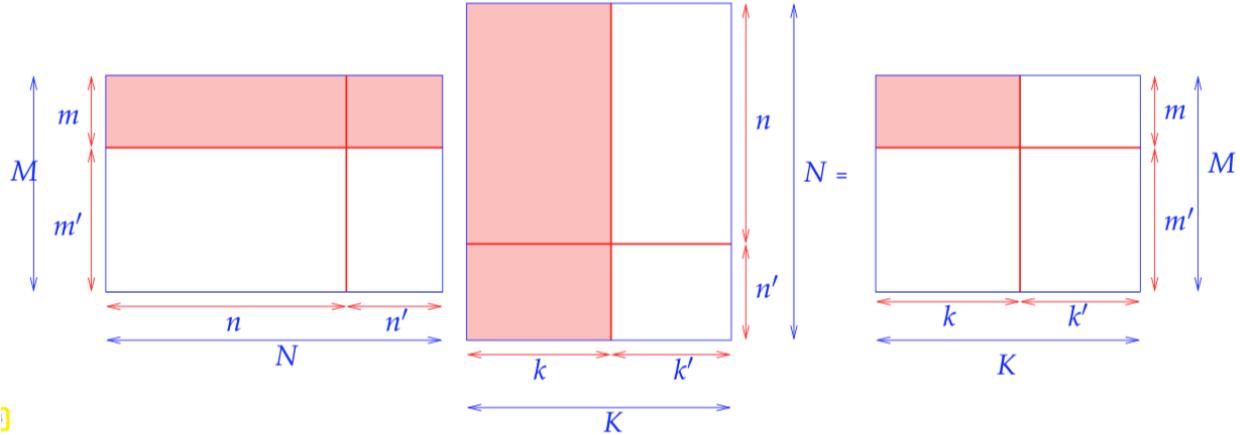
$$\begin{array}{lll} \mathbf{A}_{11} \in \mathbb{K}^{m,n} & \mathbf{A}_{12} \in \mathbb{K}^{m,n'} & \mathbf{B}_{11} \in \mathbb{K}^{n,k} & \mathbf{B}_{12} \in \mathbb{K}^{n,k'} \\ \mathbf{A}_{21} \in \mathbb{K}^{m',n} & \mathbf{A}_{22} \in \mathbb{K}^{m',n'} & \mathbf{B}_{21} \in \mathbb{K}^{n',k} & \mathbf{B}_{22} \in \mathbb{K}^{n',k'} \end{array}.$$

This matrices serve as sub-matrices or matrix blocks and are assembled into larger matrices

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} \in \mathbb{K}^{M,N}, \quad \mathbf{B} = \begin{bmatrix} \mathbf{B}_{11} & \mathbf{B}_{12} \\ \mathbf{B}_{21} & \mathbf{B}_{22} \end{bmatrix} \in \mathbb{K}^{N,K}.$$

It turns out that the matrix product \mathbf{AB} can be computed by the same formula as the product of simple 2×2 -matrices:

$$\blacktriangleright \quad \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{B}_{11} & \mathbf{B}_{12} \\ \mathbf{B}_{21} & \mathbf{B}_{22} \end{bmatrix} = \begin{bmatrix} \mathbf{A}_{11}\mathbf{B}_{11} + \mathbf{A}_{12}\mathbf{B}_{21} & \mathbf{A}_{11}\mathbf{B}_{12} + \mathbf{A}_{12}\mathbf{B}_{22} \\ \mathbf{A}_{21}\mathbf{B}_{11} + \mathbf{A}_{22}\mathbf{B}_{21} & \mathbf{A}_{21}\mathbf{B}_{12} + \mathbf{A}_{22}\mathbf{B}_{22} \end{bmatrix}. \quad (1.3.16)$$



Bottom line: one can compute with block-structured matrices in *almost* (*) the same ways as with matrices with real/complex entries, see [12, Sect. 1.3.3].



(*): you must not use the commutativity of multiplication (because matrix multiplication is not commutative).

1.3.2 BLAS – Basic Linear Algebra Subprograms

BLAS (Basic Linear Algebra Subprograms) is a specification (API) that prescribes a set of low-level routines for performing common linear algebra operations such as vector addition, scalar multiplication, dot products, linear combinations, and matrix multiplication. They are the de facto low-level routines for linear algebra libraries ([Wikipedia](#)).

The BLAS API is standardised by the [BLAS technical forum](#) and, due to its history dating back to the 70s, follows conventions of FORTRAN 77, see the [Quick Reference Guide](#) for examples. However, wrappers for other programming languages are available. CPU manufacturers and/or developers of operating systems usually supply highly optimised implementations:

- **OpenBLAS**: open source implementation with some general optimisations, available under BSD license.
- **ATLAS** (Automatically Tuned Linear Algebra Software): open source BLAS implementation with auto-tuning capabilities. Comes with C and FORTRAN interfaces and is included in Linux distributions.
- **Intel MKL** (Math Kernel Library): commercial highly optimised BLAS implementation available for all Intel CPUs. Used by most proprietary simulation software and also MATLAB.

Experiment 1.3.17 (Multiplying matrices in MATLAB)

MATLAB-code 1.3.18: Timing different implementations of matrix multiplication in MATLAB

```
1 % MATLAB script for timing different implementations of matrix
   multiplications
2 nruns = 3; times = [];
3 for n=2.^2:10 % n = size of matrices
4 fprintf('matrix size n = %d\n',n);
```

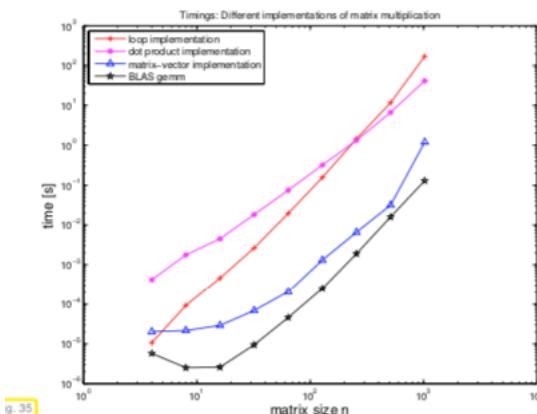
```

5 A = rand(n,n); B = rand(n,n); C = zeros(n,n);
6 t1 = realmax;
7 % loop based implementation (no BLAS)
8 for l=1:nruns
9     tic;
10    for i=1:n, for j=1:n
11        for k=1:n, C(i,j) = C(i,j) + A(i,k)*B(k,j); end
12    end, end
13    t1 = min(t1,toc);
14 end
15 t2 = realmax;
16 % dot product based implementation (BLAS level 1)
17 for l=1:nruns
18     tic;
19     for i=1:n
20         for j=1:n, C(i,j) = dot(A(i,:),B(:,j)); end
21     end
22     t2 = min(t2,toc);
23 end
24 t3 = realmax;
25 % matrix-vector based implementation (BLAS level 2)
26 for l=1:nruns
27     tic;
28     for j=1:n, C(:,j) = A*B(:,j); end
29     t3 = min(t3,toc);
30 end
31 t4 = realmax;
32 % BLAS level 3 matrix multiplication
33 for l=1:nruns
34     tic; C = A*B; t4 = min(t4,toc);
35 end
36 times = [ times; n t1 t2 t3 t4];
37 end

38
39 figure('name','mmtiming');
40 loglog(times(:,1),times(:,2),'r+-',...
41         times(:,1),times(:,3),'m*-',...
42         times(:,1),times(:,4),'b^-',...
43         times(:,1),times(:,5),'kp-');
44 title('Timings: Different implementations of matrix
45       multiplication');
46 xlabel('matrix size n','fontsize',14);
47 ylabel('time [s]','fontsize',14);
48 legend('loop implementation','dot product implementation',...
49         'matrix-vector implementation','BLAS gemm (MATLAB *)',...
50         'location','northwest');

51 print -depsc2 '../PICTURES/mvtiming.eps';

```



Platform:

- ◆ Mac OS X 10.6
- ◆ Intel Core 7, 2.66 GHz
- ◆ L2 256 kB, L3 4 MB, Mem 4 GB
- ◆ MATLAB 7.10.0 (R 2010a)

In MATLAB we can achieve a tremendous gain in execution speed by relying on compact matrix/vector operations that invoke efficient BLAS routine.

Advise: avoid loops in MATLAB and replace them with vectorised operations.

To some extent the same applies to EIGEN code, a corresponding timing script is given here:

C++11 code 1.3.19: Timing different implementations of matrix multiplication in EIGEN

→ GITLAB

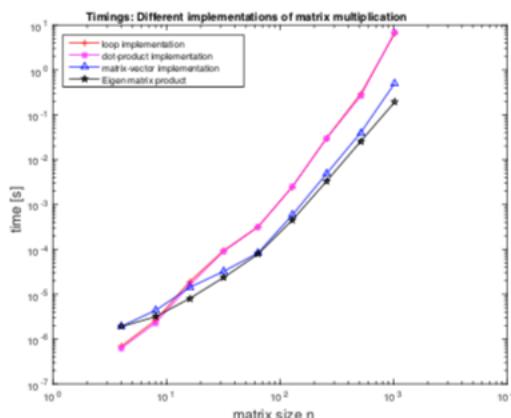
```

2  ///!  script for timing different implementations of matrix
3  ///!  multiplications
4  ///!  no BLAS is used in Eigen!
5  void mmtiming(){
6      int nruns = 3, minExp = 2, maxExp = 10;
7      MatrixXd timings(maxExp-minExp+1,5);
8      for(int p = 0; p <= maxExp-minExp; ++p){
9          Timer t1, t2, t3, t4; // timer class
10         int n = std::pow(2, minExp + p);
11         MatrixXd A = MatrixXd::Random(n,n);
12         MatrixXd B = MatrixXd::Random(n,n);
13         MatrixXd C = MatrixXd::Zero(n,n);
14         for(int q = 0; q < nruns; ++q){
15             // Loop based implementation no template magic
16             t1.start();
17             for(int i = 0; i < n; ++i)
18                 for(int j = 0; j < n; ++j)
19                     for(int k = 0; k < n; ++k)
20                         C(i,j) += A(i,k)*B(k,j);
21             t1.stop();
22             // dot product based implementation little template magic
23             t2.start();
24             for(int i = 0; i < n; ++i)
25                 for(int j = 0; j < n; ++j)
26                     C(i,j) = A.row(i).dot( B.col(j) );
27             t2.stop();
28             // matrix-vector based implementation middle template magic
29             t3.start();
30             for(int j = 0; j < n; ++j)
31                 C.col(j) = A * B.col(j);
32             t3.stop();
33             // Eigen matrix multiplication template magic optimized

```

```

33     t4.start();
34     C = A * B;
35     t4.stop();
36   }
37   timings(p,0)=n; timings(p,1)=t1.min(); timings(p,2)=t2.min();
38   timings(p,3)=t3.min(); timings(p,4)=t4.min();
39 }
40 std::cout << std::scientific << std::setprecision(3) << timings <<
41   std::endl;
42 //Plotting
43 mgl::Figure fig;
44 fig.setFontSize(4);
45 fig.setlog(true, true);
46 fig.plot(timings.col(0),timings.col(1), " +r-").label("loop
47   implementation");
48 fig.plot(timings.col(0),timings.col(2), " *m-").label("dot-product
49   implementation");
50 fig.plot(timings.col(0),timings.col(3), " ^b-").label("matrix-vector
51   implementation");
52 fig.plot(timings.col(0),timings.col(4), " ok-").label("Eigen matrix
53   product");
54 fig.xlabel("matrix size n"); fig.ylabel("time [s]");
55 fig.legend(0.05,0.95); fig.save("mmtiming");
56 }
```



Platform:

- ◆ ubuntu 14.04 LTS
- ◆ i7-3517U CPU @ 1.90GHz × 4
- ◆ L1 32 KB, L2 256 KB, L3 4096 KB, Mem 8 GB
- ◆ gcc 4.8.4, -O3

In EIGEN we can achieve some gain in execution speed by relying on compact matrix/vector operations that invoke efficient EIGEN routine. Notice that loops are not as punished as in MATLAB!

The same applies to PYTHON code, a corresponding timing script is given here:

PYTHON-code 1.3.20: Timing different implementations of matrix multiplication in PYTHON

```

1 # script for timing different implementations of matrix multiplications
2 import numpy as np
3 from matplotlib import pyplot as plt
4 import timeit
5
6 def mm_loop_based(A, B, C):
7     m, n = A.shape
```

```

8     _, p = B.shape
9     for i in range(m):
10        for j in range(p):
11            for k in range(n):
12                C[i, j] += A[i, k] * B[k, j]
13
14
15 def mm_blas1(A, B, C):
16     m, n = A.shape
17     _, p = B.shape
18     for i in range(m):
19         for j in range(p):
20             C[i, j] = np.dot(A[i, :], B[:, j])
21
22
23 def mm_blas2(A, B, C):
24     m, n = A.shape
25     _, p = B.shape
26     for i in range(m):
27         C[i, :] = np.dot(A[i, :], B)
28
29
30 def mm_blas3(A, B, C):
31     C = np.dot(A, B)
32
33
34 def main():
35     nrungs = 3
36     res = []
37     for n in 2**np.mgrid[2:11]:
38         print('matrix size n = {}'.format(n))
39         A = np.random.uniform(size=(n, n))
40         B = np.random.uniform(size=(n, n))
41         C = np.random.uniform(size=(n, n))
42
43         tloop = min(timeit.repeat(lambda: mm_loop_based(A, B, C),
44                                     repeat=nrungs, number=1))
45         tblas1 = min(timeit.repeat(lambda: mm_blas1(A, B, C),
46                                     repeat=nrungs, number=1))
47        tblas2 = min(timeit.repeat(lambda: mm_blas2(A, B, C),
48                                     repeat=nrungs, number=1))
49        tblas3 = min(timeit.repeat(lambda: mm_blas3(A, B, C),
50                                     repeat=nrungs, number=1))
51         res.append((n, tloop, tblas1, tblas2, tblas3))
52
53 ns, tloops, tblas1s, tblas2s, tblas3s = np.transpose(res)
54 plt.figure()
55 plt.loglog(ns, tloops, 'o', label='loop implementation')
56 plt.loglog(ns, tblas1s, '+', label='dot product
57             implementation')
```

```

57 plt.loglog(ns, blas2s, '*', label='matrix-vector
58     implementation')
59 plt.loglog(ns, blas3s, '^', label='BLAS gemm (np.dot)')
60 plt.legend(loc='upper left')
61 plt.savefig('../PYTHON_PICTURES/mvtiming.eps')
62 plt.show()
63
64 if __name__ == '__main__':
65     main()

```

BLAS routines are grouped into “levels” according to the amount of data and computation involved (asymptotic complexity, see Section 1.4.1 and [5, Sect. 1.1.12]):

- **Level 1:** `vector` operations such as scalar products and vector norms.
asymptotic complexity $O(n)$, (with $n \triangleq$ vector length),
e.g.: dot product: $\rho = \mathbf{x}^T \mathbf{y}$
- **Level 2:** `vector-matrix` operations such as matrix-vector multiplications.
asymptotic complexity $O(mn)$, (with $(m, n) \triangleq$ matrix size),
e.g.: matrix \times vector multiplication: $\mathbf{y} = \alpha \mathbf{A} \mathbf{x} + \beta \mathbf{y}$
- **Level 3:** `matrix-matrix` operations such as matrix additions or multiplications.
asymptotic complexity often $O(nmk)$, (with $(n, m, k) \triangleq$ matrix sizes),
e.g.: matrix product: $\mathbf{C} = \mathbf{AB}$

Syntax of BLAS calls:

The functions have been implemented for different types, and are distinguished by the first letter of the function name. E.g. `sdot` is the dot product implementation for single precision and `ddot` for double precision.

- ◆ **BLAS LEVEL 1:** vector operations, asymptotic complexity $O(n)$, $n \triangleq$ vector length

- dot product $\rho = \mathbf{x}^T \mathbf{y}$

$\text{xDOT}(N, X, INCX, Y, INCY)$

- $x \in \{S, D\}$, scalar type: $S \triangleq$ type `float`, $D \triangleq$ type `double`
- $N \triangleq$ length of vector (modulo stride `INCX`)
- $X \triangleq$ vector \mathbf{x} : array of type X
- $INCX \triangleq$ stride for traversing vector X
- $Y \triangleq$ vector \mathbf{y} : array of type X
- $INCY \triangleq$ stride for traversing vector Y

- vector operations $\mathbf{y} = \alpha \mathbf{x} + \mathbf{y}$

$\text{xAXPY}(N, ALPHA, X, INCX, Y, INCY)$

- $x \in \{S, D, C, Z\}$, $S \triangleq$ type `float`, $D \triangleq$ type `double`, $C \triangleq$ type `complex`

- $N \triangleq$ length of vector (modulo stride INCX)
 - ALPHA \triangleq scalar α
 - $X \triangleq$ vector x : array of type x
 - INCX \triangleq stride for traversing vector X
 - $Y \triangleq$ vector y : array of type x
 - INCY \triangleq stride for traversing vector Y
- ◆ **BLAS LEVEL 2:** matrix-vector operations, asymptotic complexity $O(mn)$, $(m, n) \triangleq$ matrix size
- matrix \times vector multiplication $y = \alpha Ax + \beta y$

**XGEMV (TRANS, M, N, ALPHA, A, LDA, X,
INCX, BETA, Y, INCY)**

- $x \in \{S, D, C, Z\}$, scalar type: S \triangleq type float, D \triangleq type double, C \triangleq type complex
- $M, N \triangleq$ size of matrix A
- ALPHA \triangleq scalar parameter α
- $A \triangleq$ matrix A stored in *linear array* of length $M \cdot N$ (*column major arrangement*)

$$(A)_{i,j} = A[N * (j - 1) + i].$$
- LDA \triangleq “leading dimension” of $A \in \mathbb{K}^{n,m}$, that is, the number n of rows.
- $X \triangleq$ vector x : array of type x
- INCX \triangleq stride for traversing vector X
- BETA \triangleq scalar parameter β
- $Y \triangleq$ vector y : array of type x
- INCY \triangleq stride for traversing vector Y

- **BLAS LEVEL 3:** matrix-matrix operations, asymptotic complexity $O(mnk)$, $(m, n, k) \triangleq$ matrix sizes

- matrix \times matrix multiplication $C = \alpha AB + \beta C$

**XGEMM (TRANS A, TRANS B, M, N, K,
ALPHA, A, LDA, X, B, LDB,
BETA, C, LDC)**

(☞ meaning of arguments as above)

Remark 1.3.21 (BLAS calling conventions)

The BLAS calling syntax seems queer in light of modern object oriented programming paradigms, but it is a legacy of FORTRAN77, which was (and partly still is) the programming language, in which the BLAS

routines were coded.

It is a very common situation in scientific computing that one has to rely on old codes and libraries implemented in an old-fashioned style.

Example 1.3.22 (Calling BLAS routines from C/C++)

When calling BLAS library functions from C, all arguments have to be passed by reference (as pointers), in order to comply with the argument passing mechanism of FORTRAN77, which is the model followed by BLAS.

C++-code 1.3.23: BLAS-based SAXPY operation in C++

```
1 #define daxpy_ daxpy
2 #include <iostream>
3
4 // Definition of the required BLAS function. This is usually done
5 // in a header file like blas.h that is included in the EIGEN3
6 // distribution
7 extern "C" {
8     int daxpy_(const int* n, const double* da, const double* dx,
9                 const int* incx, double* dy, const int* incy);
10 }
11
12 using namespace std;
13
14 int main(){
15     const int      n      = 5; // length of vector
16     const int      incx   = 1; // stride
17     const int      incy   = 1; // stride
18     double alpha = 2.5;      // scaling factor
19
20     // Allocated raw arrays of doubles
21     double* x  = new double [n];
22     double* y  = new double [n];
23
24     for (size_t i=0; i<n; i++){
25         x[i] = 3.1415 * i;
26         y[i] = 1.0 / (double)(i+1);
27     }
28
29     cout << "x=["; for (size_t i=0; i<n; i++) cout << x[i] << ' ';
30     cout << "]" << endl;
31     cout << "y=["; for (size_t i=0; i<n; i++) cout << y[i] << ' ';
32     cout << "]" << endl;
33
34     // Call the BLAS library function passing pointers to all arguments
35     // (Necessary when calling FORTRAN routines from C
36     daxpy_(&n, &alpha, x, &incx, y, &incy);
37
38     cout << "y = " << alpha << " * x + y = ]";
39     for (int i=0; i<n; i++) cout << y[i] << ' '; cout << "]" << endl;
40     return(0);
41 }
```

When using EIGEN in a mode that includes an external BLAS library, all this calls are wrapped into EIGEN methods.

Example 1.3.24 (Using Intel Math Kernel Library (Intel MKL) from EIGEN)

The [Intel Math Kernel Library](#) is a highly optimized math library for Intel processors and can be called directly from EIGEN ([Using Intel® Math Kernel Library from Eigen](#)) using the correct compiler flags.

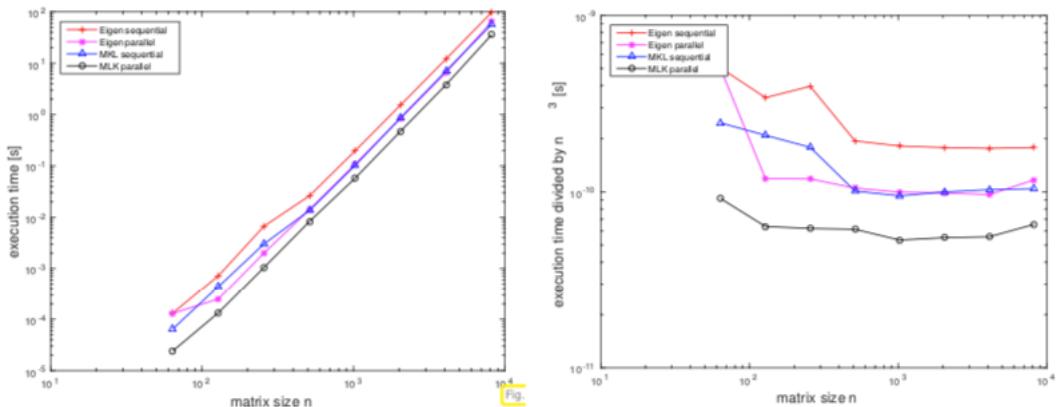
C++-code 1.3.25: Timing of matrix multiplication in EIGEN for MKL comparison → [GITLAB](#)

```

2 //! script for timing different implementations of matrix
   multiplications
3 void mmeigenmkl(){
4     int nruns = 3, minExp = 6, maxExp = 13;
5     MatrixXd timings(maxExp-minExp+1,2);
6     for(int p = 0; p <= maxExp-minExp; ++p){
7         Timer t1; // timer class
8         int n = std::pow(2, minExp + p);
9         MatrixXd A = MatrixXd::Random(n,n);
10        MatrixXd B = MatrixXd::Random(n,n);
11        MatrixXd C = MatrixXd::Zero(n,n);
12        for(int q = 0; q < nruns; ++q){
13            t1.start();
14            C = A * B;
15            t1.stop();
16        }
17        timings(p,0)=n; timings(p,1)=t1.min();
18    }
19    std::cout << std::scientific << std::setprecision(3) << timings <<
      std::endl;
20 }
```

Timing results:

<i>n</i>	EIGEN sequential [s]	EIGEN parallel [s]	MKL sequential [s]	MKL parallel [s]
64	1.318e-04	1.304e-04	6.442e-05	2.401e-05
128	7.168e-04	2.490e-04	4.386e-04	1.336e-04
256	6.641e-03	1.987e-03	3.000e-03	1.041e-03
512	2.609e-02	1.410e-02	1.356e-02	8.243e-03
1024	1.952e-01	1.069e-01	1.020e-01	5.728e-02
2048	1.531e+00	8.477e-01	8.581e-01	4.729e-01
4096	1.212e+01	6.635e+00	7.075e+00	3.827e+00
8192	9.801e+01	6.426e+01	5.731e+01	3.598e+01



Timing environment:

- ◆ ubuntu 14.04 LTS
- ◆ i7-3517U CPU @ 1.90GHz × 4
- ◆ L1 32 KB, L2 256 KB, L3 4096 KB, Mem 8 GB
- ◆ gcc 4.8.4, -O3

1.4 Computational effort

Large scale numerical computations require immense resources and execution time of numerical codes often becomes a central concern. Therefore, much emphasis has to be put on

1. designing algorithms that produce a desired result with (nearly) minimal computational effort (defined precisely below),
2. exploit possibilities for parallel and vectorised execution,
3. organising algorithms in order to make them fit memory hierarchies,
4. implementing codes that make optimal use of hardware resources and capabilities,

While Item 2–Item 4 are out of the scope of this course and will be treated in more advanced lectures, Item 1 will be a recurring theme.

The following definition encapsulates what is regarded as a measure for the “cost” of an algorithm in computational mathematics.

Definition 1.4.1. Computational effort

The **computational effort** required by a numerical code amounts to the number of **elementary operations** (additions, subtractions, multiplications, divisions, square roots) executed in a run.

(1.4.2) What computational effort does not tell us

Fifty years ago counting elementary operations provided good predictions of runtimes, but nowadays this is no longer true.

(1.4.2) What computational effort does not tell us

Fifty years ago counting elementary operations provided good predictions of runtimes, but nowadays this is no longer true.

“Computational effort $\not\propto$ runtime”



The computational effort involved in a run of a numerical code is only loosely related to overall execution time on modern computers.

This is conspicuous in Exp. 1.2.29, where algorithms incurring exactly the same computational effort took different times to execute.

The reason is that on today's computers a key bottleneck for fast execution is latency and bandwidth of memory, *cf.* the discussion at the end of Exp. 1.2.29 and [9]. Thus, concepts like **I/O-complexity** [1, 6] might be more appropriate for gauging the efficiency of a code, because they take into account the pattern of memory access.

1.4.1 (Asymptotic) complexity

The concept of computational effort from Def. 1.4.1 is still useful in a particular context:

Definition 1.4.4. (Asymptotic) complexity

The **asymptotic complexity** of an algorithm characterises the worst-case dependence of its computational effort on one or more **problem size parameter(s)** when these tend to ∞ .

- *Problem size parameters* in numerical linear algebra usually are the lengths and dimensions of the vectors and matrices that an algorithm takes as inputs.
- *Worst case* indicates that the maximum effort over a set of admissible data is taken into account.

When dealing with asymptotic complexities a mathematical formalism comes handy:

Definition 1.4.5. Landau symbol [2, p. 7]

We write $F(n) = O(G(n))$ for two functions $F, G : \mathbb{N} \rightarrow \mathbb{R}$, if there exists a constant $C > 0$ and $n_* \in \mathbb{N}$ such that

$$F(n) \leq C G(n) \quad \forall n \geq n_* .$$

More generally, $F(n_1, \dots, n_k) = O(G(n_1, \dots, n_k))$ for two functions $F, G : \mathbb{N}^k \rightarrow \mathbb{R}$ implies the existence of a constant $C > 0$ and a threshold value $n_* \in \mathbb{N}$ such that

$$F(n_1, \dots, n_k) \leq C G(n_1, \dots, n_k) \quad \forall n_1, \dots, n_k \in \mathbb{N}, \quad n_\ell \geq n_*, \quad \ell = 1, \dots, k .$$

Remark 1.4.6 (Meaningful “ O -bounds” for complexity)

Of course, the definition of the Landau symbol leaves ample freedom for stating meaningless bounds; an algorithm that runs with linear complexity $O(n)$ can be correctly labelled as possessing $O(\exp(n))$ complexity.

Yet, whenever the Landau notation is used to describe asymptotic complexities, the bounds have to be **sharp** in the sense that no function with slower asymptotic growth will be possible inside the \mathcal{O} . To make this precise we stipulate the following.

Sharpness of a complexity bound

Whenever the asymptotic complexity of an algorithm is stated as $\mathcal{O}(n^\alpha \log^\beta n \exp(\gamma n^\delta))$ with non-negative parameters $\alpha, \beta, \gamma, \delta \geq 0$ in terms of the problem size parameter n , we take for granted that choosing a smaller value for any of the parameters will no longer yield a valid (or provable) asymptotic bound.

In particular

- ◆ complexity $\mathcal{O}(n)$ means that the complexity is not $\mathcal{O}(n^\alpha)$ for any $\alpha < 1$,
- ◆ complexity $\mathcal{O}(\exp(n))$ excludes asymptotic complexity $\mathcal{O}(n^p)$ for any $p \in \mathbb{R}$.

Terminology: If the asymptotic complexity of an algorithm is $\mathcal{O}(n^p)$ with $p = 1, 2, 3$ we say that it is of “linear”, “quadratic”, and “cubic” complexity, respectively.

Remark 1.4.8 (Relevance of asymptotic complexity)

§ 1.4.2 warned us that computational effort and, thus, asymptotic complexity, of an algorithm for a concrete problem on a particular platform may not have much to do with the actual runtime (the blame goes to memory hierarchies, internal pipelining, vectorisation, etc.).

Then, why do we pay so much attention to asymptotic complexity in this course?

To a certain extent, the asymptotic complexity allows to predict the *dependence of the runtime* of a particular implementation of an algorithm *on the problem size* (for large problems).

For instance, an algorithm with asymptotic complexity $O(n^2)$ is likely to take $4\times$ as much time when the problem size is doubled.

(1.4.9) Concluding polynomial complexity from runtime measurements

Available: “Measured runtimes” $t_i = t_i(n_i)$ for different values $n_1, n_2, \dots, n_N, n_i \in \mathbb{N}$, of the problem size parameter

Conjectured: power law dependence $t_i \approx Cn_i^\alpha$ (also “algebraic dependence”), $\alpha \in \mathbb{R}$

How can we glean evidence that supports or refutes our conjecture from the data? Look at the data in **doubly logarithmic scale!**

$$t_i = Cn_i^\alpha \Rightarrow \log(t_i) \approx \log C + \alpha \log(n_i), \quad i = 1, \dots, N.$$

- If the conjecture holds true, then the points (n_i, t_i) will approximately lie on a *straight line* with **slope** α in a doubly logarithmic plot (which can be created in MATLAB by the **loglog** plotting command and in EIGEN with the **Figure-command** `fig.setlog(true, true);`).
- quick “visual test” of conjectured asymptotic complexity

More rigorous: Perform linear regression on $(\log n_i, \log t_i), i = 1, \dots, N$ (\rightarrow Chapter 3)

1.4.2 Cost of basic operations

Performing elementary BLAS-type operations through simple (nested) loops, we arrive at the following obvious complexity bounds:

operation	description	#mul/div	#add/sub	asympt. complexity
dot product	$(\mathbf{x} \in \mathbb{R}^n, \mathbf{y} \in \mathbb{R}^n) \mapsto \mathbf{x}^H \mathbf{y}$	n	$n - 1$	$O(n)$
tensor product	$(\mathbf{x} \in \mathbb{R}^m, \mathbf{y} \in \mathbb{R}^n) \mapsto \mathbf{x}\mathbf{y}^H$	nm	0	$O(mn)$
matrix product ^(*)	$(\mathbf{A} \in \mathbb{R}^{m,n}, \mathbf{B} \in \mathbb{R}^{n,k}) \mapsto \mathbf{AB}$	mnk	$mk(n - 1)$	$O(mnk)$

Remark 1.4.10 (“Fast” matrix multiplication)

(*): The $O(mnk)$ complexity bound applies to “straightforward” matrix multiplication according to (1.3.1).

For $m = n = k$ there are (sophisticated) variants with better asymptotic complexity, e.g., the **divide-and-conquer Strassen algorithm** [13] with asymptotic complexity $O(n^{\log_2 7})$:

Start from $\mathbf{A}, \mathbf{B} \in \mathbb{K}^{n,n}$ with $n = 2\ell$, $\ell \in \mathbb{N}$. The idea relies on the block matrix product (1.3.16) with $\mathbf{A}_{ij}, \mathbf{B}_{ij} \in \mathbb{K}^{\ell,\ell}$, $i, j \in \{1, 2\}$. Let $\mathbf{C} := \mathbf{AB}$ be partitioned accordingly: $\mathbf{C} = \begin{pmatrix} \mathbf{C}_{11} & \mathbf{C}_{12} \\ \mathbf{C}_{21} & \mathbf{C}_{22} \end{pmatrix}$. Then tedious elementary computations reveal

$$\begin{aligned}\mathbf{C}_{11} &= \mathbf{Q}_0 + \mathbf{Q}_3 - \mathbf{Q}_4 + \mathbf{Q}_6, \\ \mathbf{C}_{21} &= \mathbf{Q}_1 + \mathbf{Q}_3, \\ \mathbf{C}_{12} &= \mathbf{Q}_2 + \mathbf{Q}_4, \\ \mathbf{C}_{22} &= \mathbf{Q}_0 + \mathbf{Q}_2 - \mathbf{Q}_1 + \mathbf{Q}_5,\end{aligned}$$

where the $\mathbf{Q}_k \in \mathbb{K}^{\ell,\ell}$, $k = 1, \dots, 7$ are obtained from

$$\begin{aligned}\mathbf{Q}_0 &= (\mathbf{A}_{11} + \mathbf{A}_{22}) * (\mathbf{B}_{11} + \mathbf{B}_{22}), \\ \mathbf{Q}_1 &= (\mathbf{A}_{21} + \mathbf{A}_{22}) * \mathbf{B}_{11}, \\ \mathbf{Q}_2 &= \mathbf{A}_{11} * (\mathbf{B}_{12} - \mathbf{B}_{22}), \\ \mathbf{Q}_3 &= \mathbf{A}_{22} * (-\mathbf{B}_{11} + \mathbf{B}_{21}), \\ \mathbf{Q}_4 &= (\mathbf{A}_{11} + \mathbf{A}_{12}) * \mathbf{B}_{22}, \\ \mathbf{Q}_5 &= (-\mathbf{A}_{11} + \mathbf{A}_{21}) * (\mathbf{B}_{11} + \mathbf{B}_{12}), \\ \mathbf{Q}_6 &= (\mathbf{A}_{12} - \mathbf{A}_{22}) * (\mathbf{B}_{21} + \mathbf{B}_{22}).\end{aligned}$$

Beside a considerable number of matrix additions (computational effort $O(n^2)$) it takes only 7 multiplications of matrices of size $n/2$ to compute \mathbf{C} ! Strassen’s algorithm boils down to the *recursive application* of these formulas for $n = 2^k$, $k \in \mathbb{N}$.

A refined algorithm of this type can achieve complexity $O(n^{2.36})$, see [3].

1.4.3 Reducing complexity in numerical linear algebra: Some tricks

In computations involving matrices and vectors complexity of algorithms can often be reduced by performing the operations in a particular order:

Example 1.4.11 (Efficient associative matrix multiplication)

We consider the multiplication with a **rank-1-matrix**. Matrices with rank 1 can always be obtained as the tensor product of two vectors, that is, the matrix product of a column vector and a row vector. Given $\mathbf{a} \in \mathbb{K}^m$, $\mathbf{b} \in \mathbb{K}^n$, $\mathbf{x} \in \mathbb{K}^n$ we may compute the vector $\mathbf{y} = \mathbf{ab}^\top \mathbf{x}$ in two ways:

$$\mathbf{y} = (\mathbf{ab}^\top) \mathbf{x} . \quad (1.4.12) \qquad \qquad \mathbf{y} = \mathbf{a}(\mathbf{b}^\top \mathbf{x}) . \quad (1.4.13)$$

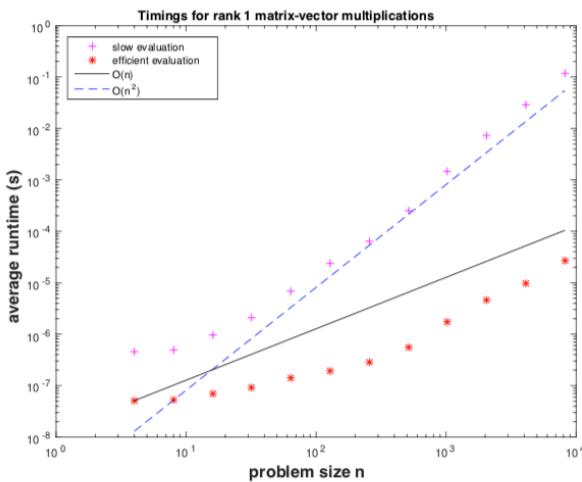
- $\mathbf{T} = (\mathbf{a} * \mathbf{b}.transpose()) * \mathbf{x};$ $\mathbf{t} = \mathbf{a} * \mathbf{b}.dot(\mathbf{x});$
- complexity $O(mn)$ ► complexity $O(n + m)$ ("linear complexity")

Visualization of evaluation according to (1.4.12):

$$\left[\begin{array}{|c|} \hline \end{array} \right] \cdot \left[\begin{array}{|c|} \hline \end{array} \right] = \left[\begin{array}{|c|} \hline \end{array} \right] \left[\begin{array}{|c|} \hline \end{array} \right]$$

Visualization of evaluation according to (1.4.13):

$$\left[\begin{array}{|c|} \hline \end{array} \right] \left[\begin{array}{|c|} \hline \end{array} \right] \cdot \left[\begin{array}{|c|} \hline \end{array} \right] = \left[\begin{array}{|c|} \hline \end{array} \right] [\square]$$



▷ average runtimes for efficient/inefficient matrix×vector multiplication with rank-1 matrices , see § 1.4.9 for the rationale behind choosing a *doubly logarithmic plot*.

Platform:

- ◆ ubuntu 14.04 LTS
- ◆ i7-3517U CPU @ 1.90GHz × 4
- ◆ L1 32 KB, L2 256 KB, L3 4096 KB,
- ◆ 8 GB main memory
- ◆ gcc 4.8.4, -O3

C++11 code 1.4.14: EIGEN code for Ex. 1.4.11 → GITLAB

```

2  /// This function compares the runtimes for the multiplication
3  /// of a vector with a rank-1 matrix  $\mathbf{ab}^T$ ,  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$ 
4  /// using different associative evaluations.
5  /// Runtime measurements consider minimal time for
6  /// several (nruns) runs
7  MatrixXd dottenstming() {
8      const int nruns = 3, minExp = 2, maxExp = 13;
9      // Matrix for storing recorded runtimes
10     MatrixXd timings(maxExp-minExp+1,3);
11     for(int i = 0; i <= maxExp-minExp; ++i){
12         Timer tfool, tsmart; // Timer objects
13         const int n = std::pow(2, minExp + i);
14         VectorXd a = VectorXd::LinSpaced(n,1,n);
15         VectorXd b = VectorXd::LinSpaced(n,1,n).reverse();
16         VectorXd x = VectorXd::Random(n,1), y(n);
17         for(int j = 0; j < nruns; ++j){
18             // Grossly wasteful evaluation
19             tfool.start(); y = (a*b.transpose())*x;      tfool.stop();
20             // Efficient implementation
21             tsmart.start(); y = a * b.dot(x); tsmart.stop();
22         }
23         timings(i,0)=n;
24         timings(i,1)=tsmart.min(); timings(i,2)=tfool.min();
25     }
26     return timings;
27 }
```

Complexity can sometimes be reduced by reusing intermediate results.

Example 1.4.15 (Hidden summation)

The asymptotic complexity of the EIGEN code

```
Eigen::MatrixXd AB = A*B.transpose();
y = AB.triangularView<Eigen::Upper>() *x;
```

when supplied with two low-rank matrices $\mathbf{A}, \mathbf{B} \in \mathbb{K}^{n,p}$, $p \ll n$, in terms of $n \rightarrow \infty$ obviously is $O(n^2)$, because an intermediate $n \times n$ -matrix \mathbf{AB}^T is built.

First, consider the case of a tensor product (= rank-1) matrix, that is, $p = 1$, $\mathbf{A} \leftrightarrow \mathbf{a} = [a_1, \dots, a_n]^\top \in \mathbb{K}^n$,

$\mathbf{B} \leftrightarrow \mathbf{b} = [b_1, \dots, b_n] \in \mathbb{K}^n$. Then

$$\begin{aligned} \mathbf{y} = \text{triu}(\mathbf{ab}^T)\mathbf{x} &= \begin{bmatrix} a_1 b_1 & a_1 b_2 & \dots & \dots & a_1 b_n \\ 0 & a_2 b_2 & a_2 b_3 & \dots & \dots & a_2 b_n \\ \vdots & \ddots & \ddots & \ddots & & \vdots \\ & & \ddots & \ddots & \ddots & \vdots \\ \vdots & & & \ddots & \ddots & \vdots \\ 0 & \dots & \dots & 0 & a_n b_n \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \\ &= \underbrace{\begin{bmatrix} a_1 & & & & \\ & \ddots & & & \\ & & \ddots & & \\ & & & \ddots & \\ & & & & a_n \end{bmatrix}}_{\mathbf{T}} \underbrace{\begin{bmatrix} 1 & 1 & \dots & \dots & 1 \\ 0 & 1 & 1 & \dots & \dots & 1 \\ \vdots & \ddots & \ddots & \ddots & & \vdots \\ & \ddots & \ddots & \ddots & \ddots & \vdots \\ \vdots & & \ddots & \ddots & \ddots & \vdots \\ 0 & \dots & \dots & 0 & 1 \end{bmatrix}}_{\mathbf{T}} \underbrace{\begin{bmatrix} b_1 & & & & \\ & \ddots & & & \\ & & \ddots & & \\ & & & \ddots & \\ & & & & b_n \end{bmatrix}}_{\mathbf{T}} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}. \end{aligned}$$

Thus, the core problem is the fast multiplication of a vector with an upper triangular matrix \mathbf{T} described in EIGEN syntax by `Eigen::MatrixXd::Ones(n, n).triangularView<Eigen::Upper>()`. Note that multiplication of a vector \mathbf{x} with \mathbf{T} yields a vector of partial sums of components of \mathbf{x} starting from last component. This can be achieved by invoking the special C++ command `std::partial_sum` in the numeric header ([documentation](#)). We also observe that

$$\mathbf{AB}^T = \sum_{\ell=1}^p (\mathbf{A})_{:, \ell} ((\mathbf{B})_{:, \ell})^\top,$$

so that the computations for the special case $p = 1$ discussed above can simply be reused p times!

C++11 code 1.4.16: Efficient multiplication with the upper diagonal part of a rank- p -matrix in EIGEN → GITLAB

```

2  //! Computation of y = triu(AB^T)x
3  //! Efficient implementation with backward cumulative sum
4  //! (partial_sum)
5  template<class Vec, class Mat>
6  void Irtrimulteff(const Mat& A, const Mat& B, const Vec& x, Vec& y){
7      const int n = A.rows(), p = A.cols();
8      assert( n == B.rows() && p == B.cols()); // size mismatch
9      for(int l = 0; l < p; ++l){
10          Vec tmp = (B.col(l).array() * x.array()).matrix().reverse();
11          std::partial_sum(tmp.data(), tmp.data() + n, tmp.data());
12          y += (A.col(l).array() * tmp.reverse().array()).matrix();
13      }
14 }
```

This code enjoys the obvious complexity of $O(pn)$ for $p, n \rightarrow \infty$, $p < n$. The code offers an example of a function templated with its argument types, see § 0.2.5. The types **Vec** and **Mat** must fit the concept of EIGEN vectors/matrices.

The next concept from linear algebra is important in the context of computing with multi-dimensional arrays.

Definition 1.4.17. 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

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

Example 1.4.18 (Multiplication of Kronecker product with vector)

The function $(\mathbf{A} \otimes \mathbf{B})\mathbf{x}$ when invoked with two matrices $\mathbf{A} \in \mathbb{K}^{m,n}$ and $\mathbf{B} \in \mathbb{K}^{l,k}$ and a vector $\mathbf{x} \in \mathbb{K}^{nk}$, will suffer an asymptotic complexity of $\mathcal{O}(m \cdot n \cdot l \cdot k)$, determined by the size of the intermediate dense matrix $\mathbf{A} \otimes \mathbf{B} \in \mathbb{K}^{ml,nk}$.

Using the partitioning of the vector \mathbf{x} into n equally long sub-vectors

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}^1 \\ \mathbf{x}^2 \\ \vdots \\ \mathbf{x}^n \end{bmatrix}, \quad \mathbf{x}^j \in \mathbb{K}^k,$$

we find the representation

$$(\mathbf{A} \otimes \mathbf{B})\mathbf{x} = \begin{bmatrix} (\mathbf{A})_{1,1}\mathbf{Bx}^1 + (\mathbf{A})_{1,2}\mathbf{Bx}^2 + \cdots + (\mathbf{A})_{1,n}\mathbf{Bx}^n \\ (\mathbf{A})_{2,1}\mathbf{Bx}^1 + (\mathbf{A})_{2,2}\mathbf{Bx}^2 + \cdots + (\mathbf{A})_{2,n}\mathbf{Bx}^n \\ \vdots \\ \vdots \\ (\mathbf{A})_{m,1}\mathbf{Bx}^1 + (\mathbf{A})_{m,2}\mathbf{Bx}^2 + \cdots + (\mathbf{A})_{m,n}\mathbf{Bx}^n \end{bmatrix}.$$

The idea is to form the products \mathbf{Bx}^j , $j = 1, \dots, n$, once, and then combine them linearly with coefficients given by the entries in the rows of \mathbf{A} :

**C++ code 1.4.19: Efficient multiplication of Kronecker product with vector in EIGEN
→ GITLAB**

```

2 //! @brief Multiplication of Kronecker product with vector y = (A ⊗ B)x.
3 //! Elegant way using reshape
4 //! WARNING: using Matrix::Map we assume the matrix is in ColMajor
5 //!           format, *beware* you may incur bugs if matrix is in RowMajor instead
6 //! @param[in] A Matrix m × n
7 //! @param[in] B Matrix l × k
8 //! @param[in] x Vector of dim nk
9 //! @param[out] y Vector y = kron(A,B) *x of dim ml
10 template <class Matrix, class Vector>

```

```

9  void kronmultv(const Matrix &A, const Matrix &B, const Vector &x,
10 Vector &y){
11     unsigned int m = A.rows(); unsigned int n = A.cols();
12     unsigned int l = B.rows(); unsigned int k = B.cols();
13     // 1st matrix mult. computes the products  $Bx^j$ 
14     // 2nd matrix mult. combines them linearly with the coefficients
15     // of A
16     Matrix t = B * Matrix::Map(x.data(), k, n) * A.transpose(); //
y = Matrix::Map(t.data(), m*l, 1);
}

```

Recall the reshaping of a matrix in EIGEN in order to understand this code: ??.

The asymptotic complexity of this code is determined by the two matrix multiplications in Line 14. This yields the asymptotic complexity $O(lkn + mnl)$ for $l, k, m, n \rightarrow \infty$.

PYTHON-code 1.4.20: Efficient multiplication of Kronecker product with vector in PYTHON

```

1 def kronmultv(A, B, x):
2     n, k = A.shape[1], B.shape[1]
3     assert x.size == n * k, 'size mismatch'
4     xx = np.reshape(x, (n, k))
5     Z = np.dot(xx, B.T)
6     yy = np.dot(A, Z)
7     return np.ravel(yy)

```

Note that different reshaping is used in the PYTHON code due to the default row major storage order.

1.5 Machine Arithmetic

1.5.1 Experiment: Loss of orthogonality

(1.5.1) Gram-Schmidt orthogonalisation

From linear algebra [10, Sect. 4.4] or Ex. 0.2.41 we recall the fundamental algorithm of **Gram-Schmidt orthogonalisation** of an ordered finite set $\{a^1, \dots, a^k\}$, $k \in \mathbb{N}$, of vectors $a^\ell \in \mathbb{K}^n$:

Input: $\{\mathbf{a}^1, \dots, \mathbf{a}^k\} \subset \mathbb{K}^n$

```

1:  $\mathbf{q}^1 := \frac{\mathbf{a}^1}{\|\mathbf{a}^1\|_2}$  % 1st output vector
2: for  $j = 2, \dots, k$  do
   { % Orthogonal projection
3:    $\mathbf{q}^j := \mathbf{a}^j$ 
4:   for  $\ell = 1, 2, \dots, j-1$  do (GS)
5:     {  $\mathbf{q}^j \leftarrow \mathbf{q}^j - \mathbf{a}^j \cdot \mathbf{q}^\ell \mathbf{q}^\ell$  }
6:   if ( $\mathbf{q}^j = 0$ ) then STOP
7:   else {  $\mathbf{q}^j \leftarrow \frac{\mathbf{q}^j}{\|\mathbf{q}^j\|_2}$  }
8:   }

```

Output: $\{\mathbf{q}^1, \dots, \mathbf{q}^j\}$

☞ Notation: $\|\cdot\|_2 \triangleq$ Euclidean norm of a vector $\in \mathbb{K}^n$

The following code implements the Gram-Schmidt orthonormalization of a set of vectors passed as the columns of a matrix $\mathbf{A} \in \mathbb{R}^{n,k}$.

In linear algebra we have learnt that, if it does not **STOP** prematurely, this algorithm will compute **orthonormal vectors** $\mathbf{q}^1, \dots, \mathbf{q}^k$ satisfying

$$\text{Span}\{\mathbf{q}^1, \dots, \mathbf{q}^\ell\} = \text{Span}\{\mathbf{a}^1, \dots, \mathbf{a}^\ell\}, \quad (1.5.2)$$

for all $\ell \in \{1, \dots, k\}$.

More precisely, if $\mathbf{a}^1, \dots, \mathbf{a}^\ell, \ell \leq k$, are linearly independent, then the Gram-Schmidt algorithm will not terminate before the $\ell+1$ -th step.

C++11 code 1.5.3: Gram-Schmidt orthogonalisation in EIGEN → GITLAB

```

2 template <class Matrix>
3 Matrix gramschmidt( const Matrix & A ) {
4     Matrix Q = A;
5     // First vector just gets normalized, Line 1 of (GS)
6     Q.col(0).normalize();
7     for(unsigned int j = 1; j < A.cols(); ++j) {
8         // Replace inner loop over each previous vector in Q with fast
9         // matrix-vector multiplication (Lines 4, 5 of (GS))
10        Q.col(j) -= Q.leftCols(j) * (Q.leftCols(j).transpose() *
11            A.col(j));//
12        // Normalize vector, if possible.
13        // Otherwise columns of A must have been linearly dependent
14        if( Q.col(j).norm() <= 10e-9 * A.col(j).norm() ) { //
15            std::cerr << "Gram-Schmidt failed: A has lin. dep
16            columns." << std::endl;
17            break;
18        } else { Q.col(j).normalize(); } // Line 7 of (GS)
19    }
20    return Q;
21 }
```

PYTHON-code 1.5.4: Gram-Schmidt orthogonalisation in PYTHON

```

1 def gramschmidt(A):
2     _, k = A.shape
3     Q = A[:, [0]] / np.linalg.norm(A[:, 0])
4     for j in range(1, k):
5         q = A[:, j] - np.dot(Q, np.dot(Q.T, A[:, j]))
6
7         nq = np.linalg.norm(q)
8         if nq < 1e-9 * np.linalg.norm(A[:, j]):
9             break
10        Q = np.column_stack([Q, q / nq])
11
12    return Q

```

Note the different loop range due to the zero-based indexing in PYTHON.

Experiment 1.5.5 (Unstable Gram-Schmidt orthonormalization)

If $\{\mathbf{a}^1, \dots, \mathbf{a}^k\}$ are linearly independent we expect the output vectors $\mathbf{q}^1, \dots, \mathbf{q}^k$ to be orthonormal:

$$(\mathbf{q}^\ell)^\top \mathbf{q}^m = \delta_{\ell,m}, \quad \ell, m \in \{1, \dots, k\}. \quad (1.5.6)$$

This property can be easily tested numerically, for instance by computing $\mathbf{Q}^\top \mathbf{Q}$ for a matrix $\mathbf{Q} = [\mathbf{q}^1, \dots, \mathbf{q}^k] \in \mathbb{R}^{n,k}$.

C++11 code 1.5.7: Wrong result from Gram-Schmidt orthogonalisation EIGEN → GITLAB

```

2 void gsroundoff(MatrixXd& A){
3     // Gram-Schmidt orthogonalization of columns of A, see Code 1.5.3
4     MatrixXd Q = gramschmidt(A);
5     // Test orthonormality of columns of Q, which should be an
6     // orthogonal matrix according to theory
7     cout << setprecision(4) << fixed << "I = "
8     << endl << Q.transpose()*Q << endl;
9     // EIGEN's stable internal Gram-Schmidt orthogonalization by
10    // QR-decomposition, see Rem. 1.5.9 below
11    HouseholderQR<MatrixXd> qr(A.rows(), A.cols()); //
12    qr.compute(A); MatrixXd Q1 = qr.householderQ(); //
13    // Test orthonormality
14    cout << "I1 = " << endl << Q1.transpose()*Q1 << endl;
15    // Check orthonormality and span property (1.5.2)
16    MatrixXd R1 = qr.matrixQR().triangularView<Upper>();
17    cout << scientific << "A-Q1*R1 = " << endl << A-Q1*R1 << endl;
18}

```

We test the orthonormality of the output vectors of Gram-Schmidt orthogonalization for a special matrix $\mathbf{A} \in \mathbb{R}^{10,10}$, a so-called [Hilbert matrix](#), defined by $(\mathbf{A})_{i,j} = (i+j-1)^{-1}$. Then Code 1.5.7 produces the following output:

```
I =
 1.0000  0.0000 -0.0000  0.0000 -0.0000  0.0000 -0.0000 -0.0000 -0.0000 -0.0000
 0.0000  1.0000 -0.0000  0.0000 -0.0000  0.0000 -0.0000 -0.0000 -0.0000 -0.0000
-0.0000 -0.0000  1.0000  0.0000 -0.0000  0.0000 -0.0000 -0.0000 -0.0000 -0.0000
 0.0000  0.0000  0.0000  1.0000 -0.0000  0.0000 -0.0000 -0.0000 -0.0000 -0.0000
-0.0000 -0.0000 -0.0000 -0.0000  1.0000  0.0000 -0.0008 -0.0007 -0.0007 -0.0006
 0.0000  0.0000  0.0000  0.0000  0.0000  1.0000 -0.0540 -0.0430 -0.0360 -0.0289
-0.0000 -0.0000 -0.0000 -0.0000 -0.0008 -0.0540  1.0000  0.9999  0.9998  0.9996
-0.0000 -0.0000 -0.0000 -0.0000 -0.0007 -0.0430  0.9999  1.0000  1.0000  0.9999
-0.0000 -0.0000 -0.0000 -0.0000 -0.0007 -0.0360  0.9998  1.0000  1.0000  1.0000
-0.0000 -0.0000 -0.0000 -0.0000 -0.0006 -0.0289  0.9996  0.9999  1.0000  1.0000
```

Obviously, the vectors produced by the function `gramschmidt` fail to be orthonormal, contrary to the predictions of rigorous results from linear algebra!

However, Line 11, Line 12 of Code 1.5.7 demonstrate another way to orthonormalize the columns of a matrix using EIGEN's built-in class template **HouseholderQR**:

```
II =
 1.0000 -0.0000  0.0000 -0.0000 -0.0000 -0.0000 -0.0000 -0.0000 -0.0000  0.0000
-0.0000  1.0000 -0.0000  0.0000  0.0000 -0.0000 -0.0000 -0.0000  0.0000 -0.0000
 0.0000 -0.0000  1.0000 -0.0000 -0.0000  0.0000  0.0000  0.0000  0.0000  0.0000
-0.0000  0.0000 -0.0000  1.0000  0.0000 -0.0000 -0.0000  0.0000  0.0000  0.0000
-0.0000  0.0000 -0.0000  0.0000  1.0000 -0.0000  0.0000 -0.0000 -0.0000  0.0000
-0.0000 -0.0000  0.0000 -0.0000 -0.0000  1.0000 -0.0000 -0.0000  0.0000 -0.0000
-0.0000 -0.0000  0.0000 -0.0000 -0.0000 -0.0000  1.0000  0.0000  0.0000 -0.0000
-0.0000 -0.0000  0.0000 -0.0000 -0.0000 -0.0000 -0.0000  1.0000  0.0000  0.0000
-0.0000 -0.0000  0.0000  0.0000 -0.0000 -0.0000  0.0000  1.0000 -0.0000  0.0000
-0.0000 -0.0000  0.0000  0.0000 -0.0000 -0.0000 -0.0000 -0.0000  1.0000 -0.0000
```

Now we observe apparently perfect orthogonality (1.5.6) of the columns of the matrix `Q1` in Code 1.5.7. There is another algorithm that reliably yields the theoretical output of Gram-Schmidt orthogonalization.

“Computers cannot compute”

Computers cannot compute “properly” in **R**: numerical computations may not respect the laws of analysis and linear algebra!

This introduces an important new aspect in the study of numerical algorithms.

Remark 1.5.9 (Stable orthonormalization by QR-decomposition)

In Code 1.5.7 we saw the use of the **EIGEN** class **HouseholderQR<MatrixType>** for the purpose of Gram-Schmidt orthogonalisation. The underlying theory and algorithms will be explained later in Section 3.3.3. There we will have the following insight:

- Up to signs the columns of the matrix **Q** available from the QR-decomposition of **A** are the same vectors as produced by the Gram-Schmidt orthogonalisation of the columns of **A**.

Code 1.5.7 demonstrates a case where a desired result can be obtained by two *algebraically equivalent* computations, that is, they yield the same result in a mathematical sense. Yet, when implemented on a computer, the results can be *vastly different*. One algorithm may produce junk (“unstable algorithm”), whereas the other lives up to the expectations (“stable algorithm”)

Supplement to Exp. 1.5.5: despite its ability to produce orthonormal vectors, we get as output for **D=A-Q1*R1** in Code 1.5.7:

```
D =
2.2204e-16 3.3307e-16 3.3307e-16 1.9429e-16 1.9429e-16 5.5511e-17 1.3878e-16 6.9389e-17 8.3267e-17 9.7145e-17
0.0000e+00 1.1102e-16 8.3267e-17 5.5511e-17 0.0000e+00 5.5511e-17 -2.7756e-17 0.0000e+00 0.0000e+00 4.1633e-17
-5.5511e-17 5.5511e-17 2.7756e-17 5.5511e-17 0.0000e+00 0.0000e+00 0.0000e+00 -1.3878e-17 1.3878e-17 1.3878e-17
0.0000e+00 5.5511e-17 2.7756e-17 2.7756e-17 0.0000e+00 1.3878e-17 -1.3878e-17 -1.3878e-17 0.0000e+00 1.3878e-17
0.0000e+00 2.7756e-17 0.0000e+00 1.3878e-17 1.3878e-17 1.3878e-17 0.0000e+00 1.3878e-17 1.3878e-17 4.1633e-17
-2.7756e-17 2.7756e-17 1.3878e-17 4.1633e-17 2.7756e-17 1.3878e-17 0.0000e+00 -1.3878e-17 2.7756e-17 2.7756e-17
0.0000e+00 2.7756e-17 0.0000e+00 2.7756e-17 2.7756e-17 1.3878e-17 0.0000e+00 1.3878e-17 2.7756e-17 2.0817e-17
0.0000e+00 2.7756e-17 2.7756e-17 1.3878e-17 1.3878e-17 1.3878e-17 0.0000e+00 1.3878e-17 2.0817e-17 2.7756e-17
1.3878e-17 1.3878e-17 1.3878e-17 2.7756e-17 1.3878e-17 0.0000e+00 -1.3878e-17 6.9389e-18 -6.9389e-18 1.3878e-17
0.0000e+00 2.7756e-17 1.3878e-17 1.3878e-17 1.3878e-17 1.3878e-17 0.0000e+00 0.0000e+00 1.3878e-17 1.3878e-17
```

- ➔ The computed QR-decomposition apparently fails to meet the exact algebraic requirements stipulated by Thm. 3.3.9. However, note the tiny size of the “defect”.

Remark 1.5.10 (QR-decomposition in **EIGEN**)

The two different QR-decompositions (3.3.3.1) and (3.3.3.1) of a matrix **A** $\in \mathbb{K}^{n,k}$, $n, k \in \mathbb{N}$, can be computed in **EIGEN** as follows:

```
HouseholderQR<MatrixXd> qr(A);
// Full QR-decomposition (3.3.3.1)
Q = qr.householderQ();
// Economical QR-decomposition (3.3.3.1)
thinQ = qr.householderQ() *
MatrixXd::Identity(A.rows(), std::min(A.rows(), A.cols()));
```

The returned matrices **Q** and **R** correspond to the QR-factors **Q** and **R** as defined above. See the discussion, whether the “economical” decomposition is truly economical: some expression template magic.

Remark 1.5.11 (QR-decomposition in PYTHON)

The two different QR-decomposition (3.3.3.1) and (3.3.3.1) of a matrix $A \in \mathbb{K}^{n,k}$, $n, k \in \mathbb{N}$, can be computed in PYTHON as follows:

```
1 Q, R = np.linalg.qr(A, mode='reduced') # Economical (3.3.3.1)
2 Q, R = np.linalg.qr(A, mode='complete') # Full (3.3.3.1)
```

The returned matrices Q and R correspond to the QR-factors Q and R as defined above.

1.5.2 Machine Numbers

(1.5.12) The finite and discrete set of machine numbers

The reason, why computers must fail to execute exact computations with real numbers is clear:

Computer = finite automaton

> can handle only *finitely many* numbers, not \mathbb{R}

machine numbers, set \mathbb{M}

Essential property: \mathbb{M} is a *finite, discrete* subset of \mathbb{R} (its numbers separated by gaps)

The set of machine numbers \mathbb{M} cannot be closed under elementary arithmetic operations $+, -, \cdot, /$, that is, when adding, multiplying, etc., two machine numbers the result may not belong to \mathbb{M} .

The results of elementary operations with operands in \mathbb{M} have to be mapped back to \mathbb{M} , an operation called **rounding**.



roundoff errors (*ger.*: Rundungsfehler) are inevitable

The impact of roundoff means that mathematical identities may not carry over to the computational realm.
As we have seen above in Exp. 1.5.5

Computers cannot compute “properly” !

numerical computations \neq **analysis
linear algebra**

This introduces a *new* and *important* aspect in the study of numerical algorithms!

(1.5.13) Internal representation of machine numbers

Now we give a brief sketch of the internal structure of machine numbers $\in \mathbb{M}$. The main insight will be that

“Computers use floating point numbers (scientific notation)”

Example 1.5.14 (Decimal floating point numbers)

Some 3-digit normalized decimal floating point numbers:

valid: $0.723 \cdot 10^2$, $0.100 \cdot 10^{-20}$, $-0.801 \cdot 10^5$
invalid: $0.033 \cdot 10^2$, $1.333 \cdot 10^{-4}$, $-0.002 \cdot 10^3$

General form of an m -digit normalized decimal floating point number:

never = 0 !
 $x = \pm [0] . \underbrace{[] [] [] []}_{m \text{ digits of mantissa}} \dots [] [] \cdot 10^E$
exponent $\in \mathbb{Z}$

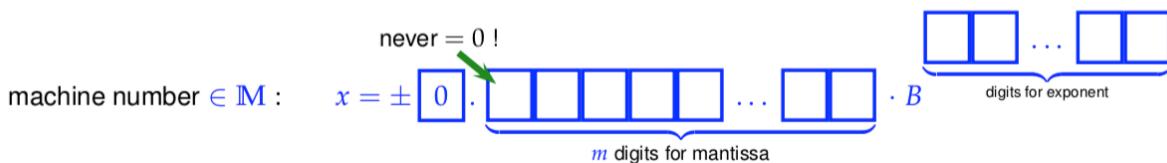
Of course, computers are restricted to a *finite range* of exponents:

Definition 1.5.15. Machine numbers/floating point numbers → [2, Sect. 2.1]

Given

- ☞ basis $B \in \mathbb{N} \setminus \{1\}$,
- ☞ exponent range $\{e_{\min}, \dots, e_{\max}\}$, $e_{\min}, e_{\max} \in \mathbb{Z}$, $e_{\min} < e_{\max}$,
- ☞ number $m \in \mathbb{N}$ of digits (for mantissa),

the corresponding set of machine numbers is

$$\mathbb{M} := \{d \cdot B^E : d = i \cdot B^{-m}, i = B^{m-1}, \dots, B^m - 1, E \in \{e_{\min}, \dots, e_{\max}\}\}$$


Remark 1.5.16 (Extremal numbers in \mathbb{M})

Clearly, there is a largest element of \mathbb{M} and two that are closest to zero. These are mainly determined by the range for the exponent E , cf. Def. 1.5.15.

- Largest machine number (in modulus) : $x_{\max} = \max |\mathbb{M}| = (1 - B^{-m}) \cdot B^{e_{\max}}$
 Smallest machine number (in modulus) : $x_{\min} = \min |\mathbb{M}| = B^{-1} \cdot B^{e_{\min}}$

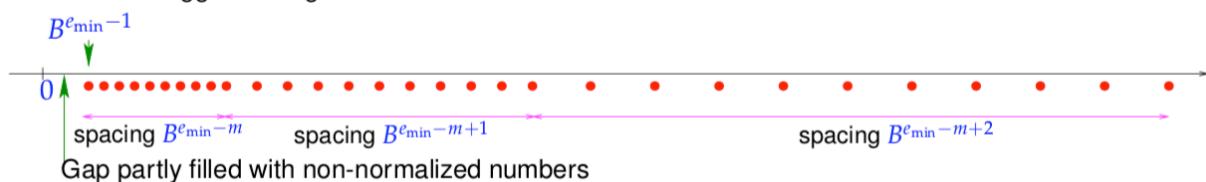
In C++ these extremal machine numbers are accessible through the

`std::numeric_limits<double>::max()`
 and `std::numeric_limits<double>::min()`

functions. Other properties of arithmetic types can be queried accordingly from the `numeric_limits` header.

Remark 1.5.17 (Distribution of machine numbers)

From Def. 1.5.15 it is clear that there are equi-spaced sections of \mathbb{M} and that the gaps between machine numbers are bigger for larger numbers:



Non-normalized numbers violate the lower bound for the mantissa i in Def. 1.5.15.

No surprise: for modern computers $B = 2$ (binary system), the other parameters of the universally implemented machine number system are

single precision : $m = 24^*, E \in \{-125, \dots, 128\} \geq 4$ bytes

double precision : $m = 53^*, E \in \{-1021, \dots, 1024\} \geq 8$ bytes

*: including bit indicating sign

The standardisation of machine numbers is important, because it ensures that the same numerical algorithm, executed on different computers will nevertheless produce the same result.

Remark 1.5.19 (Special cases in IEEE standard)

```
double x = exp(1000), y = 3/x, z = x*sin(M_PI), w = x*log(1);  
cout << x << endl << y << endl << z << endl << w << endl;
```

Output:

```
1 inf  
2 0  
3 inf  
4 -nan
```



$E = e_{\max}, M \neq 0 \doteq \text{NaN} = \text{Not a number} \rightarrow \text{exception}$

$E = e_{\max}, M = 0 \doteq \text{Inf} = \text{Infinity} \rightarrow \text{overflow}$

$E = 0 \doteq \text{Non-normalized numbers} \rightarrow \text{underflow}$

$E = 0, M = 0 \doteq \text{number 0}$

(1.5.20) Characteristic parameters of IEEE floating point numbers (double precision)

- ☞ C++ does not always fulfill the requirements of the IEEE 754 standard and it needs to be checked with `std::numeric_limits<T>::is_iec559`.

C++11-code 1.5.21: Querying characteristics of double numbers → [GITLAB](#)

```
2 #include <limits>
3 #include <iostream>
4 #include <iomanip>
5
6 using namespace std;
7
8 int main() {
9     cout << std::numeric_limits<double>::is_iec559 << endl
10    << std::defaultfloat << numeric_limits<double>::min() << endl
11    << std::hexfloat << numeric_limits<double>::min() << endl
12    << std::defaultfloat << numeric_limits<double>::max() << endl
13    << std::hexfloat << numeric_limits<double>::max() << endl;
14 }
```

Output:

```
1 true
2 2.22507e-308
3 0010000000000000
4 1.79769e+308
5 7 f e f f f f f f f f f f f f
```

1.5.3 Roundoff errors

Experiment 1.5.22 (Input errors and roundoff errors)

The following computations would always result in 0, if done in exact arithmetic.

C++11-code 1.5.23: Demonstration of roundoff errors → [GITLAB](#)

```
2 #include <iostream>
3 int main () {
4     std::cout.precision(15);
5     double a = 4.0/3.0, b = a-1, c = 3*b, e = 1-c;
6     std::cout << e << std::endl;
7     a = 1012.0/113.0; b = a-9; c = 113*b; e = 5+c;
8     std::cout << e << std::endl;
9     a = 83810206.0/6789.0; b = a-12345; c = 6789*b; e = c-1;
10    std::cout << e << std::endl;
11 }
```

Output:

1 2.22044604925031e-16
2 6.75015598972095e-14
3 -1.60798663273454e-09

Can you devise a similar calculation, whose result is even farther off zero? Apparently the rounding that inevitably accompanies arithmetic operations in M can lead to results that are far away from the true result.

For the discussion of errors introduced by rounding we need important notions.

Definition 1.5.24. Absolute and relative error → [2, Sect. 1.2]

Let $\tilde{x} \in \mathbb{K}$ be an approximation of $x \in \mathbb{K}$. Then its **absolute error** is given by

$$\epsilon_{\text{abs}} := |x - \tilde{x}|,$$

and its **relative error** is defined as

$$\epsilon_{\text{rel}} := \frac{|x - \tilde{x}|}{|x|}.$$

Remark 1.5.25 (Relative error and number of correct digits)

The **number of correct** (significant, valid) **digits** of an approximation \tilde{x} of $x \in \mathbb{K}$ is defined through the relative error:

$$\text{If } \epsilon_{\text{rel}} := \frac{|x - \tilde{x}|}{|x|} \leq 10^{-\ell}, \text{ then } \tilde{x} \text{ has } \ell \text{ correct digits, } \ell \in \mathbb{N}_0$$

(1.5.26) Floating point operations

We may think of the elementary binary operations $+, -, *, /$ in \mathbb{M} comprising two steps:

- ① Compute the exact result of the operation.
- ② Perform rounding of the result of ① to map it back to \mathbb{M} .

Definition 1.5.27. Correct rounding

Correct rounding ("rounding up") is given by the function

$$\text{rd} : \left\{ \begin{array}{ccc} \mathbb{R} & \rightarrow & \mathbb{M} \\ x & \mapsto & \max \arg\min_{\tilde{x} \in \mathbb{M}} |x - \tilde{x}| \end{array} \right.$$

(Recall that $\arg\min_x F(x)$ is the set of arguments of a real valued function F that makes it attain its (global) minimum.)

Of course, ① above is not possible in a strict sense, but the effect of both steps can be realised and yields a **floating point realization of $\star \in \{+, -, \cdot, /\}$** .

☞ Notation: write $\tilde{\star}$ for the floating point realization of $\star \in \{+, -, \cdot, /\}$:

Then ① and ② may be summed up into

$$\text{For } \star \in \{+, -, \cdot, /\}: \quad x \tilde{\star} y := \text{rd}(x \star y).$$

Remark 1.5.28 (Breakdown of associativity)

As a consequence of rounding addition $\tilde{+}$ and multiplication $\tilde{\times}$ as implemented on computers fail to be associative. They will usually be commutative, though this is not guaranteed.

(1.5.29) Estimating roundoff errors → [2, p. 23]

Let us denote by EPS the largest **relative error** (\rightarrow Def. 1.5.24) incurred through rounding:

$$\text{EPS} := \max_{x \in I \setminus \{0\}} \frac{|\text{rd}(x) - x|}{|x|}, \quad (1.5.30)$$

where $I = [\min |\mathbb{M}|, \max |\mathbb{M}|]$ is the range of positive machine numbers.

For machine numbers according to Def. 1.5.15 EPS can be computed from the defining parameters B (base) and m (length of mantissa) [2, p. 24]:

$$\text{EPS} = \frac{1}{2}B^{1-m}. \quad (1.5.31)$$

However, when studying roundoff errors, we do not want to delve into the intricacies of the internal representation of machine numbers. This can be avoided by just using a single bound for the relative error due to rounding, and, thus, also for the relative error potentially suffered in each elementary operation.

Assumption 1.5.32. “Axiom” of roundoff analysis

There is a small positive number EPS , the **machine precision**, such that for the elementary arithmetic operations $\star \in \{+, -, \cdot, /\}$ and “hard-wired” functions* $f \in \{\exp, \sin, \cos, \log, \dots\}$ holds

$$x \tilde{\star} y = (x \star y)(1 + \delta) \quad , \quad \tilde{f}(x) = f(x)(1 + \delta) \quad \forall x, y \in \mathbb{M},$$

with $|\delta| < \text{EPS}$.

*: this is an ideal, which may not be accomplished even by modern CPUs.

► relative roundoff errors of elementary steps in a program bounded by machine precision !

Example 1.5.33 (Machine precision for IEEE standard)

C++ tells the machine precision as following:

C++11-code 1.5.34: Finding out EPS in C++ → [GITLAB](#)

```
2 #include <iostream>
3 #include <limits> // get various properties of arithmetic types
4 int main() {
5     std::cout.precision(15);
6     std::cout << std::numeric_limits<double>::epsilon() << std::endl;
7 }
```

Output:

¹ 2.22044604925031 e-16

Knowing the machine precision can be important for checking the validity of computations or coding termination conditions for iterative approximations.

```
cout.precision(25);
double eps =
    numeric_limits<double>::epsilon();
cout << fixed << 1.0 + 0.5*eps << endl
    << 1.0 - 0.5*eps << endl
    << (1.0 + 2/eps) - 2/eps << endl
```

Output:

In fact, the following “definition” of EPS is sometimes used:

EPS is the smallest positive number $\in \mathbb{M}$ for which $1 + \tilde{\text{EPS}} \neq 1$ (in \mathbb{M}):

We find that $\tilde{\text{EPS}} = 1$ actually complies with the “axiom” of roundoff error analysis, Ass. 1.5.32:

$$\frac{2}{\text{EPS}} = \left(1 + \frac{2}{\text{EPS}}\right)(1 + \delta) \Rightarrow |\delta| = \left|\frac{\text{EPS}}{2 + \text{EPS}}\right| < \text{EPS}.$$



Do we have to worry about these tiny roundoff errors ?

YES

(→ Exp. 1.5.5):

- accumulation of roundoff errors
- amplification of roundoff errors

Remark 1.5.36 (Testing equality with zero)



Since results of numerical computations are almost always polluted by roundoff errors:
Tests like `if (x == 0)` are pointless and even dangerous, if `x` contains the result
of a numerical computation.

► Remedy: Test `if (abs(x) < eps*s) ...`,
 $s \doteq$ positive number, compared to which $|x|$ should be small.

We saw a first example of this practise in Code 1.5.3, Line 13.

Remark 1.5.37 (Overflow and underflow)

overflow) \doteq |result of an elementary operation| $>$ `max{M}`

\doteq IEEE standard \Rightarrow `Inf`

underflow \doteq $0 < |$ result of an elementary operation| $<$ `min{|M \ {0}|}`

\doteq IEEE standard \Rightarrow use non-normalized numbers (!)

The Axiom of roundoff analysis Ass. 1.5.32 does *not hold* once non-normalized numbers are encountered:

C++11-code 1.5.38: Demonstration of over-/underflow → [GITLAB](#)

```
2 #include <iostream>
3 #define _USE_MATH_DEFINES
```

```

4 #include <cmath>
5 #include <limits>
6 using namespace std;
7 int main () {
8     cout.precision(15);
9     double min = numeric_limits<double>::min();
10    double res1 = M_PI*min/123456789101112;
11    double res2 = res1*123456789101112/min;
12    cout << res1 << endl << res2 << endl;
13 }

```

Output:

1	5.68175492717434e-322
2	3.15248510554597



Try to avoid underflow and overflow

A simple example teaching how to avoid overflow during the computation of the norm of a 2D vector [2, Ex. 2.9]:

$$r = \sqrt{x^2 + y^2}$$

straightforward evaluation: overflow, when $|x| > \sqrt{\max |M|}$ or $|y| > \sqrt{\max |M|}$.

$$r = \begin{cases} |x| \sqrt{1 + (y/x)^2} & , \text{if } |x| \geq |y| \\ |y| \sqrt{1 + (x/y)^2} & , \text{if } |y| > |x| \end{cases}$$

➤ no overflow!

1.5.4 Cancellation

In general, predicting the impact of roundoff errors on the result of a multi-stage computation is very difficult, if possible at all. However, there is a constellation that is particularly prone to dangerous amplification of roundoff and still can be detected easily.

Example 1.5.39 (Computing the zeros of a quadratic polynomial)

The following simple EIGEN code computes the real roots of a quadratic polynomial $p(\xi) = \xi^2 + \alpha\xi + \beta$ by the discriminant formula

$$p(\xi_1) = p(\xi_2) = 0, \quad \xi_{1/2} = \frac{1}{2} (\alpha \pm \sqrt{D}), \quad \text{if } D := \alpha^2 - 4\beta \geq 0. \quad (1.5.40)$$

C++11-code 1.5.41: Discriminant formula for the real roots of $p(\xi) = \xi^2 + \alpha\xi + \beta$ [→ GITLAB](#)

```
2  ///! C++ function computing the zeros of a quadratic polynomial
3  ///!  $\xi \rightarrow \xi^2 + \alpha\xi + \beta$  by means of the familiar discriminant
4  ///! formula  $\xi_{1,2} = \frac{1}{2}(-\alpha \pm \sqrt{\alpha^2 - 4\beta})$ . However
5  ///! this implementation is vulnerable to round-off! The zeros are
6  ///! returned in a column vector
```

```
7  Vector2d zerosquadpol(double alpha, double beta){
8      Vector2d z;
9      double D = std::pow(alpha, 2) - 4*beta; // discriminant
10     if(D < 0) throw "no real zeros";
11     else{
12         // The famous discriminant formula
13         double wD = std::sqrt(D);
14         z << (-alpha-wD)/2, (-alpha+wD)/2; //
15     }
16     return z;
17 }
```

This formula is applied to the quadratic polynomial $p(\xi) = (\xi - \gamma)(\xi - \frac{1}{\gamma})$ after its coefficients α, β have been computed from γ , which will have introduced *small* relative roundoff errors (of size **EPS**).

C++11-code 1.5.42: Testing the accuracy of computed roots of a quadratic polynomial
→ GITLAB

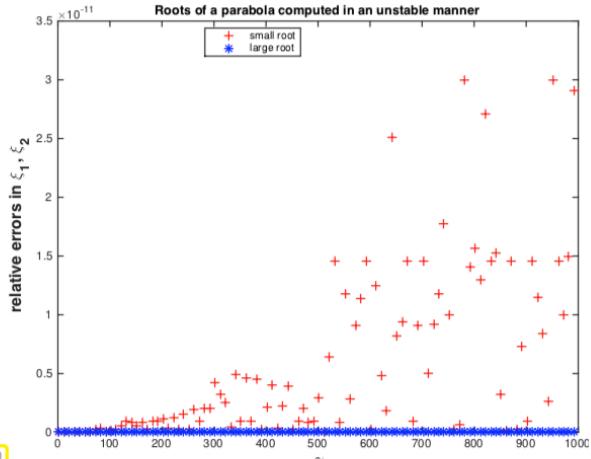
```
2 //! Eigen Function for testing the computation of the zeros of a
   parabola
3 void compzeros() {
4     int n = 100;
5     MatrixXd res(n,4);
6     VectorXd gamma = VectorXd::LinSpaced(n, -2,992);
7     for(int i = 0; i < n; ++i){
8         double alpha = -(gamma(i) + 1./gamma(i));
9         double beta = 1.;
10        Vector2d z1 = zerosquadpol(alpha, beta);
11        Vector2d z2 = zerosquadpolstab(alpha, beta);
12        double ztrue = 1./gamma(i), z2true = gamma(i);
13        res(i,0) = gamma(i);
14        res(i,1) = std::abs((z1(0)-ztrue)/ztrue);
15        res(i,2) = std::abs((z2(0)-ztrue)/ztrue);
16        res(i,3) = std::abs((z1(1)-z2true)/z2true);
17    }
18 // Graphical output of relative error of roots computed by unstable
   implementation
19 mgl::Figure fig1;
20 fig1.setFontSize(3);
21 fig1.title("Roots of a parabola computed in an unstable manner");
22 fig1.plot(res.col(0),res.col(1), " +r").label("small root");
23 fig1.plot(res.col(0),res.col(3), " *b").label("large root");
24 fig1.xlabel("\gamma");
25 fig1.ylabel("relative errors in \xi_1, \xi_2");
26 fig1.legend(0.05,0.95);
27 fig1.save("zqperrinstab");
28 // Graphical output of relative errors (comparison), small roots
29 mgl::Figure fig2;
30 fig2.title("Roundoff in the computation of zeros of a parabola");
31 fig2.plot(res.col(0), res.col(1), " +r").label("unstable");
32
33 fig2.plot(res.col(0), res.col(3), " *m").label("stable");
34 fig2.xlabel("\gamma");
35 fig2.ylabel("relative errors in \xi_1");
36 fig2.legend(0.05,0.95);
37 fig2.save("zqperr");
```

Observation:

Roundoff incurred during the computation of α and β leads to “wrong” roots.

For large γ the computed small root may be fairly inaccurate as regards its *relative error*, which can be several orders of magnitude larger than machine precision EPS .

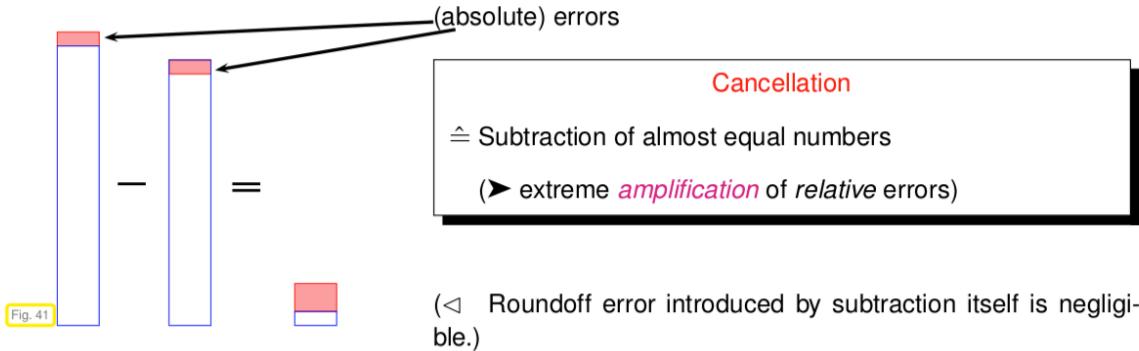
The large root always enjoys a small relative error about the size of EPS .



In order to understand why the small root is much more severely affected by roundoff, note that its computation involves the subtraction of two large numbers, if γ is large. This is the typical situation, in which **cancellation** occurs.

(1.5.43) Visualisation of cancellation effect

We look at the exact subtraction of two almost equal positive numbers both of which have small relative errors (red boxes) with respect to some desired exact value (indicated by blue boxes). The result of the subtraction will be small, but the errors may add up during the subtraction, ultimately constituting a large fraction of the result.



Example 1.5.44 (Cancellation in decimal system)

We consider two positive numbers x, y of about the same size afflicted with relative errors $\approx 10^{-7}$. This means that their seventh decimal digits are perturbed, here indicated by *. When we subtract the

two numbers the perturbed digits are shifted to the left, resulting in a possible relative error of $\approx 10^{-3}$:

$$\begin{array}{rcl}
 x & = & 0.123467* \\
 y & = & 0.123456* \\
 \hline
 x - y & = & 0.000011* = 0.11*000 \cdot 10^{-4}
 \end{array}
 \quad \begin{array}{l}
 \leftarrow \text{7th digit perturbed} \\
 \leftarrow \text{7th digit perturbed} \\
 \leftarrow \text{3rd digit perturbed}
 \end{array}$$

↑
padded zeroes

Again, this example demonstrates that cancellation wreaks havoc through **error amplification**, not through the roundoff error due to the subtraction.

Example 1.5.45 (Cancellation when evaluating difference quotients → [4, Sect. 8.2.6], [2, Ex. 1.3])

From analysis we know that the derivative of a differentiable function $f : I \subset \mathbb{R} \rightarrow \mathbb{R}$ at a point $x \in I$ is the limit of a **difference quotient**

$$f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}.$$

This suggests the following approximation of the derivative by a difference quotient with small but finite $h > 0$

$$f'(x) \approx \frac{f(x+h) - f(x)}{h} \quad \text{for } |h| \ll 1.$$

Results from analysis tell us that the **approximation error** should tend to zero for $h \rightarrow 0$. More precise quantitative information is provided by the Taylor formula for a twice continuously differentiable function [2, p. 5]

$$f(x+h) = f(x) + f'(x)h + \frac{1}{2}f''(\xi)h^2 \quad \text{for some } \xi = \xi(x, h) \in [\min\{x, x+h\}, \max\{x, x+h\}], \quad (1.5.46)$$

from which we infer

$$\frac{f(x+h) - f(x)}{h} - f'(x) = \frac{1}{2}hf''(\xi) \quad \text{for some } \xi = \xi(x, h) \in [\min\{x, x+h\}, \max\{x, x+h\}]. \quad (1.5.47)$$

We investigate the approximation of the derivative by difference quotients for $f = \exp$, $x = 0$, and different values of $h > 0$:

C++11-code 1.5.48: Difference quotient approximation of the derivative of exp → GITLAB

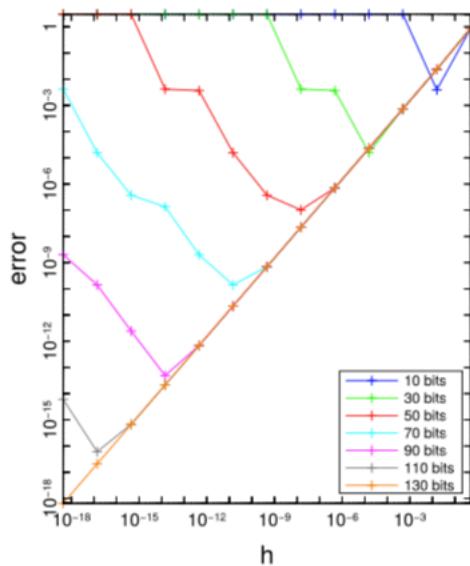
```

2 //!! Difference quotient approximation
3 //!! of the derivative of exp
4 void diffq () {
5     double h = 0.1, x = 0.0;
6     for(int i = 1; i <= 16; ++i){
7         double df = (exp(x+h)-exp(x))/h;
8         cout << setprecision(14) << fixed;
9
10        cout << setw(5) << -i
11        << setw(20) << df-1 << endl;
12        h /= 10;
13    }
}

```

Measured relative errors ▷

We observe an initial decrease of the relative approximation error followed by a steep increase when h drops below 10^{-8} .



$\log_{10}(h)$	relative error
-1	0.05170918075648
-2	0.00501670841679
-3	0.00050016670838
-4	0.00005000166714
-5	0.00000500000696
-6	0.00000049996218
-7	0.00000004943368
-8	-0.0000000607747
-9	0.00000008274037
-10	0.00000008274037
-11	0.00000008274037
-12	0.00008890058234
-13	-0.00079927783736
-14	-0.00079927783736
-15	0.11022302462516
-16	-1.000000000000000

That the observed errors are really due to round-off errors is confirmed by the following numerical results reported besides, using a variable precision floating point module of EIGEN, the [MPFRC++ Support module](#).

C++11-code 1.5.49: Evaluation of difference quotients with variable precision → GITLAB

```

2 // This module provides support for multi precision floating point
   numbers via the MPFR C++ library which itself is built upon
   MPFR/GMP.
3 #include <unsupported/Eigen/MPRealSupport>
4 //!! Numerical differentiation of exponential function with extended
   precision arithmetic

```

```

5 //! Uses the unsupported EIGEN MPFRC++ Support module
6 void numericaldifferentiation(){
7     // declare numeric type
8     typedef mpfr::mpreal numeric_t;
9     // bit number that should be evaluated
10    int n = 7, k = 13;
11    VectorXi bits(n); bits << 10,30,50,70,90,110,130;
12    MatrixXd experr(13, bits.size() + 1);
13    for(int i = 0; i < n; ++i){
14        // set precision to bits(i) bits (double has 53 bits)
15        numeric_t::set_default_prec(bits(i));
16        numeric_t x = "1.1"; // Evaluation point in extended precision
17        for(int j=0;j<k;++j) {
18            numeric_t h = mpfr::pow("2", -1-5*j); // width of
19            // difference quotient in extended precision
20            // compute (absolute) error
21            experr(j,i+1) = mpfr::abs(( mpfr::exp(x+h) -
22                mpfr::exp(x) ) / h - mpfr::exp(x)).toDouble();
23            experr(j,0) = h.toDouble();
24        }
25    }
26    // Plotting
27    // ...

```

Line 16: The literal "1.1" instead of 1.1 prevents the conversion to double.

- Obvious culprit: cancellation when computing the numerator of the difference quotient for small $|h|$ leads to a strong amplification of inevitable errors introduced by the evaluation of the transcendent exponential function.

We witness the competition of two opposite effects: Smaller h results in a better approximation of the derivative by the difference quotient, but the impact of cancellation is the stronger the smaller $|h|$.

$$\left. \begin{array}{l} \text{Approximation error } f'(x) - \frac{f(x+h) - f(x)}{h} \rightarrow 0 \\ \text{Impact of roundoff} \rightarrow \infty \end{array} \right\} \text{as } h \rightarrow 0 .$$

In order to provide a rigorous underpinning for our conjecture, in this example we embark on our first roundoff error analysis merely based on the "Axiom of roundoff analysis" Ass. 1.5.32: As in the computational example above we study the approximation of $f'(x) = e^x$ for $f = \exp, x \in \mathbb{R}$.

$$\begin{aligned} df &= \frac{e^{x+h} (1 + \delta_1) - e^x (1 + \delta_2)}{h} && \text{correction factors take into account roundoff:} \\ &= e^x \left(\frac{e^h - 1}{h} + \frac{\delta_1 e^h - \delta_2}{h} \right) && |\delta_1|, |\delta_2| \leq \text{eps} . \\ \Rightarrow |df| &\leq e^x \left(\frac{e^h - 1}{h} + \text{eps} \frac{1+e^h}{h} \right) && 1 + O(h) \quad O(h^{-1}) \quad \text{for } h \rightarrow 0 \end{aligned}$$

(Note that the estimate for the term $(e^h - 1)/h$ is a particular case of (1.5.47).)

► relative error: $\left| \frac{e^x - df}{e^x} \right| \approx h + \frac{2\text{eps}}{h} \rightarrow \min \quad \text{for } h = \sqrt{2\text{eps}} .$

In double precision: $\sqrt{2\text{eps}} = 2.107342425544702 \cdot 10^{-8}$

Remark 1.5.50 (Cancellation during the computation of relative errors)

In the numerical experiment of Ex. 1.5.45 we computed the relative error of the result by subtraction, see Code 1.5.48. Of course, massive cancellation will occur! Do we have to worry?

In this case cancellation can be tolerated, because we are *interested only* in the magnitude of the relative error. Even if it was affected by a large relative error, this information is still not compromised.

For example, if the relative error has the exact value 10^{-8} , but can be computed only with a huge relative error of 10%, then the perturbed value would still be in the range $[0.9 \cdot 10^{-8}, 1.1 \cdot 10^{-8}]$. Therefore it will still have the correct magnitude and still permit us to conclude the number of valid digits correctly.

Remark 1.5.51 (Cancellation in Gram-Schmidt orthogonalisation of Exp. 1.5.5)

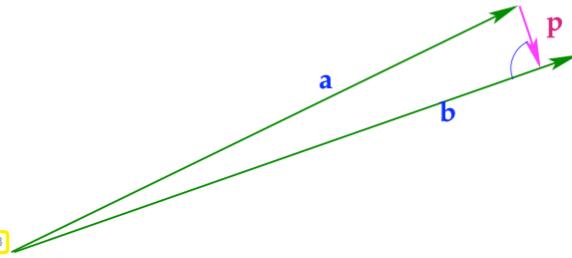
The matrix created by the MATLAB command `A = hilb(10)`, the so-called **Hilbert matrix**, has columns that are almost linearly dependent.

Cancellation when computing orthogonal projection
of vector **a** onto space spanned by vector **b** ▷

$$\mathbf{p} = \mathbf{a} - \frac{\mathbf{a} \cdot \mathbf{b}}{\mathbf{b} \cdot \mathbf{b}} \mathbf{b}.$$

If **a**, **b** point in almost the same direction, $\|\mathbf{p}\| \ll \|\mathbf{a}\|, \|\mathbf{b}\|$, so that a “tiny” vector **p** is obtained by subtracting two “long” vectors, which implies cancellation.

Fig. 43



This can happen in Line 10 of Code 1.5.3.

Example 1.5.52 (Cancellation: roundoff error analysis)

We consider a simple arithmetic expression written in two ways:

$$a^2 - b^2 = (a + b)(a - b), \quad a, b \in \mathbb{R};.$$

We evaluate this term by means of two **algebraically equivalent** algorithms for the input data $a = 1.3$, $b = 1.2$ in 2-digit decimal arithmetic with standard rounding. (“Algebraically equivalent” means that two algorithms will produce the same results in the absence of roundoff errors.

Algorithm A	Algorithm B
$x := a \tilde{+} a = 1.7$ (rounded)	$x := a \tilde{+} b = 2.5$ (exact)
$y := b \tilde{-} b = 1.4$ (rounded)	$y := a \tilde{-} b = 0.1$ (exact)
$x \tilde{-} y = 0.30$ (exact)	$x * y = 0.25$ (exact)

Algorithm B produces the exact result, whereas Algorithm A fails to do so. Is this pure coincidence or an indication of the superiority of algorithm B? This question can be answered by **roundoff error analysis**. We demonstrate the approach for the two algorithms A & B and general input $a, b \in \mathbb{R}$.

Roundoff error analysis heavily relies on Ass. 1.5.32 and dropping terms of "higher order" in the machine precision, that is terms that behave like $O(\text{EPS}^q)$, $q > 1$. It involves introducing the relative roundoff error for every elementary operation through a factor $(1 + \delta)$, $|\delta| \leq \text{EPS}$.

Algorithm A:

$$x = a^2(1 + \delta_1), y = b^2(1 + \delta_2)$$

$$\tilde{f} = (a^2(1 + \delta_1) - b^2(1 + \delta_2))(1 + \delta_3) = f + a^2\delta_1 - b^2\delta_2 + (a^2 - b^2)\delta_3 + O(\text{EPS}^2)$$



$$\frac{|\tilde{f} - f|}{|f|} \leq \text{EPS} \frac{a^2 + b^2 + |a^2 - b^2|}{|a^2 - b^2|} + O(\text{EPS}^2) = \text{EPS} \left(1 + \frac{|a^2 + b^2|}{|a^2 - b^2|} \right) + O(\text{EPS}^2). \quad (1.5.53)$$

will be neglected

► For $a \approx b$ the relative error of the result of Algorithm A will be much larger than the machine precision **EPS**. This reflects cancellation in the last subtraction step.

Algorithm B:

$$x = (a + b)(1 + \delta_1), y = (a - b)(1 + \delta_2)$$

$$\tilde{f} = (a + b)(a - b)(1 + \delta_1)(1 + \delta_2)(1 + \delta_3) = f + (a^2 - b^2)(\delta_1 + \delta_2 + \delta_3) + O(\text{EPS}^2)$$



$$\frac{|\tilde{f} - f|}{|f|} \leq |\delta_1 + \delta_2 + \delta_3| + O(\text{EPS}^2) \leq 3\text{EPS} + O(\text{EPS}^2). \quad (1.5.54)$$

► Relative error of the result of Algorithm B is always $\approx \text{EPS}$!

In this example we see a general guideline at work:

If inevitable, subtractions prone to cancellation should be done as early as possible.

The reason is that input data and initial intermediate results are usually not as much tainted by roundoff errors as numbers computed after many steps.

(1.5.55) Avoiding disastrous cancellation

The following examples demonstrate a few fundamental techniques for steering clear of cancellation by using alternative formulas that yield the same value (in exact arithmetic), but do not entail subtracting two numbers of almost size.

Example 1.5.56 (Stable discriminant formula → Ex. 1.5.39, [2, Ex. 2.10])

If ξ_1 and ξ_2 are the two roots of the quadratic polynomial $p(\xi) = \xi^2 + \alpha\xi + \beta$, then $\xi_1 \cdot \xi_2 = \beta$ (Vieta's formula). Thus once we have computed a root, we can obtain the other by simple division.



Idea:

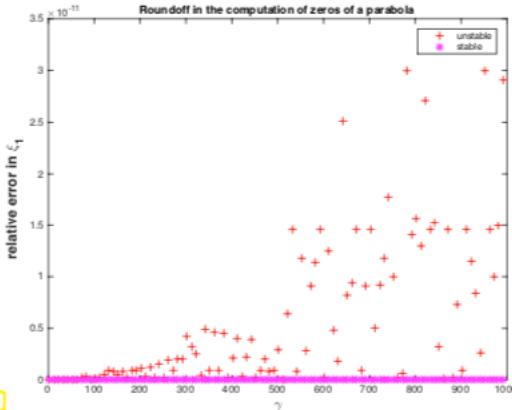
- ① Depending on the sign of α compute "stable root" without cancellation.
- ② Compute other root from Vieta's formula (avoiding subtraction)

C++11-code 1.5.57: Stable computation of real root of a quadratic polynomial → GITLAB

```

2  ///! C++ function computing the zeros of a quadratic polynomial
3  ///!  $\xi \rightarrow \xi^2 + \alpha\xi + \beta$  by means of the familiar discriminant
4  ///! formula  $\xi_{1,2} = \frac{1}{2}(-\alpha \pm \sqrt{\alpha^2 - 4\beta})$ .
5  ///! This is a stable implementation based on Vieta's theorem.
6  ///! The zeros are returned in a column vector
7  VectorXd zerosquadpolstab(double alpha, double beta){
8      Vector2d z(2);
9      double D = std::pow(alpha, 2) - 4*beta; // discriminant
10     if(D < 0) throw "no real zeros";
11     else{
12         double wD = std::sqrt(D);
13         // Use discriminant formula only for zero far away from 0
14         // in order to avoid cancellation. For the other zero
15         // use Vieta's formula.
16         if(alpha >= 0){
17             double t = 0.5*(-alpha-wD); //
18             z << t, beta/t;
19         }
20         else{
21             double t = 0.5*(-alpha+wD); //
22             z << beta/t, t;
23         }
24     }
25     return z;
26 }
```

→ Invariably, we add numbers with the same sign in Line 17 and Line 21.



Numerical experiment based on the driver code
Code 1.5.42.

Observation:

The new code can also compute the small root of the polynomial $p(\xi) = (\xi - \gamma)(\xi - \frac{1}{\gamma})$ (expanded in monomials) with a relative error $\approx \text{EPS}$.

Example 1.5.58 (Exploiting trigonometric identities to avoid cancellation)

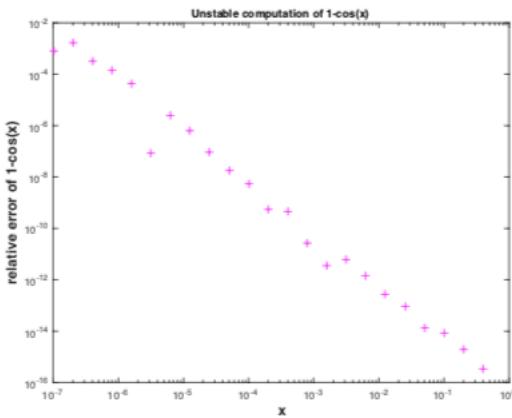
The task is to evaluate the integral

$$\int_0^x \sin t dt = \underbrace{1 - \cos x}_{I} = \underbrace{2 \sin^2(x/2)}_{II} \quad \text{for } 0 < x \ll 1, \quad (1.5.59)$$

and this can be done by the two different formulas *I* and *II*.

Relative error of expression *I* ($1 - \cos(x)$) with respect to equivalent expression *II* ($2 * \sin(x/2)^2$) ▷

Expression *I* is affected by cancellation for $|x| \ll 1$, since then $\cos x \approx 1$, whereas expression *II* can be evaluated with a relative error $\approx \text{EPS}$ for all x .

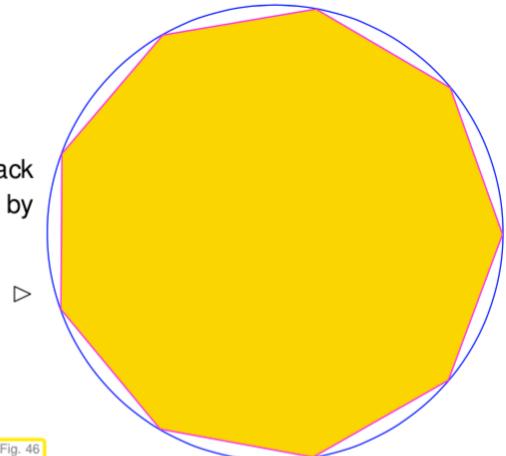


Analytic manipulations offer ample opportunity to rewrite expressions in equivalent form immune to cancellation.

Example 1.5.60 (Switching to equivalent formulas to avoid cancellation)

Now we see an example of a computation allegedly dating back to Archimedes, who tried to approximate the area of a circle by the areas of inscribed regular polygons.

Approximation of circle by regular n -gon, $n \in \mathbb{N}$



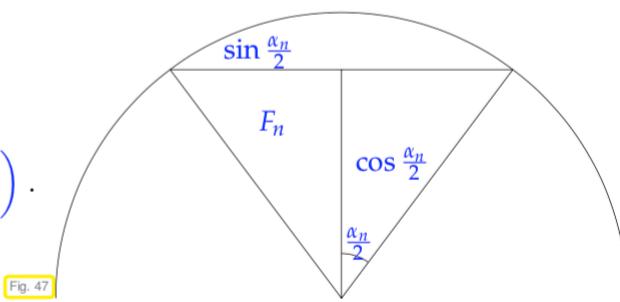
Area of n -gon:

$$A_n = n \cos \frac{\alpha_n}{2} \sin \frac{\alpha_n}{2} = \frac{n}{2} \sin \alpha_n = \frac{n}{2} \sin\left(\frac{2\pi}{n}\right).$$

Recursion formula for A_n derived from

$$\sin \frac{\alpha_n}{2} = \sqrt{\frac{1 - \cos \alpha_n}{2}} = \sqrt{\frac{1 - \sqrt{1 - \sin^2 \alpha_n}}{2}},$$

$$\text{Initial approximation: } A_6 = \frac{3}{2}\sqrt{3}.$$



C++11-code 1.5.61: Tentative computation of circumference of regular polygon → GITLAB

```

2  ///! Approximation of Pi by approximating the circumference of a
3  ///! regular polygon
4  MatrixXd ApproxPiInstable(double tol = 1e-8, int maxIt = 50){
5      double s=sqrt(3)/2.; double An=3.*s;// initialization (hexagon case)
6      unsigned int n = 6, it = 0;
7      MatrixXd res(maxIt,4); // matrix for storing results
8      res(it,0) = n; res(it,1) = An;
9      res(it,2) = An - M_PI; res(it,3)=s;
10     while( it < maxIt && s > tol ){// terminate when s is 'small enough'
11         s = sqrt((1.- sqrt(1.-s*s))/2.); // recursion for area
12         n *= 2; An = n/2.*s;           // new estimate for circumference
13         ++it;
14         res(it,0) =n; res(it,1) =An;// store results and (absolute) error
15         res(it,2) = An - M_PI; res(it,3)=s;
16     }
17     return res.topRows(it);
18 }
```

The approximation deteriorates after applying the recursion formula many times:

n	A_n	$A_n - \pi$	$\sin \alpha_n$
6	2.598076211353316	-0.543516442236477	0.866025403784439
12	3.000000000000000	-0.141592653589794	0.500000000000000
24	3.105828541230250	-0.035764112359543	0.258819045102521
48	3.132628613281237	-0.008964040308556	0.130526192220052
96	3.139350203046872	-0.002242450542921	0.065403129230143
192	3.141031950890530	-0.000560702699263	0.032719082821776
384	3.141452472285344	-0.000140181304449	0.016361731626486
768	3.141557607911622	-0.000035045678171	0.008181139603937
1536	3.141583892148936	-0.000008761440857	0.004090604026236
3072	3.141590463236762	-0.000002190353031	0.002045306291170
6144	3.141592106043048	-0.000000547546745	0.001022653680353
12288	3.141592516588155	-0.000000137001638	0.000511326906997
24576	3.141592618640789	-0.000000034949004	0.000255663461803
49152	3.141592645321216	-0.000000008268577	0.000127831731987
98304	3.141592645321216	-0.000000008268577	0.000063915865994
196608	3.141592645321216	-0.000000008268577	0.000031957932997
393216	3.141592645321216	-0.000000008268577	0.000015978966498
786432	3.141593669849427	0.0000001016259634	0.000007989485855
1572864	3.141592303811738	-0.000000349778055	0.000003994741190
3145728	3.141608696224804	0.000016042635011	0.000001997381017
6291456	3.141586839655041	-0.000005813934752	0.000000998683561
12582912	3.141674265021758	0.000081611431964	0.000000499355676
25165824	3.141674265021758	0.000081611431964	0.000000249677838
50331648	3.143072740170040	0.001480086580246	0.000000124894489
100663296	3.159806164941135	0.018213511351342	0.000000062779708
201326592	3.181980515339464	0.040387861749671	0.000000031610136
402653184	3.354101966249685	0.212509312659892	0.000000016660005
805306368	4.242640687119286	1.101048033529493	0.000000010536712
1610612736	6.000000000000000	2.858407346410207	0.000000007450581

Where does cancellation occur in Line 11 of Code 1.5.61? Since $s \ll 1$, computing $1 - s$ will not trigger cancellation. However, the subtraction $1 - \sqrt{1 - s^2}$ will, because $\sqrt{1 - s^2} \approx 1$ for $s \ll 1$:

Where does cancellation occur in Line 11 of Code 1.5.61? Since $s \ll 1$, computing $1 - s$ will not trigger cancellation. However, the subtraction $1 - \sqrt{1 - s^2}$ will, because $\sqrt{1 - s^2} \approx 1$ for $s \ll 1$:

$$\sin \frac{\alpha_n}{2} = \sqrt{\frac{1 - \cos \alpha_n}{2}} = \sqrt{\frac{1 - \sqrt{1 - \sin^2 \alpha_n}}{2}}$$

For $\alpha_n \ll 1$: $\sqrt{1 - \sin^2 \alpha_n} \approx 1$
 Cancellation here!

We arrive at an *equivalent* formula not vulnerable to cancellation essentially using the identity $(a + b)(a - b) = a^2 - b^2$ in order to eliminate the difference of square roots in the numerator.

$$\begin{aligned}
 \sin \frac{\alpha_n}{2} &= \sqrt{\frac{1 - \sqrt{1 - \sin^2 \alpha_n}}{2}} = \sqrt{\frac{1 - \sqrt{1 - \sin^2 \alpha_n}}{2} \cdot \frac{1 + \sqrt{1 - \sin^2 \alpha_n}}{1 + \sqrt{1 - \sin^2 \alpha_n}}} \\
 &= \sqrt{\frac{1 - (1 - \sin^2 \alpha_n)}{2(1 + \sqrt{1 - \sin^2 \alpha_n})}} = \frac{\sin \alpha_n}{\sqrt{2(1 + \sqrt{1 - \sin^2 \alpha_n})}}.
 \end{aligned}$$

C++11-code 1.5.62: Stable recursion for area of regular n -gon ➔ GITLAB

```
2 //! Approximation of Pi by approximating the circumference of a
3
4 //! regular polygon
5 MatrixXd apprpistable(double tol = 1e-8, int maxIt = 50){
6     double s=sqrt(3)/2.; double An=3.*s;// initialization (hexagon case)
7     unsigned int n = 6, it = 0;
8     MatrixXd res(maxIt,4); // matrix for storing results
9     res(it,0) = n; res(it,1) = An;
10    res(it,2) = An - M_PI; res(it,3)=s;
11    while( it < maxIt && s > tol ){// terminate when s is 'small enough'
12        s = s/sqrt(2*(1+sqrt((1+s)*(1-s))));// Stable recursion without
13        cancellation
14        n *= 2; An = n/2.*s;      // new estimate for circumference
15        ++it;
16        res(it,0) =n; res(it,1) =An;// store results and (absolute) error
17        res(it,2) = An - M_PI; res(it,3)=s;
18    }
19    return res.topRows(it);
20}
```

Using the stable recursion, we observe better approximation for polygons with more corners:

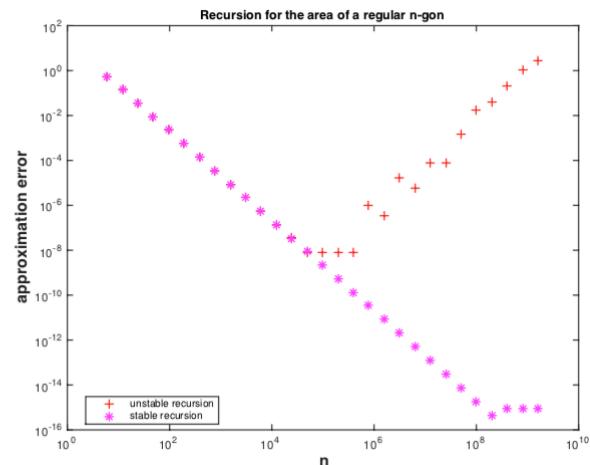
n	A_n	$A_n - \pi$	$\sin \alpha_n$
6	2.598076211353316	-0.543516442236477	0.866025403784439
12	3.000000000000000	-0.141592653589793	0.500000000000000
24	3.105828541230249	-0.035764112359544	0.258819045102521
48	3.132628613281238	-0.008964040308555	0.130526192220052
96	3.139350203046867	-0.002242450542926	0.065403129230143
192	3.141031950890509	-0.000560702699284	0.032719082821776
384	3.141452472285462	-0.000140181304332	0.016361731626487
768	3.141557607911857	-0.000035045677936	0.008181139603937
1536	3.141583892148318	-0.000008761441475	0.004090604026235
3072	3.141590463228050	-0.000002190361744	0.002045306291164
6144	3.141592105999271	-0.000000547590522	0.001022653680338
12288	3.141592516692156	-0.000000136897637	0.000511326907014
24576	3.141592619365383	-0.000000034224410	0.000255663461862
49152	3.141592645033690	-0.000000008556103	0.000127831731976
98304	3.141592651450766	-0.000000002139027	0.000063915866118
196608	3.141592653055036	-0.000000000534757	0.000031957933076
393216	3.141592653456104	-0.000000000133690	0.000015978966540
786432	3.141592653556371	-0.000000000033422	0.000007989483270
1572864	3.141592653581438	-0.000000000008355	0.000003994741635
3145728	3.141592653587705	-0.000000000002089	0.000001997370818
6291456	3.141592653589271	-0.000000000000522	0.000000998685409
12582912	3.141592653589663	-0.000000000000130	0.000000499342704
25165824	3.141592653589761	-0.00000000000032	0.000000249671352
50331648	3.141592653589786	-0.00000000000008	0.000000124835676
100663296	3.141592653589791	-0.00000000000002	0.000000062417838
201326592	3.141592653589794	0.00000000000000	0.000000031208919
402653184	3.141592653589794	0.00000000000001	0.000000015604460
805306368	3.141592653589794	0.00000000000001	0.000000007802230
1610612736	3.141592653589794	0.00000000000001	0.000000003901115

Plot of errors for approximations of π as computed by the two algebraically equivalent recursion formulas▷

Observation, cf. Ex. 1.5.45

Amplified roundoff errors due to cancellation supersedes approximation error for $n \geq 10^5$.

Roundoff errors merely of magnitude EPS in the case of stable recursion



Example 1.5.63 (Summation of exponential series)

In principle, the function value $\exp(x)$ can be approximated up to any accuracy by summing sufficiently many terms of the globally convergent exponential series.

$$\begin{aligned}\exp(x) &= \sum_{k=0}^{\infty} \frac{x^k}{k!} \\ &= 1 + x + \frac{x^2}{2} + \frac{x^3}{6} + \frac{x^4}{24} + \dots\end{aligned}$$

C++11-code 1.5.64: Summation of exponential series → GITLAB

```
2 double expeval(double x,
3                 double tol=1e-8){
4     // Initialization
5     double y = 1.0, term = 1.0;
6     long int k = 1;
7     // Termination criterion
8     while(abs(term) > tol*y) {
9         term *= x/k;      // next summand
10        y += term; // Summation
11        ++k;
12    }
13    return y;
14 }
```

Results for $\text{tol} = 10^{-8}$, $\tilde{\exp}$ designates the approximate value for $\exp(x)$ returned by the function from Code 1.5.64. Rightmost column lists relative errors, which tells us the number of valid digits in the approximate result.

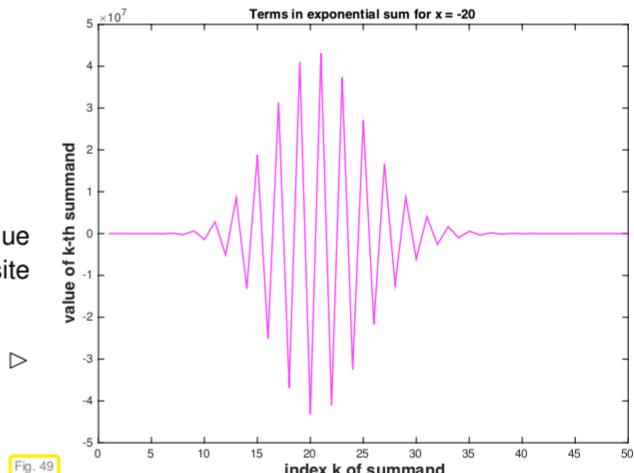
x	Approximation $\widetilde{\exp}(x)$	$\exp(x)$	$\frac{ \exp(x) - \widetilde{\exp}(x) }{\exp(x)}$
-20	6.1475618242e-09	2.0611536224e-09	1.982583033727893
-18	1.5983720359e-08	1.5229979745e-08	0.049490585500089
-16	1.1247503300e-07	1.1253517472e-07	0.000534425951530
-14	8.3154417874e-07	8.3152871910e-07	0.000018591829627
-12	6.1442105142e-06	6.1442123533e-06	0.000000299321453
-10	4.5399929604e-05	4.5399929762e-05	0.000000003501044
-8	3.3546262812e-04	3.3546262790e-04	0.000000000662004
-6	2.4787521758e-03	2.4787521767e-03	0.000000000332519
-4	1.8315638879e-02	1.8315638889e-02	0.000000000530724
-2	1.3533528320e-01	1.3533528324e-01	0.000000000273603
0	1.0000000000e+00	1.0000000000e+00	0.000000000000000
2	7.3890560954e+00	7.3890560989e+00	0.000000000479969
4	5.4598149928e+01	5.4598150033e+01	0.000000001923058
6	4.0342879295e+02	4.0342879349e+02	0.000000001344248
8	2.9809579808e+03	2.9809579870e+03	0.000000002102584
10	2.2026465748e+04	2.2026465795e+04	0.000000002143799
12	1.6275479114e+05	1.6275479142e+05	0.000000001723845
14	1.2026042798e+06	1.2026042842e+06	0.000000003634135
16	8.8861105010e+06	8.8861105205e+06	0.000000002197990
18	6.5659968911e+07	6.5659969137e+07	0.000000003450972
20	4.8516519307e+08	4.8516519541e+08	0.000000004828737

Observation:

Large relative approximation errors for $x \ll 0$.

For $x \ll 0$ we have $|\exp(x)| \ll 1$, but this value is computed by summing large numbers of opposite sign.

Terms summed up for $x = -20$



Remedy: Cancellation can be avoided by using identity

$$\exp(x) = \frac{1}{\exp(-x)} , \text{ if } x < 0 .$$



Example 1.5.65 (Combat cancellation by approximation)

In a computer code we have to provide a routine for the evaluation of

$$\int_0^1 e^{at} dt = \frac{\exp(a) - 1}{a} \quad \text{for any } a > 0 , \quad (1.5.66)$$

cf. the discussion of cancellation in the context of numerical differentiation in Ex. 1.5.45. There we observed massive cancellation.

Recall the **Taylor expansion** formula in one dimension for a function that is $m + 1$ times continuously differentiable in a neighborhood of x [14, Satz 5.5.1]

$$f(x+h) = \sum_{k=0}^m \frac{1}{k!} f^{(k)}(x) h^k + R_m(x, h) , \quad R_m(x, h) = \frac{1}{(m+1)!} f^{(m+1)}(\xi) h^{m+1} ,$$

for some $\xi \in [\min\{x, x+h\}, \max\{x, x+h\}]$, and for all sufficiently small $|h|$. Here $R(x, h)$ is called the **remainder** term and $f^{(k)}$ denotes the k derivative of f .

Cancellation in (1.5.66) can be avoided by replacing $\exp(a)$, $a > 0$, with a suitable Taylor expansion around $a = 0$ and then dividing by a :

$$\frac{\exp(a) - 1}{a} = \sum_{k=0}^m \frac{1}{(k+1)!} a^k + R_m(a) , \quad R_m(a) = \frac{1}{(m+1)!} \exp(\xi) a^{m+1} \text{ for some } 0 \leq \xi \leq a .$$

For a similar discussion see [2, Ex. 2.12].

Issue: How to choose the number m of terms to be retained in the Taylor expansion? We have to pick m large enough such that the relative approximation error remains below a prescribed threshold `tol`. To estimate the relative approximation error, we use the expression for the remainder together with the simple estimate $(\exp(a) - 1)/a > 1$ for all $a > 0$:

$$\text{rel. err.} = \frac{(e^a - 1)/a - \sum_{k=0}^m \frac{1}{(k+1)!} a^k}{(e^a - 1)/a} \leq \frac{1}{(m+1)!} \exp(\xi) a^{m+1} \leq \frac{1}{(m+1)!} \exp(a) a^m .$$

For $a = 10^{-3}$ we get

m	1	2	3	4	5
	1.0010e-03	5.0050e-07	1.6683e-10	4.1708e-14	8.3417e-18

Hence, keeping $m = 3$ terms is enough for achieving about 10 valid digits.

Relative error of unstable formula $(\exp(a) - 1.0)/a$ and relative error, when using a Taylor expansion approximation for small a \triangleright

```
if (abs(a) < 1E-3)
    v = 1.0 + (1.0/2 + 1.0/6*a)*a;
else
    v = (\exp(a)-1.0)/a;
end
```

Error computed by comparison with MATLAB's built-in function `expml` that provides a stable implementation of $\exp(x) - 1$.

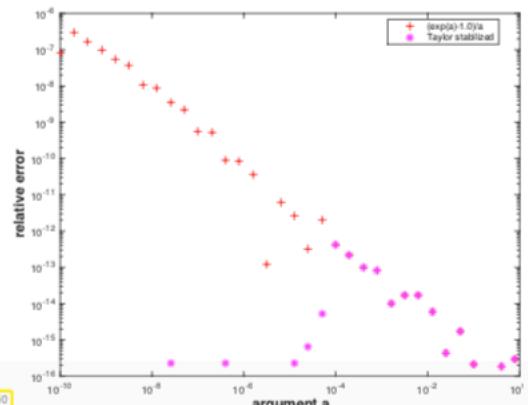


Fig. 50

1.5.5 Numerical stability

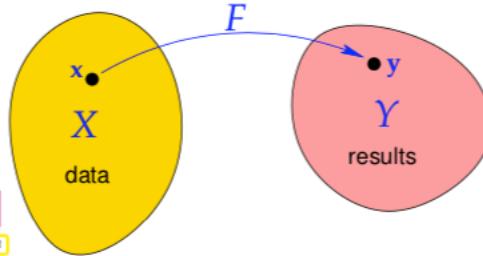
We have seen that a particular “problem” can be tackled by different “algorithms”, which produce different results due to roundoff errors. This section will clarify what distinguishes a “good” algorithm from a rather abstract point of view.

(1.5.67) The “problem”

A mathematical notion of “problem”:

- ◆ data space X , usually $X \subset \mathbb{R}^n$
- ◆ result space Y , usually $Y \subset \mathbb{R}^m$
- ◆ mapping (problem function) $F: X \mapsto Y$

A problem is a well defined *function* that assigns to each datum a result.



Note: In this course, both the data space X and the result space Y will always be subsets of finite dimensional **vector spaces**.

Example 1.5.68 (The “matrix×vector-multiplication problem”)

We consider the “problem” of computing the product \mathbf{Ax} for a given matrix $\mathbf{A} \in \mathbb{K}^{m,n}$ and a given vector $\mathbf{x} \in \mathbb{K}^n$.

-
- Data space $X = \mathbb{K}^{m,n} \times \mathbb{K}^n$ (input is a matrix and a vector)
 - Result space $Y = \mathbb{K}^m$ (space of column vectors)
 - Problem function $F: X \rightarrow Y$, $F(\mathbf{a}, \mathbf{x}) := \mathbf{Ax}$

(1.5.69) Norms on spaces of vectors and matrices

Norms provide tools for measuring errors. Recall from linear algebra and calculus [10, Sect. 4.3], [7, Sect. 6.1]:

Definition 1.5.70. Norm

X = vector space over field \mathbb{K} , $\mathbb{K} = \mathbb{C}, \mathbb{R}$. A map $\|\cdot\| : X \mapsto \mathbb{R}_0^+$ is a **norm** on X , if it satisfies

- (i) $\forall \mathbf{x} \in X: \mathbf{x} \neq 0 \Leftrightarrow \|\mathbf{x}\| > 0$ (definite),
- (ii) $\|\lambda \mathbf{x}\| = |\lambda| \|\mathbf{x}\| \quad \forall \mathbf{x} \in X, \lambda \in \mathbb{K}$ (homogeneous),
- (iii) $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\| \quad \forall \mathbf{x}, \mathbf{y} \in X$ (triangle inequality).

Examples: (for vector space \mathbb{K}^n , vector $\mathbf{x} = (x_1, x_2, \dots, x_n)^T \in \mathbb{K}^n$)

name	:	definition	EIGEN function
Euclidean norm	:	$\ \mathbf{x}\ _2 := \sqrt{ x_1 ^2 + \dots + x_n ^2}$	<code>x.norm()</code>
1-norm	:	$\ \mathbf{x}\ _1 := x_1 + \dots + x_n $	<code>x.lpNorm<1>()</code>
∞ -norm, max norm	:	$\ \mathbf{x}\ _\infty := \max\{ x_1 , \dots, x_n \}$	<code>x.lpNorm<Eigen::Infinity>()</code>

Remark 1.5.71 (Inequalities between vector norms)

All norms on the vector space \mathbb{K}^n , $n \in \mathbb{N}$, are **equivalent** in the sense that for arbitrary two norms $\|\cdot\|_1$ and $\|\cdot\|_2$ we can always find a constant $C > 0$ such that

$$\|\mathbf{v}\|_1 \leq C \|\mathbf{v}\|_2 \quad \forall \mathbf{v} \in \mathbb{K}^n. \quad (1.5.72)$$

Of course, the constant C will usually depend on n and the norms under consideration.

For the vector norms introduced above, explicit expressions for the constants " C " are available: for all $\mathbf{x} \in \mathbb{K}^n$

$$\|\mathbf{x}\|_2 \leq \|\mathbf{x}\|_1 \leq \sqrt{n} \|\mathbf{x}\|_2, \quad (1.5.73)$$

$$\|\mathbf{x}\|_\infty \leq \|\mathbf{x}\|_2 \leq \sqrt{n} \|\mathbf{x}\|_\infty, \quad (1.5.74)$$

$$\|\mathbf{x}\|_\infty \leq \|\mathbf{x}\|_1 \leq n \|\mathbf{x}\|_\infty. \quad (1.5.75)$$

The matrix space $\mathbb{K}^{m,n}$ is a vector space, of course, and can also be equipped with various norms. Of particular importance are norms *induced by vector norms* on \mathbb{K}^n and \mathbb{K}^m .

Definition 1.5.76. Matrix norm

Given vector norms $\|\cdot\|_1$ and $\|\cdot\|_2$ on \mathbb{K}^n and \mathbb{K}^m , respectively, the associated **matrix norm** is defined by

$$\mathbf{M} \in \mathbb{R}^{m,n}: \quad \|\mathbf{M}\| := \sup_{\mathbf{x} \in \mathbb{R}^n \setminus \{\mathbf{0}\}} \frac{\|\mathbf{M}\mathbf{x}\|_2}{\|\mathbf{x}\|_1}.$$

By virtue of definition the matrix norms enjoy an important property, they are **sub-multiplicative**:

$$\forall \mathbf{A} \in \mathbb{K}^{n,m}, \mathbf{B} \in \mathbb{K}^{m,k}: \quad \|\mathbf{AB}\| \leq \|\mathbf{A}\| \|\mathbf{B}\|. \quad (1.5.77)$$

☞ notations for matrix norms for *quadratic matrices* associated with standard vector norms:

$$\|\mathbf{x}\|_2 \rightarrow \|\mathbf{M}\|_2, \quad \|\mathbf{x}\|_1 \rightarrow \|\mathbf{M}\|_1, \quad \|\mathbf{x}\|_\infty \rightarrow \|\mathbf{M}\|_\infty$$

Example 1.5.78 (Matrix norm associated with ∞ -norm and 1-norm)

Rather simple formulas are available for the matrix norms induced by the vector norms $\|\cdot\|_\infty$ and $\|\cdot\|_1$

$$\begin{aligned} \text{e.g. for } \mathbf{M} \in \mathbb{K}^{2,2}: \quad & \|\mathbf{M}\mathbf{x}\|_\infty = \max\{|m_{11}x_1 + m_{12}x_2|, |m_{21}x_1 + m_{22}x_2|\} \\ & \leq \max\{|m_{11}| + |m_{12}|, |m_{21}| + |m_{22}|\} \|\mathbf{x}\|_\infty, \\ & \|\mathbf{M}\mathbf{x}\|_1 = |m_{11}x_1 + m_{12}x_2| + |m_{21}x_1 + m_{22}x_2| \\ & \leq \max\{|m_{11}| + |m_{21}|, |m_{12}| + |m_{22}|\}(|x_1| + |x_2|). \end{aligned}$$

For general $\mathbf{M} \in \mathbb{K}^{m,n}$

$$\geq \text{matrix norm} \leftrightarrow \|\cdot\|_\infty = \text{row sum norm} \quad \|\mathbf{M}\|_\infty := \max_{i=1,\dots,m} \sum_{j=1}^n |m_{ij}|, \quad (1.5.79)$$

$$\geq \text{matrix norm} \leftrightarrow \|\cdot\|_1 = \text{column sum norm} \quad \|\mathbf{M}\|_1 := \max_{j=1,\dots,n} \sum_{i=1}^m |m_{ij}|. \quad (1.5.80)$$

Sometimes special formulas for the Euclidean matrix norm come handy [5, Sect. 2.3.3]:

Lemma 1.5.81. Formula for Euclidean norm of a Hermitian matrix

$$\mathbf{A} \in \mathbb{K}^{n,n}, \mathbf{A} = \mathbf{A}^H \Rightarrow \|\mathbf{A}\|_2 = \max_{\mathbf{x} \neq 0} \frac{|\mathbf{x}^H \mathbf{A} \mathbf{x}|}{\|\mathbf{x}\|_2^2}.$$

Proof. Recall from linear algebra: Hermitian matrices (a special class of normal matrices) enjoy unitary similarity to diagonal matrices:

$$\exists \mathbf{U} \in \mathbb{K}^{n,n}, \text{diagonal } \mathbf{D} \in \mathbb{R}^{n,n}: \mathbf{U}^{-1} = \mathbf{U}^H \text{ and } \mathbf{A} = \mathbf{U}^H \mathbf{D} \mathbf{U}.$$

Since multiplication with an unitary matrix preserves the 2-norm of a vector, we conclude

$$\|\mathbf{A}\|_2 = \left\| \mathbf{U}^H \mathbf{D} \mathbf{U} \right\|_2 = \|\mathbf{D}\|_2 = \max_{i=1,\dots,n} |d_i|, \quad \mathbf{D} = \text{diag}(d_1, \dots, d_n).$$

On the other hand, for the same reason:

$$\max_{\|\mathbf{x}\|_2=1} \mathbf{x}^H \mathbf{A} \mathbf{x} = \max_{\|\mathbf{x}\|_2=1} (\mathbf{U} \mathbf{x})^H \mathbf{D} (\mathbf{U} \mathbf{x}) = \max_{\|\mathbf{y}\|_2=1} \mathbf{y}^H \mathbf{D} \mathbf{y} = \max_{i=1,\dots,n} |d_i|.$$

Hence, both expressions in the statement of the lemma agree with the largest modulus of eigenvalues of \mathbf{A} . \square

Corollary 1.5.82. Euclidean matrix norm and eigenvalues

For $\mathbf{A} \in \mathbb{K}^{n,n}$ the Euclidean matrix norm $\|\mathbf{A}\|_2$ is the square root of the largest (in modulus) eigenvalue of $\mathbf{A}^H \mathbf{A}$.

For a *normal* matrix $\mathbf{A} \in \mathbb{K}^{n,n}$ (that is, \mathbf{A} satisfies $\mathbf{A}^H \mathbf{A} = \mathbf{A} \mathbf{A}^H$) the Euclidean matrix norm agrees with the modulus of the largest eigenvalue.

(1.5.83) (Numerical) algorithm

When we talk about an “algorithm” we have in mind a *concrete code function* in MATLAB or C++; the only way to describe an algorithm is through a piece of code. We assume that this function defines another mapping $\tilde{F}: X \rightarrow Y$ on the data space of the problem. Of course, we can only feed data to the MATLAB/C++-function, if they can be represented in the set \mathbb{M} of machine numbers. Hence, implicit in the definition of \tilde{F} is the assumption that input data are subject to *rounding* before passing them to the code function proper.

Problem	Algorithm
$F: X \subset \mathbb{R}^n \rightarrow Y \subset \mathbb{R}^m$	$\tilde{F}: X \rightarrow \tilde{Y} \subset \mathbb{M}$

(1.5.84) Stable algorithm → [2, Sect. 1.3]

[Stable algorithm]

- ◆ We study a problem (\rightarrow § 1.5.67) $F : X \rightarrow Y$ on data space X into result space Y .
- ◆ We assume that both X and Y are equipped with norms $\|\cdot\|_X$ and $\|\cdot\|_Y$, respectively (\rightarrow Def. 1.5.70).
- ◆ We consider a concrete algorithm $\tilde{F} : X \rightarrow Y$ according to § 1.5.83.

We write $w(x)$, $x \in X$, for the **computational effort** (\rightarrow Def. 1.4.1) required by the algorithm for input x .

Definition 1.5.85. Stable algorithm

An algorithm \tilde{F} for solving a problem $F : X \mapsto Y$ is **numerically stable** if for all $x \in X$ its result $\tilde{F}(x)$ (possibly affected by roundoff) is the exact result for "slightly perturbed" data:

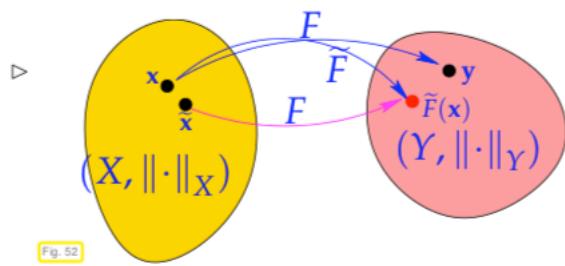
$$\exists C \approx 1: \forall x \in X: \exists \tilde{x} \in X: \|x - \tilde{x}\|_X \leq C w(x) \text{ EPS} \|x\|_X \wedge \tilde{F}(x) = F(\tilde{x}).$$

Here **EPS** should be read as machine precision according to the "Axiom" of roundoff analysis Ass. 1.5.32.

Illustration of Def. 1.5.85
($y \hat{=} \text{exact result for exact data } x$)

Terminology:

Def. 1.5.85 introduces stability in the sense of
backward error analysis



Sloppily speaking, the impact of roundoff (*) on a *stable algorithm* is of the same order of magnitude as the effect of the inevitable perturbations due to rounding the input data.

➤ For stable algorithms roundoff errors are "harmless".

(*) In some cases the definition of \tilde{F} will also involve some approximations as in Ex. 1.5.65. Then the above statement also includes approximation errors.

Example 1.5.86 (Testing stability of matrix \times vector multiplication)

Assume you are given a black box implementation of a function

```
VectorXd mvmult(const MatrixXd &A, const VectorXd &x)
```

that purports to provide a stable implementation of \mathbf{Ax} for $\mathbf{A} \in \mathbb{K}^{m,n}$, $\mathbf{x} \in \mathbb{K}^n$, cf. Ex. 1.5.68. How can we verify this claim for particular data. Both, $\mathbb{K}^{m,n}$ and \mathbb{K}^n are equipped with the Euclidean norm.

The task is, given $\mathbf{y} \in \mathbb{K}^n$ as returned by the function, to find conditions on \mathbf{y} that ensure the existence of a $\tilde{\mathbf{A}} \in \mathbb{K}^{m,n}$ such that

$$\tilde{\mathbf{A}}\mathbf{x} = \mathbf{y} \quad \text{and} \quad \|\tilde{\mathbf{A}} - \mathbf{A}\|_2 \leq C_{mn} \text{ EPS} \|\mathbf{A}\|_2, \quad (1.5.87)$$

for a small constant ≈ 1 .

In fact we can choose (easy computation)

$$\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{z}\mathbf{x}^T, \quad \mathbf{z} := \frac{\mathbf{y} - \mathbf{Ax}}{\|\mathbf{x}\|_2^2} \in \mathbb{K}^m,$$

and we find

$$\|\tilde{\mathbf{A}} - \mathbf{A}\|_2 = \|\mathbf{z}\mathbf{x}^T\|_2 = \sup_{\mathbf{w} \in \mathbb{K}^n \setminus \{\mathbf{0}\}} \frac{\mathbf{x} \cdot \mathbf{w} \|\mathbf{z}\|_2}{\|\mathbf{w}\|_2} \leq \|\mathbf{x}\|_2 \|\mathbf{z}\|_2 = \frac{\|\mathbf{y} - \mathbf{Ax}\|_2}{\|\mathbf{x}\|_2}.$$

Hence, in principle stability of an algorithm for computing \mathbf{Ax} is confirmed, if for every $\mathbf{x} \in X$ the computed result $\mathbf{y} = \mathbf{y}(\mathbf{x})$ satisfies

$$\|\mathbf{y} - \mathbf{Ax}\|_2 \leq C mn \text{EPS} \|\mathbf{x}\|_2 \|\mathbf{A}\|_2,$$

with a small constant $C > 0$ independent of data and problem size.

Remark 1.5.88 (Numerical stability and sensitive dependence on data)

A problem shows **sensitive dependence** on the data, if small perturbations of input data lead to large perturbations of the output. Such problems are also called **ill-conditioned**. For such problems stability of an algorithm is easily accomplished.



53

Example: The problem is the prediction of the position of the billiard ball after ten bounces given the initial position, velocity, and spin.

It is well known, that tiny changes of the initial conditions can shift the final location of the ball to virtually any point on the table: the billiard problem is **chaotic**.

Hence, a stable algorithm for its solution may just output a **fixed** or **random** position without even using the initial conditions!

Learning Outcomes

Principal take-home knowledge and skills from this chapter:

- Knowledge about the syntax of fundamental operations on matrices and vectors in EIGEN.
- Understanding of the concepts of computational effort/cost and asymptotic complexity in numerics.
- Awareness of the asymptotic complexity of basic linear algebra operations
- Ability to determine the (asymptotic) computational effort for a concrete (numerical linear algebra) algorithm.
- Ability to manipulate simple expressions involving matrices and vectors in order to reduce the computational cost for their evaluation.