

# Computational Mathematics 1

Lecture notes, University (of Technology) Graz  
based on the lecture by Tobias Breiten

Lukas Prokop

December 13, 2018

## Contents

## Course

↓ *This lecture took place on 2018/10/02.*

- Homepage at [uni-graz.at](http://uni-graz.at)
- Contents:
  1. Linear equation system
  2. Numerical error analysis
  3. Curve fitting (dt. "Lineare Ausgleichsrechnung")
  4. Non-linear equations
  5. Interpolation
  6. Numerical integration
- Literature:
  1. Deuffhard, Hohmann: Numerische Mathematik 1
  2. Schwarz, Köckler: Numerische Mathematik

## 1 What is Computational Mathematics?

Computational Mathematics as branch of mathematics considers the construction and analysis of algorithms for continuous mathematical problems. More specifically, based on a model and input data, we construct specific output data of interest. Formally, we associate a model using function  $f : X \rightarrow Y$  where  $X$  is the set of input data and  $Y$  the set of output data. The problem is to find output  $f(x) \in Y$  for given  $x \in X$ .

Examples:

- Addition of two numbers:  $X = \mathbb{R}^2, Y = \mathbb{R}, f : (x_1, x_2) \mapsto x_1 + x_2$
- Solution of a linear equation system:  $Ax = b$

$$X = \mathbb{R}^{n \times n} \times \mathbb{R}^n, Y = \mathbb{R}^n, f : (A, b) \mapsto A^{-1}b$$

An algorithm is defined as unambiguous sequence of steps to solve a given problem. A problem is characterized as

- the input data required for computation
- the output data that represent the solution of the algorithm
- the description of the steps to be done, including auxiliary values

We require, the algorithm is

**executable** by a machine

**terminating** after a finite number of steps

**deterministic** as the same input always leads to the same output

Furthermore we analyze the algorithm in terms of

**stability** small errors in the input lead to small errors in the output

**precision** output data should solve the problem with maximum accuracy

**efficiency** large problems can be solved in practical time

Consider the linear equation system

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n &= b \\ &\vdots \\ a_{n1}x_1 + a_{n2}x_2 + \cdots + a_{nn}x_n &= b \end{aligned}$$

The same can be specified in compact notation like  $Ax = b$  with  $A \in \mathbb{R}^{n \times n}, b \in \mathbb{R}^n$ .

**Theorem 1.1** (Cramer's Rule). Let  $A \in \mathbb{R}^{n \times n}$  with  $\det(A) \neq 0$  and  $b \in \mathbb{R}^n$ . Then there exists exactly one  $x \in \mathbb{R}^n$  such that  $Ax = b$ . We can determine  $x$  using "Cramer's Rule":

$$x_i = \frac{\det(A_i)}{\det(A)} \text{ where } A_i = \begin{bmatrix} a_{11} & \dots & a_{1,i-1} & b_1 & a_{1,i+1} & \dots & a_{1n} \\ \vdots & & & & & & \vdots \\ a_{n1} & \dots & a_{n,i-1} & b_n & a_{n,i+1} & \dots & a_{nn} \end{bmatrix}$$

**Remark** (Problem). Computation using  $\det(A) = \sum_{\sigma \in S_n} \text{sign}(\sigma) a_{1,\sigma(1)} \dots a_{n,\sigma(n)}$  requires  $n \cdot n!$  operations.

## 1.1 Solving linear eq. system by forward/backwards substitution

Consider the special case of a linear equation system of structure

$$\begin{array}{ccccccc} \gamma_{11}x_1 & +\gamma_{12}x_2 & +\dots & +\gamma_{1n}x_n & = & z_1 \\ & \gamma_{22}x_2 & +\dots & +\gamma_{2n}x_n & = & z_2 \\ & & \ddots & & = & \vdots \\ & & & \gamma_{nn}x_n & = & z_n \end{array}$$

and accordingly,  $Rx = z$  with an upper triangular matrix, hence  $\gamma_{ij} = 0$  for  $i > j$ . We retrieve  $x$  by recursive solution (by so-called "backwards substitution") starting with row  $n$ :

$$\begin{aligned} x_n &= \frac{z_n}{\gamma_{nn}} \text{ if } \gamma_{nn} \neq 0 \\ x_{n-1} &= \frac{z_{n-1} - \gamma_{n-1,n}x_n}{\gamma_{n-1,n-1}} && \text{if } \gamma_{n-1,n-1} \neq 0 \\ x_1 &= \frac{z_1 - \gamma_{12}x_2 - \dots - \gamma_{1n}x_n}{\gamma_{11}} && \text{if } \gamma_{11} \neq 0 \end{aligned}$$

Obviously, it holds that

$$\det(R) = \gamma_{11} \dots \gamma_{nn} \neq 0 \iff \gamma_{ii} \neq 0 \quad \forall i = 1, \dots, n$$

The algorithm is applicable (just like Cramer's rule), if  $\det(R) \neq 0$ . Then the existence of a solution is guaranteed and this solution is unique.

**Remark** (Computing time).

1. for the  $i$ -th row, we need  $(n - i)$  additions,  $(n - i)$  multiplications and one division
2. in total for the rows  $n$  to 1,

$$\sum_{i=1}^n (i - 1) = \frac{n(n - 1)}{2} \doteq \frac{n^2}{2}$$

Multiplication and the same amount of additions. The notation  $\doteq$  is used to denote equivalence up to terms of same degree. Analogously, the equation system  $Lx = z$  can be solved with a lower triangular matrix  $L$  by forward substitution starting with the first row.

## 1.2 Gaussian elimination process

*Basic idea:*

Beginning with a linear equation system of structure,

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n &= b_1 \\ a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n &= b_2 \\ &\vdots \\ a_{n1}x_1 + a_{n2}x_2 + \dots + a_{nn}x_n &= b_n \end{aligned} \quad (1)$$

we apply equivalence transformations to end up with an upper triangular matrix.

As a first step, we eliminate variable  $x_1$  of rows 2 to  $n$ .

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n &= b_1 \\ a'_{22}x_2 + \dots + a'_{2n}x_n &= b'_2 \\ &\vdots \\ a'_{n2}x_2 + \dots + a'_{nn}x_n &= b'_n \end{aligned} \quad (2)$$

So we apply this elimination step ( $r_i := \frac{r_i}{a_{i1}} - \frac{r_1}{a_{11}}$  where  $r_i$  denotes the  $i$ -th row) to the last  $n - 1$  rows and determine the triangular matrix recursively this way.

Consider the transformation of (1) to (2). Let  $a_{11} \neq 0$  and we recognize  $\text{row}_i := \text{row}_i - l_{i1} \cdot \text{row}_1$ . Formally,

$$\begin{aligned} \underbrace{(a_{i1} - l_{i1}a_{11})}_{=0}x_1 + \underbrace{(a_{i2} - l_{i1}a_{12})}_{a'_{i2}}x_2 + \dots + \underbrace{(a_{in} - l_{i1}a_{1n})}_{a'_{in}}x_n &= \underbrace{b_i - l_{i1}b_1}_{b'_i} \\ \implies l_{in} &= \frac{a_{i1}}{a_{11}} \end{aligned}$$

hence, the first elimination step is applicable if  $a_{11} \neq 0$  is given. If we apply this step, we retrieve a sequence of matrices:

$$A = A^{(1)} \rightarrow A^{(2)} \rightarrow \dots \rightarrow A^{(n)} = R$$

where

$$A^{(k)} = \begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & \dots & a_{1n}^{(1)} \\ 0 & a_{22}^{(2)} & \dots & a_{2n}^{(2)} \\ & & \ddots & \vdots \\ & & & a_{kk}^{(k)} & \dots & a_{kn}^{(k)} \\ & & & \vdots & & \vdots \\ & & & a_{nk}^{(k)} & \dots & a_{nn}^{(k)} \end{bmatrix}$$

with a remaining matrix of dimension  $(n - k + 1) \times (n - k + 1)$ . For every remaining matrix we can apply the elimination step

$$\begin{aligned} l_{ik} &= \frac{a_{ik}^{(k)}}{a_{kk}^{(k)}} & \text{for } i = k + 1, \dots, n \\ a_{ij}^{(k+1)} &= a_{ij}^{(k)} - l_{ik}a_{kj}^{(k)} & \text{for } i, j = k + 1, \dots, n \\ b_i^{(k+1)} &= b_i^{(k)} - l_{ik}b_k^{(k)} & \text{for } i = k + 1, \dots, n \end{aligned}$$

under the assumption that  $a_{kk}^{(k)} \neq 0$  (pivot element).

**Remark.** Every elimination step is a linear operation of the rows of  $A$ . Hence it is representable by left multiplication with  $L_k \in \mathbb{R}^{n \times n}$  according to  $A^{k+1} = L_k A^{(k)}$  and  $b^{(k+1)} = L_k b^{(k)}$ .

The following matrix is a Frobenius matrix, which has an interesting property regarding its inverse:

$$L_k = \begin{pmatrix} 1 & & & & \\ & \ddots & & & \\ & & 1 & & \\ & & -l_{k+1,k} & \ddots & \\ & & \vdots & \ddots & \\ & & -l_{n,k} & \dots & 1 \end{pmatrix} \iff L_k^{-1} = \begin{pmatrix} 1 & & & & \\ & \ddots & & & \\ & & 1 & & \\ & & l_{k+1,k} & \ddots & \\ & & \vdots & \ddots & \\ & & l_{n,k} & \dots & 1 \end{pmatrix}$$

and

$$L = L_1^{-1} \dots L_{n-1}^{-1} = \begin{pmatrix} 1 & & & & 0 \\ l_{21} & \ddots & & & \\ l_{31} & l_{32} & \ddots & & \\ \vdots & \vdots & & \ddots & \\ l_{n1} & \dots & l_{n,n-1} & 1 \end{pmatrix}$$

$$L^{-1} = L_{n-1} \dots L_1$$

This way, we retrieve an equivalent linear equation system to  $Rx = z$ ,  $R = L^{-1}A$ ,  $z = L^{-1}b$ . The representation  $A = LR$  is called *LR decomposition* of  $A$ . If they exist,  $L$  and  $R$  are unambiguously determined.

$$A = L_1 R_1 = L_2 R_2 \implies R_1 = L_1^{-1} L_2 R_2 \implies R_1 R_2^{-1} = L_1^{-1} L_2$$

$$\implies \text{diagonal representation of } L_1^{-1} L_2 \dots, R_1 R_2^{-1}$$

**Algorithm 1** (Gaussian elimination).

1.  $A = LR$  with  $R$  as upper triangular matrix,  $L$  as lower triangular matrix
2. Solve  $Lz = b$  by forward substitution
3. Solve  $Rx = z$  by backwards substitution

The entries  $l_{ik} \neq 1$  can be stored in the remaining null-entries of the matrix  $A^{(k)}$ . Thus the total memory requirements are bounded by  $n(n+1)$ .

**Remark** (Computing time).

- $\sum_{k=1}^{n-1} k^2 \doteq \frac{n^3}{3}$  for step 1
- $\sum_{k=1}^{n-1} k \doteq \frac{n^2}{2}$  for step 2 and 3

### 1.3 Pivot strategies

**Problem:** The algorithm above does not work for the simple example:

$$A = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \quad \det(A) = -1 \neq 0$$

We need row exchanges. We can evaluate, that no LU-decomposition exists.

But exchange of the rows gives

$$A = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \rightsquigarrow \hat{A} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = I = LU$$

with  $L = U = I$ .

**Example 1.2.** Consider the equation system

$$\begin{aligned} 1 \cdot 10^{-4}x_1 + 1 \cdot x_2 &= 1 \\ 1 \cdot x_1 + 1 \cdot x_2 &= 2 \end{aligned}$$

**Assumption:** Compute up to three decimal points. Then for the “accurate” solution, we get  $x_1 = 1$  and  $x_2 = 1$ .

Using the Gaussian elimination process, we get

$$\begin{aligned} l_{21} &= \frac{a_{21}}{a_{11}} = 10^4 \\ \implies 0 \cdot x_1 + (1 - 10^4 \cdot 1)x_2 &= 2 \cdot 10^4 \cdot 1 \end{aligned}$$

Therefore we get the triangular system

$$\begin{aligned} 10^{-4}x_1 + 1 \cdot x_2 &= 1 \\ -10^4x_2 &= -10^4 \end{aligned}$$

By this, we get an approximation  $x_2 = 1$  and  $x_1 = 0$ .

In contrast, if we exchange the two rows,

$$\begin{aligned} 1 \cdot x_1 + 1 \cdot x_2 &= 2 \\ 10^{-4}x_1 + 1 \cdot x_2 &= 1 \end{aligned}$$

then  $\tilde{l}_21 = \frac{10^4}{1}$  follows and therefore the triangular system

$$\begin{aligned} 1 \cdot x_1 + 1 \cdot x_2 &= 2 \\ 1 \cdot x_2 &= 1 \end{aligned}$$

with the “correct” solution is  $x_2 = 1$  and  $x_1 = 1$ .

Our exchange of the rows has given us  $|\tilde{l}_{21}| < 1$  and  $|\tilde{a}_{11}| \geq |\tilde{a}_{21}|$ .

↓ This lecture took place on 2018/10/04.

The new pivot element  $\tilde{a}_{11}$  is the largest absolute value in the first column. Thus, we apply the column pivot strategy: In every step, choose the row with the largest, absolute value in the pivot column.

**Algorithm 2** (Gaussian elimination with column pivot strategy).

1. In step  $A^{(k)} \rightarrow A^{(k+1)}$ , choose some  $p \in \{k, \dots, n\}$  such that  $|a_{pk}^{(k)}| \geq |a_{jk}^{(k)}|$  for  $j = k, \dots, n$
2. Exchange rows  $p$  and  $k$

$$A^{(k)} \rightsquigarrow \tilde{A}^{(k)} \text{ with } \tilde{a}_{ij}^{(k)} = \begin{cases} a_{kj}^{(k)} & i = p \\ a_{pj}^{(k)} & i = k \\ a_{ij}^{(k)} & \text{else} \end{cases}$$

3. Apply the elimination step  $\tilde{A}^{(k)} \rightarrow \tilde{A}^{(k+1)}$

**Remark 1.3.** We can apply the row pivot strategy with column exchange instead of the column pivot strategy with row exchange.

We use permutation matrices  $P \in \mathbb{R}^{n \times n}$  to analyze the column pivot search. By permutation  $\pi \in S_n$ , we associate the *permutation matrix*

$$P_\pi = [e_{\pi(1)}, e_{\pi(2)}, \dots, e_{\pi(n)}]$$

where  $e_j$  denotes the  $j$ -th unit vector in  $\mathbb{R}^n$ . Row and column exchanges of a matrix  $A$  can followingly be represented as

$$\pi_z : A \mapsto P_\pi A, \quad \pi_S : A \mapsto AP_\pi$$

**Remark.** It holds that  $\det(P_\pi) = \text{sign}(\pi) \in \{\pm 1\}$  and  $P_\pi^{-1} = P_\pi^T$ . [Consider  $P_\pi = e_{\pi(1)} \cdot e_1^T + e_{\pi(2)} e_2^T + \dots + e_{\pi(n)} e_n^T$ .]

$$\begin{aligned} P_\pi^T &= e_1 e_{\pi(1)}^T + e_2 e_{\pi(2)}^T + \dots + e_n e_{\pi(n)}^T \\ P_\pi P_\pi^T &= e_{\pi(1)} e_{\pi(1)}^T + e_{\pi(2)} e_{\pi(2)}^T + \dots + e_{\pi(n)} e_{\pi(n)}^T \end{aligned}$$

**Theorem 1.4.** Let  $A \in \mathbb{R}^{n \times n}$  and  $\det(A) \neq 0$ . Then there exists some permutation matrix  $P \in \mathbb{R}^{n \times n}$  such that

$$PA = LU$$

with lower and upper triangular matrices  $L, U \in \mathbb{R}^{n \times n}$ . Furthermore  $P$  can be chosen such that  $|L| := \max_{i,j} |l_{ij}| \leq 1$ .

*Proof.* We use Algorithm 2. Because  $\det(A) \neq 0$ ,  $\max_i |a_{i1}| > 0$  and therefore there exists some transposition  $\tau_1 \in S_n$  such that for  $A^{(1)} = P_{\tau_1} A$ ,

$$\max_i |a_{i1}^{(1)}| = |a_{11}^{(1)}| > 0$$

We can apply the elimination step and retrieve a matrix of structure

$$A^{(2)} = L_1 A^{(1)} = L_1 P_{\tau_1} A = \begin{bmatrix} a_{11}^{(1)} & x & x & \dots & x \\ 0 & & & & \\ \vdots & & B^{(2)} & & \\ 0 & & & & \end{bmatrix}$$

Especially, it holds that  $\max_{i,j} |l_{i,j}^{(1)}| \leq 1$  and  $\det(L_1) = 1$ . And

$$0 \neq \text{sign}(\tau_1) \cdot \det(A) = \det(A^{(2)}) = a_{11}^{(1)} \det(B^{(2)})$$

and thus  $\det(B^{(2)}) \neq 0$ . Inductively, we get  $U = A^{(n)} = L_{n-1} P_{\tau_{n-1}} \dots L_1 P_{\tau_1} A$  with  $|L_k| \leq 1$  and transposition of two numbers  $\geq k$ . If  $\pi \in S_n$  exchanges only numbers  $\geq k+1$ , then

$$L_k = \begin{bmatrix} 1 & & & & \\ & 1 & & & \\ & -l_{k+1,k} & 1 & & \\ 0 & \vdots & & \ddots & \\ & l_{n,k} & & 0 & 1 \end{bmatrix}, \hat{L}_k = P_\pi L_k P_\pi^{-1} = \begin{bmatrix} 1 & & & & \\ & 1 & & & \\ & -l_{\pi(k+1),k} & 1 & & \\ 0 & \vdots & & \ddots & \\ & l_{\pi(n),k} & & 0 & 1 \end{bmatrix}$$

By insertion of the identities  $P_\pi^{-1} P_{\tau_k}$  we get,

$$\begin{aligned} &= L_{n-1} P_{\tau_{n-1}} L_{n-2} P_{\tau_{n-1}}^{-1} P_{\tau_{n-1}} P_{\tau_{n-2}} L_{n-3} P_{\tau_{n-3}} \dots L_1 P_{\tau_1} A \\ &= (\underbrace{L_{n-1}}_{\hat{L}_{n-1}}) (P_{\tau_{n-1}} L_{n-2} P_{\tau_{n-1}}^{-1}) (P_{\tau_{n-1}} P_{\tau_{n-2}} L_{n-3} P_{\tau_{n-2}}^{-1} P_{\tau_{n-1}}^{-1} P_{\tau_{n-1}} P_{\tau_{n-2}} P_{\tau_{n-3}} \dots L_1 P_{\tau_n} A) \\ &= \hat{L}_{n-1} \hat{L}_{n-2} \hat{L}_{n-3} \dots \hat{L}_1 P_{\pi_0} A \end{aligned}$$

with  $\hat{L}_k = P_{\pi_k} L_k P_{\pi_k}^{-1}$ ,  $\pi_{n-1} = \text{id}$ ,  $\pi_k = \tau_{n-1} \dots \tau_{k+1}$ ,  $k = 0, \dots, n-2$ . Be aware that  $\pi_k$  only exchanges numbers  $\geq k+1$  and the matrices  $\hat{L}_k$  are therefore of structure previously mentioned. As a consequence, we created the decomposition  $P_{\pi_0} A = LU$  with

$$L := \hat{L}_1^{-1} \dots \hat{L}_{n-1}^{-1}, \quad L = \begin{bmatrix} 1 & & & & \\ l_{\pi_1(2),1} & 1 & & & \\ l_{\pi_2(3),1} & l_{\pi_2(3)} & \ddots & & \\ \vdots & \vdots & & \ddots & \\ l_{\pi_1(n),1} & \dots & & l_{\pi_{n-1}(n),n-1} & 1 \end{bmatrix}$$

Especially,  $|L| = 1$ . □

**Remark 1.5.** Recognize that the method mentioned in the proof can be used to compute the determinant of  $A$ .

$$\det(A) = \det(P) \cdot \det(LU) = \text{sign}(\pi_0) \cdot \gamma_{11} \dots \gamma_{nn}$$

## 1.4 Cholesky process for symmetric positive definite matrices



↓ This lecture took place on 2018/10/09.

In the following, let  $A = A^T > 0$  be symmetric positive definite.

**Remark.** By symmetric positive definite property,

1.  $\langle x, Ax \rangle > 0 \forall x \neq 0$
2. there exists an orthonormal basis  $q_1, \dots, q_n$  of eigenvectors such that  $Q^T Q = I$

**Theorem 1.6.** For every symmetric positive definite matrix  $A \in \mathbb{R}^{n \times n}$  it holds that

1.  $A$  is invertible
2.  $a_{ii} > 0$  for  $i = 1, \dots, n$
3.  $\max_{i,j=1,\dots,n} |a_{ij}| = \max_{i=1,\dots,n} a_{ii}$
4. Gaussian elimination without pivot search gives a symmetric positive definite remainder matrix in every step

*Proof.* 1. follows by  $\langle x, Ax \rangle > 0 \forall x \neq 0$ . Assume  $A$  is not invertible, then  $\exists x \in \text{kernel}(A) \implies \langle x, Ax \rangle = \langle x, 0 \rangle = 0$  which contradicts with  $A$  being symmetric positive definite.

2. also follows by  $\langle x, Ax \rangle > 0$  with choice  $x = e_i$  for  $i = 1, \dots, n$  because  $\langle e_i, Ae_i \rangle = a_{ii}$

3. Left as an exercise

4. We denote  $A = A^{(1)}$  according to

$$A^{(1)} = \begin{bmatrix} a_{11} & z^T \\ z & B^{(1)} \end{bmatrix}$$

with  $z = [a_{21}, \dots, a_{n1}]^T$ . After the first elimination step,

$$A^{(2)} = L_1 A^{(1)} = \begin{bmatrix} a_{11}^{(1)} & z^T \\ 0 & B^{(2)} \\ \vdots & \\ 0 & \end{bmatrix}$$

$$L_1 = \begin{bmatrix} 1 & & & \\ -l_{21} & \ddots & & \\ \vdots & 0 & \ddots & \\ -l_{n1} & & & 1 \end{bmatrix}$$

Right-side multiplication with  $L_1^T$  gives

$$L_1 A^{(1)} L_1^T = \begin{bmatrix} a_{11}^{(1)} & 0 & \dots & 0 \\ 0 & & & \\ \vdots & & B^{(2)} & \\ 0 & & & \end{bmatrix}$$

Furthermore with  $A^{(1)} = A > 0$  it also holds that  $L_1 A^{(1)} L_1^T > 0$ . By this, it is especially true that  $B^{(2)} > 0$ .

□

**Theorem 1.7.** *For every symmetric positive definite matrix there exists a unique decomposition  $A = LDL^T$  where  $L$  is an upper triangular matrix (with  $l_{ii} = 1$ ) and  $D$  is a positive diagonal matrix.*

*Proof.* This follows by the construction made in proof 1.4 for  $k = 2, \dots, n-1$ . Thus we get  $L = L_1^{-1} L_2^{-1} \dots L_{n-1}^{-1}$  and  $D$  as diagonal matrix of the pivot elements. □

Because all diagonal elements  $d_{ii}$  are positive definite, there exists

$$D^{\frac{1}{2}} := \text{diag} \left( \sqrt{d_{11}}, \dots, \sqrt{d_{nn}} \right)$$

**Corollary 1.8.** *There exists some  $\bar{L} := LD^{\frac{1}{2}}$  (lower triangular matrix) such that  $A = \bar{L}\bar{L}^T$ .*

**Algorithm 3** (Cholesky decomposition).

- For  $k = 1, \dots, n$ 
  - $l_{kk} := (a_{kk} - \sum_{j=1}^{k-1} l_{kj}^2)^{\frac{1}{2}}$
  - For  $i = k+1, \dots, n$ 
    - \*  $l_{ik} := \frac{a_{ik} - \sum_{j=1}^{k-1} l_{ij} l_{kj}}{l_{kk}}$

## 2 Errors and matrix conditioning

So far, we use input data  $(A, b)$  to retrieve the result  $A^{-1}b$ .

Consider an abstract problem characterized by  $(f, x)$  with given map  $f$  and input data  $x$ . Thus in theory, we have the input, apply the algorithm and retrieve the result. But in practice, we have input with some error, an error in the algorithm and an error in the result.

In the following, we investigate input errors, which we cannot avoid in the general case. In the best case we can change the problem task. This refers to *conditioning of the problem*.

The algorithm triggers errors, which we can reduce or avoid by adapting the algorithm. This refers to *stability of the algorithm*.

### 2.1 Number representation and rounding errors

**Definition 2.1.** *Let  $x$  be a real number.  $\tilde{x} \in \mathbb{R}$  is an approximating value for  $x$ . Then  $\tilde{x} - x$  is called absolute error of  $\tilde{x}$  and, assuming  $x \neq 0$ ,  $\frac{\tilde{x} - x}{x}$  is the relative error of  $x$ .*

Even if we know the input know exactly, the representation of continuous numbers yields rounding errors.

A number  $x \in M$  is represented as  $x = \text{sign}(x) \cdot a \cdot E^{e-k}$ . The number system  $M$  is defined by the following four integer parameters:

**basis**  $E \in \mathbb{N}$  with  $E > 1$ ; mostly  $E = 2$

**precision**  $k \in \mathbb{N}$

**exponent**  $e$  in domain  $e_{\min} \leq e \leq e_{\max}$  where  $e_{\min}, e_{\max} \in \mathbb{Z}$

**mantissa**  $a \in \mathbb{N}_0$  is defined as

$$0 \leq a = a_1 E^{k-1} + a_2 \cdot E^{k-2} + \dots + a_{k-1} E^1 + a_k E^0 \leq E^k - 1$$

The mantissa length  $k$  and  $a_i$  are numbers of the number system, so 0 or 1 in case  $E = 2$ . If  $x \neq 0$ , then we require, that the first number is non-zero, then

$$E^{k-1} \leq q < E^k \quad \text{if } x \neq 0$$

We call  $x$  *k-digit normalized floating point number with basis E*.

Calculating with these numbers is called *calculation with k significant digits*. This way, we get the domain of normalized floating points  $x \neq 0$ :

$$E^{e_{\min}-1} \leq |x| \leq E^{e_{\max}}(1 - E^{-k})$$

**Example 2.2.** Let  $k = 3$ ,  $E = 2$ ,  $e_{\min} = -1$  and  $e_{\max} = 3$ . Let  $x > 0$ . Then the smallest number is 0.25 and the largest number is 7. Now look at the distribution of numbers. There are 3 numbers between 1 and 2.

**Remark.** The numbers are not equidistantly distributed.

After every arithmetic operation, the result  $x$  is rounded to a uniquely defined value  $\text{rd}(x)$ . We define the *machine precision*  $\varepsilon := \frac{E}{2} E^{-k}$ .

**Lemma 2.3.** If  $x \neq 0$  is within the domain of normalized floating point values and  $\text{rd}(x) \in M$ , then

$$|\text{rd}(x) - x| \leq \frac{E^{e-k}}{2}$$

is the maximum absolute error with respect to rounding.

$$\frac{|\text{rd}(x) - x|}{|x|} \leq \frac{1}{2} E^{1-k} = \varepsilon$$

is the maximum relative error with respect to rounding.

*Proof.* Without loss of generality, let  $x > 0$ . Then we can denote,

$$x = \mu E^{e-k} \quad E^{k-1} \leq \mu \leq E^k - 1$$

So,  $x$  lies in between the neighboring floating point numbers  $x_1 = \lfloor \mu \rfloor E^{e-k}$  and  $x_2 = \lceil \mu \rceil E^{e-k}$ . By this, we either get  $\text{rd}(x) = x_1$  or  $\text{rd}(x) = x_2$ .

Rounding gives

$$|\text{rd}(x) - x| \leq \frac{x_2 - x_1}{2} \leq \frac{E^{e-k}}{2}$$

Thus, it follows that

$$\frac{|\text{rd}(x) - x|}{|x|} \leq \frac{\frac{E^{e-k}}{2}}{\mu E^{e-k}} \leq \frac{1}{2} E^{1-k} = \varepsilon$$

□

| precision | k  | $e_{\min}$ | $e_{\max}$ | $\varepsilon$                      |
|-----------|----|------------|------------|------------------------------------|
| single    | 24 | -125       | 128        | $2^{-24} \approx 6 \cdot 10^{-8}$  |
| double    | 53 | -1021      | 1024       | $2^{-53} \approx 1 \cdot 10^{-16}$ |
| extended  | 64 | -16381     | 16384      | $2^{-64} \approx 5 \cdot 10^{-20}$ |

Table 1: Examples for common number systems. Assuming basis  $E = 2$

- Lemma 2.4.**
1. It holds that  $\text{rd}(x) = x(1 + \delta)$  with  $|\delta| \leq \varepsilon$
  2. Let  $\circ$  be one of the elementary operations  $\{+, -, \cdot, /\}$ . Then  $\text{rd}(x \circ y) = (x \circ y)(1 + \delta)$  with  $|\delta| \leq \varepsilon$ .
  3. The machine precision  $\varepsilon$  is the smallest positive number  $g$  for which  $\text{rd}(1 + g) > 1$  is true.

**Definition 2.5.**

1. The significant digits of a number are mantissa with normalized floating point representation.
2. When calculating with  $k$  significant digits, all input values must be rounded to  $k$  significant digits. Followingly the results of every elementary operation will be rounded to  $k$  significant digits before continuing. Thus the result is also affected.
3. Elimination is the cancellation of leading mantissa in the subtraction of two numbers of same sign.

**Example 2.6.** Depending on the chosen expressions, we can get different results. One example is:

$$99 - 70\sqrt{2} = \sqrt{9801} - \sqrt{9800} = \frac{1}{\sqrt{9801} + \sqrt{9800}} \approx 0.05050633883346584 \dots$$

## 2.2 Condition of a problem

**Question:** What effects do deviations in input data have independent of the chosen algorithm?

|          | $99 - 70\sqrt{2}$ | $\sqrt{9801} - \sqrt{9800}$ | $\frac{1}{\sqrt{9801} - \sqrt{9800}}$ |
|----------|-------------------|-----------------------------|---------------------------------------|
| $k = 2$  | 1                 | 0                           | 0.0050                                |
| $k = 4$  | 0.02000           | 0.02000                     | 0.005051                              |
| $k = 10$ | 0.005050660000    | 0.005050660000              | 0.005050633884                        |

Table 2: Difference in result depending on the expression. Apparently elimination occurs for subtraction  $99 - 70\sqrt{2}$  yielding ill-conditioned results

We already observed, that no difference occurs between input  $x$  and all inputs  $\tilde{x}$  with an absolute error smaller than machine precision. Therefore, consider the input set  $E$ , which contains all disrupted inputs  $\tilde{x}$

$$E = \{\tilde{x} \in \mathbb{R} \mid |\tilde{x} - x| < \varepsilon \varepsilon \cdot |x|\}$$

The map  $f$  (description of the problem) maps an input set  $E$  in a resulting set  $U = f(E)$ .

**Example** (Intersection of two lines). *If we compute the intersection of lines with angles close to  $0^\circ$ , the corresponding LES is ill-conditioned. Computing intersection for almost orthogonal lines is well-conditioned.*

↓ This lecture took place on 2018/10/11.

We consider the problem  $(f, X)$  given by the map  $f : U \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$  with  $U \subset \mathbb{R}^n$  open. Let  $x \in U$  and  $\delta$  as (relative or absolute) precision of the input data be given.

Distinction between:

- $\|\tilde{x} - x\| \leq \delta$  (absolute) or  $\|\tilde{x} - x\| \leq \delta \|x\|$  (relative) for some norm  $\|\cdot\|$  on  $\mathbb{R}^n$
- $|\tilde{x}_i - x_i| \leq \delta$  (absolute) or  $|\tilde{x}_i - x_i| \leq \delta |x_i|$  (relative), hence componentwise for  $i = 1, \dots, n$ .

### 2.3 Normwise condition analysis

Assumption: The input error  $\delta$  is sufficiently small, so we use *linearized error theory* for the asymptotic behavior of  $\delta \rightarrow 0$ .

**Notation:** Two functions  $g, h : \mathbb{R}^n \rightarrow \mathbb{R}^m$  are equal in the *first approximation* or *leading approximation* for  $x \rightarrow x_0$ .

$$g(x) \doteq h(x) \text{ for } x \rightarrow x_0$$

if  $g(x) = h(x) + \mathcal{O}(\|h(x)\|)$  for  $x \rightarrow x_0$ . The Landau symbol  $o(\|h(x)\|)$  for  $x \rightarrow x_0$  denotes a function  $\varphi$  such that

$$\lim_{x \rightarrow x_0} \frac{\|\varphi(x)\|}{\|h(x)\|} = 0$$

If  $f$  is differentiable in  $x$ , we have  $f(\tilde{x}) - f(x) \doteq f'(x)(\tilde{x} - x)$  for  $\tilde{x} \rightarrow x$ . Analogously, “ $g(x) \leq h(x)$  for  $x \rightarrow x_0$ ” (componentwise).

**Definition 2.7.** The absolute normwise condition of problem  $(f, X)$  is the smallest number  $\kappa_{\text{abs}} \geq 0$ , such that

$$\|f(\tilde{x}) - f(x)\| \leq \kappa \|\tilde{x} - x\| \text{ for } \tilde{x} \rightarrow x$$

The problem  $(f, X)$  is ill-posed, if there is no such number (formally  $\kappa_{\text{abs}} = \infty$ ).

The relative normwise condition of  $(f, X)$  is the smallest number  $\kappa_{\text{rel}} \geq 0$  such that

$$\frac{\|f(\tilde{x}) - f(x)\|}{\|f(x)\|} \leq \kappa_{\text{rel}} \frac{\|\tilde{x} - x\|}{\|x\|} \text{ for } \tilde{x} \rightarrow x$$

**Remark.**  $\kappa_{\text{abs}}$  give the amplification of the absolute error,  $\kappa_{\text{rel}}$  is the relative error.

If  $f$  is differentiable in  $x$ , then

$$\kappa_{\text{abs}} = \|f'(x)\|, \kappa_{\text{rel}} = \frac{\|x\|}{\|f(x)\|} \|f'(x)\|$$

where  $\|f'(x)\|$  is the number of the Jacobian matrix  $f'(x) \in \mathbb{R}^{m \times n}$  in regards of the norm

$$\|A\| := \sup_{x \neq 0} \frac{\|Ax\|}{\|x\|} = \sup_{\|x\|=1} \|Ax\| \text{ for } A \in \mathbb{R}^{m \times n}$$

**Example 2.8.**

$$\text{Addition: } f : \mathbb{R}^2 \rightarrow \mathbb{R} \quad \begin{pmatrix} a \\ b \end{pmatrix} \mapsto f(a, b) = a + b$$

with the derivative  $f'(a, b) = (1, 1) \in \mathbb{R}^{1 \times 2}$ . For the 1-norm on  $\mathbb{R}^2$  it holds that

$$\left\| \begin{pmatrix} a \\ b \end{pmatrix} \right\| = |a| + |b|$$

The norm (associated operator norm) is

$$\|f'(a, b)\| = \sup_{\|x\|_1=1} \left| (1 \ 1) \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \right| = 1$$

We retrieve the condition numbers of addition as

$$\kappa_{\text{abs}} = 1 \text{ and } \kappa_{\text{rel}} = \frac{|a| + |b|}{|a + b|}$$

Consequence:

- $\kappa_{\text{rel}} = 1$  for addition of two numbers with same sign
- subtraction of two almost equal numbers given  $|a + b| \ll |a| + |b| \implies \kappa_{\text{rel}} \gg 1$ .

**Example 2.9** (Condition of a linear equation system  $Ax = b$ ). If we only consider  $b \in \mathbb{R}^n$  as input data, the problem is specified by the linear map  $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$  and  $b \mapsto f(b) = A^{-1}b$ . The derivative is  $f'(b) = A^{-1}$  and we retrieve the condition numbers

$$\kappa_{\text{abs}} = \|A^{-1}\| \text{ and } \kappa_{\text{rel}} = \frac{\|b\|}{\|A^{-1}b\|} \|A^{-1}\| = \frac{\|Ax\|}{\|x\|} \|A^{-1}\|$$

What about the deviations in  $A$ ? For this, consider  $A$  as input data

$$f : \mathbb{R}^{n \times n} \supset \text{GL}(n) \rightarrow \mathbb{R}^n, A \mapsto f(A) = A^{-1}b$$

for some fixed  $b \in \mathbb{R}^n$ . The map  $f$  is non-linear, but differentiable.

**Lemma 2.10.** The map  $g : \mathbb{R}^{n \times n} \supset \text{GL}(n) \rightarrow \text{GL}(n)$  with  $g(A) = A^{-1}$  is differentiable and  $g'(A) \cdot C = -A^{-1}CA^{-1}$  for all  $C \in \mathbb{R}^{n \times n}$ .

*Proof.* Consider  $I = (A + tC)(A + tC)^{-1}$  for sufficiently small  $t$  and derive by  $t$ :

$$0 = C(A + tC)^{-1} + (A + tC) \frac{d}{dt} (A + tC)^{-1}$$

For  $t = 0$ , it especially follows that  $g'(A)C = \frac{d}{dt} (A + tC)^{-1} \big|_{t=0} = -A^{-1}CA^{-1}$ .  $\square$

By Lemma 2.10, by the derivative of the solution  $f(A) = A^{-1}b$  to  $A$ ,

$$f'(A)C = -A^{-1}CA^{-1}b = -A^{-1}Cx \text{ for } C \in \mathbb{R}^{n \times n}$$

We retrieve the condition numbers

$$\kappa_{\text{abs}} = \|f'(A)\| = \sup_{\|C\|=1} \|A^{-1}Cx\| \leq \|A^{-1}\| \|x\|$$

$$\kappa_{\text{rel}} = \frac{\|Ax\|}{\|x\|} \|f'(A)\| \leq \|A\| \|A^{-1}\|$$

By the sub-multiplicativity,  $\|Ax\| \leq \|A\| \|x\|$ , of the matrix norm, it also holds for the relative condition in regards of the input  $b$ , that  $\kappa_{\text{rel}} \leq \|A\| \cdot \|A^{-1}\|$ . The value  $\kappa(A) := \|A\| \|A^{-1}\|$  is called *condition number* of matrix  $A$ .

↓ This lecture took place on 2018/10/16.

We can show that for invertible matrices, we have  $\kappa(A) = \frac{\max_{\|x\|=1} \|Ax\|}{\min_{\|x\|=1} \|Ax\|}$  depending on the choice of the norm. By this it follows that:

1.  $\kappa(A) \geq 1$
2.  $\kappa(\alpha A) = \kappa(A) \forall \alpha \in \mathbb{R}, \alpha \neq 0$   
hence the condition number is invariant under scalar transformations (in contrast to determinants).  
Remember that  $\det(\alpha A) = \alpha^n \det(A)$ .
3.  $A \in \mathbb{R}^{n \times n}$  with  $A \neq 0$  is singular iff  $\kappa(A) = \infty$

## 2.4 Componentwise condition analysis

The normwise condition analysis is inappropriate in particular special cases.

**Example 2.11.** Consider the equation system  $Ax = b$  with a diagonal matrix

$$A = \begin{pmatrix} 1 & 0 \\ 0 & \varepsilon \end{pmatrix} \quad A^{-1} = \begin{pmatrix} 1 & 0 \\ 0 & \frac{1}{\varepsilon} \end{pmatrix}$$

The problem is completely decoupled and we expect a well-conditioned problem (at least for diagonal deviations). For normwise conditioning (in regards of  $\|\cdot\|_\infty$ ), it holds that

$$\kappa_\infty(A) = \|A^{-1}\|_\infty \|A\|_\infty = \frac{1}{\varepsilon} \quad \varepsilon \rightarrow 0$$

For  $\varepsilon \rightarrow 0$  the condition number will be arbitrary large, because it permits arbitrary deviations in the matrix.

By componentwise condition analysis, we get  $(f, x) \rightsquigarrow$  smallest number  $\kappa_{\text{rel}}$  such that

$$\max_i \frac{|f_i(\tilde{x}) - f_i(x)|}{|f_i(x)|} \leq \max_i \frac{|\tilde{x}_i - x_i|}{|x_i|}$$

## 2.5 Stability of an algorithm

We now consider error, that are created, if the map  $f$  is approximated by some map  $\tilde{f}$  (algorithmic realization). It encompasses all rounding errors and approximation errors and return an approximation  $\tilde{f}(x)$  instead of  $f(x)$ .

**Question:** Is  $\tilde{f}(x)$  acceptable as an substitute of  $f(x)$ ? Compare with Figure 1.

**Definition 2.12** (Forward stability). Let  $\tilde{f}$  be the floating point implementation of an algorithm to solve problem  $(f, x)$  of relative normwise condition  $\kappa_{\text{rel}}$ . The stability indicator of a normwise forward analysis is the smallest number  $\sigma \geq 0$ , such that for all  $\tilde{x} \in E$  it holds that

$$\frac{\|\tilde{f}(\tilde{x}) - f(\tilde{x})\|}{\|f(\tilde{x})\|} \leq \sigma \kappa_{\text{rel}} \varepsilon \quad \text{for } \varepsilon \rightarrow 0$$

## 2.6 Stability of the componentwise forward analysis

$\bar{\kappa}_{\text{rel}}$  is the smallest number  $\sigma \geq 0$  such that

$$\max_i \frac{|f_i(\tilde{x}) - f_i(x)|}{|f_i(x)|} \leq \sigma \cdot \bar{\kappa}_{\text{rel}} \varepsilon \quad \text{for } \varepsilon \rightarrow 0$$

where  $\bar{\kappa}_{\text{rel}}$  is the componentwise relative conditioning of  $(f, x)$ . We call algorithm  $\tilde{f}$  stable in regards of forward analysis if  $\sigma$  is smaller than the number of consecutively executed elementary operations.

**Lemma 2.13.** For the elementary operations  $\{+, -, \cdot, /\}$  and the floating point implementations  $\{\hat{+}, \hat{-}, \hat{\cdot}, \hat{/}\}$  it holds that  $\sigma \kappa_{\text{rel}} \leq 1$ .



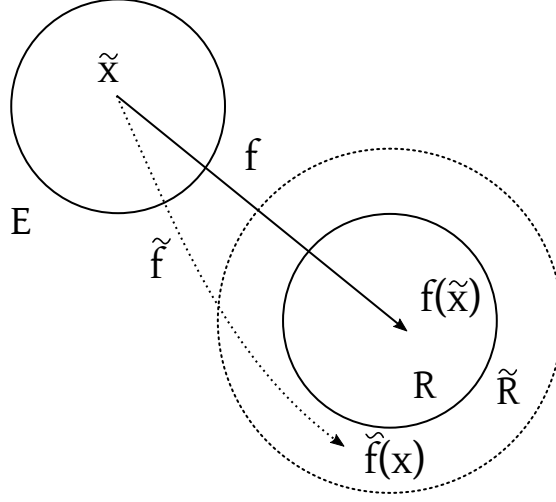


Figure 1: Map of the error with  $\tilde{f}$  instead of  $f$

*Proof.* For every floating point implementation it holds that

$$x \hat{*} y = (x * y)(1 + \delta) \quad * \in \{+, -, \cdot, /\}$$

by Lemma 2.4 for some  $\delta$  with  $|\delta| \leq \varepsilon$ . Thus,

$$\frac{|x \hat{*} y - x * y|}{x * y} = \frac{|(x * y)(1 + \delta) - x * y|}{|x * y|} = |\delta| \leq 1 \cdot \varepsilon$$

□

This approach is rather impractical (determination of the conditioning number of the problem is difficult).

**Idea:** interpret the error in the algorithm like an input error. The result  $\tilde{y} = \tilde{f}(\tilde{x})$  is considered as the exact result  $\tilde{y} = f(\hat{x})$  for distorted data  $\hat{x}$ .  $E = \{\tilde{x} \mid \|\tilde{x} - x\| \leq \varepsilon \|x\|\}$ . This is only possible if  $\tilde{y}$  is a admissible solution of  $f$  at all. In the other case, we call the algorithm *unstable*.

**Definition 2.14.** The normwise backwards error of the algorithm  $\tilde{f}$  to solve the problem  $(f, x)$  is the smallest number  $\eta \geq 0$ , for which  $\forall \tilde{x} \in E_\varepsilon$  some  $\hat{x}$  exists such that

$$\frac{\|\hat{x} - \tilde{x}\|}{\|\tilde{x}\|} \leq \eta \quad \text{for } \varepsilon \rightarrow 0$$

The componentwise backwards error is defined as

$$\max_i \frac{|\hat{x}_i - \tilde{x}_i|}{|\tilde{x}_i|} \leq \eta \quad \text{for } \varepsilon \rightarrow 0$$

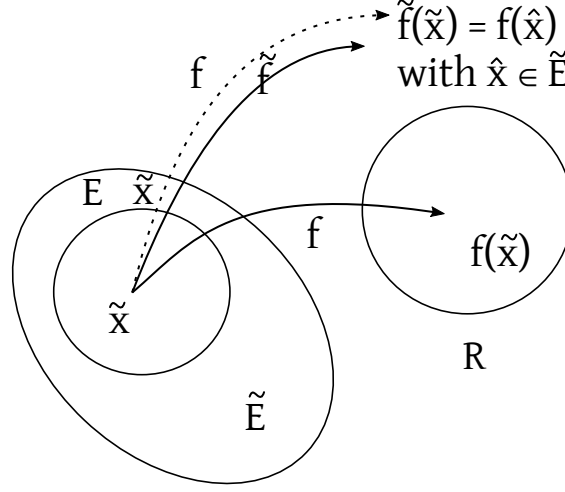


Figure 2: Backwards analysis

The algorithm is called *stable* (in regards of backwards analysis) with respect to the relative input error  $\varepsilon$  if  $\eta < \varepsilon \cdot \text{number of elementary operations}$ . For given  $\varepsilon$  we define a stability indicator of the backwards analysis as quotient

$$\sigma_R := \frac{\eta}{\varepsilon}$$

**Lemma 2.15.** For stability indicators  $\sigma$  and  $\sigma_R$ , of the forward/backward analysis respectively, it holds that  $\sigma \leq \sigma_R$ . Especially by backwards stability, forwards stability follows.

*Proof.* By definition of the backwards error, for every  $\tilde{x} \in E$  there exists some  $\hat{x}$  such that  $f(\hat{x}) = \tilde{f}(\tilde{x})$  and

$$\frac{\|\hat{x} - \tilde{x}\|}{\|\tilde{x}\|} \leq \eta = \sigma_R \cdot \varepsilon \quad \text{for } \varepsilon \rightarrow 0$$

For the relative error in the result, we have

$$\frac{\|\tilde{f}(\tilde{x}) - f(\tilde{x})\|}{\|f(\tilde{x})\|} = \frac{\|f(\hat{x}) - f(\tilde{x})\|}{\|f(\tilde{x})\|} \leq \kappa_{\text{rel}} \frac{\|\hat{x} - \tilde{x}\|}{\|\tilde{x}\|} \leq \kappa_{\text{rel}} \cdot \sigma_R \cdot \varepsilon \quad (\text{for } \varepsilon \rightarrow 0)$$

□

**Remark.**  $Ax = b$  with solution  $\tilde{x}$ , approximation  $\hat{x}$ . Error  $\|\hat{x} - \tilde{x}\|$  is not computable, because  $\tilde{x}$  is unknown. Residue  $\|b - A\hat{x}\|$ .

## 2.7 Backwards analysis for linear equation systems

**Question:** Is the LU decomposition a backwards stable algorithm?

In the following, we will use the following notation for vectors and matrices:

$$|A| \leq |B| \iff |a_{ij}| \leq |b_{ij}| \forall i, j$$

**Lemma 2.16.** *The floating point implementation of the scalar product*

$$\langle x, y \rangle := x_n y_n + \langle x^{n-1}, y^{n-1} \rangle$$

with  $x^{n-1} := (x_1, \dots, x_{n-1})^T$  and  $y^{n-1} := (y_1, \dots, y_{n-1})^T$  denotes some solution  $\langle x, y \rangle_{\text{rd}} := \text{rd}(\langle x, y \rangle)$  for some  $x, y \in \mathbb{R}^n$  such that  $\langle x, y \rangle_{\text{rd}} = \langle \hat{x}, y \rangle$  for some  $\hat{x} \in \mathbb{R}^n$  with  $|x - \hat{x}| \leq n\varepsilon |x|$ , hence the relative componentwise backwards error is  $\eta \leq n \cdot \varepsilon$  and the scalar product is stable in regards of backwards analysis (with  $2n - 1$  elementary operations by the scalar product).

*Proof.* We need to prove:

$$\forall i : |x_i - \hat{x}_i| \leq \underbrace{n\varepsilon |x_i|}_{h(\varepsilon)} + \overbrace{\sigma(\|h(\varepsilon)\|)}^{\sigma} \quad \text{for } \varepsilon \rightarrow 0 \text{ with } \lim_{\varepsilon \rightarrow 0} \frac{\|\varphi(\varepsilon)\|}{\|h(\varepsilon)\|} = 0$$

Thus, we look for a function that converges superlinear towards 0 for  $\varepsilon \rightarrow 0$  ( $\mathcal{O}(\varepsilon)$ ).

We use induction over  $n$  to prove this. For  $n = 1$ , the statement is true by Lemma 2.4. Now let  $n > 1$  and let the statement be true for  $n - 1$ . For floating point implementation of recursion it is true that

$$\langle x, y \rangle_{\text{rd}} = (x_n y_n (1 + \delta) + \langle x^{n-1}, y^{n-1} \rangle_{\text{rd}}) (1 + \tilde{\delta})$$

where  $|\delta| \leq \varepsilon$  and  $|\tilde{\delta}| \leq \varepsilon$  characterize the relative error of multiplication and addition respectively.

By the induction hypothesis, it holds that  $\langle x^{n-1}, y^{n-1} \rangle_{\text{rd}} = \langle z, y^{n-1} \rangle$  for some  $z \in \mathbb{R}^{n-1}$  with

$$|x^{n-1} - z| \leq (n-1)\varepsilon |x^{n-1}| + \mathcal{O}(\varepsilon)$$

By  $\hat{x}_n = x_n(1 + \delta)(1 + \tilde{\delta})$  and  $\hat{x}_k = z_k(1 + \tilde{\delta})$  by  $k = 1, \dots, n-1$  it follows that

$$\begin{aligned} \langle x, y \rangle_{\text{rd}} &= x_n y_n (1 + \delta)(1 + \tilde{\delta}) + \langle z(1 + \varepsilon), y^{n-1} \rangle \\ &= \hat{x}_n y_n + \langle \hat{x}^{n-1}, y^{n-1} \rangle = \langle \hat{x}, y \rangle \\ &\quad \text{with } |x_n - \hat{x}_n| \leq 2\varepsilon |x_n| + |x_n| |y_n| \delta \tilde{\delta} \\ &\quad \leq n\varepsilon |x_n| + \mathcal{O}(\varepsilon) \\ \text{and } |x_k - \hat{x}_k| &\leq |x_k - z_k| + |x_k - \hat{x}_k| \\ &\leq (n-1)\varepsilon |x_k| + \mathcal{O}(\varepsilon) + \varepsilon |z_k| \\ &\leq n\varepsilon |x_k| + \mathcal{O}(\varepsilon) \text{ for } k = 1, \dots, n-1 \end{aligned}$$

The last estimate can be made using

$$\begin{aligned} |z_k| - |x_k| &\leq |x_k - z_k| \leq (n-1)\varepsilon |x_k| + \mathcal{O}(\varepsilon) \\ \implies \varepsilon |z_k| &\leq \varepsilon |x_k| + \mathcal{O}(\varepsilon) \end{aligned}$$

□

**Theorem 2.17.** *The floating point implementation of forward substitution to solve a linear equation system  $Lx = b$  with some lower triangular matrix  $L$  computes the solution  $\hat{x}$  such that there exists some lower triangular matrix  $\hat{L}$  with  $\hat{L}\hat{x} = b$  and  $|L - \hat{L}| \leq n\varepsilon |L|$ . Hence for componentwise relative backwards errors, we have  $\eta \leq n \cdot \varepsilon$  and the forward substitution is stable in regards of backwards analysis.*

*Proof.* Again, we will use a recursive approach.

$$l_{kk}x_k = b_k - \langle l^{k-1}, x^{k-1} \rangle$$

For  $k = 1, \dots, n$ , we have

$$x^{k-1} := (x_1, \dots, x_{k-1})^T \text{ and } l^{k-1} := (l_{k,1}, \dots, l_{k,k-1})^T$$

In floating point arithmetics, the following recursion is yielded,

$$l_{kk}(1 + \delta_k)(1 + \tilde{\delta}_k)\hat{x}_k = b_k - \langle l^{k-1}, x^{k-1} \rangle_{\text{rd}} (1 + \tilde{\delta}_k)$$

where  $\delta_k, \tilde{\delta}_k$  with  $|\delta_k| \leq \varepsilon$  and  $|\tilde{\delta}_k| \leq \varepsilon$  are the relative errors of multiplication/addition. By Lemma 2.16, it follows that

$$\langle l^{k-1}, \hat{x}^{k-1} \rangle_{\text{rd}} = \langle \tilde{l}^{k-1}, \hat{x}^{k-1} \rangle$$

for some vector  $\tilde{l}^{k-1} = (\tilde{l}_{k,1}, \dots, \tilde{l}_{k,k-1})^T$  with  $|l^{k-1} - \tilde{l}^{k-1}| \leq (k-1)\varepsilon |l^{k-1}|$ . Now let  $\hat{l}_{kk} := l_{kk}(1 + \delta_k)(1 + \tilde{\delta}_k)$  and  $\hat{l}^{k-1} := (1 + \tilde{\delta}_k)\tilde{l}^{k-1}$  then

$$\begin{aligned} |\hat{l}^{k-1} - l^{k-1}| &\leq (k-1)\varepsilon |l^{k-1}| + |\tilde{\delta}_k| |\tilde{l}^{k-1}| \\ &\leq k\varepsilon |l^{k-1}| \\ |\hat{l}_{kk} - l_{kk}| &\leq 2\varepsilon \\ \implies \hat{L}\hat{x} = b \wedge |L - \hat{L}| &\leq n\varepsilon |L| \end{aligned}$$

□

↓ This lecture took place on 2018/10/18.

**Lemma 2.18.** *Let  $A$  have a LU decomposition. Then Gaussian elimination determines  $\hat{L}$  and  $\hat{U}$  such that  $\hat{L}\hat{U} = \hat{A}$  for some matrix  $\hat{A}$  with*

$$|A - \hat{A}| \leq n |\hat{L}| |\hat{U}| \varepsilon$$

**Theorem 2.19.**  *$A$  has some LU-decomposition. Then Gaussian elimination for linear equation system  $Ax = b$  gives a solution  $\hat{x}$  with  $\hat{A}\hat{x} = b$  for some matrix  $\hat{A}$  with*

$$|A - \hat{A}| \leq 2n |\hat{L}| |\hat{U}| \varepsilon$$

**Theorem 2.20.** *The Gaussian elimination with column pivot strategy for the linear equation system  $Ax = b$  determines some  $\hat{x}$  such that  $\hat{A}\hat{x} = b$  for some matrix  $\hat{A}$  with*

$$\frac{\|A - \hat{A}\|_\infty}{\|A\|_\infty} \leq 2n^3 \rho_n(A) \varepsilon \quad \text{where } \rho_n(A) := \frac{\alpha_{\max}}{\max_{i,j} |a_{ij}|}$$

and where  $\alpha_{\max}$  is the largest absolute value of an element that occurs during elimination in the remainder matrices  $A^{(1)} = A$  to  $A^{(n)} = U$ .

$$A = \begin{bmatrix} 1 & 0 & \dots & 0 & 1 \\ -1 & \vdots & \ddots & \vdots & \vdots \\ \vdots & \ddots & \ddots & \vdots & \vdots \\ -1 & \dots & \dots & -1 & 1 \end{bmatrix}$$

For matrix  $A$  it holds that  $\rho_n(A) = 2^{n-1}$ .

Consider that the previous result also provides *no* positive statement about the backwards stability of the LU decomposition.

### 3 Linear least squares

**Goal:** Solution of overdetermined linear equation systems

$$Ax = b \quad A \in \mathbb{R}^{m \times n} \quad m > n \quad b \in \mathbb{R}^m$$

by the method of linear least squares.

For this, we are going to use (numerically stable) orthogonal transformations.

#### 3.1 Gaussian process of least squares

**Problem:** Given  $m$  input data points  $(t_i, b_i)$  and  $t_i, b_i \in \mathbb{R}$  for  $i = 1, \dots, m$ . They describe the values  $b_i$  of an object at timestamp  $t_i$ .

We assume that there exists an underlying law such that the dependency of  $b$  and  $t$  can be represented by some *modelling function*  $\varphi$  with  $b(t) = \varphi(t; x_1, \dots, x_n)$ . The modelling function therefore takes  $n$  unknowns  $p_i$ .

**Example 3.1** (Ohm's law). *Consider Ohm's law  $b = xt = \varphi(t; x)$  where  $t$  denotes the current,  $b$  the voltage and  $x$  resistance.*

So want to approximate data points in the  $\mathbb{R}^2$  plane linearly.

If the model would be *exact* and if you have no errors in measurement, then we are supposed to evaluate parameters  $x_1, \dots, x_n$  such that  $b_i = b(t_i) = \varphi(t_i; x_1, \dots, x_n)$  for  $i = 1, \dots, m$ .

In reality, measurements are subject to errors and modelling functions are also just approximations of reality. We therefore require that  $b_i \approx \varphi(t_i; x_1, \dots, x_n)$  for  $i = 1, \dots, m$ . We weight the individual deviations

$$\Delta i := b_i - \varphi(t_i; x_1, \dots, x_n) \quad i = 1, \dots, m$$

By Gauss, determine  $x_i$  such that  $\sum_{i=1}^m \Delta_i^2$  becomes minimal. We use the short notation  $\Delta^2 := \sum_{i=1}^m \delta_i^2 \implies \min$ . Often it is more useful to weight the individual deviations differently. For some measurement precision  $\delta b$ , introduce weights

$$\sum_{i=1}^m \left( \frac{\Delta_i}{\delta b_i} \right) \rightarrow \min$$

We consider the special case of linear (with respect to  $x$ ) modelling functions  $\varphi$ , hence  $\varphi(t_i; x_1, \dots, x_n) = a_1(t)x_1 + \dots + a_n(t)x_n$  where  $a_1, \dots, a_n : \mathbb{R} \rightarrow \mathbb{R}$  are arbitrary functions. Furthermore we choose  $\|\cdot\| = \|\cdot\|_2$ . We can denote the least square problem as

$$\|b - Ax\| \rightarrow \min \quad \text{where } b = \begin{pmatrix} b_1 \\ \vdots \\ b_m \end{pmatrix}, x = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$$

$$\text{and } A = (a_{ij}) \in \mathbb{R}^{m \times n} \text{ with } a_{ij} := a_j(t_i)$$

## Orthogonal equations

We look for some point  $z = Ax$  of the image space  $U(A)$  of  $A$ , that has the smallest distance to some given  $b$ .  $Ax$  is the orthogonal projection of  $b$  to the subspace  $U(A) = \text{image}(A)$ .

**Theorem 3.2.** Let  $V$  be a vector space of  $\dim V < \infty$  with scalar product  $\langle \cdot, \cdot \rangle$ . Let  $U \subset V$  be a subspace. Let

$$U^\perp = \{v \in V \mid \langle v, u \rangle = 0 \forall u \in U\}$$

be the orthogonal complement in  $V$ . Then for all  $v \in V$  with respect to the norm induced by  $\|v\| = \sqrt{\langle v, v \rangle}$  such that

$$\|v - u\| = \min_{\tilde{u} \in U} \|v - \tilde{u}\| \iff v - u \in U^\perp$$

Proof sketch:

- $V = U \oplus U^\perp$
- $v = u + w, u \in U, w \in U^\perp$  unique
- $\tilde{u} \in U$  arbitrary:

$$\|v - \tilde{u}\|^2 = \|(v - u) + (u - \tilde{u})\|^2 = \|v - u\|^2 + \|u - \tilde{u}\|^2 \geq \|v - u\|^2$$

with  $v - u = w$  and  $(u - \tilde{u}) \in U$ .

↓ This lecture took place on 2018/10/23.

The unique solution  $u \in U$  of  $\min \|v - u\|$  is called *orthogonal projection* of  $V$  to  $U$ . The map  $P : V \rightarrow U$  with  $v \mapsto Pv$  with  $Pv = u$  is linear and is called *orthogonal projection* of  $V$  to  $U$ .

**Remark 3.3.** We get an analogous statement if we replace  $U$  by an affine subspace  $W = w_0 + U \subseteq V$  where  $w_0 \in V$ .  $U$  is the parallel, to  $W$ , linear subspace of  $V$ .

For all  $v \in V, w \in W: \|v - w\| = \min_{\tilde{w} \in W} \|v - \tilde{w}\| \iff v - w \in U^\perp$ .

**Theorem 3.4.** The vector  $x \in \mathbb{R}^n$  is a solution of the linear least squares problems  $\min \|b - Ax\| \iff x$  satisfies the so-called orthogonal equations  $A^T Ax = A^T b$  (dt. "Normalgleichungen").

Especially, the linear least squares problem is uniquely solvable iff  $A$  has full column rank.

*Proof.* By Theorem 3.2,  $V := \mathbb{R}^n, U := \text{image}(A) = \text{rank}(A)$ .

$$\begin{aligned} \implies \max \|b - Ax\| &\iff \langle b - Ax, A\tilde{x} \rangle = 0 \forall \tilde{x} \in \mathbb{R}^n \\ &\iff 0 = \langle A^t(b - Ax), \tilde{x} \rangle \forall \tilde{x} \in \mathbb{R}^n \\ &\iff A^t(b - Ax) = 0 \\ &\iff A^t b = A^t Ax \end{aligned}$$

Thus the first statement is proven. How about the second statement?

$$A^t A \text{ regular} \iff \text{rank}(A) = n$$

□

**Remark** (Geometric interpretation).  $b - Ax$  is orthogonal to  $\text{range}(A)$  with  $A \subseteq \mathbb{R}^n$ .

**Solution of the orthogonal equation system:** Theoretically, we can create  $A^T A$  and retrieve a symmetric positive definite matrix ( $\implies$  decomposition with Cholesky process).

**Remark** (Computing time). 1. Determination of  $A^t A$  (only one half, because of symmetry)

$$(A^t A)_{ij} = (\langle a_j, a_j \rangle)_{i,j} \quad a_j = j\text{-th column of } A \rightarrow \frac{1}{2} n^2 n$$

$$2. \text{ Cholesky decomposition } \sim \frac{1}{6} n^3$$

In case of  $m \gg n$ , the first approach will take over  $\implies \sim \frac{1}{2} n^2 n. \frac{2}{3} m^3$  for  $n \approx m$ .

**Lemma 3.5.** For some matrix  $A \in \mathbb{R}^{m \times n}$  with  $m \geq n$  and full column rank ( $\text{rank}(A) = n$ ),  $\kappa_2(A^t A) = (\kappa_2(A))^2$ .

*Proof.* By definition of condition numbers:

$$\begin{aligned} (\kappa_2(A))^2 &= \frac{\max_{\|x\|=1} \|Ax\|_2^2}{\min_{\|x\|=1} \|Ax\|_2^2} \\ &= \frac{\max_{\|x\|=1} \langle Ax, Ax \rangle}{\min_{\|x\|=1} \langle Ax, Ax \rangle} \\ &= \frac{\max \langle A^t Ax, x \rangle}{\min \langle A^t Ax, x \rangle} \quad \text{"Rayleigh eigenvalue"} \\ &\stackrel{A^t \text{ Asym.}}{=} \frac{\lambda_{\max}(A^t A)}{\lambda_{\min}(A^t A)} \\ &\stackrel{A^t \text{ Asym.}}{=} \kappa_2(A^t A) \end{aligned}$$

□

Recognize that the amplification of the error for the orthogonal equation is larger. Therefore it makes sense to avoid computing  $A^T A$  and use the original matrices  $A$  and  $A^T$  instead, e.g.

$$\begin{aligned} \begin{pmatrix} I & A \\ A^T & 0 \end{pmatrix} \begin{pmatrix} r \\ 0 \end{pmatrix} &= \begin{pmatrix} b \\ 0 \end{pmatrix} \\ \implies A^T r &= 0 \wedge r \neq Ax = b \\ \implies r &= b - Ax \quad \implies 0 = A^T r = A^T(b - Ax) \end{aligned}$$

$$\|A_3\|_2^2 = \max_{\|x\|_2=1} \|Q_j x\|_2^2 = \max_{\|x\|_2=1} \langle x, Q^T Q x \rangle = 1$$

Recognize that  $Q^T Q = I$ .

### 3.2 Orthogonalization process

We can represent the elimination process for linear equation systems (e.g. Gaussian elimination) formally as

$$A \xrightarrow{f_1} B_1 A \xrightarrow{f_2} B_2 B_1 A \xrightarrow{f_3} \dots \xrightarrow{f_k} B_k \dots B_1 A \dots \rightarrow U$$

where

$$U = \begin{pmatrix} 1 & * \\ 0 & \varepsilon \end{pmatrix}$$

and where matrices  $B_j$  describe the operations applied to matrix  $A$ . Through this process, the condition numbers of the individual matrices can be amplified, such that instabilities can occur. If we instead use orthogonal transformations  $Q_j$  for elimination, we have

$$\kappa_2(Q_j) = \|Q_j\|_2 \cdot \|Q_j\|_2^{-1} = \|Q_j\|_2 \|Q_j\|_2 = 1$$

Orthogonalization process are necessarily stable, but have higher computing times compared to Gaussian elimination. If we assume we transformed a given matrix  $A \in \mathbb{R}^{m \times n}$  with  $m \geq n$  with some orthogonal matrix  $Q \in \mathbb{R}^{m \times m}$  into upper triangular matrix, then

$$Q^T A = \begin{bmatrix} * & \dots & x \\ & \ddots & \vdots \\ 0 & & x \\ \vdots & & * \\ 0 & & \vdots \\ & & 0 \end{bmatrix} = \begin{bmatrix} U \\ 0 \end{bmatrix}$$

**Theorem 3.6.** Let  $A \in \mathbb{R}^{m \times n}$ ,  $m \geq n$ , of maximum rank  $n$ ,  $b \in \mathbb{R}^m$  and  $Q \in \mathbb{R}^{m \times m}$  be an orthogonal matrix with

$$Q^T A = \begin{bmatrix} U \\ 0 \end{bmatrix} \text{ and } Q^T b = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}$$

where  $b_1 \in \mathbb{R}^n$ ,  $b_2 \in \mathbb{R}^{m-n}$  and  $U \in \mathbb{R}^{n \times n}$  is an (invertible) upper triangular matrix. Then  $x = U^{-1} b_1$  is the solution of the linear least squares problems  $\min \|b - Ax\|_2$ .



*Proof.* Because  $Q$  is orthogonal,  $\forall x \in \mathbb{R}^n$  it follows that

$$\|b - Ax\|_2^2 = \|Q^T(b - Ax)\|_2^2 = \left\| \begin{pmatrix} b_1 - Ux \\ b_2 \end{pmatrix} \right\|_2^2 = \|b_1 - Ux\|_2^2 + \|b_2\|_2^2 \geq \|b_2\|_2^2$$

Because  $\text{rank}(A) = \text{rank}(U) = n$ ,  $U$  is invertible. The term  $\|b_1 - Ux\|_2^2$  disappear for  $x = U^{-1}b_1$   $\square$

**Remark.**  $A: 0 \neq \|r\| = \|b - Ax\| = \|b_2\|$ .

What do orthogonal transformations look like? In general, these transformations are rotations and reflections.

e.g.  $x \mapsto \alpha e_1 = Qx$  with  $Q = \begin{bmatrix} \cos(\theta) & \sin(\theta) \\ -\sin(\theta) & \cos(\theta) \end{bmatrix}$ .

Alternatively,  $x \mapsto \alpha e_1 = x = 2 \frac{\langle v, x \rangle}{\langle v, v \rangle} v$  where  $v$  is collinear to the difference  $x = \alpha e_1$ .

## Givens rotations

Givens rotations are described by matrices of form

$$Q_{k,l} := \begin{bmatrix} I & & & \\ & c & & s \\ & & I & \\ & -s & & c \\ & & & & I \end{bmatrix} \in \mathbb{R}^{m \times m}$$

with proper dimension of the unit matrix and  $c^2 + s^2 = 1$ . For  $x \in \mathbb{R}^m$  it holds that

$$x \mapsto y = Q_{k,l}x \text{ with } y_i = (Q_{k,l}x)_i = \begin{cases} cx_k + sx_l & i = k \\ -sx_k + cx_l & i = l \\ x_i & \text{else} \end{cases}$$

Thus for  $A = [a_1, \dots, a_n] \in \mathbb{R}^{m \times n}$ , we get

$$Q_{k,l}A = [Q_{k,l}a_1, \dots, Q_{k,l}a_n]$$

Hence, only rows  $k$  and  $l$  of matrix  $A$  will be modified.

**Question:** How to choose  $c$  and  $s$  to eliminate a component  $x_l$  of  $x$ ? Because  $Q_{k,l}$  operates on the  $(k, l)$ -layer, we consider without loss of generality  $m = 2$ . By  $x_k^2 + x_l^2 \neq 0$  and  $s^2 + c^2 = 1$  it follows that

$$\begin{bmatrix} c & s \\ -s & c \end{bmatrix} \begin{bmatrix} x_k \\ x_l \end{bmatrix} = \begin{bmatrix} r \\ 0 \end{bmatrix}$$

$$\Leftrightarrow r = \pm \sqrt{x_k^2 + x_l^2}, c = \frac{x_k}{r}, s = \frac{x_l}{r}$$

Computing  $x_k^2$  might gives us an exponent overflow.

$$\tau := \frac{x_k}{x_l}, s := \frac{1}{\sqrt{1 + \tau^2}}, c := s\tau \text{ if } |x_l| > |x_k|$$

and accordingly,

$$\tau := \frac{x_l}{x_k}, c = \frac{1}{\sqrt{1 + \tau^2}}, s := c\tau \text{ if } |x_k| \geq |x_l|$$

Now transform  $A$  iteratively into triangular form.

For example, consider

$$A = \begin{bmatrix} * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{bmatrix} \xrightarrow{(5,4)} \begin{bmatrix} * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ 0 & * & * & * \end{bmatrix} \xrightarrow{(4,3)} \dots$$

$$\xrightarrow{(2,1)} \begin{bmatrix} * & * & * & * \\ 0 & * & * & * \\ 0 & * & * & * \\ 0 & * & * & * \\ 0 & * & * & * \end{bmatrix} \xrightarrow{(5,4)} \dots \xrightarrow{(4,3)} \dots \xrightarrow{(5,4)} \begin{bmatrix} * & * & * & * \\ 0 & * & * & * \\ 0 & 0 & * & * \\ 0 & 0 & 0 & * \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

**Remark** (Computing time). *Computing time for a dense matrix  $A \in \mathbb{R}^{m \times n}$ :*

1.  $\sim \frac{n^2}{2}$  square roots and  $\frac{4}{3}n^3$  multiplications, if  $m \approx n$ .
2.  $m \cdot n$  are square root and  $2mn^2$  multiplications if  $m \gg n$ .

This is an alternative for Gaussian elimination. The algorithm is *more stable*, but *more expensive*.

But it can be advantageous for special matrix, e.g. Hessenberg matrices, which are “almost triangular” (so the first subdiagonal also contains non-zero entries). Here, we need  $n - 1$  Givens rotations if we want to achieve upper triangular form.

## Householder reflections

↓ This lecture took place on 2018/10/25.

For given  $0 \neq v \in \mathbb{R}^m$ , we define the matrix  $Q \in \mathbb{R}^{m \times m}$ .

$$Q = I - 2 \frac{vv^T}{v^T v}$$

They describe reflections on the mirror plane, which are vertical to  $v$ . The following properties hold:

1. symmetry:  $Q = Q^T$
2. orthogonality:  $QQ^T = I$

3.  $Q^2 = \text{id}$ , thus  $Q$  is *involutional*
4.  $\text{spec}(Q) = \{\pm 1\}$ ,  $-1$  simple,  $1$   $(n-1)$ -times.

For  $y \in \mathbb{R}^n$ ,

$$y \mapsto Qy = \left( I - 2 \frac{vv^T}{v^T v} \right) y = y - 2 \cdot \frac{\langle v, y \rangle}{\langle v, v \rangle} v$$

**Question:** How can we achieve that  $Q$  maps vector  $y$  to  $\alpha \cdot e_1$ ? Hence,

$$\alpha \cdot e_1 = y - 2 \frac{\langle v, y \rangle}{\langle v, v \rangle} \cdot v \in \text{span}(e_1) = \alpha(\{e_1\})$$

Consider that

$$|\alpha| = \|\alpha e_1\| = \left\| y - 2 \frac{\langle v, y \rangle}{\langle v, v \rangle} v \right\|_2 = \|Qy\|_2 = \|y\|_2$$

and  $y \in \text{span}(y - \alpha \cdot e_1)$ . This way we can retrieve  $Q$  by,

$$v := y - \alpha e_1 \quad \alpha = \|y\|_2$$

To avoid elimination (or subtraction) in  $v = (y_1 - \alpha, y_2, \dots, y_n)^T$ , choose  $\alpha := -\text{sign}(y_1) \cdot \|y\|_2$ . Now,

$$\langle v, v \rangle = \langle y - \alpha e_1, y - \alpha e_1 \rangle = \|y\|^2 - 2\alpha \langle e_1, y \rangle + \alpha^2 = 2\alpha^2 - 2\alpha \langle e_1, y \rangle = -2\alpha(y_1 - \alpha)$$

$Qy$  for  $y \in \mathbb{R}^m$  can be determined by

$$Qy = y - 2 \frac{\langle v, y \rangle}{2(y_1 - \alpha)} \cdot v$$

In a first step choose  $Q_1$  (where  $A \in \mathbb{R}^{m \times n}$ ):

$$[a_1 \dots a_n] = A \rightarrow A' = Q_1 \cdot A = \begin{bmatrix} \alpha_1 & & \\ 0 & & \\ \vdots & a'_2 \dots a'_n & \\ 0 & & \end{bmatrix}$$

where  $Q_1 = I - 2 \frac{v_1 v_1^T}{v_1^T v_1}$  with  $v_1 := a_1 - \alpha_1 e_1$  and  $\alpha_1 := -\text{sign}(a_{11}) \cdot \|a_1\|$ .

After the  $k$ -th step, we get

$$A^{(k)} = \begin{bmatrix} * & & * \\ & \ddots & \\ & & * \\ 0 & & & T^{(k+1)} \end{bmatrix}$$

with  $T^{(k+1)} \in \mathbb{R}^{m-k \times n-k}$ . Now we create orthogonal matrix

$$Q_{k+1} := \begin{bmatrix} I & 0 \\ 0 & \tilde{Q}_{k+1} \end{bmatrix}$$

where  $\tilde{Q}_{k+1}$  modifies matrix  $T^{(k+1)}$ . The next subcolumn will be eliminated.

After  $p = \min(m, n)$  steps, we retrieve an upper triangular matrix

$$R = Q_p \dots Q_1 \cdot A$$

$$\Leftrightarrow A = QR \text{ where } Q = Q_1 \cdot \dots \cdot Q_p$$

Based on Theorem 2.6, we retrieve the following process:

1.  $A = QR$  (QR decomposition using Householder or Givens)
2.  $\begin{pmatrix} b_1 \\ b_2 \end{pmatrix} = Q^t b$  with  $b_1 \in \mathbb{R}^n$  and  $b_2 \in \mathbb{R}^{m-n}$  (transformation of  $b$ )
3.  $Rx = b_1$  (solving by backwards substitution)

**Remark (Implementation).** Usually, we store the Householder vectors  $v_1, \dots, v_p$  in the lower half of  $A$  and the diagonal elements  $v_{1i} = \alpha_i$  with  $i = 1, \dots, p$  in a separate vector.

**Remark (Computing time).**

- $\sim 2n^2m$  multiplications if  $m \gg n$
- $\frac{2}{3}n^3$  if  $m \approx n$

It is important to recognize that this process is numerically stable!

### 3.3 Singular value decomposition

Literature:

- Gene H. Golub and Charles F. Van Loan: "Matrix computations"

If  $A \in \mathbb{R}^{m \times n}$  has non-full rank  $n$ , then the linear equation system of  $\min \|b - Ax\|$  is not unique, because  $\forall \tilde{x} : \|b - A\tilde{x}\| = \|b - A(\tilde{x} + v)\| \forall v \in \ker(A)$ .

Are there more decompositions of the matrix?

Recognize that the QR decomposition always exists anyways.

**Theorem 3.7.** Let  $A \in \mathbb{R}^{m \times n}$ . There  $\exists U \in \mathbb{R}^{m \times m}, V \in \mathbb{R}^{n \times n} : U^T U = I, V^T V = I$  and  $A = U \Sigma V^T$  with  $\Sigma \in \mathbb{R}^{n \times n}$  where  $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_p), p = \min(m, n), V = [v_1, \dots, v_n]$  and  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p \geq 0$ .

**Remark 3.8.**  $\sigma_i$  are called "singular values" of  $A$ .  $u_i$  and  $v_i$  are called left- and right-singular vectors of  $A$ .

↓ This lecture took place on 2018/10/30.

*Proof.* Let  $x \in \mathbb{R}^n$  and  $y \in \mathbb{R}^m$  such that  $\|x\|_2 = 1 = \|y\|_2$  and  $Ax = \sigma y$  where  $\sigma = \|A\|_2$ . Recognize that  $x, y$  exist because  $\|A\|_2 = \sup_{\|\tilde{x}\|=1} \|A\tilde{x}\|_2 = \max_{\|\tilde{x}\|_2=1} \|A\tilde{x}\|_2$ .

Hence, there exists  $x$  with  $\|x\|_2 = 1$  and  $\|Ax\|_2 = \|A\|_2$ . Followingly, we define  $y := \frac{Ax}{\|A\|_2}$  and it is true that  $\sigma y = Ax$ .

$$\|y\|_2 = \frac{\|Ax\|_2}{\|A\|_2} = 1$$

We construct orthonormal bases  $\{x, v_2, \dots, v_n\}$  and  $\{y, u_2, \dots, u_m\}$ . The corresponding matrices  $U = [y \ u_2 \ \dots \ u_m]$  and  $V = [x \ v_2 \ \dots \ v_n]$  are orthogonal.  $V^T V = I, U^T U = I$ . Let  $U_2 := (u_2, \dots, u_m)$  and  $V_2 := (v_2, \dots, v_n)$ . Furthermore,

$$U^T A V = \begin{bmatrix} y^T \\ U_2^T \end{bmatrix} A \begin{bmatrix} x & V_2 \end{bmatrix} = \begin{bmatrix} y^T A x & y^T A V_2 \\ U_2^T A x & U_2^T A V_2 \end{bmatrix} = \begin{bmatrix} \sigma & w^T \\ 0 & B \end{bmatrix} := A_1$$

where  $w \in \mathbb{R}^{n-1}$  and  $B \in \mathbb{R}^{m-1 \times n-1}$  because  $\left\| A_1 \begin{pmatrix} \sigma \\ w \end{pmatrix} \right\|_2^2 \geq (\sigma^2 + w^T w)$  it holds that  $\|A_1\|_2^2 \geq (\sigma^2 + w^T w)$ . On the other hand,  $\sigma^2 = \|A\|_2^2 = \|A_1\|_2^2$ . Thus  $w = 0$ . The statement then follows inductively.  $\square$

## Properties of the singular value decomposition

**Corollary 3.9.** If  $A = U \Sigma V^T \in \mathbb{R}^{m \times n}$  and  $m \geq n$  is the singular value decomposition of  $A$ , then  $A v_i = \sigma_i u_i$  and  $A^T u_i = \sigma_i v_i$  for all  $i = 1, \dots, n$ .

*Proof.* Consider columnwise the equality  $AV = U\Sigma$  and accordingly,  $A^T U = V\Sigma^T$ .  $\square$

By the corollary, it immediately follows that

$$A^T A v_i = \sigma_i^2 v_i \quad A A^T u_i = \sigma_i^2 u_i$$

hence the squares of singular values are the eigenvalues of the matrices  $A^T A$ , and accordingly  $A A^T$ .

**Illustrative interpretation:** The singular values of a matrix  $A$  give the length of the semi-axes of the ellipsis defined by  $E = \{Ax \mid \|x\|_2 = 1\}$ . The directions of semi-axes are denoted by  $u_i$ .

$$\|A\|_2^2 = \sup_{x \neq 0} \frac{\langle Ax, Ax \rangle}{\langle x, x \rangle} = \sup_{x \neq 0} \frac{\langle A^T A x, x \rangle}{\langle x, x \rangle} = \lambda_{\max}(A^T A)$$

**Corollary 3.10.** Let  $A \in \mathbb{R}^{m \times n}$ . Then,

$$\|A\|_2 = \sigma_1 \quad \|A\|_F = \sqrt{\text{trace}(A^T A)} = \sqrt{\sigma_1^2 + \dots + \sigma_p^2}$$

*Proof.* Consider that for spectral norm and Frobenius norm, it holds that  $\|A\| = \|U^T A V\| = \|\Sigma\|$ .  $\square$

**Corollary 3.11.** *If  $A$  has exactly  $r$  positive singular values, then  $\text{rank}(A) = r$  and  $\ker(A) = \text{span}\{v_{r+1}, \dots, v_n\}$  and  $\text{range}(A) = \text{span}\{u_1, \dots, u_r\}$ .*

*Proof.* Because  $U$  and  $V$  are regular, it holds that  $\text{rank}(A) = \text{rank}(U\Sigma V^T) = \text{rank}(\Sigma) = r$ .

The statements regarding the kernel and image of  $A$  result by Corollary 3.9.  $\square$

**Corollary 3.12.** *Let  $A \in \mathbb{R}^{m \times n}$  and  $\text{rank}(A) = r$ . Then  $A = \sum_{i=1}^r \sigma_i u_i v_i^T$ .*

*Proof.* Again, consider the matrix notation of the singular value decomposition:

$$A = (U\Sigma)V^T = \begin{bmatrix} \sigma_1 u_1 & \sigma_2 u_2 & \dots & \sigma_r u_r & 0 & \dots & 0 \end{bmatrix} \begin{bmatrix} v_1^T \\ \vdots \\ v_n^T \end{bmatrix} = \sum_{i=1}^r \sigma_i u_i v_i^T$$

$\square$

**Theorem 3.13** (Eckart-Young-Mirsky). *If  $k < r = \text{rank}(A)$  and  $A_k = \sum_{i=1}^k \sigma_i u_i v_i^T$ , then  $\min_{\text{rank}(B)=k} \|A - B\|_2 = \|A - A_k\|_2 = \sigma_{k+1}$ .*

*Proof.* Because  $U^T A_k V = \text{diag}(\sigma_1, \dots, \sigma_k, 0, \dots, 0)$ ,  $\text{rank}(A_k) = k$ . Furthermore,

$$\|U^T(A - A_k)V\|_2 = \text{diag}(0, \dots, 0, \sigma_{k+1}, \dots, \sigma_p)$$

and therefore  $\|A - A_k\|_2 = \sigma_{k+1}$ . Now let  $B \in \mathbb{R}^{m \times n}$  with  $\text{rank}(B) = k$  arbitrary. Then we can find an orthonormal basis  $x_1, \dots, x_{n-k}$  such that  $\ker(B) = \text{span}\{x_1, \dots, x_{n-k}\}$ . Furthermore it holds that  $\text{span}\{x_1, \dots, x_{n-k}\} \cap \text{span}\{v_1, \dots, v_{k+1}\} \neq \{0\}$ . Let  $z$  with  $\|z\|_2 = 1$  of this intersection. Then  $Bz = 0$  and  $Az = A \sum_{i=1}^{k+1} \alpha_i v_i = \sum_{i=1}^{k+1} \sigma_i \alpha_i u_i$ . Because  $\|z\|_2 = 1$ ,  $1 = \|z\|_2^2 = \left\| \sum_{i=1}^{k+1} \alpha_i v_i \right\|_2^2 = \sum_{i=1}^{k+1} \alpha_i^2$ . Followingly,

$$\|A - B\|_2^2 \geq \|(A - B)z\|_2^2 = \|Az\|_2^2 = \sum_{i=1}^{k+1} \sigma_i^2 \alpha_i^2 \geq \sigma_{k+1}^2 \underbrace{\sum_{i=1}^{k+1} \alpha_i^2}_{=1} = \sigma_{k+1}^2$$

$\square$

We go back to the topic of  $\min_x \|b - Ax\|_2$ .

Consider the equation  $b \cdot Ax = d$  and we followingly want to minimize the norm of the residue  $\|d\|_2$ . Left-sided multiplication with  $U^T$  gives

$$U^T(b - Ax) = U^T d \implies U^T b - U^T A V V^T x = U^T d = \tilde{d}$$

Because of orthogonality of  $U$ ,  $\|\tilde{d}\|_2 = \|d\|_2$  and thus we can minimize  $\|\tilde{d}\|_2$ . For this we introduce auxiliary vectors  $\tilde{x} = V^T x \in \mathbb{R}^n$  and  $\tilde{b} = U^T b \in \mathbb{R}^m$ . Because of the singular value decomposition, we retrieve the special form of  $U^T b - U^T A V V^T x = U^T d = \tilde{d}$

$$\begin{aligned} -\sigma_i \tilde{x}_i + \tilde{b}_i &= \tilde{d}_i & i &= 1, \dots, r \\ \tilde{b}_i &= \tilde{d}_i & i &= r+1, \dots, m \end{aligned} \tag{3}$$

The last  $m - r$  components are independent of  $\tilde{x}_i$ . The term  $\|\tilde{\alpha}\|_2^2 = \sum_{i=1}^m \tilde{d}_i^2$  will become minimal if  $\tilde{d}_i = 0$  for  $i = 1, \dots, r$ . In this case it holds that

$$\min \|\tilde{d}\|_2^2 = \sum_{i=r+1}^m \tilde{d}_i^2 = \sum_{i=r+1}^m \tilde{b}_i^2 = \sum_{i=r+1}^m (u_i^T b)^2$$

The first  $r$  unknowns are given by

$$\tilde{x}_i = \frac{\tilde{b}_i}{\sigma_i} \quad i = 1, \dots, r$$

because of equations 3 and simultaneously the last  $n - r$  unknowns can be chosen freely. The solution vector  $x$  therefore has the following representation:

$$x = V\tilde{x} = \sum_{i=1}^n \tilde{x}_i v_i = \sum_{i=1}^r \frac{\tilde{b}_i}{\sigma_i} v_i + \sum_{i=r+1}^n \tilde{x}_i v_i = \sum_{i=1}^r \frac{u_i^T b}{\sigma_i} v_i + \sum_{i=r+1}^n \tilde{x}_i v_i$$

with the free parameters  $\tilde{x}_{r+1}, \dots, \tilde{x}_n$ . If matrix  $A$  does not have full column rank, then the generic solution is representable as sum of a particulate solution in the subspace of  $r$  right-singular vectors  $v_i$  to positive singular values  $\sigma_i$  and an arbitrary vector of  $\text{kernel}(A)$ .

Often, the specific solution  $x^*$  with minimal Euclidean norm is used. Because of orthogonality, it is characterized by  $\tilde{x}_i = 0$  for  $i = r + 1, \dots, n$  and we get

$$x^* = \sum_{i=1}^r \frac{u_i^T b}{\sigma_i} v_i \quad \|x^*\|_2 = \min_{\|b - Ax\|_2^2} \|x\|_2$$

with  $\|b - Ax\|_2^2 = \sum_{i=r+1}^n (u_i^T b)^2$ .

## Generalized inverse

Using the singular value decomposition, we can also give the inverse of a regular matrix  $A \in \mathbb{R}^{n \times n}$ , because  $A^{-1} = (U\Sigma V^T)^{-1} = V\Sigma^{-1}U^T$ . For arbitrary matrices  $A \in \mathbb{R}^{m \times n}$  with  $\text{rank}(A) = r$ , we can define the *generalized inverse* by

$$A^\dagger = [v_1 \quad \dots \quad v_r] \text{diag} \left( \frac{1}{\sigma_1}, \dots, \frac{1}{\sigma_r} \right) \begin{bmatrix} u_1^T \\ \vdots \\ u_r^T \end{bmatrix}$$

This gives the so-called *Moore-Penrose inverse*.

## 4 Interpolation

**Problem setting:** For an unknown function  $f : \mathbb{R} \rightarrow \mathbb{R}$ , construct from given data  $f^{(j)}(t_i)$  for  $i = 0, \dots, n$  and  $j = 0, \dots, c_i$  a (efficiently computable) function  $\varphi$  which approximates  $f$ . We require the following properties:

**Interpolation property** at points  $t_i$ ,  $\varphi$  and  $f$  correspond.

$$\varphi^{(j)}(t_i) = f^{(j)}(t_i) \quad \forall i, j$$

We also call  $f^{(j)}(t_i)$  *supporting points of the function*

**Approximation property** in a (yet unknown) function space it holds that  $\|\varphi - f\| \approx 0$

#### 4.1 Polynomial interpolation

First, we assume that only function values  $f_i := f(t_i)$  for  $t = 0, \dots, n$  at pairwise different points  $t_0, \dots, t_n$  are given.

**Goal:** Find a polynomial  $P \in \mathcal{P}_n$  of degree  $\deg(P) \leq n$ .

$$P(t) = a_n t^n + a_{n-1} t^{n-1} + \dots + a_1 t + a_0$$

with  $a_0, \dots, a_n \in \mathbb{R}$  such that  $f$  is interpolated at  $t_0, \dots, t_n$ , thus  $P(t_i) = f_i$  for  $i = 0, \dots, n$ .

#### Uniqueness and condition

The values  $a_0, \dots, a_n$  are uniquely determined by  $P(t_i) = f_i$  for  $i = 0, \dots, n$ , because if  $P, Q \in \mathcal{P}_n$  satisfy  $P(t_i) = f_i$  then,

$$P(t_i) = Q(t_i) \text{ for } i = 0, \dots, n$$

and thus  $R := P - Q \in \mathbb{P}_n$  for some polynomial with  $n + 1$  zeros. Thus  $R$  is a zero polynomial, thus  $P = Q$ .

Now consider the map

$$\mathcal{P}_n \rightarrow \mathbb{R}^{n+1}, p \mapsto (p(t_0), \dots, p(t_n))$$

The map is linear and by property  $\dim(\mathcal{P}_n) = n + 1 = \dim(\mathbb{R}^{n+1})$  with injectivity, surjectivity follows (e.g. by the rank theorem).

↓ This lecture took place on 2018/11/06.

**Theorem 4.1.** For  $n + 1$  supporting points  $(Z_i, f_i)$  for  $i = 0, \dots, n$  with pairwise different points  $t_0, \dots, t_n$ , there exists a unique interpolation polynomial  $P$  of degree  $n$  such that  $P(t_i) = f_i \forall i$ .

This unique polynomial  $P$  is called interpolation polynomial of  $f$  for points  $p_0, \dots, p_n$  with  $p = P(f \mid t_0, \dots, t_n)$ .

A proof for Theorem 4.1 is not provided.

To compute  $p$ , we need a basis of  $\mathcal{P}_n$ . If  $P$  is given in *coefficient representation*  $(a_n t^n + \dots + a_1 t + a_0)$  with respect to monomial basis  $\{1, t, \dots, t^n\}$ , by  $p(t_i) = f_i$  we get the



linear equation system

$$\underbrace{\begin{bmatrix} t_i^j \end{bmatrix}}_{=: v_n} \cdot \begin{bmatrix} a_0 \\ \vdots \\ a_n \end{bmatrix} = \begin{bmatrix} f_0 \\ \vdots \\ f_n \end{bmatrix} \quad A = \begin{bmatrix} 1 & t_0^1 & t_0^2 & \dots \\ 1 & t_1^1 & t_1^2 & \dots \\ \vdots & \vdots & \vdots & \ddots \end{bmatrix}$$

$V_n$  is called *Vandermonde matrix* and it holds:

**Theorem 4.2.** Let  $V_n \in \mathbb{R}^{n+1 \times n+1}$  be Vandermonde's matrix. Then

$$\det(V_n) = \prod_{i=0}^n \prod_{j=i+1}^n (t_j - t_i)$$

**Remark 4.3.** Obviously:  $\det(V_n) \neq 0 \iff t_i \neq t_j \forall i = j$

Alternative basis: Consider the Lagrange basis  $L_0, \dots, L_n$  where these are the uniquely determined interpolation polynomials  $L_i \in \mathcal{P}_n$  with  $L_i(t_j) = \delta_{ij}$ .

This leads to an explicit representation:

$$L_i(t) = \prod_{\substack{j=0 \\ j \neq i}}^n \frac{t - t_j}{t_i - t_j}$$

The interpolation polynomial  $P$  for arbitrary supporting points  $f_0, \dots, f_n$  can be constructed by superposition of polynomials.

$$P(t) = \sum_{i=0}^n f_i L_i(t)$$

Thus,

$$P(t_j) = \sum_{i=0}^n f_i \cdot L_i(t_j) = \sum_{i=0}^n f_i \cdot \delta_{ij} = f_j$$

**Remark 4.4** (Alternative characterization). The Lagrange polynomials create an orthonormal basis of  $\mathcal{P}_n$  with respect to scalar product  $\langle P, Q \rangle = \sum_{i=0}^n P(t_i)Q(t_i)$ . Because  $\langle P, L_i \rangle = P(t_i)$ ,

$$P = \sum_{i=0}^n \langle P, L_i \rangle \cdot L_i = \sum_{i=0}^n P(t_i) \cdot L_i = \sum_{i=0}^n f_i \cdot L_i$$

**Theorem 4.5.** Let  $a \leq t_0 < \dots < t_n \leq b$ . Let  $L_{i,n}$  be the Lagrange polynomials. Then the absolute conditioning number  $\kappa_{\text{abs}}$  of polynomial interpolation  $\Phi = P(\cdot \mid t_0, \dots, t_n) \in C[a, b] \rightarrow \mathcal{P}_n$  with respect to supremum norm is the so-called Lebesgue constant  $\kappa_{\text{abs}} := \Lambda_n = \max_{t \in [a, b]} \sum |L_{i,n}(t)|$  for points  $t_0, \dots, t_n$ .

*Proof.* Polynomial interpolation is linear

$$\implies \Phi'(f)(g) = \Phi(g)$$

Thus, show:  $\|\Phi'\| = \Lambda_n$ .

$$\begin{aligned}
\forall f \in C[a, b] : |\varphi(f)(t)| &= \left| \sum f(t_i) L_{i,n}(t) \right| \\
&\leq \sum |f(t_i)| \cdot |L_{i,n}(t)| \\
&\leq \|f\|_\infty \cdot \max_{t \in [a, b]} \sum |L_{i,n}(t)| \\
\Rightarrow \|\Phi\|_\infty &= \kappa_{\text{abs}} \leq \Lambda_n
\end{aligned}$$

In the other direction: We construct  $g \in C[a, b]$ :

$$|\Phi(g(\tau))| = \|g\|_\infty \cdot \max_{t \in [a, b]} \sum_{i=0}^n |L_{i,n}(t)|$$

Let  $\tau$  be the specific point which takes this maximum.

Let  $g \in C[a, b] : \|g\|_\infty = 1$  and  $g(t_i) = \text{sign}(L_{i,n}(\tau))$ .  $g$  solves the interpolation problem  $(t_i, \text{sign}(L_{i,n}(t_i)))$ .

$$\begin{aligned}
\Rightarrow |\Phi(g)(\tau)| &= \sum |L_{i,n}(\tau)| = \|g\|_\infty \cdot \max_{t \in [a, b]} \sum_{i=0}^n |L_{i,n}(t)| \\
\Rightarrow \kappa_{\text{abs}} &= \Lambda_n
\end{aligned}$$

□

| $n$ | $\Lambda_n$ for equidistant points | $\Lambda_n$ for Chebyshev points |
|-----|------------------------------------|----------------------------------|
| 5   | 3.11                               | 2.104                            |
| 10  | 29.9                               | 2.48                             |
| 15  | 512.05                             | 2.727                            |
| 20  | 10986.53                           | 2.9                              |

Chebyshev points:

$$t_i = \cos\left(\frac{2i+1}{2n+2} \cdot \pi\right)$$

## 4.2 Polynomial evaluation by Aitken-Neville

**Motivation:** Often we are interested in  $P(t)$  for a few  $t$  and not in  $t$  itself.

**Lemma 4.6** (Aitken's lemma). *For the interpolation polynomial  $P = P(f \mid t_0, \dots, t_n)$ , the following recursion formula holds:*

$$P(f \mid t_0, \dots, t_n) = \frac{(t_0 - t)P(f \mid t_1, \dots, t_n) - (t_n - t)P(f \mid t_0, \dots, t_{n-1})}{t_0 - t_n}$$

*Proof.* Let us denote the function of the right with  $\varphi$ .  $\varphi$  is a polynomial of degree  $n$ . Thus  $\varphi \in \mathcal{P}_n$ . Furthermore,

$$\phi(t_i) = \frac{(t_0 - t_i) \cdot f_i - (t_n - t_i)f_i}{t_0 - t_n} = f_i \forall i \in \{1, \dots, n-1\}$$

For  $i = 0$  or  $i = n$ :

$$\begin{aligned}\varphi(t_0) &= -\frac{t_n - t_0}{t_0 - t_n} f(t_0) \\ \varphi(t_n) &= \frac{t_0 - t_n}{t_0 - t_n} f(t_n)\end{aligned}$$

The interpolation polynomial is unique. Thus we can conclude the claim.  $\square$

Aitken's lemma provides a recursive method for evaluation of  $P(f \mid t_0 \dots t_n)$  at point  $t$ .

**Notation:**  $P(f(t_i)) = f_i$ ,  $P_{i,k} := P(f, t_{i-k} \dots t_i)(t)$ ,  $i \geq k$ .

*Neville's scheme:* where  $P_{i0} = f_i$  and

|          | $P_{i0}$ | $P_{i1}$ | $P_{i2}$ | $\dots$ | $P_{in}$ |
|----------|----------|----------|----------|---------|----------|
| $t_0$    | $P_{00}$ |          |          |         |          |
| $t_1$    | $P_{10}$ | $P_{11}$ |          |         |          |
| $\vdots$ |          |          |          |         |          |
| $t_n$    | $P_{n0}$ | $P_{n1}$ | $\dots$  |         | $P_{nn}$ |

$$P_{ik} = \frac{t - t_{i-k}}{t_i - t_{i-k}} P_{i,k-1} + \frac{t_i - t}{t_i - t_{i-k}} P_{i-1,k-1} = P_{i,k-1} + \frac{t - t_i}{t_i - t_{i-k}} (P_{i,k-1} - P_{i-1,k-1})$$

How many operations are necessary?

$k = 1$ :  $6n$

$k = 2$ :  $6(n-1)$

$k = n$ :  $6$

In total, we get

$$6 \cdot \sum_{i=1}^n i = 6 \frac{n(n-1)}{2} = \mathcal{O}(n^2)$$

If there is one additional supporting point, the computing time is increased by  $\mathcal{O}(n)$ .

### 4.3 Hermitian interpolation and Newton basis

Now we allow, that points repeat, thus  $a \leq t_0 \leq \dots \leq t_n \leq b$ .

$$\triangle := \{t_i\}_{i=0}^n$$

If the derivatives  $f'(t_i), \dots, f^{(k)}(t_i)$  are given,  $t_i$  is expected to occur  $k+1$  times.

$$d_i := \max \{j \mid t_i = t_{i-j}\}$$

For example:

$$\begin{array}{c|ccccccc} t_i & t_0 & = t_1 & < t_2 & = t_3 & = t_4 & < t_5 & < t_6 \\ d_i & 0 & 1 & 0 & 1 & 2 & 0 & 0 \end{array}$$

Once we define  $\forall i : \mu_i : C^n[a, b] \rightarrow \mathbb{R}, \mu_i(f) = f^{(d_i)}(t_i)$ , we can define the goal of Hermitian interpolation:

We construct a polynomial  $P \in \mathcal{P}_n$  satisfying  $\mu_i(P) = \mu_i(f) \forall i = 0, \dots, n$ .

↓ This lecture took place on 2018/11/08.

**Remark** (revision of Hermitian interpolation).  $i = 0, \dots, n$  with  $\mu_i : C^n([a, b]) \rightarrow \mathbb{R}^n$  with  $\mu_i(f) := f^{(d_i)}(t_i)$

The solution  $p = P(f(t_0, \dots, t_n)) \in \mathcal{P}_n$  is called Hermitian interpolant of  $f$  at  $t_0, \dots, t_n$ .

**Theorem 4.7.** For every function  $f \in C^n([a, b])$  and every monotonic sequence  $a \leq t_0 \leq t_1 \leq \dots \leq t_n \leq b$  of (not necessarily different) points there exists exactly one polynomial  $P \in \mathcal{P}_n$  such that  $\forall i \in \{0, \dots, n\} : \mu_i(P) = \mu_i(f)$ .

*Proof.* A map  $\mu : \mathcal{P}_n \rightarrow \mathbb{R}^{n+1}$  with  $P \mapsto (\mu_0 P, \dots, \mu_n P)$  is linear and injective. By  $\mu(P) = 0$ ,  $P$  has at least  $n + 1$  roots (with multiplicities); so its the zero polynomial. Because  $\dim(\mathcal{P}_n) = \dim(\mathbb{R}^{n+1}) = n + 1$ , existence follows.  $\square$

If all points are pairwise different, we retrieve the Lagrange interpolant

$$P(f \mid t_0, \dots, t_n) = \sum_{i=0}^n f(t_i) L_i$$

However, if  $t_0 = t_1 = \dots = t_n$ , the interpolant is the aborted Taylor series around  $t = t_0$

$$P(f \mid t_0, \dots, t_n)(t) = \sum_{j=0}^n \frac{(t - t_0)^j}{j!} f^{(j)}(t_0)$$

is also called *Taylor interpolant*.

**Lemma 4.8.** Assuming  $t_i \neq t_j$ , the Hermitian interpolation polynomial  $P = P(f(t_0, \dots, t_n))$  satisfies

$$P = \frac{(t_i - t)P(f \mid t_0, \dots, \hat{t}_i, \dots, t_n) - (t_j - t)P(f \mid t_0, \dots, \hat{t}_j, \dots, t_n)}{t_i - t_j}$$

where  $\hat{t}_i$  means that the  $i$ -th point is skipped.

*Proof.* We prove the interpolation properties by insertion of the definition.

If many evaluations of polynomials are required, it is advantageous to determine polynomial  $P$  explicitly. Therefore, choose the Newton basis  $w_0, \dots, w_n$  of  $\mathcal{P}_n : w_0(t) := 1$  and  $w_i(t) := \prod_{j=0}^{i-1} (t - t_j), i \geq 1$ .  $\square$

**Definition 4.9.** The leading coefficient  $a_n$  of the interpolation polynomial  $P(f \mid t_0, \dots, t_n) = a_n t^n + a_{n-1} t^{n-1} + \dots + a_0$  of  $f$  for points  $t_0 \leq t_1 \leq \dots \leq t_n$  is called  $n$ -th divided difference of  $f$  at  $t_0, \dots, t_n$  and is denoted as  $[t_0, \dots, t_n]f := a_n$ .

**Theorem 4.10.** For every function  $f \in \mathcal{C}^n([a, b])$  and points  $t_0, \dots, t_n$ ,  $P := \sum_{i=0}^n [t_0, \dots, t_i]f \cdot w_i$  defines the interpolation polynomial  $P(f \mid t_0, \dots, t_n)$  of  $f$  at  $t_0, \dots, t_n$ . If  $f \in \mathcal{C}^{n+1}([a, b])$ , then  $f(t) = P(t) + [t_0, \dots, t_n, t]f \cdot w_{n+1}(t)$ .

*Proof.* Proof by induction over  $n$ .

For  $n = 0$ , we get value  $[t_0]f = f(t_0)$  as leading coefficient (at points  $t_0$ ). The constant interpolation polynomial  $P$  satisfies  $p := [t_0]f \cdot w_0 = f(t_0) \cdot 1$ . This is shown by the first statement for  $n = 0$ . To determine  $[t_0, t]f$ , compute the leading coefficients of the interpolation polynomial at points  $t_0$  and  $t$ . For this polynomial it holds true that

$$\begin{aligned} P(f \mid t_0, t) &= a + a_1 t =: p(t) \\ \implies f(t_0) &= p(t_0) = a_0 + a_1 t_0 \\ f(t) &= p(t) = a_0 + a_1 t \\ \implies a_1 &= [t_0, t]f = \frac{f(t) - f(t_0)}{t - t_0} \end{aligned}$$

This implies

$$P(t) = [t_0]f \cdot w_0 + [t_0, t]f \cdot w_1(t) = f(t_0) \cdot 1 + \frac{f(t) - f(t_0)}{t - t_0}(t - t_0) = f(t)$$

So the statement for case  $n = 0$  follows.

Now let  $n > 0$  and

$$P_{n-1} := P(f \mid t_0, \dots, t_{n-1}) = \sum_{i=0}^{n-1} [t_0, \dots, t_i]f \cdot w_i$$

the interpolation polynomial of  $f$  at  $t_0, \dots, t_{n-1}$ . For  $P_n = P(f(t_0, \dots, t_n))$ , we have

$$P_n(t) = [t_0, \dots, t_n]f \cdot t^n + a_{n-1} t^{n-1} + \dots + a_0 = [t_0, \dots, t_n]f \cdot w_n(t) + Q_{n-1}(t)$$

with some polynomial  $Q_{n-1} \in \mathcal{P}_{n-1}$ . Because  $Q_{n-1} = P_n - [t_0, \dots, t_n]f \cdot w_n$ ,  $Q_{n-1}$  satisfies the interpolation requirements for  $t_0, \dots, t_{n-1}$ . Therefore

$$Q_{n-1} = P_{n-1} = \sum_{i=0}^{n-1} [t_0, \dots, t_i]f \cdot w_i$$

This shows the first statement. Especially,  $P_n + [t_0, \dots, t_n, t]f \cdot w_{n+1}$  interpolates the function  $f$  at points  $t_0, \dots, t_n, t$ . So the second statement follows.  $\square$

The following properties are given:

**Theorem 4.11.** Dividing differences  $[t_0, \dots, t_n]f$  satisfying (assuming  $f \in \mathcal{C}^n([a, b])$ )

1.  $[t_0, \dots, t_n]P = 0$  for all  $P \in \mathcal{P}_{n-1}$

2. For collapsing points  $t_0 = \dots = t_n$ , we have  $[t_0, \dots, t_n]f = \frac{f^{(n)}}{n!}$
3. For  $i < j$ , we have the recursive formula  $[t_i, \dots, t_j]f = \frac{[t_{i+1}, \dots, t_j]f - [t_i, \dots, t_{j-1}]f}{t_j - t_i}$

↓ This lecture took place on 2018/11/13.

**Theorem 4.12** (Hermite-Gnocco Formula). For the  $n$ -th dividing difference of an  $n$ -times differentiable function  $f \in \mathcal{C}^n$ , we have

$$\sum^n := \left\{ s = (s_0, \dots, s_n) \in \mathbb{R}^{n+1} \mid \sum s_i = 1, s_i \geq 0 \right\}$$

$$[t_0, \dots, t_n]f = \int_{\sum^n} f^{(n)} \left( \sum_{i=0}^n s_i t_i \right) ds$$

No proof is given of this theorem.

#### 4.3.1 Approximation error

**Theorem 4.13.** If  $f \in \mathcal{C}^{n+1}$ , for the interpolation polynomial  $P(f(t_0, \dots, t_n))$  and  $t_i, t \in [a, b]$  we have:

$$f(t) - P(f \mid t_0, \dots, t_n)(t) = \frac{f^{(n+1)}(\xi)}{(n+1)!} \cdot w_{n+1}(t) \text{ with } \xi = \xi(t) \in (a, b)$$

*Proof.* Let  $t \in [a, b]$ . For  $t = t_i$ , both sides are zero. Thus, we assume  $t \neq t_i \forall i$ .

For  $\alpha \in \mathbb{R}$  we define that  $g(\tau) := f(\tau) - P(f \mid t_0, \dots, t_n)(\tau) - \alpha w_{n+1}(\tau)$ .  $g$  is  $(n+1)$ -times continuous differentiable and has roots  $t_0, \dots, t_n$ . Choose

$$\alpha = \frac{f(t) - P(f \mid t_0, \dots, t_n)(t)}{(t - t_0) \cdot \dots \cdot (t - t_n)}$$

Thus,  $\alpha$  is a root of  $g$ . This implies there are  $n+2$  roots.

By Rolle's Theorem,  $g'$  has  $n+1$  roots and consider this iteratively, we have exactly one root  $\xi \in [a, b]$  for  $g^{(n+1)}$ . For  $\xi$ , we have:

$$\begin{aligned} 0 &= g^{(n+1)}(\xi) \\ &= f^{(n+1)}(\xi) - \frac{f(t) - P(f \mid t_0, \dots, t_n)(t)}{(t - t_0) \dots (t - t_n)} \cdot (n+1)! \end{aligned}$$

Because  $P$  is a polynomial of degree  $n+1$  and  $w_{n+1}$  is a polynomial with leading coefficient 1.

If we resolve this by  $P(f \mid \dots)(t)$  shows our claim. □

**Remark 4.14.** In case of Taylor interpolation ( $t_0 = \dots = t_n$ ) we have the Lagrange remainder term of the Taylor expansion

$$f(t) - P(f \mid t_0 \dots t_n)(t) = \frac{f^{(n+1)}(\xi)}{(n+1)!} (t - t_0)^{n+1}$$

Consider the following class of functions:

$$\mathcal{F} := \left\{ f \in \mathcal{C}^{n+1}[a, b] \mid \sup_{x \in [a, b]} |f^{(n+1)}(x)| \leq M(n+1)! \right\}$$

for  $M > 0$  constant, the approximation error depends on the choice of point  $t$  with respect to  $w_{n+1}(t) = (t - t_0) \cdot \dots \cdot (t - t_n)$ .

#### 4.3.2 Numeric differentiation

Approach:

1. Choose points  $t_0, \dots, t_n$ .
2. Let  $f$  approximate this polynomial  $P \in \mathcal{P}_n$
3. Differentiate  $P \rightarrow P'$
4. Can be determined accurately such that the interpolation error can determine the accuracy of the approximation

Consider  $P(f \mid t_i)$  in Newton form:

$$P = f(t_0) + [t_0, t_1]f w_1 + \dots + [t_0 \dots t_n]f w_n$$

The approximation error is given by

$$r(t) = w_{n+1}(t) \cdot \frac{f^{(n+1)}(\xi(t))}{(n+1)!}$$

If  $f \in \mathcal{C}^{n+2}$ , we can differentiate both equation systems.

$$t = t_0 : P'(t_0) = [t_0, t_1]f + [t_0, t_1, t_2]r \cdot (t_0 - t_2) + \dots + [t_0 \dots t_n]f(t_0 - t_1) \dots (t_0 - t_{n-1})$$

$$r'(t_0) = (t_0 - t_n) \frac{f^{(n+1)}(\xi(t_0))}{(n+1)!} \quad w_{n+1}^{(k)}(t_0) = 0 \forall k$$

Thus, we get:

$$f'(t_0) = P'(t_0) + r'(t_0) \quad (\text{analogously for } t_j \neq t_0)$$

If the supporting points are distributed with distance  $k > 0$  next to  $t_0$ , we get the classical differential quotient:

**Example 4.15. Forward differential quotient** for the first derivative:  $t_1 := t_0 + k$

$$\begin{aligned} f'(t_0) &= P'_1(t_0) + r'_1(t_0) \\ &= [t_0, t_1]f + (t_0 - t_1) \cdot \frac{1}{2}f''(\xi) \\ &= \frac{f(t_0 + h) - f(t_0)}{h} - \frac{n}{2} \cdot f''(\xi) \end{aligned}$$

The resulting error is in  $\mathcal{O}(h)$  for  $h \rightarrow 0$

**Central differential quotient**  $t_0 = t_1 - k, t_1, t_2 = t_1 + k$ .

The three supporting points define a polynomial of second degree, we can use to determine  $f'(t_2)$ .

$$\begin{aligned} \Rightarrow f'(t_1) &= \underbrace{P'(t_1)}_{\deg=2} + r_2(t_1) = [t_0, t_1]f + [t_0, t_1, t_2]f \cdot \underbrace{(t_1 - t_0)}_{=h} + (t_1 - t_0)(t_1 - t_2) \cdot \frac{1}{6} \cdot f'''(\xi) \\ &= \frac{f(t_1) \cdot f(t_0)}{h} + \frac{1}{2h} (f(t_2) - 2f(t_1) + f(t_0)) \\ &= \frac{1}{2h} (f(t_1 + h) - f(t_1 - h)) - \frac{h^2}{6} f'''(\xi) \end{aligned}$$

where

$$[t_0, t_1, t_2]f = \frac{[t_1, t_2]f - [t_0, t_1]f}{t_2 - t_0} = \frac{1}{2h} \left( \frac{f(t_2) - f(t_1)}{h} - \frac{f(t_1) - f(t_0)}{h} \right)$$

The error is in  $\mathcal{O}(h^2)$  for  $h \rightarrow 0$ .

Thus, symmetric differential quotients usually give a better approximation than the one-sided quotient.

#### 4.3.3 Minmax property of the Chebyshev polynomial

(see page 214 in the book)

Our goal is to show that  $\max_{t \in [a, b]} |w_{n+1}(t)|$  is minimized using Chebyshev points. We consider the *Chebyshev polynomials*

$$T_n = \cos(n \cdot \arccos(x)) \text{ for } x \in [-1, 1]$$

For  $x \in \mathbb{R}$ , we define recursively,

$$T_k(x) = 2x \cdot T_{k-1}(x) - T_{k-2}(x) \quad T_0(x) = 1 \quad T_1(x) = x$$

- The Chebyshev polynomials have integral coefficients.
- The largest coefficient of  $T_n$  is  $a_n = 2^{n-1}$
- $T_n$  is even if  $n$  is even, otherwise it is odd
- $T_n(1) = 1, T_n(-1) = (-1)^n$
- $|T_n(x)| \leq 1 \forall x \in [-1, 1]$
- $|T_n(x)|$  takes the value 1 at the so-called *Chebyshev x-coordinates* (dt. "Tschebyscheff Abszissen")  $\tilde{x}_k = \cos(\frac{k\pi}{2})$

$$\Rightarrow |T_n(x)| = 1 \iff x = \tilde{x}_k = \cos(\frac{k\pi}{2}) \text{ for } k = 1, \dots, n$$

- $T_k$  create an orthogonal basis over  $[-1, 1]$  to the weight function  $w(x) = \frac{1}{\sqrt{1-x^2}}$



- Global representation:

$$T_k(x) = \frac{1}{2} \left( (x + \sqrt{x^2 - 1})^k + (x - \sqrt{x^2 - 1})^k \right) \quad \text{for } x \in \mathbb{R}$$

In terms of the approximation error, we look for  $q \in \mathcal{P}_{n+1}$  of degree  $n+1$  with leading coefficient 1 and the smallest supremum norm over some interval  $[a, b]$ .

$$\implies \max_{t \in [a, b]} |q(t)| \rightarrow \min$$

$\implies$  reduce the problem to interval  $[-1, 1]$ .

We use the following map:

$$x : [a, b] \rightarrow [-1, 1] \quad t \mapsto x(t) = 2 \frac{t - a}{b - a} - 1$$

Inverse function:

$$t(x) = \frac{1 - x}{2}a + \frac{1 + x}{2}b$$

If  $q \in \mathcal{P}_{n+1}$  with  $\deg(q) = n + 1$  and leading coefficient 1 is the solution of the min-max problem

$$\max_{x \in [-1, 1]} |q(x)| \rightarrow \min$$

then  $\hat{q}(t) = q(x(t))$  is the solution of the original problem with leading coefficient  $\frac{2^{n+1}}{(b-a)^{n+1}}$ .

**Theorem 4.16.** Every polynomial  $q_n \in \mathcal{P}_n$  with leading non-zero coefficient takes up in interval  $[-1, 1]$  a value with absolute value  $\geq \frac{|a_n|}{2^{n-1}}$ .

In particular, the Chebyshev polynomials  $T_n(x)$  are minimal with respect to maximum norm  $\|f\|_\infty = \max_{x \in [-1, 1]} |f(x)|$  and the polynomials of degree  $n$  with leading coefficient  $2^{n-1}$ .

*Proof.* Let  $q_n \in \mathcal{P}_n$  with leading coefficient  $a_n = 2^{n-1}$ . Remark:  $|q_n(x)| < 1 \forall x \in [-1, 1]$ .

$T_n - q_n$  has degree  $\leq n - 1$ . At the Chebyshev x-coordinates, we have  $T_n(\tilde{x}_{2k}) = 1$  and  $q_n(\tilde{x}_{2k}) < 1$ .  $T_n(\tilde{x}_{2k+1}) = -1$  and  $q_n(\tilde{x}_{2k+1}) > -1$ .

$T_n - q_n$  at the  $n + 1$  Chebyshev x-coordinates is alternatingly positive and negative, thus there are  $\geq n$  roots.

But  $\deg \leq n - 1 \implies T_n - q_n \equiv 0$  is a contradiction.

$$\implies \forall q_n \in \mathcal{P}_n \text{ with } a_n = 2^{n-1}$$

$$\exists x \in [-1, 1] : |q_n(x)| \geq 1$$

For  $q_n \in \mathcal{P}_n$  with  $a_n \neq 0$ , the following claim follows:  $\tilde{q}_n = q_n \cdot \frac{2^{n-1}}{a_n}$  is a polynomial with leading coefficient  $2^{n-1}$ .  $\square$

**Remark 4.17.** When minimizing the approximation error of the polynomial interpolation, we look for the points  $t_0, \dots, t_n \in [a, b]$  that solve the min-max problem  $\max_{t \in [a, b]} |w(t)| = \max |(t - t_0) \cdots (t - t_n)| \rightarrow \min$ .

We look for the normed polynomial  $w_n(\cdot) \in \mathcal{P}_n$  with real roots  $t_0, \dots, t_n$ :  $\max_{t \in [a, b]} |w(t)| \rightarrow \min$ .

$\Rightarrow$  for  $[a, b]$ , this is the  $n + 1$ -th Chebyshev polynomial with roots  $t_i = \cos\left(\frac{2i+1}{2n+2} \cdot \pi\right)$  with  $i = 0, \dots, n$ .

↓ This lecture took place on 2018/11/15.

#### 4.4 Trigonometric interpolation

*Motivation:* The interpolation of periodic processes (e.g. waves) it is useful to use periodic functions. Thus for this section, we assume  $f(x) = f(x \cdot 2\pi)$ .

*Idea:* Consider interpolation polynomials on the complex unit circle by variable transformation  $z : [0, 2\pi] \rightarrow \mathbb{C}, t \mapsto e^{it}$ .

Consider polynomial  $P_n$ ,

$$\begin{aligned} P_n(z(t)) &= \sum_{k=0}^{n-1} c_k (z(t))^k \\ &= \sum_{k=0}^{n-1} c_k \cdot e^{itk} =: p(t) \end{aligned}$$

We call  $P$  *trigonometric polynomial* and consider the following problem:

Given  $t_j = 2\pi \frac{j}{n}$  and points  $f_j \in \mathbb{C} : 0 \leq j < n$ . Find  $c_k$  with  $0 \leq k < n$  :  $\sum_{k=0}^{n-1} c_k \cdot e^{2\pi i k \frac{j}{n}} = f_j$  with  $j = 0, \dots, n-1$ .

In matrix representation, we get  $M \cdot c := [c_0 \ \dots \ c_{n-1}] = [f_0 \ \dots \ f_{n-1}]$  with  $M \in \mathbb{C}^{n \times n}$ ,

$$M = (m_{kj})_{k,j=0}^{n-1} \quad m_{kj} = e^{2\pi i k \frac{j}{n}}$$

Pay attention that a different indexing is used here.

The interpolation polynomial has a unique solution if  $M$  is invertible and

$$M^{-1} = \frac{1}{n} \cdot M^* = \frac{1}{n} \cdot \overline{M}^t$$

Recognize that  $M_{kj}^* = \overline{M}_{jk} = e^{-2\pi i j \frac{k}{n}}$ .

Now consider the entries  $M * M$ :

$$\begin{aligned}
(M^* M)_{k,j} &= \sum_{l=0}^{n-1} \overline{M}_{lk} \cdot M_{lj} \\
&= \sum_{l=0}^{n-1} e^{-2\pi i l \frac{k}{n}} \cdot e^{2\pi i l \frac{j}{n}} \\
&= \sum_{l=0}^{n-1} e^{-2\pi i l \frac{k-j}{n}} \\
&= \begin{cases} \sum_{l=0}^{n-1} 1 & \text{if } k = j \\ \sum_{l=0}^{n-1} \left( e^{-2\pi i \frac{k-j}{n}} \right)^l & \text{else} \end{cases} = \frac{1 - e^{-2\pi i (k-j)}}{1 - e^{-2\pi i \frac{k-j}{n}}}
\end{aligned}$$

which is a geometric sum.

For the enumerator, we have:

$$\begin{aligned}
\exp(-2\pi i (k-j)) &= \exp(-\pi i)^{2(k-j)} = 1 - ((-1)^2)^{k-j} = 0 \\
\implies M * M &= n \cdot I
\end{aligned}$$

The unique solution of the interpolation problem is given by

$$c = \frac{1}{n} \cdot M^* \cdot f \quad \text{or} \quad c_k = \sum_{j=0}^{n-1} f_j \cdot e^{-2\pi i \frac{j-k}{n} \frac{1}{n}}$$

This process is called discrete Fourier Transform (dFT)

The inverse Fourier transform corresponds to the polynomial evaluation  $f_i = \sum \dots$

#### Fast Fourier transform

Disadvantage of discrete Fourier transform: The computational requirements  $\mathcal{O}(n^2)$  are quite high. The computation can be sped up with fast Fourier Transform.

In the following we assume that  $n$  is even. For the unit root  $w_n^k := \exp(-2\pi i \frac{k}{n})$ , we have  $w_{2n}^{2k} = w_n^k = w_{\frac{n}{2}}^{\frac{k}{2}}$ . Decompose the components of the Discrete Fourier Transform in its even and odd parts:

$$\begin{aligned}
n \cdot c_k &= \sum_{j \text{ even}} w_n^{kj} f_j + \sum_{j \text{ odd}} w_n^{kj} f_j \\
&= \sum_{l=0}^{\frac{n}{2}-1} w_n^{k2l} f_{2l} + \sum_{l=0}^{\frac{n}{2}-1} w_n^{k(2l+1)} f_{2l+1} \\
&= \sum_{l=0}^{\frac{n}{2}-1} w_{\frac{n}{2}}^{kl} f_{2l} + w_n^k \sum_{l=0}^{\frac{n}{2}-1} w_{\frac{n}{2}}^{kl} f_{2l+1} \\
&=: \tilde{c}_k + w_n^k \tilde{\tilde{c}}_k
\end{aligned}$$

For  $k \leq \frac{n}{2} - 1$ ,  $\tilde{c}_k$  give the Discrete Fourier Transform for the  $\frac{n}{2}$  values  $f_0, f_2, \dots, f_{n-2}$ .  $\tilde{\tilde{c}}_k$  give the Discrete Fourier Transform for the  $\frac{n}{2}$  values  $f_1, f_3, \dots, f_{n-1}$ .

For those  $k$  we can  $c_k$  compute by two DFT by length  $\frac{n}{2}$ . The remaining  $c_k$  ( $k > \frac{n}{2} - 1$ ) can also be expressed with  $\tilde{c}_k$  and  $\tilde{\tilde{c}}_k$ , because for  $k \geq \frac{n}{2}$ , and accordingly  $k = \tilde{k} + \frac{n}{2}$ , with  $0 \leq \tilde{k} \leq \frac{n}{2} - 1$ , we have:

$$w_n^{2lk} = w_n^{2l\tilde{k}+ln} = w_n^{2l\tilde{k}} \cdot \underbrace{(w_n^n)^l}_{=1}$$

$$\begin{aligned} w_n^k &= w_n^{\tilde{k} + \frac{n}{2}} \\ &= w_n^{\tilde{k}} \cdot w_n^{\frac{n}{2}} \\ &= w_n^{\tilde{k}} \cdot e^{-\pi i} \\ &= -w_n^{\tilde{k}} \end{aligned}$$

So for  $k = \tilde{k} + \frac{n}{2} \geq \frac{n}{2} - 1$ :

$$nc_k = \tilde{c}_k - w_n^{\tilde{k}} \tilde{\tilde{c}}_k$$

If the unit roots are precomputed, the computational effort is given by  $A(n)$  for some DFT of length  $n$  by 2 DFTs of length  $\frac{n}{2}$  and  $n$  additions and  $\frac{n}{2}$  multiplications:

$$A(n) = 2A\left(\frac{n}{2}\right) + \frac{3}{2}n$$

If  $\frac{n}{2}$  is even again, we can solve the DFT by two recursions. For  $n = 2^n$ :

$$A(n) = 2n \cdot \log_2(n)$$

**Algorithm 4.** Input:  $f_0, \dots, f_{n-1}$ ,  $n = 2^m$

1. If  $n = 1$ :  $c_0 = f_0$  end
2.  $(a_0, \dots, a_{\frac{n}{2}}) = (f_0, f_2, \dots, f_{n-2})$
3.  $(b_0, \dots, b_{\frac{n}{2}}) = (f_1, f_3, \dots, f_{n-1})$
4. for  $k = 0, \dots, \frac{n}{2} - 1$ 
  - (a)  $c_k = a_k + w_n^k b_k$
  - (b)  $c_{k+\frac{n}{2}} = a_k - w_n^k b_k$

Output:  $c_0, \dots, c_{n-1}$

↓ This lecture took place on 2018/11/20.

*Problem setting:*  $f_j = f(t_j) \in \mathbb{R} \forall j$ . Can we create a pure interpolation polynomial?

We restrict ourself to  $n = 2p + 1$  and reenumerate the coefficients of the complex interpolation polynomial from  $c_0, \dots, c_{n-1}$  to  $c_0, \dots, c_p, c_{-p}, \dots, c_{-1}$ . Recognize that

$$e^{2\pi i \frac{n-1}{n}} = e^{-2\pi i \frac{1}{n}} \quad \dots \quad e^{2\pi i \frac{p+1}{n}} = e^{-2\pi i \frac{p}{n}}$$

Thus, the interpolation polynomial can be written as

$$f(t_j) = \sum_{k=-p}^p c_k \cdot e^{ikt_j}$$

Because  $f(t_j) \in \mathbb{R}$  it follows that  $f(t_j) = \overline{f(t_j)}$  and thus

$$\sum_{k=-p}^p c_k \cdot e^{ikt_j} = \sum_{k=-p}^p \overline{c_k} e^{-ikt_j} \implies \overline{c_k} = c_{-k} \text{ for } k = 1, \dots, p \text{ and } c_0 \in \mathbb{R}$$

If we write  $c_k = c_k^r + ic_k^i$ , we get

$$\begin{aligned} f(t_j) &= c_0 + \sum_{k=1}^p (c_k e^{ikt_j} + \overline{c_k} e^{-ikt_j}) \\ &= c_0 + \sum_{k=1}^p [(c_k^r + ic_k^i)(\cos(kt_j) + i \sin(kt_j)) + (c_k^r - ic_k^i)(\cos(kt_j) + i \sin(kt_j))] \\ &= c_0 + \sum_{k=1}^p (2c_k^r \cos(kt_j) - 2c_k^i \sin(kt_j)) \end{aligned}$$

Therefore, the unique interpolation polynomial is a linear combination of  $\{1, \cos(kt), \sin(kt) : k = 1 : \frac{n-1}{2}\}$  and the coefficients can be computed by complex Fourier transform.

## 4.5 Splines

*Motivation:* For high polynomial degrees, usually high oscillations occur. This is inappropriate for interpolation.

*Idea:* Split up the interval in smaller intervals, solve the problem locally.

*Problem:* In general, only  $\mathcal{C}^0$ , but not smooth.

Spline functions satisfy the following three conditions:

- Piecewise polynomials of highest degree  $k$
- $p$ -times continuously differentiable on the entire interval
- The curvature is minimized

In practice, *cubic splines* will occur most often.

**Problem setting:** Let  $n + 1$  pairwise different points  $a = t_0 < t_1 < \dots < t_b = b$  with function values  $y_0, \dots, y_n$  be given. Find a function  $g : [a, b]$  that satisfies:

1.  $s(t_i) = y_i \forall i$
2.  $s \in \mathcal{C}^1([a, b])$
3.  $s \in \mathcal{C}^k([t_i, t_{i+1}])$   
For the derivative there exist left- and rightsided limits in the points  $t_i$ , that must not necessarily correspond.
4.  $s$  minimizes the functional  $J[s] := \frac{1}{2} \int_a^b (s'')^2 dt$

**Theorem 4.17.1.**

$\exists! s : s$  satisfies the four properties

*Proof.* First, we assume the existence of  $s$  and then we derive necessary optimality conditions that lead to a unique solvable linear equation system.

Let  $s$  be a function that satisfies 1–4, let  $s_1$  be a function satisfying 1–3.

$$\implies J[s] \leq J[s_1]$$

So we write  $s_1(t) = s(t) + \varepsilon h(t)$  with parameter  $\varepsilon \in \mathbb{R}$ . Because  $\varepsilon h = s_1 - s$ ,  $h$  satisfies 2 and 3,  $h(t_i) = 0 \forall i$ .

Every  $h$  that satisfies this condition, is considered *reliable*, because it gives some  $s_1$  with condition 1–3.

Now consider  $J[s] := J[s_1] = \frac{1}{2} \cdot \int_a^b (s''(t) + \varepsilon h''(t))^2 dt$  with reliable  $h$ .

Therefore  $J(\varepsilon)$  is a quadratic polynomial in  $\varepsilon$ , that takes up its minimum in  $\varepsilon = 0$ .

Thus our necessary optimality condition is  $J'(0) = \left. \frac{dJ}{d\varepsilon} \right|_{\varepsilon=0} = 0$  for every reliable function  $h$ . Furthermore it is true that

$$J'(0) = \int_a^b s''(t) \cdot h''(t) dt$$

Two times partial integration over the subintervals  $(t_i, t_{i+1})$  gives

$$\begin{aligned} \int_{t_i}^{t_{i+1}} s''(t) h''(t) dt &= s''(t) h'(t) \Big|_{t_i}^{t_{i+1}} - \int_{t_i}^{t_{i+1}} s^{(3)}(t) \cdot h'(t) dt \\ &= s''(t) h'(t) \Big|_{t_i}^{t_{i+1}} - s^{(3)}(t) \cdot h(t) \Big|_{t_i}^{t_{i+1}} \\ &\quad + \int_{t_i}^{t_{i+1}} s^{(k)}(t) h(t) dt \end{aligned}$$

with  $h(t_i) = 0 \forall i$ .

After summing up all subsums, we get

$$0 = J'(0) = -s''(t_0^+) h'(t_0) + s''(t_1^-) h'(t_1) + s''(t_2^-) - s''(t_2^+) h'(t_2) + \dots + s''(t_n^-) h'(t_n) + \int_a^b s^{(k)}(t) h(t) dt$$

We want to achieve that  $s^{(k)}(t) = 0 \forall t$ . Thus for  $\xi : s^{(k)}(\xi) \neq 0$  construct  $h : h(\xi) \neq 0$  and  $ht \neq 0$

By  $J'(0) = 0$  follows that  $s^{(k)}(t) = 0 \forall t \neq t_i \forall i$ . In the other case, if  $s^{(k)}(\xi) > 0$  for  $\xi \in (t_i, t_{i+1})$ , we construct some reliable  $h$  with  $h(x)$  only in some neighborhood of  $\xi$ .

$$h(x) = \begin{cases} \exp\left(\frac{\delta^2}{x^2 - \delta^2}\right) & |x| < \delta \\ 0 & \text{otherwise} \end{cases} \quad \text{shifted by } \xi$$

$g^{(k)}$  is continuous, so also in some neighborhood of  $\xi + 0$ , so the integral  $> 0$ . This gives a contradiction. Furthermore it follows that  $s''(t_0^+) = s''(t_0^-) = 0$  and  $s''(t_i^-) = s''(t_i^+) \forall i = 1, \dots, n-1$ . Thus we get  $0 = J'(0)$  like above. Otherwise we create an appropriate  $h$ -function.

By the first conclusion, we get that  $s$  can be represented as polynomial of third degree in every interval  $[t_i, t_{i+1}]$ . By property 1, 2 and the second conclusion we get

$$\begin{aligned} s_i(t_i) &= s_{i-1}(t_i) \forall i = 1, \dots, n-1 & s'_i(t_i) &= s'_{i-1}(t_i) & s''_i(t_i) &= s''_{i-1}(t_i) \\ s'_0(t_0) &= s'_{n-1}(t_n) = 0 \end{aligned}$$

The spline function is a cubic, piecewise polynomial in  $C^2([a, b])$ .

It remains to show the uniqueness of  $s$ .

In every subinterval there are four freedoms to uniquely set  $s_i$ . Two conditions result by property  $s_i(t_i) = y_i$ ,  $s_i(t_{i+1}) = y_{i+1}$ . For the remaining two conditions we define  $c_i := s'_i(t_i)$ .

We now transform the interval  $[t_i, t_{i+1}]$  onto  $[0, 1]$  using  $x := \frac{t-t_i}{\tau_i}$ ,  $\tau_i = t_{i+1} - t_i$  and therefore  $s_i$  on  $q_i(x) = s_i(t) = s_i(t_i + \tau_i x)$ .

Also  $q_i$  are cubic polynomials and we get  $q_i(0) = y_i$ ,  $q_i(1) = y_{i+1}$ ,  $q'_i(0) = \tau_i c_i$ ,  $q'_i(1) = \tau_i c_{i+1}$ . This is a Hermitian interpolation problem with solution

$$q_i(x) = (1-x')^2 y_i + x^2 y_i t + 1 + x(1-x)[(1-x)(2y_i + \tau_i c_i) + x(2y_{i+1} - \tau_i c_{i+1})]$$

By this interpolation the continuity of  $s$  and  $s$  is guaranteed.

$$q''_i(0) = 6(y_{i+1} - y_i) - 2\tau_i(2c_i + c_{i+1}) \quad q''_i(1) = -6(y_{i+1} - y_i) + 2\tau_i(c_i + 2c_{i+1})$$

Because  $s''(t) = \frac{q''_i(x)}{\tau_i^2}$ , we get the conditions:

$$\begin{aligned} \frac{2c_0 + c_1}{\tau_0} &= \frac{3(y_1 - y_0)}{\tau_0^2} \\ \frac{c_{n-1} + 2c_n}{\tau_{n-1}} &= \frac{3(y_n - y_{n-1})}{\tau_{n-1}^2} \\ \frac{c_{i-1} + 2c_i}{\tau_{i-1}} + \frac{2c_i + c_{i+1}}{\tau_i} &= \frac{3(y_i - y_{i-1})}{\tau_{i-1}^2} + \frac{3(y_{i+1} - y_i)}{\tau_i^2} \forall i \neq 0, \dots, n \end{aligned}$$

Thus we retrieve a linear equation system:

$$\begin{bmatrix} \frac{2}{\tau_0} & \frac{1}{\tau_0} & & \\ \frac{1}{\tau_0} & \frac{2}{\tau_0} + \frac{2}{\tau_1} & & \\ 0 & \frac{1}{\tau_1} & & \\ & 0 & & \ddots \end{bmatrix}$$

We want to show that this linear equation system is invertible. It is strictly diagonal dominant.

**Theorem 4.17.2** (Gershgorin theorem). For  $B \in \mathbb{R}^{n \times n}$ ,  $\text{spec } B \subseteq \bigcup_{k=1}^n \left\{ z \in \mathbb{C} \mid |z - b_{kk}| \leq \sum_{l \neq k} |b_{kl}| \right\}$ .

For some arbitrary row, this circle lies strictly in the positive complex plane.

$$\Rightarrow 0 \notin \text{spec } A \iff A \text{ regular} \Rightarrow \exists! [c_0, \dots, c_m]^t : A \cdot c = \text{right-hand side}$$

Thus, a unique spline function  $s_i$  created by piecewise polynomials on  $[t_i, t_{i+1}]$  exists.

□

↓ This lecture took place on 2018/11/22.

$$s''(t_0) = 0 \quad s''(t_n) = 0$$

**Theorem 4.18** (Gershgorin circle theorem). Let  $B \in \mathbb{R}^{n \times n}$ .

$$\sigma(b) = \bigcup_{k=1}^n \left\{ z \in \mathbb{C} \mid |z - b_{kk}| \leq \sum_{l \neq k} |b_{kl}| \right\}$$

Applied to  $A$ , we can conclude that all eigenvalues of  $A$  lie in  $\mathbb{C}_{>0}$ .

$$\Rightarrow 0 \notin \text{spec}(A) \iff A \text{ is invertible} \Rightarrow \exists! : [c_0, \dots, c_n]^T : Ac = \text{right-hand side}$$

Thus, the existence of a unique spline function  $s$  follows which is realized by piecewise cubic polynomials on  $[t_i, t_{i+1}]$ .

**Remark.** The function  $s$  is called natural spline interpolant and the marginal conditions, created by minimizing the functional  $y[s]$ , are called natural marginal conditions.

Alternatively, instead we can use  $s''(t_0) = s''(t_n) = 0$ :

**Complete spline**  $s'(t_0) = y'_0, s'(t_n) = y'_n$

**Periodic spline**  $s'(t_0) = s'(t_n), s''(t_0) = s''(t_n)$

## 4.6 Spline spaces and B-splines

This is an extended spline concept of arbitrary order  $k$ .

**Definition 4.19.** Let  $\Delta := \{t_0, \dots, t_n\}$  is a grid of pairwise different points  $a = t_0 < t_1 < \dots < t_n = b$ . A spline of degree  $k - 1$  (of order  $k$ ) with respect to  $\Delta$  is a function  $s \in C^{k-2}([a, b])$ , that corresponds with a polynomial  $s_i \in \mathcal{P}_{k-1}$  in every interval  $[t_i, t_{i+1}]$  with  $i = 0, \dots, n - 1$ .

The space of splines of degree  $k - 1$  with respect to  $\Delta$  is denoted with  $S_{k,\Delta}$ .



**Remark.**  $S_{k,\Delta}$  is a real vector space and it is true that  $\mathcal{P}_{k-1} \subset S_{k,\Delta}$ . Furthermore the aborted powers of degree  $k-1$

$$(t - t_i)_+^{k-1} = \begin{cases} (t - t_i)^{k-1} & \text{if } t \geq t_i \\ 0 & \text{if } t < t_i \end{cases}$$

are contained in  $S_{k,\Delta}$ .

**Theorem.** The monomials and aborted powers create a basis

$$B := \{1, t, \dots, t^{k-1}, (t - t_1)_+^{k-1}, \dots, (t - t_{n-1})_+^{k-1}\}$$

of the spline space  $S_{k,\Delta}$ . Especially, it holds that

$$\dim(S_{k,\Delta}) = k + n - 1$$

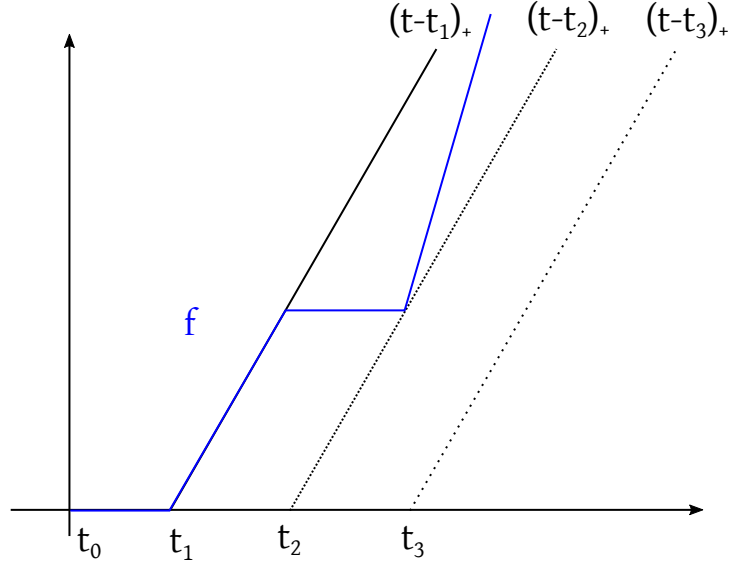


Figure 3: Linear combination of aborted powers of order 2

The basis  $B$  is not appropriate to represent a spline  $s \in S_{k,\Delta}$ , because

- the support of monomials  $t^i$  is entire  $\mathbb{R}$ ,  $\mathbb{R} \setminus \{0\}$ .
- for points with  $t_i \approx t_{i+1}$ , the aborted powers are *almost* linear dependent.

$$s(t) = \sum_{i=0}^{k-1} a_i t^i + \sum_{i=1}^{n-1} c_i (t - t_i)_+^{k-1}$$