

1. Computing with Matrices and Vectors

Machine Numbers

- $M \times M \rightarrow M$ (finite)
- Standard for machine numbers: floating-point representation

Fixed point representation

- Fixed decimal point
- Less accurate

Floating point representation

- $0,73125 \cdot 10^{12}$

$\underbrace{0,73125}_{\text{Mantissa}}$ $\underbrace{10}_{\text{Base}}$ $^{12} \text{Exponent}$

- Normalized, when Mantissa starts with 1
- Good for quickly changing the scale ($KM \rightarrow m$)
- Decimal point isn't fixed - defined through exponent

Error of rounding

- Must round to next nearest machine-number
- Absolute error: $|x - \tilde{x}|$ - Relative Error: $\frac{|x - \tilde{x}|}{|x|}$
- Machine precision "EPS": $10^{-\text{EPS}} = 10$
- Solutions for computing relative Error x
 1. Forward error: $x_{\text{exact}} - x_{\text{approximate}}$ (cannot be computed)
 2. Backward error: $b_{\text{exact}} - b_{\text{approximate}}$
 - $A x_{\text{exact}} = b \Rightarrow \text{compute } x_{\text{approximate}}$
 - $b_{\text{approximation}} = A x_{\text{approximate}}$
 - Backward error: $b_{\text{exact}} - b_{\text{approximate}}$

S.p.d matrix

- Symmetric positive definite matrix

Eigen

- MatrixXd (Dynamic double matrix)
- Matrix3d (3x3 double matrix)
- Matrix<double, dynamic, 5> (dynamic x 5 double matrix)
- MatrixXd y(6,8); (6x8 dynamic double matrix) (can be resized)
- Row-major: $\boxed{1 \ 2 \ 3 \ 4 \ 5 \ 6}$
$$\begin{array}{|c|c|c|} \hline 1 & 2 & 3 \\ \hline 4 & 5 & 6 \\ \hline \end{array}$$
- Column-major: $\boxed{1 \ 4 \ 2 \ 5 \ 3 \ 6}$

Computational effort

(= Number of basic operations that the algorithm is performing)

- Computational effort \neq runtime

(\rightarrow e.g. same number of operations, but row/column-major

- Dot-product / Skalarprodukt / Inneres-Produkt: $\Theta(n)$ ($x^T \cdot y$)
- Tensor-product / äusseres-Produkt: $\Theta(m \cdot n)$ ($x \cdot y^T$)
- Matrix-Vector-Product: $\Theta(n^2)$
- Matrix-product: $\Theta(m \cdot n \cdot k)$ ($A \cdot B$)
- Multiplying by a diagonal matrix: $\Theta(n)$

Hidden summation

- A, B matrices, we want to compute $y = \text{upper-triangular}(A \cdot B^T) \cdot x$

1. Approach: $y = \text{upper-triangular}(A \cdot B^T) \cdot x \Rightarrow \Theta(n^2 \cdot p)$

2. Approach: Decompose and regroup first: for # columns = 1

$$y = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nn} \end{bmatrix} \cdot \underbrace{\begin{bmatrix} 1 & & & & \\ & 1 & & & \\ & & 1 & & \\ & & & \ddots & \\ & & & & 1 \end{bmatrix}}_{\Theta(n)} \cdot \underbrace{\begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix}}_{\Theta(n)} \cdot \underbrace{\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}}_{\Theta(n)}$$

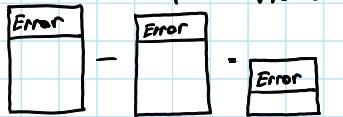
$\underbrace{\quad}_{\Theta(n) \text{ (Diagonal-Matrix-Vector)}} \quad \underbrace{\quad}_{\Theta(n) \text{ (Matrix-Vector)}} \quad \underbrace{\quad}_{\Theta(n) \text{ (Matrix-Vector)}} = z \in \Theta(n) \text{ (Matrix-Vector)}$

(\rightarrow For 1 column: $\Theta(n)$)

For p columns: $\Theta(p \cdot n)$

Cancellation

- Subtraction of two numbers, which are approximately the same size



- Errors get amplified

- Stable computation of roots:

$$\cdot p(\xi) = 1 \xi^2 + \alpha \xi + \beta$$

$$1. \xi_1 = \frac{-\alpha - \sqrt{\alpha^2 - 4\beta}}{2 \cdot 1}$$

$$2. \beta = \xi_1 \cdot \xi_2 \Rightarrow \xi_2 = \beta / \xi_1 \quad \left. \begin{array}{l} \text{Vieta's formula: } \\ x^2 + bx + c = (x-p)(x-q) \\ x^2 + bx + c = x^2 - (\underbrace{p+q})x + \underbrace{pq} \end{array} \right\}$$

Numerical stability

- When an algorithm is numerically stable, if you give me a x , the computer converts x to a slightly other \tilde{x} , because x couldn't be in the machine-numbers
- The error is close to EPS

- Algorithm \tilde{F} is outputting the exact solution to nearly the right answer, as if you rounded your input data

Mixed stability

$$\cdot \frac{\|x - \tilde{x}\|}{x} \leq O(\text{EPS})$$

- Numerical stability is only telling us, it is good at implementing $F(x)$, but $F(x)$ could be a bad function

Condition-number

- How an output value can change for a small change in the input (=sensitivity)
- Ratio of smallest/largest Eigenvalue of Matrix

2. Direct methods for solving linear systems of equations

Gauss Elimination

1. Elimination: $\mathcal{O}(n^3)$

2. Backsubstitution: $\mathcal{O}(n^2)$

$\left. \begin{array}{l} \\ \end{array} \right\} \mathcal{O}(n^3)$

- Solving triangular matrix only via Backsubstitution: $\mathcal{O}(n^2)$

- For each right-hand-side, the whole computation must be done another time

LU-Decomposition

$$- P \cdot A = L \cdot U$$

Permutation Matrix

Lower triangular with 0s

Upper triangular

1. LU-Decomposition: $A = LU \quad \mathcal{O}(n^3)$

2. Forward-Substitution: $Lz = b \quad \mathcal{O}(n^2)$

3. Backward-Substitution: $Ux = z \quad \mathcal{O}(n^2)$

- Overall complexity is the same as Gauss

- But we can solve for different right-hand sides, and only compute $A = LU$ once \Rightarrow New RHS $\mathcal{O}(n^2)$

- Low-rank-modification of A of an LSE:

• After solving $Ax = b$, we solve $\tilde{A}\tilde{x} = b$ with different A

• General rank-1-modification: $\tilde{A} = A + u \cdot v^T$

• Trick: use block-elimination: $\begin{bmatrix} A & u \\ v^T & -1 \end{bmatrix} \cdot \begin{bmatrix} \tilde{x} \\ \xi \end{bmatrix} = \begin{bmatrix} b \\ 0 \end{bmatrix}$

$$\textcircled{1} \quad \begin{cases} \tilde{A}\tilde{x} + u\xi = b \\ v^T \cdot \tilde{x} - \xi = 0 \end{cases} \quad \begin{array}{l} \xrightarrow{\quad \quad \quad} \\ \xrightarrow{\quad \quad \quad} \end{array} \quad \tilde{A}\tilde{x} + u \cdot v^T \cdot \tilde{x} = b \Leftrightarrow (A + u v^T) \cdot \tilde{x} = b$$

$$\textcircled{2} \quad \begin{cases} \tilde{x} = A^{-1} \cdot (b - u\xi) \\ v^T \cdot \tilde{x} - \xi = 0 \end{cases} \quad \begin{array}{l} \xrightarrow{\quad \quad \quad} \\ \xrightarrow{\quad \quad \quad} \end{array} \quad (v^T)^{-1} + \xi = A^{-1} \cdot (b - u\xi) \Leftrightarrow v^T \cdot A^{-1} \cdot (b - u\xi) - \xi = 0$$

$$\textcircled{3} \quad \begin{cases} \tilde{x} = A^{-1} \cdot b - A^{-1} \cdot u \cdot \xi \\ \Leftrightarrow \tilde{x} = A^{-1} \cdot b - A^{-1} \cdot u \cdot \frac{v^T \cdot A^{-1} \cdot b}{v^T \cdot A^{-1} \cdot u + 1} \end{cases} \quad \begin{array}{l} \xrightarrow{\quad \quad \quad} \\ \xrightarrow{\quad \quad \quad} \end{array} \quad \begin{aligned} &\Leftrightarrow (v^T A^{-1} u + 1) \cdot \xi = v^T \cdot A^{-1} \cdot b \\ &\Leftrightarrow \xi = \frac{v^T \cdot A^{-1} \cdot b}{v^T \cdot A^{-1} \cdot u + 1} \end{aligned}$$

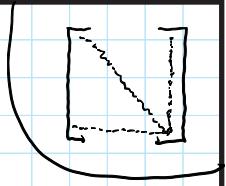
Assuming LU-Decomposition of A already available.

Arrow-Matrices

$$A_{\text{Arrow}} = \begin{bmatrix} D & C \\ b^T & \alpha \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ z \end{bmatrix} = \begin{bmatrix} b_1 \\ B \end{bmatrix} \iff \begin{bmatrix} D & C \\ b^T & \alpha \end{bmatrix} \begin{bmatrix} x_1 \\ z \end{bmatrix} = \begin{bmatrix} b_1 \\ B \end{bmatrix}$$

scalar scalar scalar

$$\begin{bmatrix} x_1 \\ z \end{bmatrix}$$



- A_{11} is diagonal, good for solving LSE
- We can also use block-matrices to compute $A \cdot A \cdot x$
 - Normal: Matrix-Matrix-Vector $O(n^3)$
 - Block-Matrix-Multiplication: Vector-Vector $O(n^2)$
- Block-LU-Decomposition: $A = L \cdot U = \begin{bmatrix} I & 0 \\ b^T D^{-1} & 1 \end{bmatrix} \cdot \begin{bmatrix} D & C \\ 0 & -b^T D^{-1} \cdot c \end{bmatrix}$
- Block-Gaussian-Decomposition:

$$\begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \rightarrow \begin{bmatrix} A_{11} & A_{12} \\ 0 & A_{22} - A_{21} \cdot A_{11}^{-1} \cdot A_{12} \end{bmatrix} \left| \begin{array}{c} b_1 \\ b_2 - A_{21} \cdot A_{11}^{-1} \cdot b_1 \end{array} \right.$$

$$\rightarrow \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \left| \begin{array}{c} A_{11}^{-1} (b_1 - A_{12} \cdot S^{-1} \cdot b_S) \\ S^{-1} \cdot b_S \end{array} \right. \quad \begin{array}{l} S = A_{22} - A_{21} \cdot A_{11}^{-1} \cdot A_{12} \\ b_S = b_2 - A_{21} \cdot A_{11}^{-1} \cdot b_1 \end{array}$$

$$\begin{cases} x_1 = A_{11}^{-1} (b_1 - A_{12} \cdot S^{-1} \cdot b_S) \\ x_2 = S^{-1} \cdot b_S \end{cases}$$

Sparse Matrices

- Most entries of matrix = 0
- Number of nonzero-Elements: $\text{nnz}(A)$
- Required memory should only depend on $\text{nnz}(A)$
- Example: Arrow-Matrix
- COO: · Triplet Format
 - (rows, columns, value)
 - if two locations are the same \rightarrow summation
- CRS: · Compressed-Row-Storage
- CCS: · Compressed-Column-Storage

- Example CRS:

| A | = | $\begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 \\ \hline 1 & 10 & 0 & 0 & -2 & 0 \\ 2 & 3 & 9 & 0 & 0 & 3 \\ 3 & 0 & 7 & 8 & 7 & 0 \\ 4 & 3 & 0 & 8 & 7 & 5 \\ 5 & 0 & 8 & 0 & 9 & 9 \\ 6 & 0 & 4 & 0 & 0 & 2 \end{matrix}$ | val-vector: 10 -2 3 9 3 7 8 7 3...9 13 4 2 -1 | col_ind-array: 1 5 1 2 6 2 3 4 1...5 6 2 5 6 | row_ptr_array: 1 3 6 9 13 17 20 |
|---|---|--|--|---|------------------------------------|
|---|---|--|--|---|------------------------------------|

- Build Sparse-Matrix

1. Use triplet format and at the end compress to CRS/CSS
2. Reserve enough storage in each row
↳ Because inserting into CRS/CSS is cost-intense, because all 3 arrays must be moved

3. Direct methods for linear least-squares problem

- Linear regression

- When we have n parameters and n measurements
- Problematic: Overfitting = When parameters fit good for the less measurements
↳ Get more measurements \Rightarrow Overdetermined system

- Normal equation

$$\begin{aligned} \cdot A \cdot x = b &\quad / \cdot A^T \\ \Leftrightarrow A^T \cdot A \cdot x = A^T \cdot b &\quad / (A^T \cdot A)^{-1} \\ \Leftrightarrow x = (A^T \cdot A)^{-1} \cdot b \end{aligned}$$

- QR-Decomposition

$$A = \underbrace{Q}_{\text{Orthogonal Unit vectors}} : R \quad \left. \begin{array}{l} \text{upper diagonal matrix} \\ \text{orthogonal unit vectors} \end{array} \right\}$$

$$\begin{aligned} \cdot \text{Full QR-decomposition: } [A] &= [Q] \cdot [R] \\ \cdot \text{Thin QR-decomposition: } [A] &= [Q] \cdot [R] \end{aligned}$$

1. Approach: Gram-Schmidt+

2. Approach: Householder-Transformation

- Singular-Value-decomposition

$$\cdot A = U \cdot \Sigma \cdot V^H$$

$$\cdot U = A \cdot \Sigma^{-1} \cdot V^{H-1}$$

• Σ : Singular values in ascending, diagonal form

• V : Normalized Eigenvectors

• Singular Value = $\sqrt{\text{Eigenvalue}}$

Approximation/Compression of Matrix via SVD

1. Compute SVD of A

2. Take only the first k rows

- If we ignore the small singular-values, it's O.K. because they don't contribute much to the matrix

Principle Component analysis

- Shows the trend of data, e.g. Stock-prices

1. Get some data $\begin{bmatrix} d_1 & d_2 & d_3 & \dots \end{bmatrix}$

2. Compute the mean (=Durchschnitt) and subtract it from each data-point $A = \begin{bmatrix} d_1 - \text{mean} & d_2 - \text{mean} & d_3 - \text{mean} & \dots \end{bmatrix}$

3. Calculate the covariance matrix $A \cdot A^T$

4. Compute the Eigenvectors $A \cdot A^T$ via SVD U -Matrix

Total least squares

- Also system Matrix A can be perturbed
- Because A and B is perturbed, we use the SVD-k-approximation trying to eliminate the measurement errors

- $\begin{bmatrix} \tilde{A} & \tilde{b} \end{bmatrix} \approx [A \ b]$
↳ Best n-rank-approximation

$$1. [A \ b] = U \cdot \Sigma \cdot V^T$$

$$2. [\tilde{A} \ \tilde{b}] = \sum_{i=1}^n \sigma_i \cdot (U_{:,i}) \cdot (V_{:,i})^T$$

$$3. [\tilde{A} \ \tilde{b}] \cdot (V)_{:,n+1} = 0 \quad / (V)_{:,n+1} \text{ is orthogonal}$$

$$\Leftrightarrow \tilde{A} \cdot (V)_{1:n, n+1} + \tilde{b} (V)_{n+1, n+1} = 0 \quad / -\tilde{b} (V)_{n+1, n+1}$$

$$\Leftrightarrow \tilde{A} \cdot (U)_{1:n, n+1} = -\tilde{b} (U)_{n+1, n+1} \quad / : (V)_{n+1, n+1}$$

$$-\tilde{b} = \tilde{A} \cdot \frac{(U)_{1:n, n+1}}{(U)_{n+1, n+1}} \quad / \tilde{A}^{-1} \quad / (-1)$$

$$\tilde{A}^{-1} \cdot \tilde{b} = -\frac{(U)_{1:n, n+1}}{(U)_{n+1, n+1}}$$

$$4. x = -\tilde{A}^{-1} \cdot \tilde{b}$$

$$x = -\frac{(V)_{1:n, n+1}}{(V)_{n+1, n+1}}$$

Constrained least-squares / Lagrange-Multipliers

- Least-Squares with side-condition

- $Ax = b$ in least squares }
- $Cx = d$ exact $\left. \begin{matrix} \\ \end{matrix} \right\} x \text{ is the same}$

- Solve via Lagrange-Multipliers:

$$1. L(y, m) := \frac{1}{2} \cdot \|Ay - b\|_2^2 + m^T(Cy - d)$$

• $x = \arg \min \max L(y, m) \Rightarrow$ Saddlepoint problem \Rightarrow Derivatives vanish

$$2. \frac{\partial f}{\partial x}(x, m) = A^T(Ax - b) + C^T \cdot m = 0$$

$$3. \frac{\partial f}{\partial m}(x, m) = Cx - d = 0$$

$$4. \text{Augmented normal equation: } \begin{bmatrix} A^T A & C^T \\ C & 0 \end{bmatrix} \begin{bmatrix} x \\ m \end{bmatrix} = \begin{bmatrix} A^T b \\ d \end{bmatrix}$$

5. Solve augmented normal equation, just as normal

4. Filtering Algorithms

- Impulse response:
 - $h = (h_0, \dots, h_{n-1})^T$
 - Output when the filter is fed with a single impulse of strength one, corresponding to input e_1
- Properties of a filter:
 - Finite length input \Rightarrow finite length output
 - Time invariant
 - Linearity
 - Output depends only on past and present
- n-th order filter: Impulse response is n long
- Given input signal x , output y can be computed as follows

$$\hookrightarrow y_k = \sum_{j=0}^{n-1} \underbrace{h_{k-j} \cdot x_j}_{\text{2. Variant}} = \sum_{j=0}^{n-1} \underbrace{h_j \cdot x_{k-j}}_{\text{2. Variant}}, h_j = 0 \text{ for } j < 0$$
- Discrete-convolution:
 - Given Input-signal $x = (x_0, \dots, x_{n-1})^T$
 - Given filter $h = (h_0, \dots, h_{n-1})^T$
 - $y = h * x \Leftrightarrow y_k = \sum_{j=0}^{n-1} h_{k-j} \cdot x_j$
- n-periodic:
 - When the Input-signal repeats itself after n inputs
 - When input is periodic, so is the output
- Discrete-periodic-convolution:
 - Output is also periodic
 - Satisfy $x_j = x_{j+n}$
 - We must zero-pad to match:

$$\hookrightarrow \text{length of } x = \text{length of } h$$
 - $$\begin{aligned} y_k &= p_k * x_k \\ &= \sum_{j=0}^{n-1} p_{k-j} \cdot x_j = \sum_{j=0}^{n-1} x_{k-j} \cdot p_j \end{aligned}$$

Periodic impulse response
 - Periodic impulse response: $p_j = \sum_{i \in \mathbb{N}} h_{j+i \cdot n} \mid j \in \{0, \dots, n-1\}$

$\hookrightarrow L\text{-periodic extension of } h$

- Periodic Signal can be represented as following

$$\begin{bmatrix} y_0 \\ \vdots \\ y_{n-1} \end{bmatrix} = \begin{bmatrix} p_0 & p_{n-1} & p_{n-2} & \cdots & \cdots & p_1 \\ p_1 & p_0 & p_{n-1} & & & \vdots \\ p_2 & p_1 & p_0 & \ddots & & \\ \vdots & \ddots & \ddots & \ddots & & \\ & & \ddots & \ddots & & \\ \vdots & & & \ddots & \ddots & \\ p_{n-1} & \cdots & & & p_1 & p_0 \end{bmatrix} \begin{bmatrix} x_0 \\ \vdots \\ x_{n-1} \end{bmatrix}$$

- Circulant-Matrix:
 - Have constant main/sub/super diagonals
 - Used to represent r_x
 - Can be determined by a vector

Discrete Fourier Transformation (DFT)

Fourier - Matrix

- All Circulant-matrices with the same dimensions have the same set of eigenvectors (=orthogonal trigonometric basis of \mathbb{C}^n)
 - Discrete convolution \Rightarrow discrete circular convolution \Rightarrow discrete fourier transformation
 - Computes the spectrum of a periodic signal
 - $W_n = e^{-\frac{j2\pi i}{n}}$ / $W_n^{ij} = \left(e^{-\frac{j2\pi i}{n}}\right)^{ij} = e^{-j2\pi ij}$
 - The scaled Fourier matrix $\frac{1}{\sqrt{n}} \cdot F_n$ is unitary

Fourier Transformation

- We get the spectrum by calculating: ($= F_n(y)$) ($y = \text{signal}$)
 $C = F_n \cdot Y$

Discrete convolution via DFT

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|
| 0 | $W^{0.0}$ | $W^{0.1}$ | $W^{0.2}$ | $W^{0.3}$ | $W^{0.4}$ | $W^{0.5}$ | $W^{0.6}$ | $W^{0.7}$ |
| 1 | $W^{1.0}$ | $W^{1.1}$ | $W^{1.2}$ | $W^{1.3}$ | $W^{1.4}$ | $W^{1.5}$ | $W^{1.6}$ | $W^{1.7}$ |
| 2 | $W^{2.0}$ | $W^{2.1}$ | $W^{2.2}$ | $W^{2.3}$ | $W^{2.4}$ | $W^{2.5}$ | $W^{2.6}$ | $W^{2.7}$ |
| 3 | $W^{3.0}$ | $W^{3.1}$ | $W^{3.2}$ | $W^{3.3}$ | $W^{3.4}$ | $W^{3.5}$ | $W^{3.6}$ | $W^{3.7}$ |
| 4 | $W^{4.0}$ | $W^{4.1}$ | $W^{4.2}$ | $W^{4.3}$ | $W^{4.4}$ | $W^{4.5}$ | $W^{4.6}$ | $W^{4.7}$ |
| 5 | $W^{5.0}$ | $W^{5.1}$ | $W^{5.2}$ | $W^{5.3}$ | $W^{5.4}$ | $W^{5.5}$ | $W^{5.6}$ | $W^{5.7}$ |
| 6 | $W^{6.0}$ | $W^{6.1}$ | $W^{6.2}$ | $W^{6.3}$ | $W^{6.4}$ | $W^{6.5}$ | $W^{6.6}$ | $W^{6.7}$ |
| 7 | $W^{7.0}$ | $W^{7.1}$ | $W^{7.2}$ | $W^{7.3}$ | $W^{7.4}$ | $W^{7.5}$ | $W^{7.6}$ | $W^{7.7}$ |

- DC can be written as multiplication with circulant matrix
 - We can use Fourier-transform instead

$$-c_k = \sum_{j=0}^{n-1} y_j \cdot w_n^{n \cdot j} \quad / F_n^{-1} = \frac{1}{n} \cdot \tilde{f_n}$$

$$-\gamma_i = \frac{1}{n} \cdot \sum_{k=0}^{n-1} c_k \cdot w_n^{-n \cdot j}$$

$$\begin{bmatrix} C_1 \\ C_2 \\ C_3 \\ C_4 \end{bmatrix} = \begin{bmatrix} w^0 & w^0 & w^0 & w^0 \\ w^0 & w^1 & w^2 & w^3 \\ w^0 & w^2 & w^4 & w^6 \\ w^0 & w^3 & w^5 & w^9 \end{bmatrix} \begin{bmatrix} S_1 \\ S_2 \\ S_3 \\ S_4 \end{bmatrix}$$

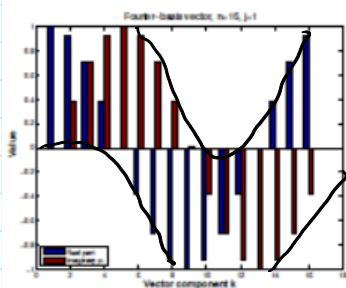
Fast-Fourier-Transformation tells you, how much periodicity you have in your signal, the more peaks, the more periodic.

- Normal DFT: $\Theta(n^2)$

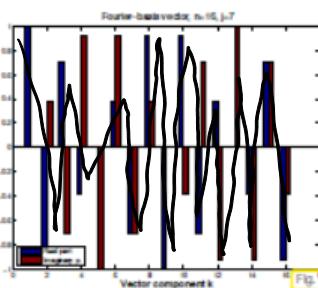
- Because $\omega_n^{2j} = \omega_m^j$ ($n=2m$), DFT can be implemented using Divide and Conquer in $\Theta(n \cdot \log(n))$

Frequency Filtering

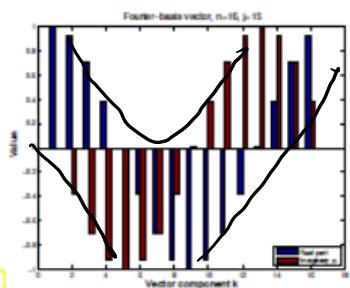
- k -th row of F_n : $(F_n x)_k = V_k^T \cdot x$ (V = trigonometric basis)
- If k is close to $0 \bmod(n) \Rightarrow$ slow frequency
- If k is in the middle $\bmod(n) \Rightarrow$ high frequency
- If k is close to $n \bmod(n) \Rightarrow$ slow frequency



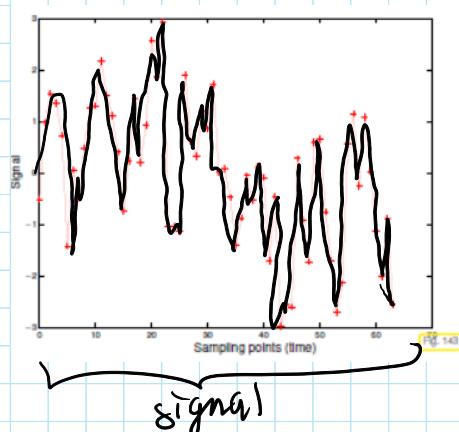
"slow oscillation/low frequency"



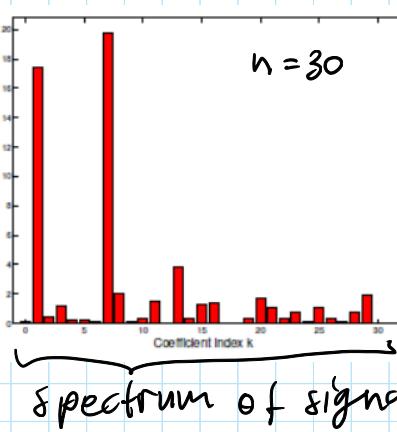
"fast oscillation/high frequency"



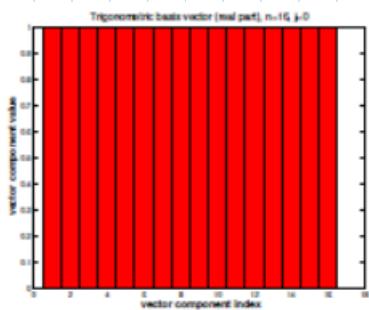
"slow oscillation/low frequency"



signal

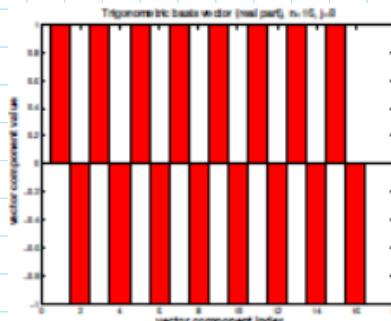


spectrum of signal



$\text{Re}(F_{16})_{:,0}$

$k \bmod(n) \Rightarrow$ slow



$\text{Re}(F_{16})_{:,8}$

$k \bmod(n) \Rightarrow$ fast

} "Low pass filter"
} Filter high frequency / noise out
} Cut them off

5. Data interpolation and Data fitting in 1D

- Given a set of data points (t_i, y_i) , find an interpolant function $f(x)$ with $f(t_i) = y_i$



Piecewise linear interpolation

- Just connect the points

Global polynomial interpolation

Lagrange interpolation

- We use Lagrange-Polynomials: $L_i(t) = \prod_{j=0, j \neq i}^n \frac{t - t_j}{t_i - t_j}$
- Interpolation: $p(t) = \sum_{i=0}^n y_i \cdot L_i(t)$
- ↳ Polynomial for a set of nodes is unique
- Polynomial of degree n is determined by ($n+1$)-data-points
- Evaluating $L_i(t)$ is $\Theta(n)$ with Horner-Scheme
- Evaluating $p(t)$ is $O(n^2)$, because we have to evaluate n polynomials
- Reducing computation for different y -values \Rightarrow Barycentric-Formula

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

$$\lambda_i = \frac{1}{(t_i - t_0) \dots (t_i - t_{i-1}) (t_i - t_{i+1}) \dots (t_i - t_n)}$$

$$\lambda = \prod_{j=0}^n (t - t_j) \cdot \sum_{i=0}^n \frac{\lambda_i}{t - t_i} \quad /: \sum_{i=0}^n \frac{\lambda_i}{t - t_i}$$

$$\prod_{j=0}^n (t - t_j) = \frac{1}{\sum_{i=0}^n \frac{\lambda_i}{t - t_i}}$$

$$p(t) = \frac{\sum_{i=0}^n \frac{\lambda_i}{t - t_i} \cdot y_i}{\sum_{i=0}^n \frac{\lambda_i}{t - t_i}}$$

$$p(t) = \prod_{j=0}^n (t - t_j) \cdot \sum_{j=0}^n \frac{d_j}{t - t_j} \cdot y_j \quad /: \sum_{j=0}^n \frac{d_j}{t - t_j} \cdot y_j$$

$$\prod_{j=0}^n (t - t_j) = \frac{p(t)}{\sum_{j=0}^n \frac{d_j}{t - t_j} \cdot y_j}$$

- Advantage: λ_i is independent of t and $y_i \Rightarrow$ precompute
- Computation of $\lambda_i, i = 0, \dots, n = O(n^2)$ (only once!)
- Cost of $\Theta(n)$ for every subsequent evaluation of p

- Update-friendly for adding/changing points: Aitken-Neville
 - We do interpolation on pieces of data, and gradually get to the whole interpolation

$$\bullet p_{k,k}(x) \equiv y_k \text{ ("constant polynomial")}, \quad k = 0, \dots, n,$$

$$\bullet p_{k,\ell}(x) = \frac{(x - t_k)p_{k+1,\ell}(x) - (x - t_\ell)p_{k,\ell-1}(x)}{t_\ell - t_k}$$

| $n =$ | 0 | 1 | 2 | 3 | |
|---------------------|---------------------|-----------------------------|-----------------------------|-----------------------------|--------------|
| t_0 | $y_0 =: p_{0,0}(x)$ | $\xrightarrow{} p_{0,1}(x)$ | $\xrightarrow{} p_{0,2}(x)$ | $\xrightarrow{} p_{0,3}(x)$ | $p_{0,3}(x)$ |
| t_1 | $y_1 =: p_{1,1}(x)$ | $\xrightarrow{} p_{1,2}(x)$ | $\xrightarrow{} p_{1,3}(x)$ | $\xrightarrow{} p_{1,4}(x)$ | $p_{1,4}(x)$ |
| t_2 | $y_2 =: p_{2,2}(x)$ | $\xrightarrow{} p_{2,3}(x)$ | $\xrightarrow{} p_{2,4}(x)$ | | |
| t_3 | $y_3 =: p_{3,3}(x)$ | $\xrightarrow{} p_{3,4}(x)$ | | | |
| + ↑ new point | $y_4 =: p_{4,4}(x)$ | | | | |

Newton Basis

- Good for many interpolation-points

$$- N_0(t) = 1$$

$$- N_1(t) = (t - t_0)$$

$$- N_n(t) = \prod_{i=0}^{n-1} (t - t_i)$$

- n = amount of points

$$- \text{Find interpolant } p(t) = \sum_{i=0}^n a_i N_i(t)$$

$$\cdot p(t_0) = a_0 + 0 + 0 + \dots + 0$$

$$\cdot p(t_1) = a_0 + a_1 N_1(t_1) + 0 + 0 + \dots + 0$$

$$\left. \begin{array}{l} \\ \\ \end{array} \right\} a_0 N_0(t_i) + a_1 N_1(t_i) + \dots + a_n N_n(t_i) = y_i$$

coefficients need to be determined

Lower triangular system

of equations

$$\left[\begin{array}{cccc|c} 1 & 0 & 0 & \dots & 0 \\ 1 & N_1(t_1) & 0 & \dots & 0 \\ 1 & N_2(t_2) & N_2(t_2) & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & N_n(t_n) & N_2(t_n) & \dots & N_n(t_n) \end{array} \right] \cdot \begin{bmatrix} a_0 \\ \vdots \\ a_n \end{bmatrix} = \begin{bmatrix} y_0 \\ \vdots \\ y_n \end{bmatrix}$$

\Rightarrow You are computing the same polynomial, only with different ways, because polynomial is unique.

Spline-Interpolation

mesh

- Each subinterval $[t_{i-1}, t_i]$ has an own polynomial of degree d
- Match the first $(d-1)$ derivatives at the nodes of the previous polynomial
- The first 3 spline-spaces
 - $d=1$: M -piecewise constant discontinuous function
 - $d=2$: M -piecewise linear continuous function
 - $d=3$: continuously differentiable M -piecewise quadratic function

→ These are relevant, because C^2 -function is smooth

Cubic spline interpolation ($d=3$)

- For each subinterval: $s_j = a_j + b_j \cdot t + c_j \cdot t^2 + d_j \cdot t^3$
 - $\hookrightarrow 4n$ coefficients to determine!
- 1. Interpolating condition: $s_j(t_{j-1}) = y_{j-1}$ & $s_j(t_j) = y_j$
 - $a_j + b_j \cdot t_{j-1} + c_j \cdot t_{j-1}^2 + d_j \cdot t_{j-1}^3 = y_{j-1}$
 - $a_j + b_j \cdot j + c_j \cdot j^2 + d_j \cdot j^3 = y_j$
- 2. Interpolating condition: Smoothness condition I (Match 1. derivative)
 - $s'_j(t_j) = s'_{j+1}(t_j) \quad (j=1, \dots, n-1)$
 - $b_j + 2c_j \cdot j + 3d_j \cdot j^2 = b_{j+1} + 2c_{j+1} \cdot j + 3d_{j+1} \cdot j^2$
- 3. Interpolating Condition: Smoothness condition II (Match 2. derivative)
 - $s''_j(t_j) = s''_{j+1}(t_j) \quad (j=1, \dots, n-1)$
 - $2c_j + 6d_j \cdot j = 2c_{j+1} + 6d_{j+1} \cdot j$
- 4. Interpolating condition: Natural / Simple Cubic-Spline
 - $s''_1(t_0) = 0$
 - $s''_n(t_n) = 0$

→ Altogether: LSE for coefficients

- Economical / Smart implementation of spline-interpolation

$$\begin{aligned} \cdot a_i &= y_{i-1} \\ \cdot b_j &= \frac{y_j - x_{j-1}}{h_j} - \frac{h_j(2 \cdot \sigma_{j-1} + \sigma_j)}{6} \\ \cdot c_j &= \frac{\sigma_{j-1}}{2} \\ \cdot d_j &= \frac{\sigma_j - \sigma_{j-1}}{6h_j} \end{aligned}$$

$$\hookrightarrow \sigma_{j-1} \cdot \frac{h_j}{6} + \sigma_j \cdot \frac{h_j + h_{j+1}}{3} + \sigma_{j+1} \cdot \frac{h_{j+1}}{6} = r_j$$

$$\boxed{h_j := t_j - t_{j-1}}$$

$\sigma_j := \text{unknown}$

$$\left[\begin{array}{ccc|c} \frac{h_1 + h_2}{3} & \frac{h_2}{6} & & \\ \frac{h_2}{6} & \frac{h_2 + h_3}{3} & \frac{h_3}{6} & \\ \vdots & \vdots & \vdots & \\ \frac{h_{n-1}}{6} & \frac{h_{n-1} + h_n}{3} & & \end{array} \right] \cdot \begin{bmatrix} \sigma_1 \\ \vdots \\ \vdots \\ \sigma_{n-1} \end{bmatrix} = \begin{bmatrix} r_1 \\ \vdots \\ \vdots \\ r_{n-1} \end{bmatrix}$$

$$\cdot \sigma_0 = 0 = \sigma_n$$

1. Solve tridiagonal system for σ_j

2. Compute a_j with y_{j-1}

3. Insert σ_j / h_j into $b_j / c_j / d_j$

6. Approximation of functions

- Given a function f , find a "simple"-approximation-fnc. \tilde{f}
- Minimize error $\|f - \tilde{f}\|$
- General idea:
 1. Get samples from the function f
 2. Interpolate points
 3. Get an approximation function \tilde{f}
- Estimate of error-norm $\|f - \tilde{f}\|$ for some norm (e.g. L^∞)
 - $\|f - \tilde{f}\| \leq C \cdot T(n)$ for $n \rightarrow \infty$ (greatest error)
 - Algebraic convergence with rate $p > 0$: $\exists p > 0: T(n) \leq n^{-p}$
 - Exponential convergence : $\exists \theta < p < 1: T(n) \leq g^n$
- (\Rightarrow We can get the convergence by measuring (n_i, ϵ_i) , where n_i is the polynomial degree and ϵ_i is the measured error)
- Algebraic convergence: $\log(\epsilon_i) \approx \log(c) - p \cdot \log(n_i) \Rightarrow$ linear in log-log-scale
- Exponential convergence: $\log(\epsilon_i) \approx \log(c) - \beta \cdot n_i \Rightarrow$ linear in lin-log-scale

$$\|f - L_T f\|_{L^\infty(I)} \leq \frac{\|f^{(n+1)}\|_{L^\infty(I)}}{(n+1)!} \max_{t \in I} |(t - t_0) \cdots (t - t_n)|$$

} Approximation of error

Taylor-Polynomial-Approximation

- $f(t) \approx \sum_{j=0}^k \frac{f^{(j)}(t_0)}{j!} \cdot (t - t_0)^j$
- Approximates $f(t)$ in a neighbourhood of t_0
(\Rightarrow Only provides approximation in a small interval)

Uniform B-spline Approximation

- Approximates $f(t)$ on a closed interval
- $p_n(t) = \sum_{j=0}^n f\left(\frac{j}{n}\right) \cdot \binom{n}{j} \cdot t^j (1-t)^{n-j}$
- As $n \rightarrow \infty$, the error $\|f - p_n\|_\infty \rightarrow 0$
- We need a high-degree-polynomial to precisely estimate $f(x)$
(\Rightarrow Inefficient!)

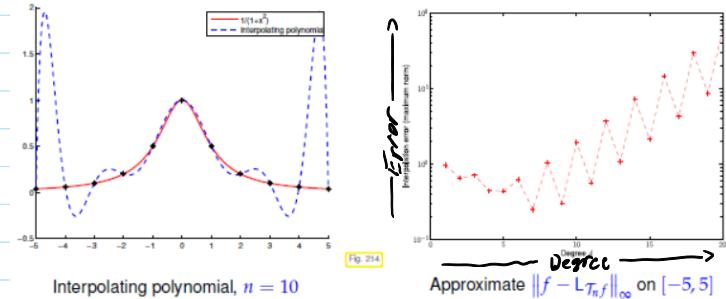
Equidistant Lagrange-Interpolation Approximation

- Take some nodes with the same distance to each other
- Use these nodes for Lagrange Interpolation
- As more points we use, the more precision we get ↗

↳ Runges - Experiment

- Strong oscillation of Lagrange-Interpolation near the endpoints

of the interval



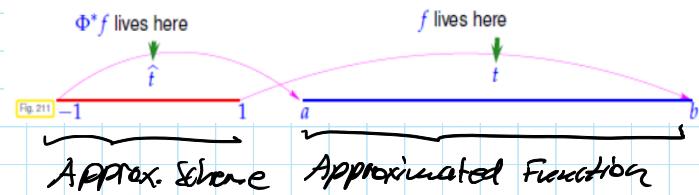
Transformation of polynomial approximation schemes

- If approx. scheme is only defined on $[-1, 1]$, then by the following trick, we can transfer any interval $[a, b]$ to $[-1, 1]$

$$\Phi: [-1, 1] \rightarrow [a, b]$$

$$\bar{\Phi}(t) = a + \frac{1}{2} (t+1) \cdot (b-a)$$

chebychev - interpolation



- Optimal choice of interpolation nodes independent of interpoland

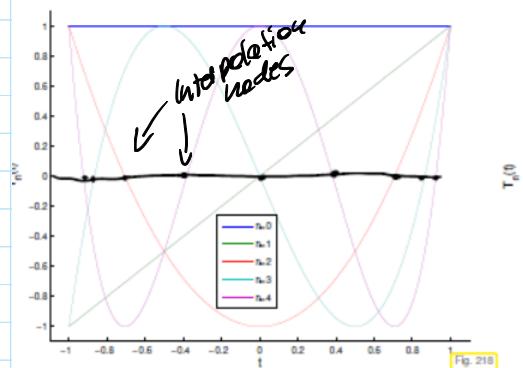
- The n^{th} chebychev polynomial is $T_n(t) = \cos(n \cdot \arccos(t))$ ($-1 \leq t \leq 1$)

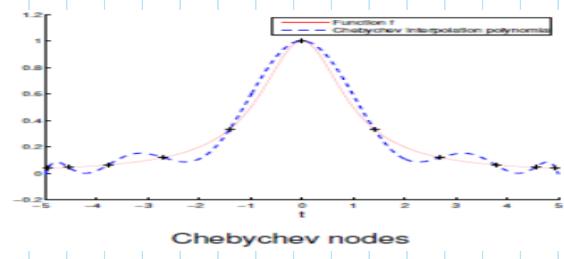
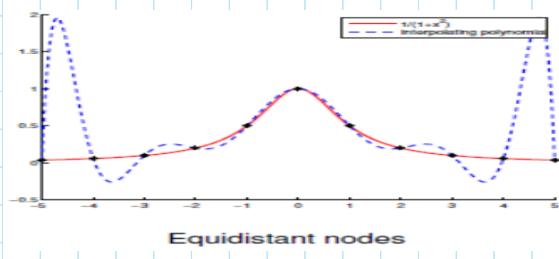
- Recursive definition: $T_0 = 1 / T_1(t) = t / T_{n+1}(t) = 2t \cdot T_n(t) - T_{n-1}(t)$

- The polynomials T_n minimize the supremum/ L^∞ -Norm

- The zeros of T_n are $t_k = \cos\left(\frac{(2k+1)}{2n} \cdot \pi\right)$

- When we use Chebychev nodes for polynomial interpolation, we call the resulting Lagrangian-approx. scheme "Chebychev interpolation"

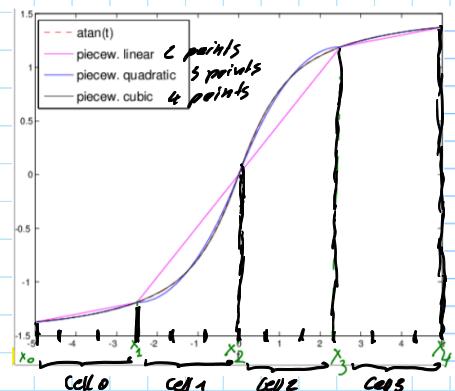




- We observe, that the Chebyshev nodes cluster at the endpoints of the interval, which prevents oscillation
- Because $\{T_0, \dots, T_n\}$ is a basis of P_n , we can write the polynomial as a Chebyshev-expansion: $p = \sum_{i=0}^{n+1} \alpha_i \cdot T_i$, $\alpha \in \mathbb{R}$
- Error approximation: $\|f - L_j(f)\|_{L^\infty(\Omega)} \leq \frac{2 \cdot 2^{n+1}}{(n+1)!} \cdot \|f^{(n+1)}\|_{L^\infty(\Omega)}$

Piecewise polynomial Lagrange-Interpolation

- We make a mesh M out of the interval $[a, b]$
 - x_j : nodes of the mesh M
 - $[x_{j-1}, x_j]$: cells of the mesh M
 - h_M : Maximal mesh width
 - n_j : local degree of the interpolating polynomial for each cell
 - T^j : set of local interpolating points (we need $n_j + 1$ of them)
- We need to satisfy: $f_{n_j}^j = f_{n_j}^{j+1}$ for continuity
- We can improve error by decreasing mesh width
- Error: $\|f - S\|_{L^\infty([x_0, x_m])} \leq \frac{h_M^{n+1}}{(n+1)!} \cdot \|f^{(n+1)}\|_{L^\infty}$



7. Numerical Quadrature

- Approximate $\int_a^b f(t) \cdot dt$ using only point evaluations of $f(t)$
 - General quadrature formula: $\int_a^b f(t) \cdot dt \approx Q_n(f) := \sum_{j=1}^n w_j^n \cdot f(c_j^n)$
- Affine Pullback

$$\int_a^b f(t) \cdot dt \approx \frac{1}{2}(b-a) \sum_{j=1}^n \widehat{w}_j \widehat{f}(\widehat{c}_j) = \sum_{j=1}^n w_j f(c_j) \quad \text{with} \quad c_j = \frac{1}{2}(1-\widehat{c}_j)a + \frac{1}{2}(1+\widehat{c}_j)b, \quad w_j = \frac{1}{2}(b-a)\widehat{w}_j.$$

weight Nodes G
EIR $\{a, b\}$

Quadrature by approximation scheme

$$\int_a^b f(t) \cdot dt \approx \int_a^b I[f(t_0), \dots, f(t_n)](t) \cdot dt = \sum_{j=1}^n f(t_j) \cdot \int_a^b I[c_j](t) \cdot dt = \sum_{j=1}^n f(t_j) \cdot w_j^n$$

Rechte Funktion Approximierte-Funktion mittels n -terter und Interpolation

- Interpolation \Rightarrow Approximation \Rightarrow Quadrature

Polynomial quadrature by approximation scheme

$$\int_a^b f(t) \cdot dt \approx \int_a^b P_{n-1}(t) \cdot dt = \sum_{i=0}^{n-1} f(t_i) \cdot \int_a^b L_i(t) \cdot dt = \sum_{i=0}^{n-1} f(t_i) \cdot w_i$$

Polynomial Nodes Weights

Lagrange
interpolant
of degree $(n-1)$

Newton-Cotes formulas

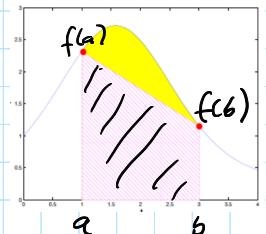
- Lagrange-Interpolation with equidistant nodes

$$t_j := a + \frac{b-a}{n-1} \cdot j$$

$$- n=2: \text{Trapezoidal-Rule: } \int_a^b f(t) \cdot dt \approx \frac{b-a}{2} (f(a) + f(b))$$

$$- n=3: \text{Simpson-Rule}$$

\Rightarrow Lagrange-Interpolation is unstable with equidistant nodes



Gauss Quadrature

- Order:
 - Quality measurement, for how good the rule is
 - Max. degree +1 of polynomials, for which the quadrature rule is guaranteed to be correct
 - Depends on the node set $T = \{t_1, \dots, t_n\}$
 - The maximal order of an n-point quad.-rule is $2n$
 ↳ With n points we can integrate a polynomial, using Lagrange-Interpolation, with max. degree of $2n-1$
- Legendre-Polynomials: $\int_{-1}^1 P_n(t) \cdot q(t) \cdot dt = 0 \quad / \quad P_n(1) = 1$
- The zeros of the Legendre-Polynomials form the Gauss-Points

Definition 7.3.29. Gauss-Legendre quadrature formulas

The n -point Quadrature formulas whose nodes, the **Gauss points**, are given by the zeros of the n -th Legendre polynomial (\rightarrow Def. 7.3.27), and whose weights are chosen according to Thm. 7.3.5, are called **Gauss-Legendre quadrature formulas**.

- Evaluation points from Legendre: $t_i, i=1, \dots, n$
- Weights for points from Lagrange: $w_j = \int_a^b L_{j-1}(t) \cdot dt$ $\left\{ \int_a^b f(t) dt = \sum_{j=1}^n w_j \cdot f(t_j) \right.$
 $\underbrace{\text{Polynom}}_{\text{Positive Constant}}$
- The Gauss-points and Gauss-weights can be computed with the Golub-Welsch-algorithm.
Legendre-roots Lagrange-Poly.

Lemma 7.3.42. Quadrature error estimates for C^r -Integrands

For every n -point quadrature rule Q_n as in (7.1.2) of order $q \in \mathbb{N}$ with weights $w_j \geq 0, j = 1, \dots, n$ we find that the quadrature error $E_n(f)$ for an integrand $f \in C^r([a, b]), r \in \mathbb{N}_0$, satisfies

$$\text{in the case } q \geq r: \quad E_n(f) \leq C q^{-r} |b-a|^{r+1} \|f^{(r)}\|_{L^\infty([a,b])}, \quad (7.3.43)$$

$$\text{in the case } q < r: \quad E_n(f) \leq \frac{|b-a|^{q+1}}{q!} \|f^{(q)}\|_{L^\infty([a,b])}, \quad (7.3.44)$$

with a constant $C > 0$ independent of n, f , and $[a, b]$.

} Error-estimation

Composite Quadrature

- Divide into small intervals
- Gauss at least as good as composite QF

8. Iterative methods for nonlinear Systems of equations

- Iterative methods for finding approximations to the solution of a nonlinear system

⇒ How to solve $f(x^*) = 0$ (=root)

Bisection algorithm

```

while( $n \leq n_{\max}$ ) do           / iterate  $n_{\max}$  times
     $b = \frac{t_r + t_r}{2}$           / compute the middle  $x^{(k)}$ 
    if ( $|f(b)| < T\delta_1$  ||  $|t_r - t_r| < T\delta_2$ ) / check if approximation
        return  $x^* = b$            / is sufficient

```

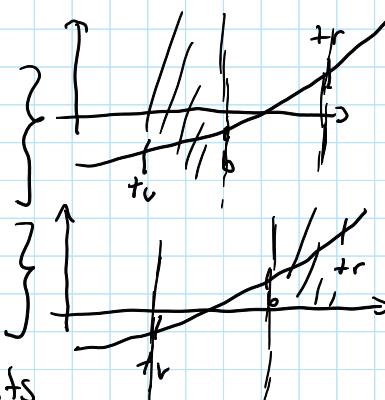
$n++$ / increase counter

if ($\text{sign}(f(t_r)) == \text{sign}(f(b))$)

$$t_r = b$$

if ($\text{sign}(f(t_r)) \neq \text{sign}(f(b))$)

$$t_r \approx b$$



- We always search, where a root exists

- Iterative error: $e^{(k)} = x^{(k)} - x^*$

- Rate of convergence: $x^{(k)} \rightarrow x^*$

• For bisection method: $|e^{(1)}| \leq \frac{1}{2} |a-b|$ } "linear type"
 $|e^{(2)}| \leq \frac{1}{2^2} |a-b|$ } convergence by factor $\frac{1}{2}$

Fixed point convergence in 1D

- $f(x) = 0 \Leftrightarrow \Phi(x) := f(x) + x$ (=finding fixed point)

- x^* is a fixed point of Φ : • $\Phi(x^*) = x^*$

$$\bullet f(x^*) = 0$$

- Requirement: Φ is Lipschitz stetig on $[a, b]$

$$\hookrightarrow \exists L > 0 \quad \forall x, y \in [a, b]: |\Phi(x) - \Phi(y)| \leq L \cdot |x-y|$$

• $L < 1$ (Φ is a contractive mapping)

- General construction of fixed-point iteration

1. Rewrite equivalently: $f(x) = 0 \Leftrightarrow \Phi(x) = x$

$$\hookrightarrow x \cdot e^{x-1} = y \quad x \cdot e^{x-1} = y$$

$$x \cdot e^{x-1} = 0 \quad x \cdot e^{x-1} = 0$$

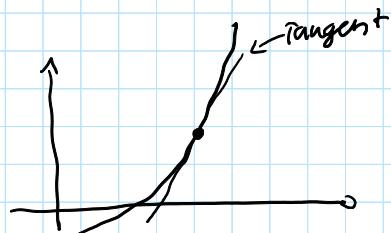
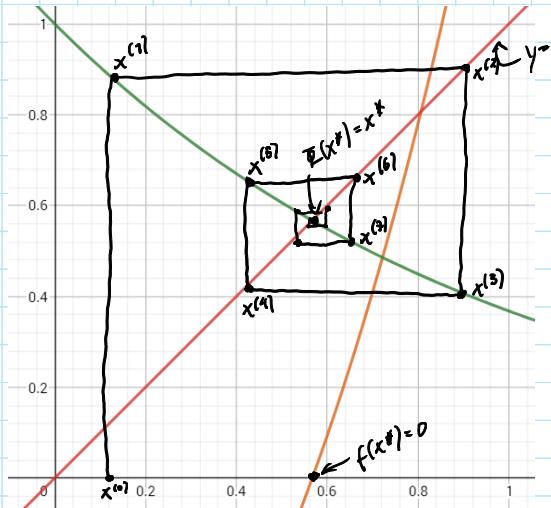
$$\frac{x}{y} = e^{-x} \quad \frac{x}{y} = \frac{1+x}{1+e^x}$$

$$(y = \frac{1+x}{1+e^x})$$

2. Start with initial guess $x^{(0)} = 0.1$

3. Use the fixed point iteration: $x^{(k+1)} = \Phi(x^{(k)})$

$$\hookrightarrow \underline{x^{(k+1)}} = e^{-x^{(k)}} \quad \underline{x^{(k+1)}} = \frac{1+x^{(k)}}{1+e^{x^{(k)}}}$$



Newton Iteration in 1D

- Gives us quadratic convergence

- Requirement: $f \in C^1$

- Intuition: Approx. around $x^{(k)}$: $f(x) \approx f(x^{(k)}) + (x - x^{(k)}) \cdot f'(x^{(k)})$

$$0 = f(x^{(k)}) + (x^{(k+1)} - x^{(k)}) \cdot f'(x^{(k)})$$

$$x^{(k+1)} = x^{(k)} - \frac{f(x^{(k)})}{f'(x^{(k)})}$$

- Newton's method converges quadratically, if $f'(x^*) \neq 0$

- Requires computation of $f'(x^{(k)})$ for each iteration \rightarrow costly!

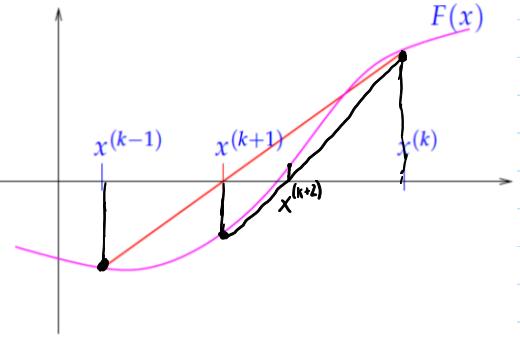
- Alternative: Secant method

- Replace $f'(x^{(k)})$ by an approximation

$$\cdot f'(x^{(k)}) \approx \frac{f(x^{(k)}) - f(x^{(k-1)})}{x^{(k)} - x^{(k-1)}}$$

• It's called a 2-point-method

↳ depends on the most 2 recent computed points



Nonlinear systems of equations

- n equations, n unknowns

- convergence in \mathbb{R}^n independent of choice of norm

• Local convergence: Only convergence in a neighbourhood \rightarrow good 1st guess

• Global convergence: 1st guess isn't important

- Fixed point iteration in higher dimensions \rightarrow Lipschitz-stability

• Requirements: - As in 1D, the mapping must be contractive

- $\bar{\Phi}$ converges at least local linearly, if

$$\|D\bar{\Phi}(x^*)\| < 1 \quad (D\bar{\Phi}(x) = \text{Jacobian-matrix})$$

- $\bar{\Phi}$ converges at least linearly, if

$$\text{everywhere } \|D\bar{\Phi}(x^*)\| < 1$$

- Termination criteria:

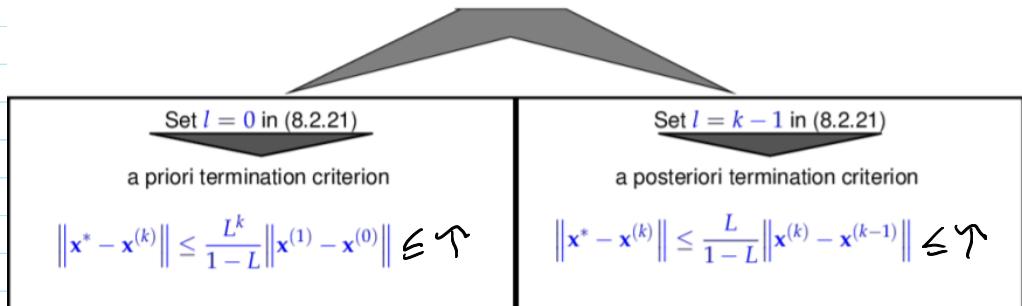
$$\cdot \|F(x^{(k)})\| \leq \gamma$$

$$\cdot \|x^{(k+1)} - x^{(k)}\| \leq \gamma \cdot \|x^{(k+1)}\|$$

• But: $\|F(x^{(k)}) - F(x^*)\| \text{ small} \not\Rightarrow \|x^{(k)} - x^*\| \text{ small}$

↳ ultimate goal: guarantee of form $\|x^{(k)} - x^*\| \leq \gamma$

$$\|x^* - x^{(k)}\| \leq \frac{L^{k-l}}{1-L} \|x^{(l+1)} - x^{(l)}\| \quad \{ \text{Approximation (8.2.21)}$$



compute, how many steps

compute at every step

- Newton's method in higher dimensions

- 1D: $f'(x^{(k)}) \neq 0$

- Higher dimensions: invertibility of Jacobian-Matrix

- $x^{(k+1)} := x^{(k)} - \underbrace{DF(x^{(k)})^{-1}}_{\text{Jacobian}} \cdot F(x^{(k)})$

- Newton correction: $= -DF(x^{(k)})^{-1} \cdot F(x^{(k)})$

↳ To compute Newton-correction: solve LSE: $DF(x^{(k)}) \cdot y = -F(x^{(k)})$

- Stopping criteria: $\|DF(x^{(k-1)})^{-1} \cdot F(x^{(k)})\| \leq \epsilon \cdot \|x^{(k)}\|$

- No need to compute $x^{(k)}$, even if $x^{(k-1)}$

- was a good enough approximation

- $DF(x^{(k)}) \approx DF(x^{(k+1)})$ at the best steps

• Problem: Newton correction is costly because of Jacobian

- Secant method in higher dimensions

- Broyden's quasi Newton Method

$$x^{(k+1)} = x^{(k)} - \left[J_{k-1} + \frac{F(x^{(k)}) \cdot (x^{(k)} - x^{(k-1)})^T}{\|x^{(k)} - x^{(k-1)}\|_2^2} \right]^{-1} \cdot F(x^{(k)})$$

$J_k^{-1} \approx DF(x^{(k)})^{-1}$

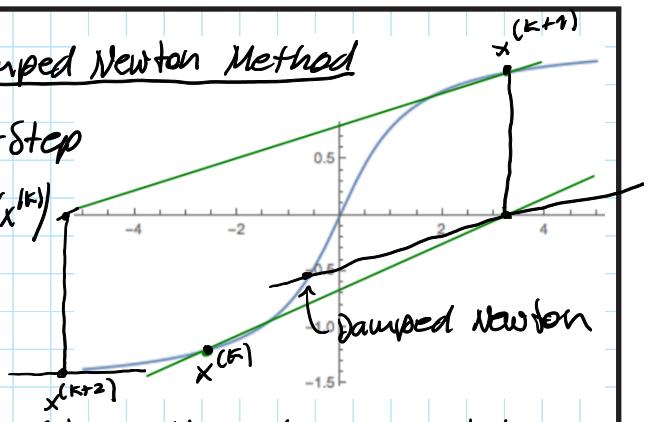
- Problem: Overshooting / solution: Damped Newton Method

- Idea: Damp the Newton-Step

$$x^{(k+1)} := x^{(k)} - \lambda^{(k)} \cdot Df(x^{(k)})^{-1} f(x^{(k)})$$

- $\lambda^{(k)}$: Damping factor of the
 k^{th} -Iteration

- Strategy: Largest possible $\lambda^{(k)}$ so that distance between iterate points decrease



Affine invariant damping strategy

Choice of damping factor: affine invariant natural monotonicity test [?, Ch. 3]:

$$\text{choose "maximal" } 0 < \lambda^{(k)} \leq 1: \quad \|\Delta\bar{x}(\lambda^{(k)})\| \leq (1 - \frac{\lambda^{(k)}}{2}) \|\Delta x^{(k)}\|_2 \quad (8.4.57)$$

where $\Delta x^{(k)} := D F(x^{(k)})^{-1} F(x^{(k)})$ → current Newton correction ,

$\Delta\bar{x}(\lambda^{(k)}) := D F(x^{(k)})^{-1} F(x^{(k)} + \lambda^{(k)} \Delta x^{(k)})$ → tentative simplified Newton correction .

1. set $\lambda^{(k)} = 1$

2. check natural monotonicity test

3. repeatedly take $\lambda^{(k)} = \frac{\lambda^{(k)}}{2}$ until NMT passes for

the first time and calculate $x^{(k)}$ with damping-factor

8.5 Unconstrained Optimization

- General question: How to find min/max of $f(x)$?
- Maximizing $f(x) \Leftrightarrow$ minimizing $-f(x)$

↳ Therefore only consider minimization problems

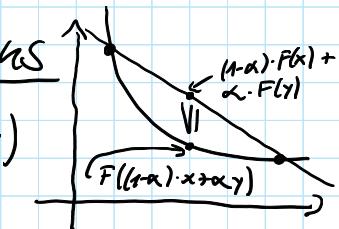
Minimization/Optimization with differentiable functions

- $F: \mathbb{R}^n \rightarrow \mathbb{R}$ differentiable
 - ∇F : direction of greatest increase
 - $-\nabla F$: direction of greatest descent
 - $\nabla F(x) = 0$: local/global max/min or saddle point
 - If F is twice differentiable check Hessian-Matrix at stationary point where $F(x) = 0$, because

$$F(x) \approx F(\bar{x}) + \underbrace{\nabla F(\bar{x})^\top (x - \bar{x})}_{=0 \text{ at stationary point}} + \frac{1}{2} (x - \bar{x})^\top \cdot H_F(\bar{x}) \cdot (x - \bar{x})$$
 - At stationary point: $F(x) \approx F(\bar{x}) + \frac{1}{2} \cdot (x - \bar{x})^\top \cdot H_F(\bar{x}) \cdot (x - \bar{x})$
 - $H_F(\bar{x})$ negative definite: \bar{x} local maximum
 - $H_F(\bar{x})$ indefinite: \bar{x} saddle point
 - $H_F(\bar{x})$ positive definite: \bar{x} local minimum
- ↳ Check by Cholesky factorization / Eigenvalues

Minimization/Optimization with convex functions

- Convex: $F((1-\alpha) \cdot x + \alpha \cdot y) \leq (1-\alpha) \cdot F(x) + \alpha \cdot F(y)$
- ↳ Local minimum $\stackrel{?}{=} \text{global minimum}$



Methods in 10

Newton's method for C^2 functions

- Applied to $f'(x)$ to search for $f'(x) = 0 \Rightarrow \text{min/max}$
- $x_{k+1} = x_k - \frac{f'(x)}{f''(x)}$

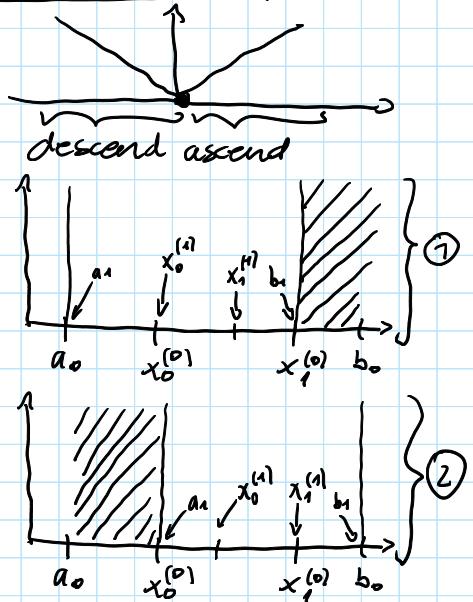
- Golden-Section-Search for non-differentiable functions

- Function must be unimodal:
- In each iteration, interval size is reduced by factor λ
 \hookrightarrow Linear-type-convergence
- In each iteration, only one point is computed
- while $|b-a| > \text{tolerance}$
- if $f_1 > f_0$ ①

$$\begin{aligned} b &= x_1 \\ x_1 &= x_0 \\ f_1 &> f_0 \\ x_0 &= a + (1-\lambda)(b-a) \\ f_0 &= f(x_0) \end{aligned}$$

if $f_0 \geq f_1$ ②

$$\begin{aligned} a &= x_0 \\ x_0 &= x_1 \\ f_0 &= f_1 \\ x_1 &= a + \lambda(b-a) \\ f_1 &= f(x_1) \end{aligned}$$



Methods in higher dimensions

- Gradient descend

- Greatest descend direction: $-\nabla F$

$$x^{(k+1)} = x^{(k)} - t^{(k)} \cdot \nabla f(x^{(k)})$$

1. Start with initial guess $x^{(0)}$

2. While $\|\nabla F\|_2 < \text{tol}$

3. Find stepsize $t^{(k)}$ through line search

$$x^{(k+1)} = x^{(k)} - t^{(k)} \cdot \nabla F(x^{(k)})$$

- How to find stepsize $t^{(k)}$ in each iteration?

\hookrightarrow with backtracking!

• Problem: line search in each iteration

- Helper: Backtracking Line search

$$\cdot F(x - \alpha \cdot \nabla F(x)) \leq F(x) - \alpha \cdot \|\nabla F(x)\|^2$$

$$\leq F(x) - \alpha \cdot \epsilon \cdot \|\nabla F(x)\|^2$$

1. iterate until good decrease is reached

2. start with $\epsilon = 1$, $\{x | x \in (0, 0.5], \beta \in (0, 1)$

3. while $F(x - \alpha \cdot \nabla F(x)) > F(x) - \alpha \cdot \epsilon \cdot \|\nabla F(x)\|^2$

$$4. \quad \epsilon = \beta \cdot \epsilon$$

- Newton's method

• As in 1D: if F is twice differentiable

$$\cdot x^{(k+1)} \approx F(x^{(k)}) + \nabla F(x^{(k)})^T \cdot (x - x^{(k)}) + \frac{1}{2} \cdot (x - x^{(k)}) \cdot H_F(x^{(k)}) \cdot (x - x^{(k)})$$

• Problem: computing H_F in each iteration

- BFGS-Method (=Quasi Newton)

• Approximate $H_F(x^{(k)})$ by B_k

• B_{k+1} is obtained from simple update of B_k

$$\cdot B_{k+1} \cdot \underbrace{(x^{(k+1)} - x^{(k)})}_{\delta^{(k)}} = \underbrace{\nabla F(x^{(k+1)}) - \nabla F(x^{(k)})}_{\gamma^{(k)}}$$

$$\left[\begin{array}{l} \cdot B_{k+1} \cdot \delta^{(k)} = \gamma^{(k)} \\ \cdot B_{k+1} = B_k + \alpha \cdot u \cdot u^T + \beta \cdot v \cdot v^T \\ \cdot u = \gamma^{(k)} \\ \cdot v = B_k \cdot \delta^{(k)} \end{array} \right]$$

$$\hookrightarrow B_{k+1} = B_k + \frac{\gamma^{(k)} \cdot \gamma^{(k)T}}{\gamma^{(k)T} \cdot \delta^{(k)}} - \frac{B_k \cdot \delta^{(k)} \cdot \delta^{(k)T} \cdot B_k^T}{\delta^{(k)T} \cdot B_k \cdot \delta^{(k)}}$$