

Friedrich-Alexander-Universität Erlangen-Nürnberg  
Technische Fakultät  
Department Informatik  
Lehrstuhl für Systemsimulation 10

Skript der Veranstaltung

# Algorithmik kontinuierlicher Systeme (AlgoKS)

gehalten im Sommersemester 2018  
von Prof. Dr. Ulrich Rüde

# VORBEMERKUNG

Diese(s) Skript enthält den Inhalt der Vorlesung „Algorithmik kontinuierlicher Systeme“ des Sommersemester 2018 bei Prof. Dr. Ulrich Rüde. Es wurde anhand von Mitschriften in  $\text{\LaTeX}$  gesetzt und erhebt daher weder Anspruch auf Korrektheit noch Vollständigkeit und ist *offensichtlich inoffiziell*. Bei Unstimmigkeiten und evtl. vorhandenen Fehlern bitte ich um eine Email an untenstehende Adresse. Dieses Skript stellt damit insbesondere **keine** offizielle Veröffentlichung des Lehrstuhl für Systemsimulation 10 am Department Informatik der Friedrich-Alexander-Universität Erlangen-Nürnberg dar.

**Florian Frank** — [florian.ff.frank@fau.de](mailto:florian.ff.frank@fau.de)  
**Philipp Hutterer** — [philipp.hutterer@fau.de](mailto:philipp.hutterer@fau.de)  
Version vom 14. April 2019

# INHALTSVERZEICHNIS

	Seite
<b>1 Direkte Verfahren für Lineare Gleichungssysteme</b>	<b>3</b>
1.1 LR-Zerlegung . . . . .	3
1.1.1 Gauß-Scherungen . . . . .	3
1.1.2 LR-Zerlegung mit Spaltenpivotsuche . . . . .	5
1.1.3 Vorwärts- und Rückwärtselimination . . . . .	6
1.1.4 Anwendungen . . . . .	6
1.1.5 Sonderfälle . . . . .	7
1.2 QR-Zerlegung . . . . .	7
1.2.1 Orthogonale Matrizen . . . . .	8
1.2.2 QR-Zerlegung mittels Givens-Rotationen . . . . .	8
1.2.3 QR-Zerlegung mittels Householderspiegelungen . . . . .	9
1.2.4 Orthonormalisierung . . . . .	11
1.3 Aufgaben . . . . .	11
<b>2 Lineare Ausgleichsprobleme</b>	<b>16</b>
<b>3 Matrizen - Datenstrukturen und praktische Verfahren</b>	<b>20</b>
3.1 Row And Column Storage (RaCS) . . . . .	20
3.2 Compressed Row Storage (CRS) . . . . .	20
3.3 Compressed Column Storage (CCS) . . . . .	21
3.4 Matrizen und Graphen . . . . .	21
3.5 Blockmatrizen . . . . .	22
3.6 Aufgaben . . . . .	22
<b>4 Diskretisierung und Quantisierung</b>	<b>24</b>
4.1 Diskretisierung . . . . .	24
4.1.1 Faltung . . . . .	24
4.1.2 Fourier-Transformation . . . . .	25
4.1.3 Abtasttheorem . . . . .	25
4.1.4 Aliasing und Anti-Aliasing . . . . .	25
4.2 Quantisierung . . . . .	26
4.2.1 Gleitpunktzahlen und Maschinenzahlen . . . . .	26
4.2.2 Quantisierung von Farbwerten mit dem <i>median-cut-algorithm</i> . . . . .	28
4.2.3 Vektorquantisierung . . . . .	29
4.3 Aufgaben . . . . .	31
<b>5 Rekonstruktion kontinuierlicher Daten — Interpolation und Approximation</b>	<b>37</b>
5.1 Interpolation im eindimensionalen Raum . . . . .	37
5.1.1 Lokale Verfahren . . . . .	37
5.1.2 Globale Verfahren — Polynominterpolation . . . . .	39
5.2 Multivariate Interpolation . . . . .	41
5.2.1 Lokale Interpolationsverfahren . . . . .	41

5.2.2	Globale Interpolationsverfahren im zweidimensionalen Raum . . . . .	43
5.3	Freiformmodellierung und Bézierkurven . . . . .	44
5.4	Aufgaben . . . . .	46
<b>6</b>	<b>Singulärwertzerlegung und Hauptkomponentenanalyse</b>	<b>49</b>
6.1	Elemente der linearen Algebra . . . . .	49
6.2	Singular-value decomposition – SVD . . . . .	51
6.3	Principal component analysis – PCA . . . . .	52
6.4	Aufgaben . . . . .	53
<b>7</b>	<b>Iterative Verfahren</b>	<b>55</b>
7.1	Allgemeine Theorie von iterativen Verfahren . . . . .	55
7.1.1	Newtonverfahren . . . . .	57
7.1.2	Sekantenverfahren . . . . .	58
7.1.3	Bisektionsverfahren . . . . .	58
7.1.4	Regula falsi . . . . .	58
7.2	Jacobi-Verfahren . . . . .	58
7.3	Gauss-Seidel-Verfahren . . . . .	59
7.4	SOR-Verfahren . . . . .	59
<b>8</b>	<b>Nichtlineare Optimierung</b>	<b>61</b>
8.1	Optimierung mittels Newtonverfahren . . . . .	61
8.2	Abstiegsverfahren . . . . .	61
8.2.1	Abstiegsverfahren des goldenen Schnittes . . . . .	62
8.2.2	Auffinden der Suchrichtung . . . . .	62
<b>9</b>	<b>Numerische Integration</b>	<b>64</b>
9.1	Monte-Carlo . . . . .	64
9.2	Newton-Cotes-Formeln . . . . .	64
9.2.1	Mittelpunktsregel . . . . .	64
9.2.2	Trapezregel . . . . .	65
9.2.3	Simpsonregel . . . . .	65
<b>A</b>	<b>Tabelle der Algorithmenaufwände und Fehlerabschätzungen</b>	<b>66</b>
A.1	Aufwände . . . . .	66
A.2	Fehlerabschätzungen . . . . .	67
<b>B</b>	<b>Python 3 — Zusammenfassung und Kodeschnipsel</b>	<b>68</b>
B.1	Zusammenfassung der wichtigsten Befehle . . . . .	68

# DIREKTE VERFAHREN FÜR LINEARE GLEICHUNGSSYSTEME

Wie seit den einführenden Mathematikvorlesungen bekannt ist, lassen sich lineare Gleichungssysteme im Allgemeinen als eine Gleichung der Form

$$Ax = b,$$

mit  $A \in \mathbb{K}^{n \times m}$ ,  $b \in \mathbb{K}^n$  und  $x \in \mathbb{K}^m$  darstellen. Vorerst wollen wir den Standardfall einer  $n \times n$ -Matrix betrachten.

Wie der gemeine Informatikstudierende bereits seit letztem Semester weiß, unterteilt man Algorithmen zum Lösen von linearen Gleichungssystemen in zwei Kategorien, **direkte** und *iterative* Verfahren. Letztere werden dann zusammen mit einer Daseinsberechtigung in Kapitel 7 näher besprochen. Wir beschäftigen uns damit jetzt mit den direkten Verfahren zum Lösen von linearen Gleichungssystemen, wobei hier die Lösung – bei Existenz – nach *endlich* vielen Schritten erreicht wird, vorausgesetzt man rechnet exakt<sup>1</sup>. Wir klassifizieren diese Verfahren nun weiter und schränken uns dabei zunächst auf (voll besetzte) **nicht singuläre** Systemmatrizen mit einer rechten Seite ein. Eine oft verwendete Strategie dabei ist es die Systemmatrix  $A$  so zu faktorisieren, dass sich die daraus entstehenden Teilprobleme leicht lösen lassen. Zu den wichtigsten und hier besprochenen Verfahren zählen die **LR-Zerlegung** ( $\rightarrow$  Kapitel 1.1), sowie die **QR-Zerlegung** ( $\rightarrow$  Kapitel 1.2).

## 1.1 LR-Zerlegung

Ziel der LR-Zerlegung ist eine Faktorisierung der Matrix  $A$  in eine linke (untere) und eine rechte (obere) Dreiecksmatrix. Der klassische Gauss-Eliminations-Algorithmus, welcher sich bemüht die Matrix sukzessive in Dreiecksgestalt zu überführen, baut dabei auf Transformationen über *Zeilenoperationen* und *Permutationen* auf. Zeilenoperationen ergeben sich dabei durch die Multiplikation mit geeigneten Matrizen, wir wollen diese dann auch Gauß-Scherung nennen.

### 1.1.1 Gauß-Scherungen

Bevor wir Gauß-Scherungen definieren können, müssen wir uns Matrizen definieren, die zu einer Zeile  $i$  einer Matrix das  $\alpha$ -fache einer anderen Zeile  $j$  addieren. Wir definieren also

$$N_{ij}(\alpha) = (\eta_{ij})_{rc} = \begin{cases} 1 & , \text{ wenn } r = c \\ \alpha & , \text{ wenn } r = i \text{ und } c = j. \\ 0 & \text{ sonst} \end{cases}$$

Ebenfalls erkennen wir leicht, dass aufgrund der Gestalt von  $N_{ij}(\alpha)$  sich das Inverse der Matrix leicht berechnen lässt mit

$$N_{ij}^{-1}(\alpha) = N_{ij}(-\alpha).$$

<sup>1</sup>Das das „exakte“ Rechnen zum Problem werden kann, haben wir bereits am einführenden Beispiel gesehen ... (genauerer siehe später)



## 1.1.2 LR-Zerlegung mit Spaltenpivotsuche

Wir erkennen selbstverständlich das größte Manko an dem eben entwickelten Verfahren, denn es funktioniert nur, wenn die Pivotelemente  $a_{ii} \neq 0$  sind. Sollte dem nicht so sein, so müssen zwingend Zeilenvertauschungen durchgeführt werden. Diese werden durch Multiplikationen mit Permutationsmatrizen dargestellt. Man sollte die Pivotsuche allerdings vor allem wegen der numerischen Stabilität der Verfahren immer durchführen, man kann auf sie verzichten, sofern symmetrisch positiv definite Matrizen vorliegen.

### Eine mögliche Strategie der Pivotsuche

Man sucht in der Restspalte das **betragsmäßig größte** Element und versucht dies durch das Vertauschen von Zeilen auf die **Diagonale** zu bringen.

Sei  $A$  eine  $n \times n$ -Matrix und  $(e_i)_{1 \leq i \leq n}$  die Standardbasisvektoren, so folgen die Permutationsmatrizen, welche wir mit  $P_{i,j}$  – für das Vertauschen der  $i$ -ten mit der  $j$ -ten Zeile – bezeichnen wollen, dem Schema

$$P_{i,j} = \begin{pmatrix} e_1, & \dots, & e_j, & \dots, & e_i, & \dots, & e_n \\ \text{1. Spalte} & \dots & \text{i. Spalte} & \dots & \text{j. Spalte} & \dots & \text{n. Spalte} \end{pmatrix}.$$

Wir erkennen dann, dass sich Permutationen und Gauß-Scherungen gegenseitig abwechseln, aufgrund der einfachen Struktur der  $P_{i,j}$  diese sich aber „vorziehen“ lassen. Im Allgemeinen gilt somit, dass sich die rechte obere Dreiecksmatrix durch

$$R = L_{n-1} \cdot P_{n-1,a} \cdots L_1 \cdot P_{1,i} \cdot A$$

berechnet und damit dann

$$A = P_{1,i} \cdot L_1^{-1} \cdots P_{n-1,a} \cdot L_{n-1}^{-1} = \prod_{i=1}^{n-1} P_{i,r_i} \cdot \prod_{i=1}^{n-1} \tilde{L}_i \cdot R = P \cdot L \cdot R,$$

wobei sich dann  $\tilde{L}_i$  aus  $L_i^{-1}$  durch das Vertauschen der jeweiligen Elemente in der Spalte hervorgeht, sobald man ein  $P$  vorzieht.

### Verfahren 1.1 (Allgemeines Vorgehen)

Für  $i = 1, \dots, n - 1$  mache:

- ① Finde das betragsmäßig größte Element in der  $i$ -ten (Rest-)spalte. Wir wollen dann die Zeile des Elements mit  $\ell$  bezeichnen.
- ② Vertausche nun in  $A$  und  $b$  die beiden (Rest-)zeilen  $i$  und  $\ell \rightsquigarrow$  Merke Permutation in Permutationsmatrix  $P_{i,\ell}$
- ③ Schreibe die  $i$ -te Zeile von  $R$  und  $i$ -te Spalte von  $L$  und führe die Gaußelimination durch.  $\rightsquigarrow$  Erhalte dadurch ein aktualisiertes  $A$ .

Man löst das System  $Ax = b$  dann mittels  $Lx = P^{-1}b$ , wobei die Inverse von  $P$  schnell bestimmt ist, man dreht dazu einfach das Produkt um.

### 1.1.3 Vorwärts- und Rückwärtselimination

#### 1.1.3.1 Vorwärtselimination

$$x_i = \frac{c_i - \sum_{j=1}^{i-1} l_{ij} \cdot x_j}{l_{ii}} \quad \text{für } 1 \leq i \leq n$$

```

1 def forward_substitute(L, b):
2     """Return y such that L @ y = b"""
3     ...
4     for j in ...:
5         ...
6         for i in range(j):
7             ...
8             ...
9         ...
10    pass

```

#### 1.1.3.2 Rückwärtselimination

$$x_i = \frac{c_i - \sum_{j=i+1}^n r_{ij} \cdot x_j}{r_{ii}} \quad \text{für } 1 \leq i \leq n$$

```

1 def backward_substitute(R, b):
2     """Return x such that R @ x = y"""
3     ...
4     for j in ...:
5         ...
6         for i in range(j):
7             ...
8             ...
9         ...
10    pass

```

### 1.1.4 Anwendungen

#### 1.1.4.1 Determinanten

Benötigt man die Determinante einer  $n \times n$ -Matrix (und nur einer solchen, da Determinanten nur für quadratische Matrizen definiert sind), so löst man dies über eine LR-Zerlegung (ohne Pivotsuche). Man nutzt dazu die Formeigenschaften von L und R aus, indem man feststellt, dass

$$\det(A) = \det(L \cdot R) = \det(L) \cdot \det(R) = 1 \cdot \prod_{i=1}^n r_{ii}.$$

Für die Determinante mit Pivotsuche ergibt sich ein Vorfaktor:

$$\det(A) = \det(P \cdot L \cdot R) = \det(P) \cdot \det(L) \cdot \det(R) = \pm 1 \cdot 1 \cdot \prod_{i=1}^n r_{ii}.$$

Die Determinante von P ist  $\pm 1$  je nach gerader (+1) oder ungerader (-1) Anzahl der Vertauschungen (Zeile- oder Spalte), die benötigt werden, um daraus eine Einheitsmatrix zu machen ( $\rightarrow$  siehe C1).

#### 1.1.4.2 Effizientes Lösen mehrerer Gleichungssysteme mit gleicher Koeffizientenmatrix A

Der allgemeine Ansatz ist es, nur **einmal** die LR-Zerlegung von A zu bestimmen (im Aufwand  $\mathcal{O}(n^3)$ ) und danach für jede rechte Seite eine Lösung mittels Vorwärts- und Rückwärtssubstitution zu bestimmen. Sind jedoch **alle rechten Seiten** zu Beginn bereits bekannt, so werden sie in einer m-spaltigen Matrix B zusammengefasst. Man löst dann  $AX = B$ , wobei  $X, B \in \mathbb{R}^{n \times m}$  durch die LR-Zerlegung der „erweiterten Matrix“  $(A \mid B)$  und erhält somit dann L und  $(R \mid B^*)$ . Damit ist X bestimmbar durch Rückwärtssubstitution von  $RX = B^*$ , die Vorwärtssubstitution entfällt hierbei.



### Berechnung der Inversen

Das Berechnen der Inversen mag zwar im Allgemeinen ein Kunstfehler und Entlassungsgrund sein, manchmal möchte man diese jedoch wirklich ohne größere Umstände ausrechnen. Im Normalfall gelingt dies ohne Probleme mit obiger Formel, setzt man für B einfach die Einheitsmatrix ein. X ist dann die inverse Matrix zu A. Aufgrund von sehr sehr starken Rundungsfehlern, welche beim Berechnen der inversen durchaus auftreten können, ist dieser Ansatz allerdings nicht als Zwischenschritt zu empfehlen!

## 1.1.5 Sonderfälle

### 1.1.5.1 LR-Zerlegungen für nicht vollbesetzte Matrizen — Bandmatrizen

Bei Bandmatrizen beschränken wir die Elimination auf diejenigen Elemente, welche **innerhalb** des Bandes und **unterhalb** der Diagonalen liegen. Es ergibt sich somit ein Aufwand bei Bandbreite  $b$  und Matrixordnung  $n$  von  $\mathcal{O}(b^2n)$ .

Bei tridiagonalen Matrizen ist die LR-Zerlegung mit  $\mathcal{O}(n)$  Operationen lösbar – mit insgesamt jeweils  $n - 1$  Divisionen, Additionen und Multiplikationen -.

$$\begin{aligned} r_1 &= b_1 & l_1 &= \frac{a_1}{r_1} \\ r_i &= b_i - l_{i-1} \quad \text{für } 1 \leq i \leq n \cdot c_{i-1} & l_i &= \frac{a_i}{r_i} \quad \text{für } 1 \leq i < n \end{aligned}$$

### Probleme bei dünnbesetzten Matrizen

Allgemein erhalten wir bei dünnbesetzten Matrizen einen sogenannten „Fill-In“, bei dem bereits vorhandene Nullen in der Matrix durch die Elimination zerstört werden. Dies treibt selbstverständlich neben der Zahl der Operationen auch den Speicherbedarf in die Höhe. Mittels sog. „*sparse matrix solvers*“ soll der Fill-In durch eine geeignete Permutation der Matrix minimiert werden. Gleichzeitig ist das Thema längst noch nicht abgeforscht, insbesondere die Entwicklung solcher Algorithmen für Parallelrechner ist weiterhin ein wichtiges Thema.

### 1.1.5.2 Cholesky-Zerlegung für positiv definite Matrizen

Ist A symmetrisch und positiv definit, so kann man A auch faktorisieren in

$$A = L \cdot D \cdot L^T,$$

wobei das L aus der LR-Zerlegung und das D die Diagonalanteile des R darstellt. Aufgrund der positiven Definitheit kann man schlussfolgern, dass D in der Diagonale nur positive Elemente hat, in Alternative lässt sich somit die Faktorisierung als

$$A = L \cdot D \cdot L^T = L \cdot D^{1/2} \cdot D^{1/2} \cdot L^T = \tilde{L} \cdot \tilde{L}^T$$

schreiben. Der schlussendliche Algorithmus nutzt dann aus, dass die beiden Dreiecksfaktoren gleich sind, man braucht damit später nur einen der beiden explizit berechnen und spart damit grob die Hälfte an Operationen.

Ein weiterer Vorteil ist das Nichtauftreten von 0-Pivots und die damit einhergehende Stabilität des Verfahrens ohne Pivotsuche, was eine Pivotsuche überflüssig macht.

Der einzige Nachteil ist die starke Einschränkung auf die Matrizenfamilie.

## 1.2 QR-Zerlegung

Ziel der QR-Zerlegung ist eine Faktorisierung der Matrix A in eine orthogonale Matrix Q und eine rechte obere Dreiecksmatrix R.

### 1.2.1 Orthogonale Matrizen

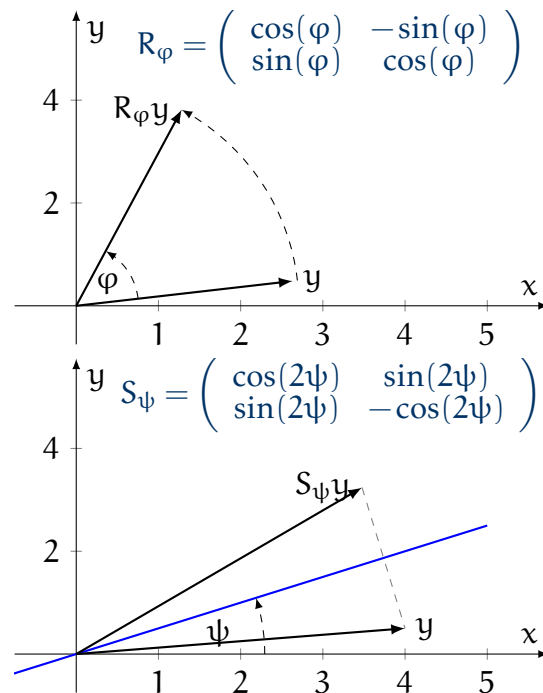
#### Definition 1.1 (Orthogonale Matrix)

Eine Matrix  $Q \in \mathbb{R}^{n \times n}$  heißt **orthogonal**, falls eine der folgenden äquivalenten Bedingungen erfüllt ist:

- $Q^T Q = E_n = Q Q^T$
- $Q^{-1} = Q^T$
- Die Spalten oder Zeilen von  $Q$  bilden eine Orthonormalbasis
- Die Abbildung  $Q$  ist winkel- und längentreu.
- $Q$  erhält das Skalarprodukt, es gilt also:

$$Qx \circ Qy = x \circ y$$

Die Determinante einer orthonormalen Matrix ist immer  $\pm 1$ .



orthogonale Matrizen im zweidimensionalen Fall: Rotationsmatrizen oder Spiegelungen

### 1.2.2 QR-Zerlegung mittels Givens-Rotationen

Während wir uns Rotationsmatrizen im zweidimensionalen bereits im obigen Beispiel überlegt haben, wollen wir nun den Aufbau einer solchen Rotationsmatrix auf den  $n$ -dimensionalen Raum verallgemeinern. Im Allgemeinen betrachten wir dabei allerdings weiterhin bloß endlichdimensionale Räume. Wir definieren die Givensrotationsmatrix an der Achse  $ij$  mit dem Winkel  $\varphi$  mit

$$J_{ij}(\varphi) = (g_{kl})_{k,l \leq n} = \begin{pmatrix} 1 & \dots & 0 & \dots & 0 & \dots & 0 \\ \vdots & \ddots & \vdots & & \vdots & & \vdots \\ 0 & \dots & c & \dots & -s & \dots & 0 \\ \vdots & & \vdots & \ddots & \vdots & & \vdots \\ 0 & \dots & s & \dots & c & \dots & 0 \\ \vdots & & \vdots & & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & \dots & 0 & \dots & 1 \end{pmatrix} = \begin{cases} 1 & \text{für } k = l, k \neq i, k \neq j \\ c = \cos(\varphi) & \text{für } k = l = j, k = l = i \\ s = \sin(\varphi) & \text{für } k = j, l = i \\ -s = -\sin(\varphi) & \text{für } k = i, l = j \\ 0 & \text{sonst} \end{cases}.$$

**Nutzen der Givensrotationsmatrizen** Mit einer Givensrotation  $J_{ij}$  kann man – unter dem richtigen Winkel  $\varphi$  – Nullen in eine Matrix an gewissen Stellen einführen. Dabei ist die Rotationsmatrix  $J_{ij}(\varphi_{\text{korrekt}})$  für die Einführung der 0 an der Stelle  $a_{ij}$  der Matrix  $A$  zuständig.  $c$  und  $s$  berechnen sich dann über

$$c = \frac{\text{sgn}(a_{jj}) \cdot a_{jj}}{\sqrt{a_{jj}^2 + a_{ij}^2}} \quad \text{und} \quad s = -\frac{\text{sgn}(a_{jj}) \cdot a_{ij}}{\sqrt{a_{jj}^2 + a_{ij}^2}}.$$

Für Givensrotationsmatrizen gilt im Allgemeinen:

$$J_{ij}^{-1}(\varphi) = J_{ij}(-\varphi) = J_{ij}^T(\varphi)$$

**Umsetzung** Wichtig bei der Umsetzung der QR-Zerlegung mit Givensrotationsmatrizen ist es eine Reihenfolge zu schaffen, in der einmal geschaffene Nullen *nicht wieder zerstört* werden. Sei  $A \in \mathbb{R}^{m \times n}$  und  $n^* = \min\{m-1, n\}$ , so ist eine QR-Zerlegung der Matrix  $A$  durch

$$\begin{aligned} A &= \underbrace{J_{2,1}^T \cdot \dots \cdot J_{m,n^*}^T}_{=Q} \cdot \underbrace{J_{m,n^*} \cdot \dots \cdot J_{2,1}}_{=R} \cdot A \\ &= Q \cdot R \end{aligned}$$

gegeben. Dabei ist das Produkt aller Rotationen selbst eine orthogonale  $m \times m$ -Matrix, was meistens nicht explizit mitberechnet werden muss.

```

1 def givens_rotation(A, i, j):
2     # You will do this in exercise 3 ...
3     pass
4
5
6 def qr_decompose(A):
7     # You will do this in exercise 3 ...
8     pass

```

**Eigenschaften der Givensrotation** Das Verfahren ist im Vergleich zur LR-Zerlegung auch ohne Pivotsuche numerisch stabil und hat insbesondere eine bessere Konditionierung als die LR-Zerlegung mit

$$\kappa_2(A) = \kappa_2(R) \quad \text{und} \quad \kappa_2(Q) = 1.$$

Ebenfalls ist es einfach auf Matrizen  $\in \mathbb{R}^{m \times n}$  und  $m > n$  anwendbar, wird also auch oft zur Ausgleichsrechnung verwendet.

Das Verfahren mit Givensrotationen ist per se ebenfalls leichter parallelisierbar als das Verfahren der Householderspiegelungen, aber dafür etwas schwerer zu implementieren.

#### Aufwandsverbesserung

Eines der Hauptprobleme der QR-Zerlegung mit Givensrotationsmatrizen ist der deutlich höhere Aufwand im Vergleich zur LR-Zerlegung. (4-mal so hoch)

- Ein Verbesserungsvorschlag ist es die Rotationen durch Skalierungen so zu modifizieren, dass die Diagonalelemente alle gleich 1 sind. Die so erfolgten Skalierungen werden dann separat aufmultipliziert, wodurch man dann insgesamt in etwa die Hälfte des Rechenaufwands spart. Dieses Verfahren ist auch unter dem Namen der *rationalen Givensrotationen* bekannt.

### 1.2.3 QR-Zerlegung mittels Householderspiegelungen

Wir betrachten nun eine andere Art der QR-Zerlegung durch Spiegelungen. Wir wollen durch eine solche Spiegelung eine ganze Zeile mit den notwendigen Nullen füllen. Wir überlegen uns dazu zuerst, was eigentlich eine Householdertransformation darstellen soll.

**Householdertransformationen** Eine Householdertransformation beschreibt im Allgemeinen die Spiegelung eines Vektors an einer Hyperebene durch die Null im euklidischen Raum, welche durch ihren Normalenvektor  $v$  mit Einheitslänge  $v^T v = \|v\|^2 = 1$  repräsentiert wird.

Die Householdertransformation eines Vektors  $x$  findet man mit

$$x' = x - 2vv^T x,$$

mit Überlegungen stellt man dann fest, dass

$$x' = (E_n - 2vv^T)x.$$

Wir wollen dann die Householdermatrix  $H$  definieren als

$$H = E_n - 2vv^T.$$

Wir erkennen, dass  $H$  sowohl symmetrisch als auch zu sich selbst invers ist, woraus unmittelbar folgt, dass  $H$  orthogonal ist. Da ebenfalls ein jeder Vektor  $u$  durch

$$v = \frac{u}{\|u\|}$$

normalisiert werden kann, ist  $H$  ebenfalls durch

$$H = E_n - 2 \cdot \frac{uu^T}{\|u\|^2}$$

darstellbar. Mit einem speziell definierten Vektor

$$u = x - \|x\| \cdot e_1,$$

wobei  $e_1$  der erste Standardbasisvektor des  $\mathbb{R}^n$  ist, können wir die Spiegelung  $x'$  von  $x$  einfach angeben mit

$$x' = Hx = \|x\| \cdot e_1.$$

Das heißt, dass durch die Spiegelung alle Positionen bis auf die erste ausgenullt werden. Damit lässt sich leicht eine obere Dreiecksmatrix erzeugen.

## QR-Zerlegung durch Householdertransformationen

### Verfahren 1.2 (QR-Zerlegung mit Householdertransformationen)

Sei  $A \in \mathbb{R}^{n \times m}$  und  $n < m$ , so gilt:

- ① Zu Beginn des Verfahrens wird die erste Householderspiegelung  $H_1$  mittels

$$H_1 = E_n - 2 \cdot \frac{uu^T}{\|u\|^2} \quad \text{und} \quad u = x - \|x\| \cdot e_1,$$

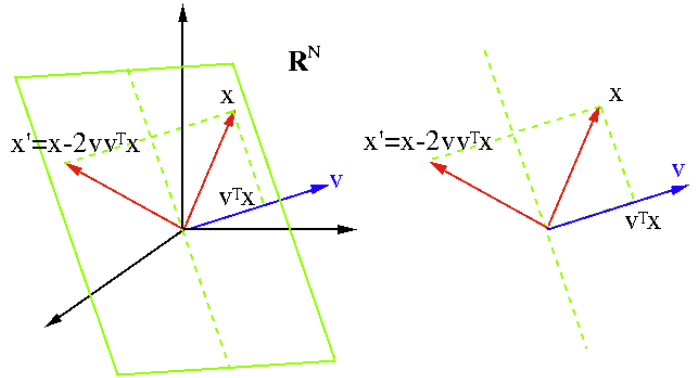
wobei sich  $u$  durch

$$u = c_1 - \|c_1\| \cdot e_1,$$

wobei  $c_1$  der erste Spaltenvektor der Matrix  $A$  ist, darstellen lässt.

- ② Sei  $H_{(k-1):1}A = \prod_{i=0}^{k-2} H_{(k-1)-i}$  mit

$$H_{(k-1):1}A = \left( \begin{array}{c|c} R_{k-1} & * \\ \hline 0 & A^{(k-1)} \end{array} \right)$$



die multiplikative Verknüpfung der bereits konstruierten Householderspiegelungen mit der Ursprungsmatrix  $A$ . Dann konstruieren wir wie in Schritt ① eine Householderspiegelung  $H'_k$  für die Matrix  $A^{(k-1)}$  und erstellen dadurch dann die Householdermatrix  $H_k$  mit

$$H_k = \begin{pmatrix} E_{k-1} & 0 \\ 0 & H'_k \end{pmatrix}.$$

③ Nach  $n - 1$  Schritten (*falls*  $n > m$ :  $n$  Schritten) setzen wir dann  $R$  und  $Q$  als

$$R = H_{(n-1):1} \cdot A \quad \text{und} \quad Q = H_{(n-1):1}^{-1} = H_{1:(n-1)}.$$

Der Aufwand der QR-Zerlegung durch Householderspiegelungen ist mit  $\mathcal{O}(n^3)$  beschrieben, er ist in etwa doppelt so groß wie der einer LR-Zerlegung.

### 1.2.4 Orthonormalisierung

#### Orthonormalisierung

Das Gram-Schmidt-Verfahren (*Standardverfahren der linearen Algebra für die Orthonormalisierung von linear unabhängigen Vektoren*) ist in der klassischen Variante **numerisch instabil**. Zur Verbesserung nehme man das modifizierte Gram-Schmidt-Verfahren oder aber eine QR-Zerlegung.

#### Verfahren 1.3 (QR-Zerlegung mit Householdertransformationen)

Seien mit  $b_i$  ( $1 \leq i \leq k$ )  $k$  linear unabhängige Vektoren gegeben und mit  $B = (b_1 \ \cdots \ b_k)$  die zugehörige Matrix, deren QR-Zerlegung mit

$$B = Q \cdot R = (q_1 \ \cdots \ q_k) \cdot R$$

bekannt ist.

Dann ist ein System von *orthonormalen* Vektoren mit den  $q_i$  gegeben.

## 1.3 Aufgaben

Wir wollen nunmehr das Stoffgebiet durch einige (*weiterführende*) Aufgaben aufarbeiten.

### Aufgabe 1.1

Bestimmen Sie die LR-Zerlegung der folgenden Matrix  $A_1$  mit

$$A_1 = \begin{pmatrix} 2 & -1 & 1 & -1 \\ -2 & 2 & -2 & 2 \\ 2 & -2 & 4 & -3 \\ 0 & 1 & -3 & 5 \end{pmatrix}$$

### Aufgabe 1.2

Von der Matrix  $A_2$  ist die QR-Zerlegung bekannt, es gilt

$$A_2 = \begin{pmatrix} -1 & 1 & 1 & 1 \\ 1 & 1 & -1 & 1 \\ -1 & -1 & 2 & 0 \\ 1 & -1 & 0 & 2 \end{pmatrix} = \frac{1}{2} \cdot \underbrace{\begin{pmatrix} -1 & 1 & 1 & 1 \\ 1 & 1 & 1 & -1 \\ -1 & -1 & 1 & -1 \\ 1 & -1 & 1 & 1 \end{pmatrix}}_{=: Q} \cdot \underbrace{\begin{pmatrix} 2 & 0 & -2 & 1 \\ 0 & 2 & -1 & 0 \\ 0 & 0 & 1 & 2 \\ 0 & 0 & 0 & 1 \end{pmatrix}}_{=: R}$$

**Teilaufgabe 1.2a**

Lösen Sie das folgende lineare Gleichungssystem (unter Verwendung der QR-Zerlegung)

$$Ax = b \quad \text{für} \quad b = (4 \ 2 \ 2 \ 4)^T.$$

**Teilaufgabe 1.2b**

Bestimmen Sie einen Vektor  $w$ , so dass die zugehörige Householderspiegelung  $H_w$  bei Anwendung auf  $A$  (von Teilaufgabe 1.2a) Nullen in der ersten Spalte erzeugt, das heißt

$$H_w A = \begin{pmatrix} * & * & * & * \\ 0 & * & * & * \\ 0 & * & * & * \\ 0 & * & * & * \end{pmatrix}.$$

**Teilaufgabe 1.2c**

Berechnen Sie den Betrag der Determinante von  $A$ , also  $|\det A|$

**Aufgabe 1.3**

Bestimmen Sie die QR-Zerlegung der Matrix  $A_3$  mit

$$A_3 = \begin{pmatrix} 4 & 1 \\ 3 & 2 \end{pmatrix}$$

**Aufgabe 1.4**

Gegeben ist eine pentadiagonale  $n \times n$ -Matrix  $A = (a_{ij})$ , das heißt  $a_{ij} = 0$  falls  $|i - j| > 2$ . Für ebendiese Matrix soll mit Hilfe von GIVENSROTATIONEN die QR-Zerlegung bestimmt werden. Wieviele solcher GIVENSROTATIONEN werden dazu benötigt?

$$C = \begin{pmatrix} * & * & * & 0 & 0 & \dots & \dots & 0 \\ * & * & * & * & 0 & & & & & & \\ * & * & * & * & * & \ddots & & & & & \vdots \\ 0 & * & * & * & * & \ddots & \ddots & & & & \vdots \\ 0 & 0 & * & * & * & \ddots & \ddots & \ddots & & & \\ & & \ddots & \ddots & \ddots & \ddots & \ddots & \ddots & 0 & 0 & \\ \vdots & & & \ddots & \ddots & \ddots & * & * & * & 0 & \\ \vdots & & & & \ddots & \ddots & * & * & * & * & \\ 0 & \dots & \dots & & 0 & 0 & * & * & * & & \end{pmatrix}$$

Struktur einer pentadiagonalen Matrix

**Aufgabe 1.5****Teilaufgabe 1.5a**

Bestimmen Sie die Householdertransformation  $S$ , die den Vektor  $a = (4 \ 2 \ 5 \ 2)^T$  auf ein Vielfaches des zweiten Einheitsvektors  $e_2 = (0 \ 1 \ 0 \ 0)^T$  spiegelt.

**Teilaufgabe 1.5b**

Bestimmen Sie die QR-Zerlegung der Matrix  $A$  mit

$$A = \begin{pmatrix} 1 & 0 & 3 \\ 0 & 3/4 & -3 \\ 0 & -1 & -2 \end{pmatrix}$$

mit einer Householdertransformation.

**Teilaufgabe 1.5c**

Sei dann die partielle Zerlegung  $Z$  der Matrix  $A = Q_1 Z$  mittels einer Givensrotation  $Q_1$  durch

$$Q_1 = \begin{pmatrix} 0 & 1 & 0 \\ -1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} \quad \text{und} \quad Z = \frac{1}{5} \cdot \begin{pmatrix} 4 & 8 & 9 \\ 0 & 15 & 20 \\ 3 & 6 & 13 \end{pmatrix}$$

bekannt. Beenden Sie die begonnene QR-Zerlegung mit Hilfe von Givensrotationen.

**Aufgabe 1.6**

Gegeben sei die Matrix  $A_5$  sowie deren LR-Zerlegung und ein Vektor  $b$

$$A_5 = \begin{pmatrix} 2 & 1 & 1 \\ 6 & 4 & 4 \\ 4 & 1 & 3 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 3 & 1 & 0 \\ 2 & -1 & 1 \end{pmatrix} \cdot \begin{pmatrix} 2 & 1 & 1 \\ 0 & 1 & 1 \\ 0 & 0 & 2 \end{pmatrix} \quad \text{und} \quad b = \begin{pmatrix} -1 \\ 2 \\ 1 \end{pmatrix}.$$

Lösen Sie das lineare Gleichungssystem  $A_5 x = b$  unter Zuhilfenahme der LR-Zerlegung.

**Aufgabe 1.7**

Bestimmen Sie die LR-Zerlegung folgender Matrix mit Pivotsuche:

$$A = \begin{pmatrix} 0 & 4 & 2 \\ 4 & 6 & 4 \\ -2 & -1 & 3 \end{pmatrix}$$

**Aufgabe 1.8**

Bestimmen Sie die Komplexität der folgenden Operationen. Sie dürfen annehmen, dass die Länge der Spaltenvektoren  $n$  und die Größe von Matrizen  $n \times n$  ist.

Szenario	Komplexität
Bestimmung der LR-Zerlegung einer $k$ -diagonalen Matrix	$\mathcal{O}(\quad)$
Bestimmung der LR-Zerlegung einer 6-diagonalen Matrix	$\mathcal{O}(\quad)$
Lösen von $Rx = b$ für untere Dreiecksmatrizen $R$	$\mathcal{O}(\quad)$
Lösen von $Ax = b$ für vollbesetzte symmetrische Matrizen $A$	$\mathcal{O}(\quad)$
Bestimmung der QR-Zerlegung einer vollbesetzten Matrix	$\mathcal{O}(\quad)$
Komplexität des Lösen des LGSs $Ax = b$ mittels Gauß-Elimination	$\mathcal{O}(\quad)$
Sei $A$ eine nichtsinguläre Matrix und $b$ ein Vektor. Welche Komplexität hat dann die Bestimmung der LR-Zerlegung $A = PLR$ ?	$\mathcal{O}(\quad)$
Komplexität des Lösen eines Gleichungssystems, wenn die LR-Zerlegung bereits vorliegt?	$\mathcal{O}(\quad)$
Wie viele Jacobirotationen benötigt die Bestimmung der QR-Zerlegung für eine tridiagonale Matrix?	
Bestimmung der LR-Zerlegung einer vollbesetzten Matrix	$\mathcal{O}(\quad)$
Komplexität des Bestimmens der Determinante bei bereits gegebener QR-Zerlegung	$\mathcal{O}(\quad)$

Szenario	Komplexität
Komplexität des Bestimmens der Determinante bei bereits gegebener LR-Zerlegung	$\mathcal{O}(\quad)$
Komplexität des Bestimmens des Betrags der Determinante bei bereits gegebener QR-Zerlegung	$\mathcal{O}(\quad)$
Welche Komplexität hat das Lösen von $Ax = b$ bei einer $k$ -diagonalen Matrix $A$ ?	$\mathcal{O}(\quad)$

### Aufgabe 1.9

#### Teilaufgabe 1.9a

Welche der folgenden Verfahren können zum Lösen von linearen Gleichungssystemen verwendet werden?

- | ja                       | nein                     |                       |
|--------------------------|--------------------------|-----------------------|
| <input type="checkbox"/> | <input type="checkbox"/> | Gauß-Seidel-Verfahren |
| <input type="checkbox"/> | <input type="checkbox"/> | Quicksort             |
| <input type="checkbox"/> | <input type="checkbox"/> | cg-Verfahren          |
| <input type="checkbox"/> | <input type="checkbox"/> | QR-Zerlegung          |
| <input type="checkbox"/> | <input type="checkbox"/> | Romberg-Verfahren     |

#### Teilaufgabe 1.9b

Welche der folgenden Methoden ist zum Invertieren von Matrizen sinnvoll?

- | ja                       | nein                     |  |
|--------------------------|--------------------------|--|
| <input type="checkbox"/> | <input type="checkbox"/> | LR-Zerlegung   |
| <input type="checkbox"/> | <input type="checkbox"/> | Newton-Cotes-Regel   |
| <input type="checkbox"/> | <input type="checkbox"/> | SVD  |
| <input type="checkbox"/> | <input type="checkbox"/> | Dijkstra-Algorithmus   |
| <input type="checkbox"/> | <input type="checkbox"/> | Ich versuche das Invertieren von Matrizen wenn möglich zu vermeiden. |

### Aufgabe 1.10

Das Modul `SoLver` stellt grundlegende Funktionen zum Lösen von linearen Gleichungssystemen bereit. Dabei sollen Sie die entsprechenden Methoden implementieren.

#### Teilaufgabe 1.10a

Implementieren Sie die Funktion `decomposeLU(M)`. Diese soll eine LR-Zerlegung einer übergebenen Matrix  $M$  (dargestellt als NumPy-Array) berechnen und die linke untere Dreiecksmatrix  $L$  und rechte obere Dreiecksmatrix  $R$  der LR-Zerlegung zurückgeben.

**Hinweis:** Sie sollen explizit überprüfen, ob eine LR-Zerlegung der übergebenen Matrix überhaupt möglich ist!

#### Teilaufgabe 1.10b

Implementieren Sie nun die Funktion `forward_substitute(L, b)`, die eine untere Dreiecksmatrix  $L$  (als  $n \times n$  NumPy-Array) und einen  $n$ -Vektor  $b$  (als  $n \times 1$  NumPy-Array) übergeben bekommt und das Gleichungssystem  $Ly = b$  nach  $y$  löst und diesen zurückgibt. Nutzen Sie dabei die besondere Struktur der Matrix  $L$ .

#### Teilaufgabe 1.10c



Implementieren Sie nun die Funktion `solveSystem(M, B)`. Diese soll das Gleichungssystem  $M \cdot X = B$  nach  $X$  lösen. Verwenden Sie dazu bereits implementierte Methoden des Moduls. *Sie können davon ausgehen, dass das Modul `So1ver` bereits die Funktion `backward_substitute(R, y)` kennt und die übergebenen Matrizen zueinander kompatibel sind.*

#### Teilaufgabe 1.10d

Implementieren Sie nun die Funktion `calcDeterminant(M)`. Diese soll die Determinante der übergebenen Matrix  $M$  berechnen und zurückgeben. Verwenden Sie dazu bereits implementierte Methoden des Moduls.

### Aufgabe 1.11

Wir betrachten nun noch einmal das Modul `So1ver`, welches nun für die QR-Zerlegung erweitert werden soll. In dieser Aufgabe können Sie der Einfachheit halber davon ausgehen, dass die übergebenen Parameter **immer** kompatibel sind, Sie brauchen also *keine Fehlerbehandlung*.

#### Teilaufgabe 1.11a

Implementieren Sie die Funktion `givens_rotation(A, i, j)`. Diese soll die Givensrotationsmatrix  $J$  zur Matrix  $A$  erstellen, so dass  $J \cdot A$  an der Stelle  $(i, j)$  eine 0 besitzt.

#### Teilaufgabe 1.11b

Implementieren Sie anschließend die Funktion `householder(A, j)`, welche die Householderspiegelungsmatrix einer übergebenen Matrix  $A$  in Spalte  $j$  berechnet und zurückgibt.

#### Teilaufgabe 1.11c

Implementieren Sie anschließend die Funktion `decomposeQR(A)`, welche die QR-Zerlegung einer übergebenen Matrix  $A$  berechnet und zurückgibt. Zum Erzeugen der Nullen sollen Sie dabei Givensrotationsmatrizen verwenden.

#### Teilaufgabe 1.11d

Implementieren Sie nun die Funktion `backward_substitute(R, y)`, die eine obere Dreiecksmatrix  $R$  (als  $n \times n$  NumPy-Array) und einen  $n$ -Vektor  $y$  (als  $n \times 1$  NumPy-Array) übergeben bekommt und das Gleichungssystem  $Rx = y$  nach  $x$  löst und diesen zurückgibt. *Nutzen Sie dabei die besondere Struktur der Matrix  $R$ .*

## LINEARE AUSGLEICHSPROBLEME

### Definition 2.1 (*Über- und unterbestimmte lineare Gleichungssysteme*)

Sei  $A \in \mathbb{R}^{n \times m}$ ,  $b \in \mathbb{R}^m$ , so heißt das Gleichungssystem ...

- ... **unterbestimmt**, falls das Gleichungssystem *weniger* Gleichungen als Unbekannte hat.
- ... **überbestimmt**, falls das Gleichungssystem *mehr* Gleichungen als Unbekannte hat.

Dabei gehen linear abhängige Zeilen nur **einfach** in die Wertung als Gleichung mit ein.

Man stellt schnell fest, dass überbestimmte Gleichungssysteme im Allgemeinen **keine Lösung** haben, weil sich einzelne Gleichungen widersprechen könnten. Vor allem bei aus mehrfachen Messwerten generierten Gleichungssystemen kommt dies wegen unvermeidlicher Ungenauigkeiten relativ häufig vor.

Wenn das Gleichungssystem nicht lösbar ist, so ist das Residuum  $r(x) = b - Ax$  nichtverschwindend für allgemeine  $x$ , naheliegender ist deswegen die Suche nach einem  $\operatorname{argmin} r(x)$ . Wir betrachten deswegen folgendes Problem:

### Problem 2.1 (*Lineares Ausgleichsproblem*)

Gegeben seien  $A \in \mathbb{R}^{n \times m}$  und  $b \in \mathbb{R}^m$ . Gesucht ist ein  $x \in \mathbb{R}^n$ , so dass das Residuum

$$r(x) = b - Ax$$

minimiert wird.

Wir suchen also ein  $\hat{x}$  mit

$$\hat{x} = \operatorname{argmin}_{x \in \mathbb{R}^n} r(x).$$

### Verfahren 2.1 (*Methode der „kleinsten Quadrate“*)

Wir lösen das lineare Ausgleichsproblem, indem wir das Residuum  $r(x)$  bzgl. der euklidischen Norm minimieren. Das heißt das lineare Ausgleichsproblem wird umformuliert zu:

Gegeben seien  $A \in \mathbb{R}^{n \times m}$  und  $b \in \mathbb{R}^m$ . Gesucht ist ein  $x \in \mathbb{R}^n$ , so dass die Norm des Residuums

$$\|r(x)\|_2 = \|b - Ax\|_2 = \|Ax - b\|_2$$

minimiert wird.

Das Ausgleichsproblem heißt linear, weil der Parameter  $x$  nur linear in den Defekt  $Ax - b$  eingeht. Während das Residuum immer ein Minimum hat, ist der Minimierer  $x$  nur dann eindeutig bestimmt, wenn  $\operatorname{rang}(A) = m$ , also die Spalten von  $A$  linear unabhängig sind.

**Lösen mittels QR-Zerlegungen** Ist die QR-Zerlegung bereits bekannt, so gilt:

**Satz 2.1 (Lineare Ausgleichsprobleme bei bekannter QR-Zerlegung)**

Ist die QR-Zerlegung einer Matrix  $A \in \mathbb{R}^{n \times m}$  bekannt, so ist zum Verfahren der „kleinsten Quadrate“ äquivalent:

$$\arg \min \|Ax - b\|_2 = \arg \min \|Rx - Q^T b\|_2$$

Bei der QR-Zerlegung ist das Beste, das man erreichen kann, die oberen  $m$  Gleichungen zu erfüllen, der Rest spielt dabei dann „keine Rolle“ mehr. Man löst dann das „verkleinerte System“ per Rückwärtssubstitution. Der Rest bestimmt das Residuum, es verschwindet für kein  $x$ , ist aber zeitgleich ein Maßstab für die Erfüllbarkeit des Systems.

**Normalgleichung**

**Satz 2.2 (Lineare Ausgleichsprobleme bei bekannter QR-Zerlegung)**

Problem 2.1 ist äquivalent zum folgenden Optimierungsproblem. Bestimme das Minimum des Funktionals

$$\eta : x \mapsto \|b - Ax\|_2^2 = \sum_{i=1}^n \left( b_i - \sum_{j=1}^m a_{ij} x_j \right)^2.$$

Mathematisch gesehen ist ein Lösungsverfahren alle partiellen Ableitungen der Null gleichzusetzen, wir erhalten somit dann die Normalgleichung des linearen Ausgleichsproblems mit

$$A^T A x = A^T b,$$

dies entspricht einem quadratischen  $m \times m$ -Gleichungssystem, wobei  $A^T A$  eine symmetrische  $m \times m$ -Matrix ist, die positiv definit und invertierbar ist, sofern  $A$  vollen Rang hat ( $\text{rang}(A) = m$ ). Zu Lösen ist dieses Gleichungssystem dann über Verfahren wie Cholesky.

**Ausgleichsgerade** Gegeben sei eine Reihe von Messwerten  $(x_i, y_i)$  mit  $1 \leq i \leq n$ , deren Ausgleichsgerade  $y = a_1 \cdot x + a_0$  gesucht ist. Wir wissen, dass die Gleichung

$$\begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} \cdot \begin{pmatrix} a_0 \\ a_1 \end{pmatrix} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

ideal ist. Mit der Normalgleichung ergibt sich dann als Lösung:

$$A^T A = \begin{pmatrix} n & \sum_{i=1}^n x_i \\ \sum_{i=1}^n x_i & \sum_{i=1}^n x_i^2 \end{pmatrix} = n \cdot \begin{pmatrix} 1 & \bar{X} \\ \bar{X} & \bar{X}^2 \end{pmatrix}$$

$$A^T b = n \cdot \begin{pmatrix} \bar{Y} \\ \bar{XY} \end{pmatrix}.$$

Mittels Cramer'scher Regel folgt dann die Lösung als

$$a_0 = \frac{\bar{Y} \cdot \bar{X}^2 - \bar{X} \cdot \bar{XY}}{\bar{X}^2 - \bar{X}^2} \quad \text{und} \quad a_1 = \frac{\bar{Y} \cdot \bar{X}^2 - \bar{X} \cdot \bar{XY}}{\bar{X}^2 - \bar{X}^2}.$$

## Anwendung 1 – Bestimmung einer Ausgleichsfunktion

### Problem 2.2 (Bestimmung einer Ausgleichsfunktion)

Gegeben seien  $n$  Messwerte  $(x_i, y_i)$ . Suche dann unter den Polynomen  $f_a(x) = \sum_{i=0}^m a_i \cdot x^i$  – mit  $a \in \mathbb{R}^m$  – dasjenige, welches zu den Messwerten am besten passt.

## Anwendung 2 – Bestimmung einer Ausgleichsebene

### Problem 2.3 (Bestimmung einer Ausgleichsebene)

Gegeben seien  $n$  Messwerte  $(x_i, y_i)$  mit jeweiligen Messwerten  $z_i$ . Suche dann unter den Ebenen  $z = a + bx + cy$ , welche die Messdaten am besten *im quadratischen Mittel* approximiert.

## Aufgaben

### Aufgabe 2.1

Gegeben sei die QR-Zerlegung der Matrix  $A_4$  mit

$$A_4 = \frac{1}{3} \cdot \underbrace{\begin{pmatrix} 2 & 0 & 1 & -2 \\ 2 & 0 & -2 & 1 \\ 1 & 0 & 2 & 2 \\ 0 & 3 & 0 & 0 \end{pmatrix}}_{=: Q_4} \cdot \underbrace{\begin{pmatrix} 3 & -1 \\ 0 & 1 \\ 0 & 0 \\ 0 & 0 \end{pmatrix}}_{=: R_4}.$$

Bestimmen Sie die Lösung des linearen Ausgleichsproblem der kleinsten Quadrate ( $\min_x \|A_4 x - b\|_2$ ) mit Hilfe der QR-Zerlegung für  $b = (2 \ -2 \ 3 \ -1)^T$ .

### Aufgabe 2.2

#### Teilaufgabe 2.2a

Wie lautet die Normalengleichung des linearen Ausgleichsproblems?

#### Teilaufgabe 2.2b

Gegeben sind die Datenpunkte  $x_i$  und zugehörige Datenwerte  $y_i$  mit

$$\begin{array}{l|llll} x_i & -1 & 0 & 2 & 3 \\ y_i & 0 & 0 & 2 & 2 \end{array}$$

Bestimmen Sie die Gerade  $y = m \cdot x + t$ , die diese Daten im *least squares* Sinne am besten approximiert.

### Aufgabe 2.3

Gegeben seien die folgenden 2D-Datenpunkte  $p_i = (x_i \ y_i)^T$ :

$$p_0 = \begin{pmatrix} -1 \\ 1 \end{pmatrix}, p_1 = \begin{pmatrix} 0 \\ -2 \end{pmatrix}, p_2 = \begin{pmatrix} 1 \\ -1 \end{pmatrix}, p_3 = \begin{pmatrix} 2 \\ 2 \end{pmatrix}$$

#### Teilaufgabe 2.3a

Bestimmen Sie das lineare überbestimmte Gleichungssystem  $Au = b$ , welches die Bedingungen  $f(x_i) = y_i$  für eine quadratische Funktion  $f(x) = ax^2 + bx + c$  repräsentiert. Geben Sie dazu  $A$ ,  $b$  und  $u$  explizit an.

**Teilaufgabe 2.3b**

Nennen Sie ein weiteres Verfahren neben der Gaußschen Normalgleichung um überbestimmte lineare Gleichungssysteme im *least squares* Sinn zu lösen.

**Aufgabe 2.4**

Bestimmen Sie ob folgende Gleichungssysteme über-, unterbestimmt oder eindeutig lösbar sind:

unterbestimmt    eindeutig lösbar    überbestimmt

$$\begin{pmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{pmatrix} \cdot x = \begin{pmatrix} 3 \\ 2 \\ 1 \end{pmatrix}$$

$$\begin{pmatrix} 1 & 2 \\ 0 & 2 \\ 0 & 3 \end{pmatrix} \cdot x = \begin{pmatrix} 3 \\ 2 \\ 1 \end{pmatrix}$$

$$\begin{pmatrix} 1 & 2 & 0 \\ 0 & 2 & 3 \\ 2 & 4 & 0 \end{pmatrix} \cdot x = \begin{pmatrix} 2 \\ 1 \\ 4 \end{pmatrix}$$

$$\begin{pmatrix} 1 & 2 \\ 0 & 2 \\ 2 & 4 \end{pmatrix} \cdot x = \begin{pmatrix} 1 \\ 4 \\ 2 \end{pmatrix}$$

$$\begin{pmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} \cdot x = \begin{pmatrix} 3 \\ 2 \\ 1 \end{pmatrix}$$

$$\begin{pmatrix} 1 & 3 \\ 2 & 0 \\ 1 & 0 \end{pmatrix} \cdot x = \begin{pmatrix} 2 \\ 3 \\ 4 \end{pmatrix}$$

# MATRIZEN - DATENSTRUKTUREN UND PRAKTISCHE VERFAHREN

Wir überlegen uns nun wie man besonders *dünn besetzte Matrizen* möglichst effizient speichert. Dazu wollen wir möglichst **keine Nullwerte** speichern. Bezeichne nun im Folgenden  $A$  eine dünnbesetzte  $m \times n$ -Matrix auf den reellen Zahlen. Wir definieren dann die Anzahl an Nichtnullwerten einer Matrix  $A$  als

$$nn_A = \#\{x \mid \exists 1 \leq i \leq m, 1 \leq j \leq n. x = A_{ij} \wedge x \neq 0\},$$

die Anzahl an Nichtnullwerten einer Zeile  $z$  der Matrix  $A$  als

$$nnr_z = \#\{x \mid \exists 1 \leq j \leq n. x = A_{zj} \wedge x \neq 0\}$$

und die Anzahl an Nichtnullwerten einer Spalte  $s$  der Matrix  $A$  als

$$nnc_s = \#\{x \mid \exists 1 \leq i \leq m. x = A_{is} \wedge x \neq 0\}.$$

! Im Folgenden gilt, dass die Indizierung im Allgemeinen bei 1 beginnt.

## 3.1 Row And Column Storage (RaCS)

Wir speichern in einem  $nn_A$  großen Feld  $val$  die Nichtnullwerte der Matrix  $A$  und in den Feldern  $col\_ind$  und  $row\_ind$  die zugehörigen Spalten- respektive Zeilenindizes. Allgemein gilt also:

$$A_{ij} = val_k \Leftrightarrow row\_ind_k = i \wedge col\_ind_k = j.$$

Der Speicheraufwand beträgt nur noch  $3 \cdot nn_A \ll m \cdot n$ , da  $nn_A \ll m \cdot n$ .

## 3.2 Compressed Row Storage (CRS)

Wir verbessern nun das obige Verfahren, indem wir statt  $row\_ind$  nur noch Pointer speichern, die die jeweiligen Zeilenanfänge markieren. Statt  $row\_ind$  gibt es dann ein Feld  $row\_ptr$ , welches die Zeiger enthält. Die Elemente aus  $row\_ptr$  zeigen dann auf die Elemente der anderen Arrays. Allgemein gilt hier also:

$$A_{ij} = val_k \Leftrightarrow (row\_ptr_i \leq k < row\_ptr_{i+1}) \wedge col\_ind_k = j.$$

**Wie signalisiert man das Ende?** Der letzte Eintrag in  $row\_ptr$  ist immer  $nn_A + 1$ .

**Wie signalisiert man leere Zeilen?** Die leeren Zeilen machen kein Problem mit der induktiven Definition des Feldes  $row\_ptr$ :

- $row\_ptr_1 = 1$
- $row\_ptr_i = row\_ptr_{i-1} + nnr_{i-1}$

**Vorteile** Der Speicheraufwand beläuft sich auf  $2 \cdot nn_A + m + 1$ . Ebenfalls erlaubt das CRS einen besonders effizienten Zeilenzugriff, sowie eine besonders effiziente Matrix-Vektor-Multiplikation:

```

1 def matMulVek_CRS(val, col_ind, row_ptr, x):
2     n = x.shape()[0]
3     y = np.zeros((n,1))
4     for i in range(0, n):
5         y[i] = 0 #Eigentlich irrelevant, da np.zeros oben
6         for j in range(row_ptr[i], row_ptr[i+1]):
7             y[i] = y[i] + val[j] * x[col_ind[j]]

```

Mit CRS gestaltet sich allerdings das Einfügen und Löschen von Werten als besonders trickreiche Aufgabe. Die Arrays müssten dafür nämlich umstrukturiert werden, eventuell könnte hier eine Implementierung mit verketteten Listen Abhilfe schaffen, allerdings ist das dann mit einem höheren Suchaufwand verbunden.

Ebenfalls gibt es keinen wahlfreien Zugriff auf einzelne Elemente, denn obwohl sich der Zeilenzugriff als besonders effizient herausstellt, ist der Spaltenzugriff sehr ineffizient zu gestalten. Man sucht deswegen auf sortiertem Teilarray in der jeweiligen Zeile, was dann im besten Fall logarithmischen Aufwand mit sich zieht (*binary search*).

**Block Compressed Row Storage (BCRS)** Als Unterart der CRS gibt es die BCRS-Methodik, bei der dicht besetzte Gebiete in einzelnen Blöcken zusammengefasst werden. Jeder nichtnull Block wird dann intern normal gespeichert, da bei dicht besetzten Blöcken eine normale Speicherung effizienter ist. Diese Blöcke werden dann als Elemente für die CRS verwendet, womit dann auch null-Gebiete nicht mehr gespeichert werden.

Vor allem für große Blockgrößen ist das BCRS deutlich effizienter als das reine CRS.

### 3.3 Compressed Column Storage (CCS)

Ein Manko, welches wir erkennen können, ist, dass vor allem für ein  $m \gg n$  die CRS ineffizienter wäre, als wenn man über die Zeilen adressiert. Wir definieren deswegen das CCS-Verfahren analog zu CRS, erhalten damit dieselben Vor- wie auch Nachteile. Statt dem `row_ptr` gibt es jetzt damit dann ein Feld `col_ptr`. Allgemein gilt dann also:

$$A_{ij} = \text{val}_k \Leftrightarrow (\text{col\_ptr}_j \leq k < \text{col\_ptr}_{j+1}) \wedge \text{row\_ind}_k = i.$$

**Wie signalisiert man das Ende?** Der letzte Eintrag in `col_ptr` ist immer  $nn_A + 1$ .

**Wie signalisiert man leere Spalten?** Die leeren Spalten machen kein Problem mit der induktiven Definition des Feldes `col_ptr`:

- $\text{col\_ptr}_1 = 1$
- $\text{col\_ptr}_i = \text{col\_ptr}_{i-1} + \text{nnc}_{i-1}$

### 3.4 Matrizen und Graphen

Im Allgemeinen kann die Besetzungsstruktur einer allgemeinen  $n \times n$ -Matrix auch als Graph dargestellt werden. Sei  $A \in \mathbb{R}^{n \times n}$ , so definieren wir mit

$$V = \{P_i \mid 1 \leq i \leq n\} \quad \text{und} \quad E = \{(P_i, P_j) \mid a_{ij} \neq 0\}$$

einen **gerichteten** Graphen  $G = (V, E)$ . Per Zusammenhangsanalyse des Graphen kann man dann feststellen, ob sich das Gleichungssystem in mehrere **unabhängige Teilprobleme** unterteilen lässt, welche sich dann separat lösen lassen.

## 3.5 Blockmatrizen

Sei  $M$  eine Matrix der Größe  $m \times n$ . Wir zerlegen nun die Zahl der Spalten und Zeilen ganzzahlig in  $m = \sum_{i=1}^q m_i$  und  $n = \sum_{i=1}^r n_i$ , wobei  $q$  und  $r$  die Anzahl an Summanden beschreiben. Dann lässt sich  $M$  darstellen als

$$M = \begin{pmatrix} M_{11} & \cdots & M_{1r} \\ \vdots & \ddots & \vdots \\ M_{q1} & \cdots & M_{qr} \end{pmatrix}, \text{ mit } M_{ij} \in \mathbb{R}^{m_i \times n_j}.$$

### Gründe für Blockmatrizen

Besonders gerne wird eine Matrix in Blockmatrizen zerlegt, wenn ein Problem *nicht mehr komplett in den Speicher passt*. Man passt dann die Blockgröße an den verfügbaren Speicher an, ein Analogon findet sich in der Optimierung für Cachearchitekturen.

**i** Ein weiterer Grund für Blockmatrizen ist die *effizientere Bereitstellung der Daten*, man passt dabei dann die Blockgröße dem gegebenen Cache an, was das *Nutzen von Pipelines* erlaubt. Besonders bei GPGPUs ist diese Herangehensweise entscheidend, man erhält somit bis zu mehreren Größenordnungen an Performancegewinn.

## 3.6 Aufgaben

### Aufgabe 3.1

Speichern Sie die folgende Matrix  $A$  im CRS-Format ab. Die Indizierung beginnt dabei bei 1.

$$A = \begin{pmatrix} 1 & 2 & 7 & 0 \\ -2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 3 & 0 & 0 & 4 \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

### Aufgabe 3.2

Sei ein Vektor  $b = (1 \ 1 \ -1 \ 1 \ 0)^T$  und eine Matrix  $B$  im CCS-Format mit

$$\begin{aligned} \text{val} &= [1, 4, 7, -2, 5] \\ \text{col\_ind} &= [2, 3, 4, 2, 3] \\ \text{row\_ptr} &= [1, 1, 2, 4, 6] \end{aligned}$$

gegeben. Bestimmen Sie den Vektor  $b^T B$  ohne die Matrix zu rekonstruieren. *Hinweis: Zwischenschritte müssen erkennbar sein! Die Indizierung soll bei 1 beginnen.*

### Aufgabe 3.3

*Hinweis: In dieser Aufgabe gilt stets, dass die Indizierung bei 0 beginnen soll. Teilaufgabe 3.3a*

Wie ergibt sich die Dimension einer  $m \times n$  Matrix aus dem CRS-Format?

#### Teilaufgabe 3.3b

Wie ergibt sich die Dimension einer  $m \times n$  Matrix aus dem CCS-Format?

#### Teilaufgabe 3.3c

Wie kann man eine  $n \times n$ -Tridiagonalmatrix möglichst effizient speichern?



**Teilaufgabe 3.3d**

Sei eine Matrix  $A$  im CCS-Format gegeben. Wie lässt sich die Transponierte von  $A$  unter Ausnutzung der Struktur möglichst effizient berechnen und abspeichern?

**Aufgabe 3.4**

Sei nun eine Matrix  $M$  mit Zeilenanzahl 4 im CCS-Format

```

val = [4, 3, 7, -2, 8, 4]
col_ind = [1, 2, 3, 1, 2, 3]
row_ptr = [0, 0, 1, 3, 5, 6]

```

gegeben. Rekonstruieren Sie die Matrix. Die Indizierung beginnt weiterhin bei 0.

**Aufgabe 3.5**

Im Folgenden soll je eine Klasse für die CCS/CRS-Matrixspeicherformen erstellt werden. Fehlerbehandlungen sind **nicht erforderlich!** Für beide Klassen sollten folgende Funktionen implementiert werden:

`__init__(self, other)` – **Konstruktor** Der Konstruktor soll eine gegebene Matrix `other` (dargestellt als  $m \times n$  NumPy-Array) in das jeweilige Format übertragen. Die zu füllenden Instanzvariablen sind:

- `_rowPtr` oder `_colPtr`
- `_colIndices` oder `_rowIndices`
- `_values`
- `_width` und `_height`: Die Breite und Höhe der Matrix
- `_nonZeroElements`: Die Anzahl an Elementen, die nicht 0 sind.

`toMatrix(self)` Erstellt aus einer im jeweiligen Format gespeicherten Matrix wieder eine „normale“ Matrix (Darzustellen als NumPy-Array).

`getEntry(self, i, j)` Bestimmt den Eintrag der Matrix in Zeile  $i$  und Spalte  $j$

`mul(self, other)` Multipliziert eine Matrix im gegebenen Format mit einer anderen „normalen“ Matrix. *Achtung: Bei CRS ist `self @ other`, bei CCS `other @ self` zu berechnen. Sie können davon ausgehen, dass die Ihnen übergebenen Matrizen zueinander kompatibel sind!*

`transpose(self)` Berechnet das Transponierte der Matrix. Sie können hier auch die Klassenmethode `fromFinal(cls, hoehe, breite, values, rowIndices, colPtr)` oder `fromFinal(cls, hoehe, breite, values, colIndices, rowPtr)` verwenden. Diese erstellt eine Matrix mit den gegebenen Einträgen.

`toCCS(self)` und `toCRS(self)` Wandelt eine gegebene Matrix in das jeweils komplementäre Schema um. *Achtung: Es ist **nicht gewünscht** die Matrix in irgendeinem Zwischenschritt zu rekonstruieren.*

# DISKRETISIERUNG UND QUANTISIERUNG

## 4.1 Diskretisierung

### 4.1.1 Faltung

Zielsetzung ist die Beantwortung der Frage, wie sich ein Signal bei der Übertragung über ein System verhält. Hierzu werden einfache, in ihrem Verlauf vollständig bekannte Signale angegeben und der zeitliche Verlauf der Ausgangssignale berechnet.

Seien also Funktionen  $f, g : D \subseteq \mathbb{R}^n \rightarrow \mathbb{C}$  gegeben, so definieren wir die Faltung als

$$(f * g)(x) := \int_D f(\tau) \cdot g(x - \tau) d\tau,$$

wobei das Integral auf der rechten Seite für fast alle  $x$  wohldefiniert sein soll. Im besonderen Fall, dass  $f, g \in \mathcal{L}^1(\mathbb{R}^n)$  – womit  $f$  und  $g$  integrierbare Funktionen mit endlichem uneigentlichem Betragsintegral sind – ist diese Bedingung immer erfüllt.

**Faltung von periodischen Funktionen** Seien  $f, g$  *periodische* Funktionen der Periode  $T$  **einer Veränderlichen**, so definiert man die Faltung von  $f$  und  $g$  als

$$(f * g)(x) := \frac{1}{T} \int_a^{a+T} f(\tau) \cdot g(x - \tau) d\tau,$$

wobei  $f * g$  wieder eine  $T$ -periodische Funktion und die Faltung auf einem beliebigen Intervall mit Periodenlänge  $T$  definiert ist.

**Diskrete Faltung** Seien  $f, g : D \rightarrow \mathbb{C}$  Funktionen mit dem diskreten Definitionsbereich  $D \subseteq \mathbb{Z}$ , so ist die diskrete Faltung von  $f$  und  $g$  definiert als

$$(f * g)(n) = \sum_{k \in D} f(k) g(n - k).$$

**Rechengesetze** Die Menge  $\mathcal{L}^1(\mathbb{R}^n)$  bildet zusammen mit der Addition und Faltung einen kommutativen Ring. En detail heißt das also:

- Kommutativität —  $f * g = g * f$
- Assoziativität —  $f * (g * h) = (f * g) * h$
- Distributivität —  $f * (g + h) = (f * g) + (f * h)$
- Assoziativität mit der skalaren Multiplikation —  $\lambda(f * g) = (\lambda f) * g = f * (\lambda g)$

**Glättungsfilter** Tiefpass-, Bartlettfilter oder die diskrete Gaussfunktion

**Kantendetektion** Solbeloperator

**Separabilität** Ein Filter heißt **separierbar** genau dann, wenn seine Matrixdarstellung – die *Filtermaske* — den Rang 1 hat.

**4.1.2 Fourier-Transformation**

Die Fouriertransformation baut auf der harmonischen Schwingungslehre in Zusammenhang mit einem mathematischen Theorem auf.

Im Allgemeinen werden harmonische Schwingungen mittels

$$\exp(i \cdot kx) = \cos(kx) + i \cdot \sin(kx)$$

dargestellt, wobei  $k$  die Wellenzahl,  $\lambda = \frac{2\pi}{k}$  die Wellenlänge und  $\nu = \frac{k}{2\pi}$  die Frequenz ist.

**Satz 4.1 (Superpositionsdarstellung)**

Eine jede (quadratisch integrierbare) Funktion  $f \in \text{Abb}(\mathbb{R}, \mathbb{C})$  ist durch Superposition von harmonischen Schwingungen darstellbar:

$$f(x) = \frac{1}{(2\pi)^n} \int_{\mathbb{R}^n} (\mathcal{F}f)(y) e^{iy \cdot x} dy .$$

Die Funktion  $\mathcal{F}f$  heißt **Spektralfunktion** oder auch **Fouriertransformierte** von  $f$  und es gilt ebenso:

$$(\mathcal{F}f)(y) = \frac{1}{(2\pi)^{n/2}} \int_{\mathbb{R}^n} f(x) e^{-iy \cdot x} dx .$$

**Satz 4.2 (Faltungssatz)**

Seien  $f, g$  faltbare Funktionen, so gilt

$$\mathcal{F}(f * g) = \sqrt{2\pi} \cdot \mathcal{F}f \cdot \mathcal{F}g$$

beziehungsweise

$$\mathcal{F}f * \mathcal{F}g = \sqrt{2\pi} \cdot \mathcal{F}(f \cdot g)$$

**4.1.3 Abtasttheorem****Satz 4.3 (Abtasttheorem)**

Ist das Spektrum  $\tilde{g}(k)$  einer kontinuierlichen Funktion  $g(x)$  bandbegrenzt, sprich gilt

$$\tilde{g}(k) = 0 \quad \text{für} \quad |k| \geq k_{\max},$$

und wählt man eine Schrittweite  $\Delta x$  mit  $\Delta x \cdot k_{\max} \leq \pi$ , so kann  $g(x)$  aus den Abtastwerten  $(g(j \cdot \Delta x))_{j \in \mathbb{Z}}$  **exakt rekonstruiert werden**.

Grobe Faustregel: Abtastrate ist die doppelte maximale Frequenz.

**4.1.4 Aliasing und Anti-Aliasing**

Aliasing tritt immer dann auf, wenn das zu messende Signal **hochfrequente Anteile** enthält, welche gemäß Abtasttheorem *nicht reproduzierbar* sind. Durch Unterabtasten der hohen Frequenzen wird dann eine nicht vorhandene, niedrigere Frequenz vorgetäuscht.

Beim Antialiasing *entfernt* man nun **vor dem Abtasten** die höheren Frequenzen im Frequenzraum durch eine Multiplikation mit einem Tiefpass, im Ortsraum durch eine Faltung mit der sinc-Funktion.

## 4.2 Quantisierung

Ziel ist es ein  $r \in \mathbb{R}$  mit  $r_{\min} \leq r \leq r_{\max}$  im Wertebereich  $W := [r_{\min}, r_{\max}]$  zu approximieren. Haben wir *nahezu gleichverteilte* Punkte, so können wir das Intervall in *äquidistante* Teilintervalle aufteilen und jeweils deren *Mittelpunkt* als Repräsentanten nehmen. (UNIFORME QUANTISIERUNG)

Sind die Punkte ungleichverteilt, so gibt es unterschiedliche Strategien von Maschinenzahlen bis zum Median-Cut-Verfahren, man spricht von einer *nicht uniformen Quantisierung*.

### 4.2.1 Gleitpunktzahlen und Maschinenzahlen

Die Menge der **normalisierten t-stelligen Gleitpunktzahlen**  $FP_{B,t}$  zur Basis  $B$  sind definiert als

$$FP_{B,t} = \left\{ \pm M \cdot B^E \mid M, E \in \mathbb{Z}, E \text{ beliebig}, B^{t-1} \leq M < B^t \text{ oder } M = 0 \right\}.$$

$M$  beschreibt hierbei die Mantisse,  $E$  den Exponenten. Die Normalisierung der Zahlen erzwingt man durch die Bedingung  $B^{t-1} \leq M$ .

Die Menge der **Maschinenzahlen** ist dann definiert als

$$MFP_{B,t,\alpha,\beta} = \left\{ g \in FP_{B,t} \mid \alpha \leq E \leq \beta \right\}.$$

Dabei sind  $B, t, \alpha$  und  $\beta$  durch die jeweilige Implementierung festgelegt und werden **niemals** explizit gespeichert. Für jede Maschinenzahl wird dann Vorzeichen, Mantisse und Exponent gespeichert, heutzutage vorzugsweise im IEEE-754-Format.

Die Gleitpunktzahlen sind dann die Repräsentanten, die Quantisierung erfolgt dann durch das Runden. Aus den vorherigen Semestern sollte ebenfalls bekannt sein, dass das Rechnen mit Gleitpunktzahlen bei naiver Implementierung **immer** zu **Rundungsfehlern** führen kann. Der **relative Fehler** ist dabei beschränkt durch die Auflösung (es gilt  $\rho = 1/B^{t-1}$ ), der *absolute Fehler* aber hingegen kann aufgrund der *logarithmischen* Verteilung sehr **groß** werden. Neben Rundungsfehlern kann es selbstverständlich ebenfalls zu *Diskretisierungsfehlern* oder *Iterationsfehlern* kommen.

Der Wunsch ist es, dass das Ergebnis einer Gleitpunktoperation  $*_{FP}$  zu dem Prozess der exakten Resultatsberechnung mit anschließender Rundung nach Modell äquivalent ist. Die Auflösung  $\rho$  ist dabei dann eine Schranke für den relativen Rundungsfehler. Wir stellen somit dann die starke Hypothese auf:

Es existiert ein  $\tilde{\varepsilon} \in \mathcal{O}(\rho)$ , so dass

$$a *_{FP} b = (a * b) \cdot (1 + \varepsilon),$$

wobei  $|\varepsilon| = |\varepsilon(*, a, b)| \leq \tilde{\varepsilon}$ .

### Wichtige Informationen

Kommutativität gilt auch weiterhin bei Fließkommaoperationen, Assoziativität und Distributivität hingegen **nicht!**

Hinsichtlich des **relativen Fehlers** sind Multiplikation und Divisionen gutartig, wohingegen die Addition und Subtraktion sehr gefährlich werden können, unter anderem wegen des Auslöschungseffektes.

! Durch die Verwendung von Gleitkommazahlen, welche dem Informatiker aus verschiedenen Veranstaltungen bereits bekannt sein sollten, werden Zahlen aufgrund der **festen** Mantisse mit einem gewissen relativen **Rundungsfehler** gespeichert. Der Fehler ist ursprünglich die Anzahl an Mantissenziffern, kann aber durch gewisse Rechenoperationen **drastisch erhöht werden**. Man nennt diesen Effekt dann **Auslöschung**. Wir betrachten dazu folgendes **Beispiel 4.1:16-ziffrige Mantisse**

$$\underbrace{0,7948236243749276}_{\text{Fehler von } 10^{-16}} - \underbrace{0,7948236243749241}_{\text{Fehler von } 10^{-16}} = \underbrace{0,35}_{\text{Fehler von } 10^{-2}} \cdot 10^{-14}$$



**Kondition und Stabilität von Verfahren** In der Numerik werden durch die Konzepte der *Kondition* und *Stabilität* die Anfälligkeit eines Verfahrens gegenüber Störungen beschrieben. Die *Kondition* eines **Problems** ist vom verwendeten Verfahren **unabhängig**, die *Stabilität* eines **Verfahrens** hängt hingegen maßgeblich vom **Verfahren** ab.

Ein Problem heißt dann **gut konditioniert**, wenn – *exakte Rechnung vorausgesetzt* – die Fehler der Eingabedaten *nicht verstärkt* werden. Damit gilt dann:

$$\frac{|F(\tilde{x}) - F(x)|}{|F(x)|} \approx \frac{|\tilde{x} - x|}{|x|}.$$

Ein Verfahren heißt (**vorwärts-)**stabil, wenn die Rechenfehler des Verfahrens sich nicht akkumulieren, wenn also gilt, dass

$$\frac{|\tilde{F}(y) - F(y)|}{|F(y)|} \ll 1.$$

### Stabilität von Operationen

! Während die arithmetischen Grundoperationen immer **stabil** sind, sind es die **Hintereinanderausführung** von *stabilen* Operationen **nicht automatisch**.

**Akzeptable Ergebnisse** Die obige Definition der Vorwärtsstabilität ist vor allem bei schlecht konditionierten Problemen verbesserungswürdig. Wir wollen deswegen ein Ergebnis  $y'$  **akzeptabel**, wenn es als exakt betrachtetes Ergebnis zu nur leicht gestörten Eingabewerten verstanden werden kann. Wenn also  $x$  die normale Eingabe,  $F(x)$  das exakte Ergebnis und  $x'$  der leicht gestörte Eingabewert mit dem gestörten Ergebnis  $F(x')$  ist, so muss der Zusammenhang  $\|x' - x\| \leq \varepsilon$  erfüllt sein mit  $\varepsilon$  sehr klein. Man nennt ein solches Verfahren dann auch *rückwärtsstabil*.

Der Nutzen von akzeptablen Ergebnissen hängt dann aber auch wieder von der Konditionierung des Problems ab, denn akzeptable Ergebnisse müssen nicht unbedingt sehr genau sein.

**Gute Kondition**

- + – Kleine Eingabestörungen sind **unkritisch** zu betrachten, da sie nur zu **kleinen** Störungen der Ergebnisse führen
- Es lohnt sich einen **guten** Algorithmus zu entwickeln, welcher die **gute Kondition** ausnutzt.

**Schlechte Kondition**

- + – Lösungen reagieren **empfindlich** auf **kleinste Störungen** der Eingabe.
- + – Eine Verbesserung des genutzten Verfahrens ist hier **sinnlos** und **verschwendete Zeit**, da das Verfahren nichts an der Kondition des Problems ändert.
- + – Aus Sicht des Anwenders ist die Aufgabenstellung hier **besonders kritisch zu hinterfragen**.

**Definition der Kondition** Wir definieren die (*relative*) **Kondition**  $\kappa$  eines Problems  $y = F(x)$ , mit  $F \in \text{Abb}(\mathbb{R}^m, \mathbb{R}^n)$  als die kleinste Zahl  $\kappa \geq 0$ , für die gilt, dass ein  $\delta > 0$  existiert, so dass für alle  $\tilde{x}$  mit  $0 < \|\tilde{x} - x\| < \delta$  gilt, dass

$$\frac{\|F(\tilde{x}) - F(x)\|}{\|F(x)\|} \leq \kappa \cdot \frac{\|\tilde{x} - x\|}{\|x\|}.$$

**Konditionszahl einer Matrix** Wir definieren die Kondition einer (*invertierbaren*) Matrix  $A \in \mathbb{R}^{n \times n}$  als

$$\kappa(A) = \frac{\max_{\|x\|=1} \|Ax\|}{\min_{\|x\|=1} \|Ax\|}.$$

Es gilt dann auch

$$\kappa(AB) \leq \kappa(A) \cdot \kappa(B).$$

**i** Von den in Kapitel 1 vorgestellten Verfahren ist **nur** die LR-Zerlegung **mit** Pivotsuche und die QR-Zerlegung stabil. Bei der QR-Zerlegung gilt dann sogar:

$$\kappa(A) = \kappa(Q \cdot R) = \kappa(R)$$

## 4.2.2 Quantisierung von Farbwerten mit dem *median-cut-algorithm*

### Verfahren 4.1 (*Verfahren der mittleren Schnitte (median-cut-algorithm)*)

**Eingabe:** Punktwolke  $\mathfrak{P}$  in Gebiet  $\mathcal{G}$ , Anzahl an Schritten  $N$ .

**Globale Variablen:**  $s$  aktuelle Anzahl an Schritten, vorinitialisiert mit 0.

- ① Finde eine am besten passende umgebende Box der Wolke  $\mathfrak{P}$  (*best-fitting bounding box*), dies liefert dann ein Rechteck oder auch einen Quader.
- ② Teil den Quader an der längsten Kante, so dass **beide Teile** (in etwa) *gleich viele* Punkte enthalten. Damit bildet man den *Median* der Werte.
- ③ Falls  $s \geq N$ , so halte an, **andernfalls** gehe zu Schritt ①.

Die Wahl der Repräsentanten erfolgt dann über beispielsweise den Mittelpunkt der Quader oder aber auch den Mittelwert der enthaltenen Punkte.

### 4.2.3 Vektorquantisierung

**Minimierer der geringsten Quadrate (*least square minimizer*) im eindimensionalen** Der *least square minimizer* ist eine Alternative bei nicht gleichverteilten Werten mit nicht äquidistanter Schrittweite und einem Repräsentanten aus einem Schrittweitenintervall.

Wir nehmen nun an, dass die Werte gemäß einer Wahrscheinlichkeitsdichtefunktion  $p(r)$  verteilt sind und dass damit dann

$$P(r \leq R) = \int_{-\infty}^R p(\rho) d\rho$$

gilt. Ist  $r$  dann eine Realisierung einer Zufallsvariablen mit Dichte  $p(\rho)$ , welche die Wahrscheinlichkeitsverteilung der Punkte beschreibt, also

$$p(\rho) = 0, \text{ wenn } r_{\max} < \rho < r_{\min}$$

gilt, so ist der mittlere quadratische Fehler durch  $\varepsilon = \sum_{j=0}^{L-1} \int_{w_j}^{w_{j+1}} (\rho - r_j)^2 p(\rho) d\rho$  gegeben.

**Satz 4.4**

Das obige Funktional  $\varepsilon$  wird minimiert, falls die Bedingung der nächsten Nachbarn (*nearest neighbor*)

$$w_j = \frac{r_{j-1} + r_j}{2}$$

und die Schwerpunktsbedingung

$$r_j = \frac{\int_{w_j}^{w_{j+1}} \rho \cdot p(\rho) d\rho}{\int_{w_j}^{w_{j+1}} p(\rho) d\rho}$$

gelten.

**Mehrdimensionale Approximation** Wir wollen nun ein  $r \in \mathbb{R}^n$  mit  $w_{i,\min} \leq r_i \leq w_{i,\max}$  und  $1 \leq i \leq n$  approximieren. Der gesamte Wertebereich spannt sich dann auf als

$$W := \prod_{i=1}^n [w_{i,\min}, w_{i,\max}]$$

Wir haben dazu dann  $L$  Vektoren  $r_i$ , welche die Repräsentanten der disjunkten Teilgebiete  $W_i$  darstellen, womit sich dann ein Quantisierungsfehler von

$$\varepsilon = \sum_{i=1}^L \int_{W_i} \|r - r_i\|^2 \cdot \rho(r) dr,$$

wobei  $\rho(r)$  eine  $r$  dimensionale Dichte darstellt. Wir stellen einen ähnlichen Satz zu Optimalitätskriterien wie oben auf und erkennen

**Satz 4.5**

Das obige Fehlerfunktional  $\varepsilon$  wird für gegebene Repräsentanten  $r_i$  mit  $1 \leq i \leq L$  minimiert, wenn für  $r \in W_i$  gilt, dass

$$\forall j \neq i. \|r - r_i\| \leq \|r - r_j\|$$

Ist eine Zerlegung von  $W$  in disjunkte  $W_i$  gegeben, so gilt dann für die Repräsentanten

$$r_i = \frac{\int_{W_i} r \cdot \rho(r) \, dr}{\int_{W_i} \rho(r) \, dr}.$$

**Iteratives Verfahren zur Bestimmung der optimalen Vektorquantisierung** Seien nun jeweils Mengen von  $M$   $n$ -dimensionalen Trainingsvektoren  $x_i$  und  $L$   $n$ -dimensionalen Kodevektoren  $r_i$  gegeben. Zudem sei eine zu den Kodevektoren gehörige Partition  $P$  des  $\mathbb{R}^n$  ebenfalls wie eine jedem Trainingsvektor einen Kodevektor zuordnende Funktion  $Q \in \text{Abb}(\mathbb{R}^n, \mathbb{R}^n)$  bekannt. Liegt ein Trainingsvektor  $x_i$  dann in  $P_j$ , so wird er dann auch durch  $Q(x_i) = r_j$  approximiert. Unser Ziel ist es den Fehler, gegeben durch

$$\varepsilon = \frac{1}{M \cdot n} \cdot \sum_{i=1}^M \|x_i - Q(x_i)\|^2,$$

zu minimieren. Auch hier gibt es wieder Optimalitätskriterien, welche denen der vorausgegangenen Sätze ähneln. So gibt es auch hier die Bedingung der nächsten Nachbarn (*nearest neighbor condition*), welche die Partition  $P$  auf Voronoiregionen aufspaltet mit

$$P_i = \left\{ x \mid \forall 1 \leq k \leq L. \|x - r_i\|^2 \leq \|x - r_k\|^2 \right\},$$

sowie die Schwerpunktsbedingung (*centroid condition*), sprich die Repräsentanten werden mit

$$r_i = \frac{1}{|M_i|} \cdot \sum_{m \in M_i} x_m,$$

wobei  $M_i = \left\{ m \mid x_m \in P_i \right\}$  gilt, gewählt.

#### Verfahren 4.2 (Lloyds-Algorithmus (*k-means clustering*))

① Berechne mit

$$r_i^* = \frac{1}{M} \cdot \sum_{i=1}^M x_i \quad \text{und} \quad \varepsilon^* = \frac{1}{M \cdot n} \cdot \sum_{i=1}^M \|x_i - r_i^*\|^2$$

unter der Annahme, dass  $P_1 = \mathbb{R}^n$  und  $M_1 = \{1, \dots, M\}$ , eine Startlösung.

② Spalte die Kodevektoren auf in

$$r_i^{(0)} = (1 + \beta) \cdot r_i^* \quad \text{und} \quad r_{i+L}^{(0)} = (1 - \beta) \cdot r_i^*.$$

③ Ordne die Samples zu den neuen Kodevektoren zu, sprich bestimme  $P_i$  und  $M_i$ .

④ Berechne die neuen Kodevektoren, bestimme also die  $r_i$ .

⑤ **Solange** der Abstand der alten zu den neuen Kodevektoren größer ist als eine *vorher festgelegte* Schranke, gehe zu Schritt ③

⑥ **Falls** die gewünschte Anzahl der Kodevektoren noch nicht erreicht ist gehe zu Schritt ②.



## 4.3 Aufgaben

### Aufgabe 4.1

Die Filterung eines diskreten Signals  $f$  mit der Filtermaske  $g$  ist definiert durch die diskrete Faltung:

$$f^{\text{filtered}}(n) = \sum_k f(k) \cdot g(n-k)$$

#### Teilaufgabe 4.1a

Die unendliche periodische Zahlenfolge

$$[\dots, -1, 0, 1, -1, 0, 1, -1, 0, 1, \dots]$$

wird (i) mit einem Tiefpassfilter der Größe 3 ( $[g_{-1}, g_0, g_1] = 1/3 \cdot [1, 1, 1]$ ) und (ii) mit einem SOBEL-Filter ( $[g_{-1}, g_0, g_1] = [-1, 0, 1]$ ) gefiltert. Beide Filterungen sind unabhängig voneinander zu betrachten! Bestimmen Sie jeweils die Ergebnisse.

#### Teilaufgabe 4.1b

Beschreiben Sie mit je einem Stichwort, welche Auswirkungen ein Sobelfilter und ein Tiefpassfilter auf ein Graustufenbild haben.

#### Teilaufgabe 4.1c

Welche der folgenden 2D-Filter sind separierbar? Geben Sie im Falle der Separierbarkeit jeweils auch die beiden 1D-Filter an.

$$F_1 = \frac{1}{9} \cdot \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix} \quad F_2 = \frac{1}{8} \cdot \begin{pmatrix} 0 & 1 & 2 \\ 1 & 0 & 1 \\ 2 & 1 & 0 \end{pmatrix} \quad F_3 = \begin{pmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{pmatrix} \quad F_4 = \frac{1}{6} \cdot \begin{pmatrix} 1 & 2 & 2 & 1 \\ 2 & 4 & 4 & 2 \\ 1 & 2 & 2 & 1 \end{pmatrix}$$

#### Teilaufgabe 4.1d

Was sind die Vorteile von separierbaren Filtern?

**Teilaufgabe 4.1e**

Filtern Sie das nachfolgende  $5 \times 5$  Bild mit der Filtermaske  $G = \frac{1}{4} \cdot \begin{pmatrix} 1 & -2 & 1 \\ 2 & -4 & 2 \\ 1 & -2 & 1 \end{pmatrix}$ . Dabei sollen die Randpunkte *nicht* berechnet werden.

$$\text{Input} = \begin{pmatrix} 12 & 8 & 4 & 8 & 12 \\ 16 & 12 & 8 & 8 & 12 \\ 20 & 16 & 12 & 12 & 12 \\ 16 & 12 & 8 & 8 & 12 \\ 12 & 8 & 4 & 8 & 12 \end{pmatrix}$$

**Teilaufgabe 4.1f**

Ein (Grau-)bild mit  $N \times N$  Pixeln soll mit einem **nicht** separierbaren Filter der Größe  $M$  gefiltert werden. Dabei werden nur die inneren  $(N-k) \times (N-k)$  Pixel berechnet. Berechnen Sie die Anzahl an arithmetischen Operationen, sowie den absoluten Anteil der Multiplikationen daran. Was ändert sich wenn ein separierbarer Filter verwendet wird?

**Aufgabe 4.2**

Eine Funktion  $f$  soll nun nacheinander mit  $h_1$  und  $h_2$  gefaltet werden.

**Teilaufgabe 4.2a**

Ist die Reihenfolge der Operanden hierbei wichtig?

**Teilaufgabe 4.2b**

Die zweimalige Faltung  $f \mapsto f * h_1 \mapsto (f * h_1) * h_2$  kann man auch durch **eine** Faltung mit  $f$  erreichen. Mit welcher Funktion gelingt das?

**Teilaufgabe 4.2c**

Was geschieht, wenn man eine Funktion  $f$  mit der Diracfunktion  $\delta_a$  faltet?

**Teilaufgabe 4.2d**

Prüfen Sie die folgenden Behauptungen auf Korrektheit und kreuzen Sie das entsprechende Kästchen an. Hierbei seien die Funktionen  $f, g, h$  jeweils so gewählt, dass die Faltungen existieren.

richtig	falsch	
<input type="checkbox"/>	<input type="checkbox"/>	$\forall f, g, h. f * (g + h) = (f * g) + (f * h)$
<input type="checkbox"/>	<input type="checkbox"/>	$\forall f, g. f * g = g * f$
<input type="checkbox"/>	<input type="checkbox"/>	$\forall f. f * f = f$
<input type="checkbox"/>	<input type="checkbox"/>	$\forall f, g, h. f * (g * h) = (h * g) * f$

**Aufgabe 4.3**

Gegeben sind die Funktionen

$$g(x) := \begin{cases} 1 & , \text{ falls } 0 \leq x \leq 1 \\ 0 & \text{sonst} \end{cases} \quad \text{und} \quad h(x) := \begin{cases} 1+x & , \text{ falls } -1 \leq x \leq 0 \\ 1-x & , \text{ falls } 0 \leq x \leq 1 \\ 0 & \text{sonst} \end{cases} .$$

Beantworten Sie die folgenden Fragen zum Faltungsprodukt  $f = g * h$ :

**Teilaufgabe 4.3a**

Auf welchem Intervall ist  $f(x)$  ungleich Null?

**Teilaufgabe 4.3b**

An welcher Stelle  $x_0$  nimmt die Funktion  $f(x)$  das Maximum an?

**Teilaufgabe 4.3c**

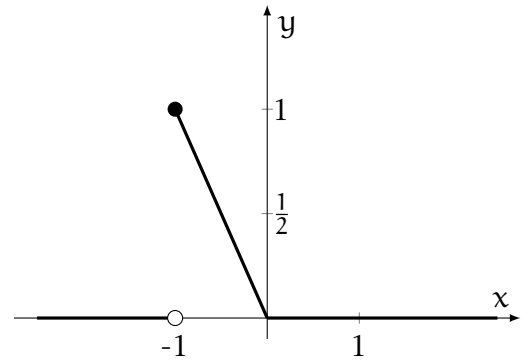
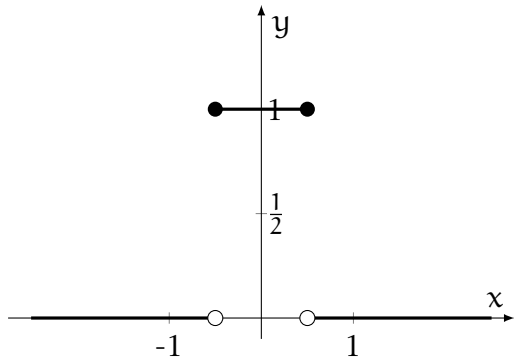
Welches ist der maximale Funktionswert  $f(x_0)$ ?

**Aufgabe 4.4**

Gegeben sind nun die Funktionen

$$h_1(x) := \begin{cases} 1 & , \text{ falls } |x| \leq 1/2 \\ 0 & \text{sonst} \end{cases} \quad \text{und} \quad h_2(x) := \begin{cases} -x & , \text{ falls } -1 \leq x \leq 0 \\ 0 & \text{sonst} \end{cases} ,$$

sowie die Funktionsgraphen der Funktionen:



Zeichnen Sie das Ergebnis der Faltung  $h_1 * h_2$ .

**Aufgabe 4.5**

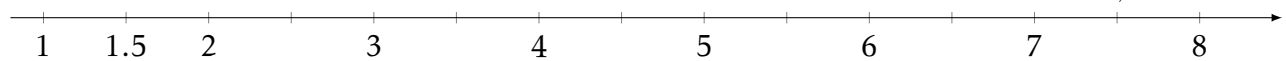
Gegeben sind die Funktionen

$$h_3(x) := \begin{cases} 0 & , \text{ falls } x < 0 \\ x & \text{sonst} \end{cases} \quad \text{und} \quad h_4(x) := \begin{cases} 2 & , \text{ falls } 0 < x < 1 \\ 0 & \text{sonst} \end{cases} .$$

Geben Sie das Ergebnis der Faltung  $h_3 * h_4$  als mathematischen Ausdruck an, indem Sie das Faltungsintegral lösen.

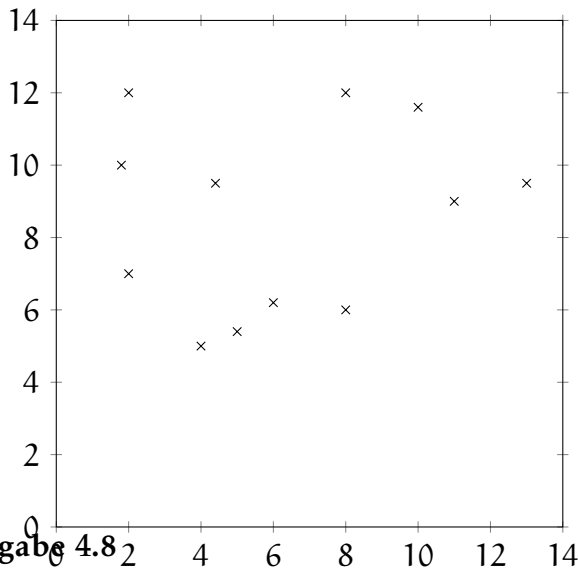
**Aufgabe 4.6**

Betrachten Sie  $FP_{2,3}$ , die Menge der Gleitpunktzahlen zur Basis  $B = 2$  mit Mantissenlänge  $t = 3$ . Markieren Sie in der Abbildung (jeweils durch ein Kreuz) alle  $a \in FP_{2,3}$  mit  $1 \leq a \leq 8$ :



**Aufgabe 4.7**

Gegeben ist eine Punktwolke mit 12 Punkten. Führen Sie zwei Schritte des Median-Cut-Verfahrens durch, wählen Sie nach dem zweiten Schritt einen jeweiligen Repräsentanten als Mittelpunkt der jeweiligen Punktwolke und benutzen Sie dabei die folgenden Vorlagen:



**Aufgabe 4.8**

Prüfen Sie die folgenden Behauptungen auf Korrektheit und kreuzen Sie das entsprechende Kästchen an.

**Teilaufgabe 4.8a**

Bei welchen Operationen kann die Auslöschung zweier positiver Gleitpunktzahlen fast gleicher Größe auftreten?

- |                          |                          |                |
|--------------------------|--------------------------|----------------|
| richtig                  | falsch                   |                |
| <input type="checkbox"/> | <input type="checkbox"/> | Addition       |
| <input type="checkbox"/> | <input type="checkbox"/> | Multiplikation |
| <input type="checkbox"/> | <input type="checkbox"/> | Division       |
| <input type="checkbox"/> | <input type="checkbox"/> | Subtraktion    |

**Teilaufgabe 4.8b**

Welche Rechengesetze sind bei Rechenoperationen mit Gleitpunktzahlen erfüllt?

- |                          |                          |                   |
|--------------------------|--------------------------|-------------------|
| richtig                  | falsch                   |                   |
| <input type="checkbox"/> | <input type="checkbox"/> | Assoziativgesetz  |
| <input type="checkbox"/> | <input type="checkbox"/> | Kommutativgesetz  |
| <input type="checkbox"/> | <input type="checkbox"/> | Distributivgesetz |

**Teilaufgabe 4.8c**

- |                          |                          |   |
|--------------------------|--------------------------|---|
| richtig                  | falsch                   |   |
| <input type="checkbox"/> | <input type="checkbox"/> | Die erreichbare Genauigkeit von Gleitpunktzahlen wird durch die Anzahl verwendeter Bits für den Exponenten, der Wertebereich durch die Anzahl Bits für die Mantisse bestimmt. |
| <input type="checkbox"/> | <input type="checkbox"/> | Die erreichbare Genauigkeit von Gleitpunktzahlen wird durch die Anzahl verwendeter Bits für die Mantisse, der Wertebereich durch die Anzahl Bits für den Exponenten bestimmt. |



# REKONSTRUKTION KONTINUIERLICHER DATEN — INTERPOLATION UND APPROXIMATION

## Interpolation

+ Bekannt seinen  $n$  ( $n \in \mathbb{N}$ ) Punkte im Raum. Soll die rekonstruierte Funktion nun **exakt** durch die Punkte verlaufen, so sprechen wir von **Interpolation**. Es ist *nicht in jedem Fall* das bessere Verfahren (*Messfehler*).

## Approximation

- Bekannt seinen  $n$  ( $n \in \mathbb{N}$ ) Punkte im Raum. Soll die rekonstruierte Funktion nun nur *näherungsweise* durch die Punkte verlaufen, so sprechen wir von **Approximation**. Vorteil: Besserer Umgang mit *Messfehlern*.

## 5.1 Interpolation im eindimensionalen Raum

### Problem 5.1 (*Interpolation*)

Von einer unbekanntem Funktion<sup>1</sup>  $f \in \text{Abb}(\mathbb{R}, \mathbb{R})$  seien die Werte an endlich vielen Stellen, sowie gewisse Grundeigenschaften bekannt. Wir wollen die Stützstellen als  $x_i$ , die Werte der Funktion an den Stützstellen (*Stützwerte*) als  $y_i$  mit  $1 \leq i \leq n$  auffassen. Wir wollen dann  $f$  aus den gegebenen endlich vielen Punkten *näherungsweise* rekonstruieren.

Wir wählen dazu eine Klasse  $K$  von Funktionen mittels derer wir die Rekonstruktion starten wollen. Wir setzen als Kriterium an  $K$  voraus, dass ...

- ...  $f$  in  $K$  gut approximierbar sein muss.
- ...  $K$  auf Computern **gut darstellbar** und *manipulierbar* sein muss.

### 5.1.1 Lokale Verfahren

Bei **lokalen Verfahren** hängt der interpolierte Wert  $p(x)$  **nur** von den unmittelbar „benachbarten“ Werten ab. Gilt beispielsweise  $x \in [x_i, x_{i+1}]$ , so gehen in die Berechnung von  $p(x)$  nur  $y_i$  und  $y_{i+1}$  mit ein.

Seien nun im Folgenden also eine Menge  $\mathcal{A}_S$  von  $n$  Stützstellen  $x_i$  und eine Menge  $\mathcal{O}_S$  von  $n$  Stützwerten  $y_i$  von  $f$  gegeben und eine kontinuierliche Repräsentation  $p(x)$  von  $f$  gesucht.

**Stückweise konstante Interpolation (*nearest neighbor*)** Um den Wert an der Stelle  $x$  anzunähern, suche den nächsten Nachbarn, sprich die nächst gelegene Stützstelle  $x_j$ . Der interpolierte Wert ist dann mit  $y_j$  gegeben. Die Repräsentation  $p$  lässt sich dann auch allgemein durch

$$p(x) = \begin{cases} y_1 & x_1 \leq x \leq \frac{1}{2}(x_1 + x_2) \\ y_i & \frac{1}{2}(x_{i-1} + x_i) < x \leq \frac{1}{2}(x_i + x_{i+1}) \\ y_n & \frac{1}{2}(x_{n-1} + x_n) < x \leq x_n \end{cases} \quad \text{für alle } 2 \leq i \leq n-1$$

<sup>1</sup>Selbstverständlich sind Daten im Allgemeinen (*in der Realität*) durch Relationen beschrieben, wir bleiben der Einfachheit halber zuerst bei Funktionen.

bestimmen. Der einzige Rechenaufwand bei dieser Methode beläuft sich auf die Suche der Nachbarn, welche bei äquidistanten Stützstellen mit Schrittweite  $h$  in  $\mathcal{O}(1)$ , bei bereits geordneten Stützstellen in  $\mathcal{O}(\log(n))$  und andernfalls in  $\mathcal{O}(n \log(n))$  bestimmen lässt.

Problematisch ist die Ungenauigkeit, welche mit einer konstanten Interpolation einhergeht. Als obere Schranke für den Fehler bei glattem  $f$  lässt sich

$$|p(x) - f(x)| \leq \frac{\overbrace{\max_{1 \leq i \leq n} \{x_{i+1} - x_i\}}^{=:h} \cdot \max_{x_1 \leq \xi \leq x_n} \{|f'(\xi)|\}}{2}$$

**Stückweise lineare Interpolation** Die stückweise lineare Interpolation ist das am weitesten verbreitete Interpolationsverfahren. Um den Wert an der Stelle  $x$  hierbei anzunähern, suchen wir zuerst die *nächst gelegene* linke und rechte Stützstelle  $x_i$  und  $x_{i+1}$ , so dass  $x_i \leq x \leq x_{i+1}$  gilt, und interpolieren dann im Intervall  $[x_i, x_{i+1}]$  linear. Wir stellen dann für  $1 \leq i < n$  die Funktion

$$p_i(x) = m_i \cdot (x - x_i) + y_i = \underbrace{\frac{y_{i+1} - y_i}{x_{i+1} - x_i}}_{=:m_i} \cdot (x - x_i) + y_i$$

auf, erhalten damit dann insgesamt die Interpolationsfunktion

$$p(x) = \sum_{i=1}^{n-1} p_i(x) \cdot \mathbb{1}_{[x_i, x_{i+1}]}$$

Alternativ kann für das  $p_i$  auch das gewichtete Mittel von  $y_i$  und  $y_{i+1}$  gewählt werden, es ergibt sich somit dann

$$p_i(x) = \omega \cdot y_i + (1 - \omega) \cdot y_{i+1} = \underbrace{\frac{x_{i+1} - x}{x_{i+1} - x_i}}_{=: \omega} \cdot y_i + \underbrace{\frac{x - x_i}{x_{i+1} - x_i}}_{=: 1 - \omega} \cdot y_{i+1}$$

**Catmull-Rom — Stückweise kubische Interpolation** Wir schätzen hierbei an jeder Stützstelle  $x_i$  die erste Ableitung  $y'_i$  und finden dann auf jedem Teilintervall  $T_i := [x_i, x_{i+1}]$  das (*eindeutige*) kubische Polynom, welches in den Endpunkten neben den Stützwerten auch noch die geschätzten Ableitungen interpoliert. Wir können Ableitungen im Allgemeinen über Vorwärts-, Rückwärts- sowie der zentralen Differenz abschätzen, dabei gilt, dass letztere den geringsten Fehler macht. Damit berechnet sich  $y'_i$  im einfachen Fall durch

$$y'_i = \frac{y_{i+1} - y_{i-1}}{x_{i+1} - x_{i-1}},$$

will man die zentrale Differenz exakt betrachten, so nimmt man das gewichtete Mittel von Vorwärts- und Rückwärtsdifferenz, damit ist dann  $y'_i$  berechenbar durch

$$y'_i = \frac{x_i - x_{i-1}}{x_{i+1} - x_{i-1}} \cdot y'_{i, \text{fw}} + \frac{x_{i+1} - x_i}{x_{i+1} - x_{i-1}} \cdot y'_{i, \text{bw}} = \frac{x_i - x_{i-1}}{x_{i+1} - x_{i-1}} \cdot \underbrace{\frac{y_{i+1} - y_i}{x_{i+1} - x_i}}_{=: y'_{i, \text{fw}}} + \frac{x_{i+1} - x_i}{x_{i+1} - x_{i-1}} \cdot \underbrace{\frac{y_i - y_{i-1}}{x_i - x_{i-1}}}_{=: y'_{i, \text{bw}}}$$

Das kubische Polynom  $p$  bestimmt sich dann mit

$$p(x) = a_0 \cdot (x_{i+1} - x)^3 + a_1 \cdot (x_{i+1} - x)^2 \cdot (x - x_i) + a_2 \cdot (x_{i+1} - x) \cdot (x - x_i)^2 + a_3 \cdot (x - x_i)^3,$$

wobei sich die  $a_i$  durch

$$a_0 = \frac{y_i}{(x_{i+1} - x_i)^3}, \quad a_1 = 3 \cdot a_0 + \frac{y'_i}{(x_{i+1} - x_i)^2}, \quad a_2 = 3 \cdot a_3 - \frac{y'_{i+1}}{(x_{i+1} - x_i)^2} \quad \text{und} \quad a_3 = \frac{y_{i+1}}{(x_{i+1} - x_i)^3}$$

bestimmen lassen.



**Fehlerabschätzungen lokaler Interpolanten** Wir wollen nun kurz annehmen, dass die Stützstellen äquidistant mit Schrittweite  $h$  verteilt sind, sowie die zu interpolierende Funktion  $f$  genügend oft differenzierbar ist.<sup>2</sup> Dann lassen sich folgende Fehlerabschätzungen erkennen:

Verfahrensname	obere Fehlerschranke	„Fehlerkomplexität“
Konstante Interpolation ( <i>nearest neighbor</i> )	$\leq \frac{ f'(\xi) }{2 \cdot h}$	$\in \mathcal{O}(h)$
Lineare Interpolation	$\leq \frac{ f''(\xi) }{8 \cdot h^2}$	$\in \mathcal{O}(h^2)$
Kubische Interpolation ( <i>catmull rom</i> )	$\leq \frac{ f'''(\xi) }{24 \cdot h^3}$	$\in \mathcal{O}(h^3)$

### 5.1.2 Globale Verfahren — Polynominterpolation

Bei **globalen Verfahren** hängt der interpolierte Wert  $p(x)$  von allen verfügbaren Werten ab. Wir rufen uns aus den Mathematikveranstaltungen in Erinnerung, dass die Polynome bis zum Grad  $n - 1$  einen  $n$ -dimensionalen Vektorraum aufspannen:

$$\mathbb{P}_n = \text{span} \left\{ x^i \mid i \in \{0, \dots, n - 1\} \right\}$$

Man nennt diese Basis auch **Monom-** oder **Taylorbasis**.

**Lagrangepolynominterpolation** Wie im lokalen Fall ist auch bei der Lagrangepolynominterpolation eine Menge  $\mathcal{A}_S$  an  $n$  Stützstellen  $x_i$ , sowie eine Menge  $\mathcal{O}_S$  an  $n$  Stützwerten  $y_i$  gegeben. Ein Polynom  $p \in \mathbb{P}_n$  interpoliert diese Werte genau dann, wenn für alle  $i$  gilt, dass  $p(x_i) = y_i$ . Setzen wir eine „Polynomvorlage“ in die  $n$  Interpolationsbedingungen ein, erhalten wir ein lineares Gleichungssystem mit  $n$  Gleichungen für die  $n$  unbekanntenen Koeffizienten:

$$\underbrace{\begin{pmatrix} 1 & x_1 & \cdots & x_1^{n-1} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & \cdots & x_n^{n-1} \end{pmatrix}}_A \cdot \underbrace{\begin{pmatrix} c_1 \\ \vdots \\ c_n \end{pmatrix}}_c = \underbrace{\begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}}_y$$

Dabei ist  $A$  die Vandermondematrix des Problems, dessen Determinante nicht 0 ist und sich über

$$\det(A) = \prod_{1 \leq j < k \leq n} (x_j - x_k)$$

berechnen lässt. Damit ist sie ebenfalls quadratisch, voll besetzt und nicht singular, solange alle Stützstellen paarweise verschieden sind. Das Problem mit der Vandermondematrix ist die schlechte Konditionieren bei schon mäßig gewähltem  $n$ . Dennoch können wir aus der Nichtsingularität der Matrix die Existenz eines eindeutig bestimmten Interpolationspolynoms vom Grad  $n - 1$  folgern. **WEGEN DER SCHLECHTEN KONDITIONIERUNG DER MATRIX IST DIESER WEG ALLERDINGS DER FALSCH!**

#### Rekursionsformel von Aitken-Neville

$$p_{i,k}(x) = \begin{cases} y_i & , \text{wenn } k = 0 \\ p_{i,k-1}(x) \cdot \frac{x_{i+k} - x}{x_{i+k} - x_i} + p_{i+1,k-1}(x) \cdot \frac{x - x_i}{x_{i+k} - x_i} & , \text{sonst} \end{cases}$$

<sup>2</sup>Die Forderung der Äquidistanz kann auch weggelassen werden, statt der globalen Schrittweite  $h$  tritt dann allerdings wie oben der maximale Abstand aufeinanderfolgender Stützstellen auf. Das neue  $h$  berechnet sich dann durch  $h = \max_{1 \leq i \leq n} \{x_{i+1} - x_i\}$ .

Wir stellen fest, dass diese Rekursionsformel das gewichtete Mittel der beiden Polynome niedrigeren Grades berechnet. Diese „Ursprungspolynome“ vom Grad  $k-1$  haben dann jeweils  $k$  Interpolationspunkte, von denen sich  $k-1$  überlappen. Das Polynom am Ende eines Schrittes hat alle Interpolationspunkte der beiden Polynome des Schrittes vorher und ist vom Grad  $k$ .

Dieses Rekursionsschema darf ebenso wenig wie das der Fibonaccizahlen „naiv“ implementiert werden, da dann für den Polynomgrad  $n$  der Aufwand in  $\mathcal{O}(2^n)$  liegt. Stattdessen verwendet man gerne einen Ansatz der dynamischen Programmierung, was dann auf ein Dreiecksschema führt. Dieses werden wir im Zusammenhang mit der Newtoninterpolation näher kennenlernen.

**Alternative Polynombasen I — Bernsteinpolynombasen** Bernsteinpolynombasen sind sehr beliebt bei Graphikern für geometrische Modellierungsaufgaben, für Polynominterpolationen eher weniger geeignet. Wir werden sie wieder bei den Bézierkurven antreffen.

$$B_i^n(x) = \binom{n}{i} \cdot (1-x)^{n-i} \cdot x^i$$

Mit obigen Funktionen gilt dann für den Raum  $\mathbb{P}_{n+1}$ :

$$\mathbb{P}_{n+1} = \text{span} \left\{ B_i^n(x) \mid 0 \leq i \leq n \right\}$$

**Alternative Polynombasen II — Lagrangepolynombasen** Zu den gegebenen Stützstellen  $x_i$  definieren wir die Polynome

$$L_i(x) = \frac{\prod_{\substack{j=1 \\ j \neq i}}^n (x - x_j)}{\prod_{\substack{j=1 \\ j \neq i}}^n (x_i - x_j)}$$

für alle  $1 \leq i \leq n$ . Dabei gilt, dass die  $L_i$  linear unabhängig sind und damit eine Basis von  $\mathbb{P}_n$  bilden, sowie dass  $L_i(x_j) = \delta_{ij}$ .

Das (eindeutige) Interpolationspolynom ist dann durch

$$p(x) = \sum_{i=1}^n y_i \cdot L_i(x)$$

gegeben.

**Alternative Polynombasen III — Newtonpolynombasen** Zu den gegebenen Stützstellen  $x_i$  definieren wir die Polynome

$$N_i(x) = \prod_{j=1}^{i-1} (x - x_j)$$

für alle  $1 \leq i \leq n$ . Das (eindeutige) Interpolationspolynom ist dann durch

$$p(x) = \sum_{i=1}^n a_{i-1} \cdot N_i(x)$$

gegeben. Die  $a_j$  berechnen sich über den Algorithmus von **Aitken-Neville**. In dieser Darstellung ist dann auch das Interpolationspolynom effizient für alle  $x$  auswertbar. Wir betrachten dazu:

```

1 def evaluateNewton(x, a, xi, n): #a, xi: np.array(n)
2     eva = a[n-1]
3     for i in range(n-2, -1, -1):
4         eva = eva * (x - xi[i]) + a[i]
5     return eva

```

### Der Algorithmus von Aitken-Neville

#### Verfahren 5.1 (Algorithmus von Aitken-Neville)

**Eingabe:** Stützstellen  $x_i \in \mathcal{A}_S$ , Stützwerte  $y_i \in \mathcal{O}_S$ .

**Ausgabe:** Koeffizienten  $a_i$  mit  $0 \leq i \leq n - 1$ .

①  $P = (p_{ik})_{\substack{1 \leq i \leq n-k \\ 0 \leq k \leq n-1}}$  und  $k \leftarrow 1$ .

② Für  $1 \leq i \leq n$  mache  $p_{i,0} \leftarrow y_i$ .

③ **Solange**  $k < n$  **tue** Für  $1 \leq i \leq n - k$  mache  $p_{i,k} = \frac{p_{i+1,k-1} - p_{i,k-1}}{x_{i+k} - x_i}$

④  $a \leftarrow p_1$ .

**Aufwand**  $\in \mathcal{O}(n^2)$

Das Verfahren lässt sich auch in Dreiecksform darstellen.

## 5.2 Multivariate Interpolation

Wir wollen nun Funktionen mehrerer Veränderlicher betrachten. Wie im univariaten Fall sind wieder Stützstellen  $x^i = (x_1^i, \dots, x_n^i) \in \mathcal{A}_S$  und Stützwerte  $f_i = f(x^i) \in \mathcal{O}_S$  gegeben und wir suchen Näherungswerte von  $f$  an den dazwischenliegenden Stellen. Wir wollen uns dabei auf maximal trivariate Funktionen einschränken, ähnliche Verfahren sind allerdings durchaus auch für allgemein multivariate Funktionen durchführbar.

Während wir im eindimensionalen Fall Punkte und Intervalle betrachtet haben, betrachten wir im zweidimensionalen Gitter mit *Kanten und Zellen* sowie **Ecken**. Die **Ecken** werden als **geometrische Eigenschaft** klassifiziert, sie geben die Lage an. *Kanten und Zellen* werden als *topologische Eigenschaften* klassifiziert, sie geben Nachbarschaftsbeziehungen an. Es gibt eine Formel für planare Gitter ohne Löcher, dass

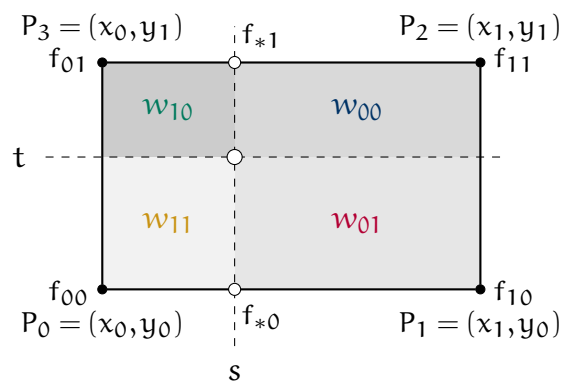
$$|\#Zellen| - |\#Kanten| + |\#Ecken| = 1.$$

### 5.2.1 Lokale Interpolationsverfahren

Auch hier berücksichtigen wir **ausschließlich** die Daten in der **Umgebung** des jeweiligen Punktes. Dies wird vorallem **zellenweise** deutlich, wenn wir **nur** Werte in den Eckpunkten der Zelle betrachten.

**Bilineare Interpolation** Wir verstehen hierunter ein zweistufiges Verfahren, bei welchem zuerst zwei Interpolationen in  $x$ -Richtung und anschließend eine Interpolation in  $y$ -Richtung mit den Ergebnissen des ersten Schrittes stattfindet. Wir erhalten somit in der bilinearen Interpolation das Ergebnis

$$\begin{aligned} BL(s, t) &= (1 - \beta) \cdot \underbrace{\left( (1 - \alpha) \cdot f_{00} + \alpha \cdot f_{10} \right)}_{=f_{*0}} \\ &+ \beta \cdot \underbrace{\left( (1 - \alpha) \cdot f_{01} + \alpha \cdot f_{11} \right)}_{=f_{*1}} \end{aligned}$$



mit

$$\alpha = \frac{s - x_0}{x_1 - x_0} \quad \text{und} \quad \beta = \frac{t - y_0}{y_1 - y_0}.$$

Eine andere Sichtweise geht über die eingezeichneten Flächenanteile, wir erhalten somit

$$\begin{aligned} \text{BL}(s, t) &= \frac{(y_1 - t) \cdot (x_1 - s)}{(y_1 - y_0) \cdot (x_1 - x_0)} \cdot f_{00} + \frac{(y_1 - t) \cdot (s - x_0)}{(y_1 - y_0) \cdot (x_1 - x_0)} \cdot f_{10} + \frac{(t - y_0) \cdot (s - x_0)}{(y_1 - y_0) \cdot (x_1 - x_0)} \cdot f_{01} \\ &+ \frac{(t - y_0) \cdot (x_1 - s)}{(y_1 - y_0) \cdot (x_1 - x_0)} \cdot f_{11} \\ &= w_{00} \cdot f_{00} + w_{10} \cdot f_{10} + w_{01} \cdot f_{01} + w_{11} \cdot f_{11} \end{aligned}$$

### Bivariate Polynome auf Tensorprodukten

i

Liegen die Stützstellen auf einem „Tensorprodukt“gitter, so kann die Interpolation zuerst in x-Richtung und dann in y-Richtung – *oder umgekehrt* – jeweils als **eindimensionale Interpolation** ausgeführt werden.

### Multilineare Interpolation im n-dimensionalen Raum

i

Es gibt

$$2^{n-1} + 2^{n-2} + \dots + 2 + 1$$

lineare Interpolationen in **alle** n Raumrichtungen (*auch hier ist die Reihenfolge egal*).

**Baryzentrische Koordinaten** Seien drei Punkte R, S, T in der Ebene gegeben, die **nicht** auf einer Geraden liegen. Dann lässt sich ein jeder Punkt P der Ebene *eindeutig* darstellen durch

$$P = \rho \cdot P + \sigma \cdot S + \tau \cdot T,$$

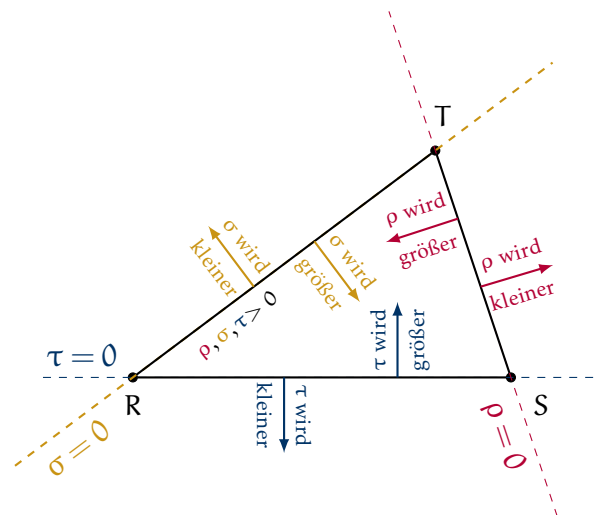
wobei  $\rho + \sigma + \tau = 1$  gelten muss. Das Tupel  $(\rho, \sigma, \tau)$  heißt dann **baryzentrische Koordinate** von P bezüglich  $\Delta(R, S, T)$ . Die Koordinaten werden über das Lösen des Gleichungssystems

$$\begin{pmatrix} 1 & 1 & 1 \\ x_R & x_S & x_T \\ y_R & y_S & y_T \end{pmatrix} \cdot \begin{pmatrix} \rho \\ \sigma \\ \tau \end{pmatrix} = \begin{pmatrix} 1 \\ x_P \\ y_P \end{pmatrix}$$

mittels LR-Zerlegung.

Zudem gilt:

$$\rho = \frac{\text{area}(\Delta(P, S, T))}{\text{area}(\Delta(P, S, T))}, \quad \sigma = \frac{\text{area}(\Delta(R, P, T))}{\text{area}(\Delta(P, S, T))} \quad \text{und} \quad \tau = \frac{\text{area}(\Delta(R, S, P))}{\text{area}(\Delta(P, S, T))}$$



### $n$ -dimensionale baryzentrische Koordinaten

Statt eines Dreiecks betrachtet man nun einen Simplex bestehend aus  $n$  linear unabhängigen Vektoren und  $n + 1$  Punkten  $X_i$ . Damit ist jeder Punkt  $P$  des Raumes darstellbar durch

$$P = \sum_{i=1}^{n+1} \xi_i \cdot X_i.$$

Die Berechnung findet genauso wie im zweidimensionalen statt, man stellt wieder die Matrix  $A = (a_{ij})$  mit

$$a_{0j} = 1, a_{ij} = (X_j)_i$$

auf und löst das so entstehende LGS durch LR-/QR-Zerlegung oder einem iterativen Verfahren.

**Lineare Interpolation** Seien nun in den Eckpunkten des Dreiecks  $\Delta(R, S, T)$  die Daten  $f_R, f_S$  und  $f_T$  gegeben, so erhält man den Wert des linearen Interpolanten an einer Stelle  $P \in \Delta(R, S, T)$  durch

$$f_P = \rho \cdot f_R + \sigma \cdot f_S + \tau \cdot f_T,$$

wobei  $\rho, \sigma, \tau$  die baryzentrischen Koordinaten von  $P$  bezüglich  $\Delta(R, S, T)$  sind.

### lineare Interpolation im $n$ -dimensionalen Raum

Gegeben seien die baryzentrischen Koordinaten  $\xi$  des Punktes  $P$  im  $n$ -dimensionalen Raum zum Simplex  $(X_i)_{1 \leq i \leq n+1}$  und die Daten der Funktion an den Randpunkten, so berechnet sich der Wert des linearen Interpolanten durch

$$f_P = \sum_{i=1}^{n+1} \xi_i \cdot f_{X_i}.$$

## 5.2.2 Globale Interpolationsverfahren im zweidimensionalen Raum

**Tensorproduktansatz** Wir setzen hierfür nunmehr voraus, dass die Stützstellen ein rechteckiges Feld bilden, also sich die Menge an Stützstellen durch

$$\mathcal{A}_S = \left\{ (x_i, y_j) \mid 1 \leq i \leq n, 1 \leq j \leq m \right\}$$

darstellen lässt. Seien dann zu jeder Stützstelle ein Funktionswert gegeben und ein Polynom

$$p(x, y) = \sum_{i=0}^{n-1} \sum_{j=0}^{m-1} c_{pq} x^i y^j$$

gesucht, welches die Interpolationsbedingungen  $p(x_i, y_j) = f_{ij}$  erfüllt.

Effiziente Berechnungen werden durch das Verwenden der Newtonpolynome garantiert, es ist das Gleichungssystem  $XDY^T = F$  zu lösen, wobei  $X$  und  $Y$  die Newtonpolynomdreiecksmatrizen ausgewertet in den jeweiligen Stützstellen  $x_i$  oder  $y_i$  sind und  $F$  die Funktionswertematrix. Man löst dann die jeweils  $n + m$  auftretenden **eindimensionalen Interpolationsprobleme** möglichst effizient durch Aitken-Neville.

### Probleme der Polynominterpolation

Im allgemeinen Fall ist eine solche Polynominterpolation allerdings nicht möglich. Denn nicht nur die Zahl der Freiheitsgrade ist problematisch (*Welchen Polynomgrad wählen wir für eine gewisse Anzahl an Punkten?*) auch die (theoretische) **Lösbarkeit** ist nicht garantiert. Ebenso können interpolierte Funktionen **mehrdeutig** sein.

**Radiale Basisfunktionen** Die Idee der *radialen Basisfunktionen* ist es eine eindimensionale Funktion  $\tilde{h}(r)$  *radialsymmetrisch* fortzusetzen. Eine Funktion  $h$  ist genau dann radialsymmetrisch, wenn gilt, dass

$$h(x) = h(\|x\|).$$

Wir setzen dann den Interpolanten an als

$$g(x, y) = \sum_{i=1}^n d_i \cdot h_i(x, y),$$

wobei die  $d_i$  die noch zu bestimmenden Koeffizienten darstellen und  $h_i$  die um die Stützstelle  $(x_i, y_i)$  *zentrierte radiale Basisfunktion* ist. Sei  $\tilde{h}(r)$  die eindimensionale Funktion, welche mit

$$h_0(x, y) = \tilde{h}(\sqrt{x^2 + y^2})$$

radialsymmetrisch fortgesetzt wurde, für die um die Stützstelle  $(x_i, y_i)$  zentrierte RBF bestimmt man dann über

$$h_i(x, y) = h_0(x - x_i, y - y_i).$$

Die Koeffizienten werden dann durch die  $n$  Interpolationsbedingungen

$$g(x_j, y_j) = \sum_{i=1}^n d_i \cdot h_i(x_j, y_j) = f_j$$

bestimmt.

**Radiale Basisfunktionen mit linearer Präzision** Der Unterschied ist hierbei, dass lineare Funktionen möglichst exakt rekonstruiert werden. Wir unterscheiden deswegen im Ansatz und setzen

$$f(x, y) = a + bx + cy + \sum_{i=1}^n d_i h_i(x, y)$$

an. Damit verschwinden die Momente aller Ordnungen kleiner gleich 1. Für das Bestimmen der Koeffizienten lösen wir daraufhin folgendes  $(n+3) \times (n+3)$  Gleichungssystem:

$$\begin{pmatrix} 0 & 0 & 0 & \sum_j h_1(x_j, y_j) & \cdots & \sum_j h_n(x_j, y_j) \\ 0 & 0 & 0 & \sum_j x_j h_1(x_j, y_j) & \cdots & \sum_j x_j h_n(x_j, y_j) \\ 0 & 0 & 0 & \sum_j y_j h_1(x_j, y_j) & \cdots & \sum_j y_j h_n(x_j, y_j) \\ 1 & x_1 & y_1 & h_1(x_1, y_1) & \cdots & h_n(x_1, y_1) \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & y_n & h_1(x_n, y_n) & \cdots & h_n(x_n, y_n) \end{pmatrix} \cdot \begin{pmatrix} a \\ b \\ c \\ d_1 \\ \vdots \\ d_n \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ f_1 \\ \vdots \\ f_n \end{pmatrix}$$

## 5.3 Freiformmodellierung und Bézierkurven

**Lokalität** Eine lokale Veränderung der Eingabedaten (*von beispielsweise einem Stützwert*) darf *nach Möglichkeit* nur **lokale** Auswirkungen haben, idealerweise ändert sich die Kurvengestalt auch nur in einem kleinen Bereich um den Stützpunkt.

**Bernsteinpolynome**

$$B_i^n(x) = \binom{n}{i} \cdot (1-x)^{n-i} \cdot x^i$$

Die oben definierten Polynome spannen dann den Raum  $\mathbb{P}_{n+1}$  auf:

$$\mathbb{P}_{n+1} = \text{span} \left\{ B_i^n(x) \mid 0 \leq i \leq n \right\}$$

Für die Funktionen gilt dann auch allgemein, dass für  $t \in [0, 1]$   $0 \leq B_i^n(t) \leq 1$  und  $B_i^n(t)$  eine  $i$ -fache Nullstelle in  $t = 0$ , sowie eine  $(n - i)$ -fache in  $t = 1$  hat. Ebenfalls gilt für alle  $t$ , dass  $\sum_{i=0}^n B_i^n(t) = 1$ .

**Bézierkurven** Gegeben seien eine Menge an **Kontrollpunkten**  $\mathcal{B} = \{b_0, \dots, b_n\} \subset \mathbb{R}^d$ , so heißt die Kurve  $C$  im  $\mathbb{R}^d$  **Bézierkurve** vom Grad  $n$ , wenn

$$C(t) = \sum_{i=0}^n b_i \cdot B_i^n(t),$$

wobei  $t \in [0, 1]$ . Verbindet man die Kontrollpunkte durch einen Polygonzug, so erhält man das der Bézierkurve zugehörige **Kontrollpolygon**.

**Formeigenschaften von Bézierkurven** Die Geometrie des Kontrollpolygons spiegelt dann auch (grob) die Geometrie der Kurve wider. Wir erhalten

(1) **Interpolation der Endpunkte** Es muss gelten, dass

$$C(0) = b_0 \quad \text{und} \quad C(1) = b_n.$$

(2) **Tangentiale Tangentenbedingung** In den Endpunkten nähert sich die Kurve tangential an das Kontrollpolygon an, es gilt also

$$C'(0) = n(b_1 - b_0) \quad \text{und} \quad C'(1) = n(b_n - b_{n-1})$$

(3) **Konvexe Hülle** Die Bézierkurve liegt in der **konvexen Hülle** der Kontrollpunkte. Die **konvexe Hülle** einer Menge  $M$  ist die kleinste Obermenge  $M'$  von  $M$ , welche ebenfalls konvex ist.

(4) **Affine Invarianz** Eine affine Abbildung ist mit  $\Phi(x) = Ax + b$  gegeben, um eine Bézierkurve affin zu transformieren müssen **bloß die Kontrollpunkte affin transformiert** werden.

(5) **Variationsreduzierend** Für jede Gerade  $g$  gilt, dass die Anzahl an Schnittpunkten von  $g$  mit der Kurve kleiner oder gleich der Anzahl an Schnittpunkten von  $g$  mit dem Kontrollpolygon ist.

```
1 | def de_casteljau(P, t):
2 |     # You will do this in exercise 7 ...
3 |     pass
```

**Auswertung** Wie oben bereits festgestellt ist die naive Auswertung der Bernsteinpolynome teuer und ineffizient; Abhilfe wird durch

**Horner-Bézier** (*etwas effizienter*) oder durch **fortgesetzte lineare Interpolation mit de Casteljau** (*stabil, siehe rechts*) geschaffen.

## 5.4 Aufgaben

### Aufgabe 5.1

Gegeben seien  $n$  Punkte  $(x_i, y_i)$  und ein Polynom  $p(x)$  vom Grad  $n - 1$ , welches die  $n$  Punkte interpoliert.

#### Teilaufgabe 5.1a

Wie lauten zu den gegebenen Stützstellen  $x_i$  gehörenden Newtonpolynome  $q_j(x)$ ?

#### Teilaufgabe 5.1b

Das Polynom  $p(x)$  kann man mit Hilfe der so angegebenen Newtonbasis darstellen. Geben Sie dabei eine Rekursionsformel zur Berechnung der Koeffizienten  $a_j$  an.

#### Teilaufgabe 5.1c

Bestimmen Sie die Koeffizienten des Newtoninterpolationspolynoms zu den folgenden Punkten und geben Sie anschließend das Newtoninterpolationspolynom an und werten Sie jenes Polynom an der Stelle  $x = 2$  *effizient* aus.

$x_i$	0	1	3
$y_i$	1	-1	7

### Aufgabe 5.2

Nennen Sie jeweils zwei Eigenschaften des Catmull-Rom- und Nearest-Neighbor-Interpolationsverfahrens.

### Aufgabe 5.3

Gegeben seien nun folgende Punkte:

$x_i$	0	1	2	4
$y_i$	1	-1	1	2

#### Teilaufgabe 5.3a

Zeichnen Sie die Funktion  $n(x) : [0,4] \rightarrow \mathbb{R}$ , welche obige Werte, gemäß der *nearest-neighbor-interpolation* stückweise konstant interpoliert.

#### Teilaufgabe 5.3b

Berechnen Sie nun die Funktion  $l(x) : [0,4] \rightarrow \mathbb{R}$ , welche obige Werte stückweise linear interpoliert.

### Aufgabe 5.4

Seien nun folgende Punkte gegeben:

$x_i$	0	2	4	6
$y_i$	-1	-3	-1	3

Geben Sie die Ableitungen  $m_1$  und  $m_2$  des Catmull-Rom-Interpolanten an den Stellen  $x_1$  und  $x_2$  an und skizzieren Sie den Interpolanten im Intervall  $[x_1, x_2]$ .

### Aufgabe 5.5

Seien nun folgende Punkte gegeben:



$i$	0	1	2	3
$x_i$	0	1	2	4
$y_i$	3	5	-1	3

Bestimmen Sie die Lagrangepolynome und geben Sie die Koeffizienten des Interpolationspolynoms an.

### Aufgabe 5.6

Gegeben sind die Eckpunkte eines Dreiecks  $\Delta(R, S, T)$  mit  $R = (3 \ 0)^T$ ,  $S = (3 \ 6)^T$  und  $T = (0 \ 3)^T$ . Die Zuordnung der baryzentrischen Koordinaten zu den Eckpunkten ist mit  $\rho \leftrightarrow R$ ,  $\sigma \leftrightarrow S$  und  $\tau \leftrightarrow T$  gegeben.

#### Teilaufgabe 5.6a

Berechnen Sie die zu den baryzentrischen Koordinaten  $\rho_i, \sigma_i$  und  $\tau_i$  gehörenden Punkte  $P_i$ :

$i$	$\rho_i$	$\sigma_i$	$\tau_i$
0	0	0	1
1	2/3	1/3	0
2	1/3	1/2	1/6

#### Teilaufgabe 5.6b

Die Funktionswerte an den Eckpunkten seien mit  $f_R = 3$ ,  $f_S = 6$  und  $f_T = 4$  bekannt. Berechnen Sie die interpolierten Funktionswerte an den in Teilaufgabe 5.6a bestimmten Punkten!

#### Teilaufgabe 5.6c

Berechnen Sie nun die baryzentrischen Koordinaten  $(\rho_Q \ \sigma_Q \ \tau_Q)^T$  des Punktes  $Q_1 = (2 \ 2)^T$  bzgl. des Dreiecks  $\Delta(R, S, T)$ .

### Aufgabe 5.7

Betrachten Sie das Prisma mit dreieckigen Deckflächen  $\Delta(R_0, S_0, T_0)$  und  $\Delta(R_1, S_1, T_1)$ , wobei

$$R_0 = \begin{pmatrix} 0 \\ 2 \\ -1 \end{pmatrix}, \quad S_0 = \begin{pmatrix} 3 \\ 0 \\ -1 \end{pmatrix}, \quad T_0 = \begin{pmatrix} 3 \\ 2 \\ -1 \end{pmatrix}, \quad R_1 = \begin{pmatrix} 0 \\ 2 \\ 2 \end{pmatrix}, \quad S_1 = \begin{pmatrix} 3 \\ 0 \\ 2 \end{pmatrix} \text{ und } T_1 = \begin{pmatrix} 3 \\ 2 \\ 2 \end{pmatrix}$$

bekannt sind. Ebenso sind die Temperaturwerte in den Eckpunkten mit

$$f_{R_0} = 12, \quad f_{S_0} = 48, \quad f_{T_0} = 36, \quad f_{R_1} = 12, \quad f_{S_1} = 72, \quad f_{T_1} = 72$$

bekannt. Diese sollen ins Innere interpoliert werden. Ermitteln Sie den interpolierten Wert im Punkt  $P = (2 \ 1 \ 0)^T$ , indem Sie zunächst entlang der senkrechten Kanten  $[R_0R_1]$ ,  $[S_0S_1]$  und  $[T_0T_1]$  linear interpolieren und dann die erhaltenen Werte  $f_R$ ,  $f_S$  und  $f_T$  mit Hilfe baryzentrischer Koordinaten linear interpolieren.

### Aufgabe 5.8

#### Teilaufgabe 5.8a

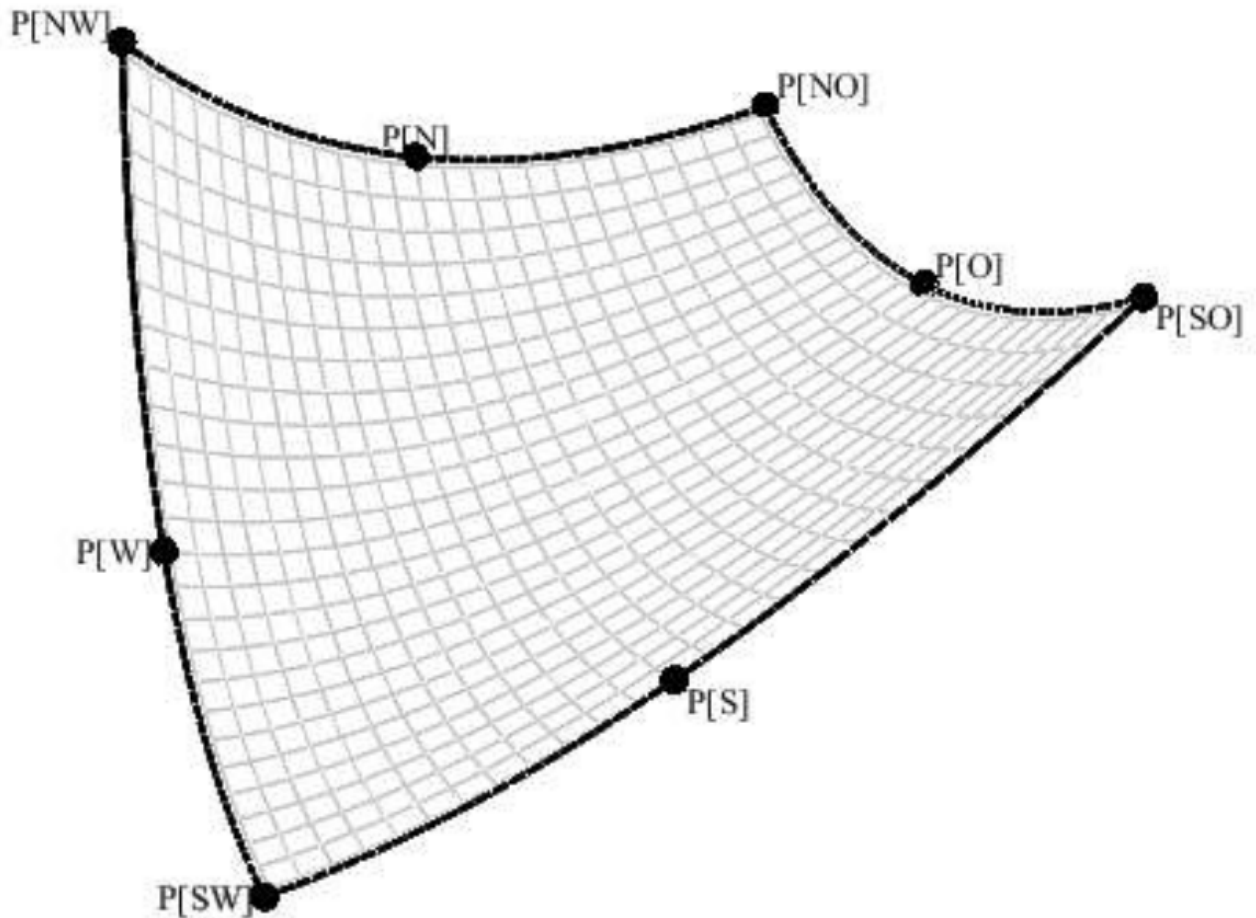
In den Ecken eines Rechtecks  $P_{0,0} = (1 \ 1)^T$ ,  $P_{1,0} = (5 \ 1)^T$ ,  $P_{1,1} = (5 \ 4)^T$  und  $P_{0,1} = (1 \ 4)^T$  sind vier Werte  $f_{00} = 10$ ,  $f_{10} = 2$ ,  $f_{11} = 3$  und  $f_{01} = 1$  vorgegeben. Wir wollen diese nun bilinear interpolieren. Die Werte des Interpolanten in einem Punkt  $P$  kann man als gewichtete Summe schreiben

$$f_P = w_{00}^P \cdot f_{00} + w_{10}^P \cdot f_{10} + w_{11}^P \cdot f_{11} + w_{01}^P \cdot f_{01}.$$

Bestimmen Sie für den Mittelpunkt  $M = (3 \ 5/2)^T$  und den Punkt  $Q = (2 \ 2)^T$  die Gewichte und den bilinearen Interpolationswert.

**Teilaufgabe 5.8b**

Man betrachte das Coonspatch zu vier gegebenen Randkurven  $C_S(s)$ ,  $C_O(t)$ ,  $C_N(s)$  und  $C_W(t)$ :



Der Mittelpunkt des Coonspatches, das heißt der Punkt  $F(0.5,0.5)$ , kann als Linearkombination der vier Kurvenmittelpunkte

$$P_S = C_S(1/2), P_O = C_O(1/2), P_N = C_N(1/2), P_W = C_W(1/2),$$

mit den vier Eckpunkten

$$P_{SW} = C_S(0) = C_W(0), P_{SO} = C_S(1) = C_O(0), P_{NO} = C_O(1) = C_N(1), P_{NW} = C_W(1) = C_N(0)$$

beschrieben werden. Es gilt also:

$$F(1/2, 1/2) = w_{SW}P_{SW} + w_S P_S + w_{SO}P_{SO} + w_O P_O + w_{NO}P_{NO} + w_N P_N + w_{NW}P_{NW} + w_W P_W$$

Bestimmen Sie nun die acht Gewichte  $w_{SW}, \dots, w_W$ .

# SINGULÄRWERTZERLEGUNG UND HAUPTKOMPONENTENANALYSE

## 6.1 Elemente der linearen Algebra

### Definition 6.1 (Norm — Wiederholung aus C1)

Sei  $V$  ein  $\mathbb{K}$ -Vektorraum. Eine Abbildung  $\|\cdot\| \in \text{Abb}(V, \mathbb{K})$  heie **Norm** genau dann, wenn fur alle  $v, u \in V$  und  $\lambda \in \mathbb{K}$  die folgenden Eigenschaften erfullt sind:

- |                                      |                                      |
|--------------------------------------|--------------------------------------|
| 1. $\ v\  \geq 0$                    | 3. $\ \lambda v\  =  \lambda  \ v\ $ |
| 2. $\ v\  = 0 \Leftrightarrow v = 0$ | 4. $\ v + u\  \leq \ v\  + \ u\ $    |

Auf dem  $\mathbb{R}^n$  gilt, dass fur  $1 \leq p < \infty$

**B**

$$\|x\|_p := \left( \sum_i |x_i|^p \right)^{1/p} \quad \text{und} \quad \|x\|_\infty := \max_i |x_i|$$

Normen definieren.

### Definition 6.2 (Matrixnorm)

Seien  $V, W$  normierte Vektorrume und die Abbildung  $F \in \text{Abb}(V, W)$  linear (also  $F \in L(V, W)$ ). Dann heie

$$\|F\| = \sup_{v \in V \setminus \{0\}} \frac{\|Fv\|_W}{\|v\|_V} = \sup_{\|v\|_V=1} \|Fv\|_W$$

die **Operator-** oder auch **Matrixnorm** von  $F$ . In endlichdimensionalen Gebieten mit kompakter Einheitskugel kann das obige Supremum auch durch ein Maximum ersetzt werden.

### Definition 6.3 (Spektralradius)

Sei  $A \in \mathbb{R}^{m \times m}$ , so heie

$$\rho(A) := \max_{i=1, \dots, m} |\lambda_i(A)|,$$

wobei  $\lambda_i(A)$  der  $i$ -te Eigenwert von  $A$  sei, der **Spektralradius** von  $A$ .

**$\|\cdot\|_1$  — Spaltensummennorm**

$$\|A\|_1 = \max_{\|x\|_1=1} \|Ax\|_1 = \max_{j=1,\dots,m} \sum_{i=1}^m |a_{ij}|$$

 **$\|\cdot\|_2$  — Spektralnorm**

$$\|A\|_2 = \sup_{\|x\|_2=1} \|Ax\|_2 = \sqrt{\rho(A^H A)}$$

Dabei ergibt sich für ein invertierbares  $A$  der Zusammenhang  $\|A\|_2 = \rho(A)$

 **$\|\cdot\|_\infty$  — Zeilensummennorm**

$$\|A\|_\infty = \max_{i=1,\dots,m} \sum_{j=1}^m |a_{ij}|$$

Wir definieren noch eine weitere Eigenschaft von Matrizen, deren Diagonaldominanz.

**Definition 6.4 (Diagonaldominanz)**

Sei  $A \in \mathbb{R}^{n \times n}$ , so heie  $A$  ...

... **(zeilenweise) diagonaldominant** genau dann, wenn fur alle  $i \in \{1, \dots, n\}$

$$|a_{ii}| > \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|$$

gilt.

... **(zeilenweise) schwach diagonaldominant** genau dann, wenn fur alle  $i \in \{1, \dots, n\}$

$$|a_{ii}| \geq \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|$$

gilt.

**Definition 6.5 (Kondition)**

Sei  $A$  eine nicht singulre Matrix, so ist die Kondition von  $A$  bzgl. einer Norm definiert als

$$\kappa(A) = \frac{\max_{\|x\|=1} \|Ax\|}{\min_{\|x\|=1} \|Ax\|}.$$

Ist  $A$  zudem noch quadratisch so folgt

$$\kappa(A) = \|A\| \cdot \|A^{-1}\|,$$

stets gilt damit, dass  $\kappa(A) \geq \left| \frac{\lambda_{\max}}{\lambda_{\min}} \right|$ .

**Satz 6.1 (Spektralsatz)**

Sei  $A \in \mathbb{R}^{n \times n}$  eine reelle **symmetrische** Matrix, so gibt es eine **Orthonormalbasis** aus **Eigenvektoren**. Gleichwertig dazu ist die Existenz einer **Faktorisierung**  $A = VDV^T$  mit  $D = \text{diag}(\lambda_1, \dots, \lambda_n)$ , den Spalten der orthogonalen Matrix  $V$  als normierte Eigenvektoren von  $A$ .

**6.2 Singular-value decomposition – SVD****Satz 6.2 (Singularwertzerlegung)**

Sei  $A \in \mathbb{R}^{m \times n}$  eine beliebige Matrix, so gibt es eine Faktorisierung  $A = U\Sigma V^T$ .  
Dabei gilt:

- $U \in \mathbb{R}^{m \times m}$  und  $V \in \mathbb{R}^{n \times n}$  sind orthogonal
- $\Sigma$  ist eine Diagonalmatrix mit  $\Sigma_{ij} = 0$  für  $i \neq j$  und  $\Sigma_{11} \geq \Sigma_{22} \geq \dots \geq 0$
- Die Spalten von  $U$  bzw.  $V$  sind die Eigenvektoren von  $AA^T$  bzw.  $A^T A$
- $\sigma_i = \Sigma_{ii} \neq 0$  sind die Singulärwerte von  $A$ , genauer  $\sigma_i = \sqrt{\lambda_i}$  wobei  $\lambda_i$  die Eigenwerte von  $AA^T$  oder  $A^T A$  sind.

Durch die Singularwertzerlegung lassen sich folgende Werte der Matrix  $A$  einfach bestimmen:

- Der **Rang** entspricht genau der Anzahl der zu Null verschiedenen Singulärwerten.  
 $\text{rang}(A) = r$  aus  $\sigma_1 \geq \dots \geq \sigma_r > 0$
- Der **Kern** entspricht dem aufgespannten Raum aus den Zeilen  $V_i$  mit  $i > r$ .  
 $\ker(A) = \text{span}\{V_{r+1}^T, \dots, V_n^T\}$
- Das **Bild** entspricht dem aufgespannten Raum aus den Spalten  $U_i$  mit  $i \leq r$ .  
 $\text{im}(A) = \text{span}\{U_1, \dots, U_r\}$
- Die **Euklidische Matrixnorm** entspricht dem größten Singulärwert.  
 $\|A\|_2 = \sigma_1 = \sigma_{\max}$
- Die **Konditionszahl** entspricht dem Quotienten aus größtem und kleinstem Singulärwert.  
 $\kappa_2 = \frac{\sigma_1}{\sigma_r} = \frac{\sigma_{\max}}{\sigma_{\min}}$

**Pseudo-Inverse** Die Pseudoinverse  $A^{-1}$  ermöglicht die Annäherung an die tatsächliche Inverse von  $A$  falls diese nicht einfach zu bestimmen ist.

$$A^{-1} = V\Sigma^{-1}U^T \quad \text{wobei} \quad \Sigma_{ii}^{-1} = 1/\Sigma_{ii}$$

**Wichtig:**  $A^{-1}$  und  $\Sigma^{-1}$  sind  $n \times m$ -Matrizen wie  $\Sigma^T$  und  $A^T$ !

$$x = A^{-1}b = \sum_{i=1}^r \frac{1}{\sigma_i} (u_i \circ b) v_i$$

$$A = \left( \begin{array}{c|c|c} u_1 & u_2 & u_3 \end{array} \right) \begin{pmatrix} \sigma_1 & 0 & 0 \\ 0 & \sigma_2 & 0 \\ 0 & 0 & \sigma_3 \end{pmatrix} \begin{pmatrix} v_1 \\ v_2 \\ v_3 \end{pmatrix}$$

## low-rank-approximation

Mit der SVD kann man eine Matrix durch Matrizen von niedrigerem Rang approximieren. Dafür wird  $A_k$  vom Rang  $k$  gesucht, mit

$$A_k = \min_{X: \text{rang}(X)=k} \|A - X\|_F$$

wobei  $A_k$  und  $X \in \mathbb{R}^{m \times n}$  Matrizen.

Für die Frobenius-Norm gilt  $\|A\|_F = \sqrt{\sum_{i=1}^r \sigma_i^2}$  ( $\sigma_i$  sind Singulärwerte von  $A$ )  
mit Hilfe der SVD:

$$A_k = U \text{diag}(\sigma_1, \dots, \sigma_k, \underbrace{0, \dots, 0}_{\substack{\text{setze } \sigma_{k+1} \\ \text{bis } \sigma_r \text{ Null}}}) V^T$$

### Was bringt die low-rank-approximation?

- Speichereffizienz: Nur  $k \cdot (m + n)$  Werte statt  $m \cdot n$
- Schnelle Berechnung von Matrix-Vektor-Multiplikation:  
 $A_k x = \sum_{i=1}^k \sigma_i (v_i \circ x) u_i$  ergibt  $k \cdot (2m + n)$  statt  $2n \cdot m$  arithmetische Operationen

**Approximationsfehler** Es ist die beste Approximation bezüglich der Frobenius-Norm  $\|\cdot\|_F$ .  
Es gelten:

$$\min_{X: \text{rang}(X)=k} \|A - X\|_F^2 = \|A - A_k\|_F^2 = \sum_{j>k} \sigma_j^2$$

$$\min_{X: \text{rang}(X)=k} \|A - X\|_2 = \|A - A_k\|_2 = \sigma_{k+1}$$

## 6.3 Principal component analysis – PCA

### Problem 6.1

*Große Datensätze sind oft schwierig analysierbar und haben nur selten hohe Aussagekraft. Die Hauptachsentransformation wirkt dem entgegen, da sie unter Beibehaltung der wesentlichen Merkmale hoch-dimensionale auf niedrig-dimensionale Probleme reduziert.*

**Schritte zur Bestimmung der PCA** Seien Daten der Form  $\{X_i\}_{1 \leq i \leq N} \in \mathbb{R}^n$  gegeben.

- Balanciere die Daten um Mittelwert  
 $\bar{X} = \frac{1}{N} \sum_i X_i \quad \bar{X}_i = X_i - \bar{X}$
- Bilde Kovarianzmatrix (symmetrisch positiv semi-definit)  
 $C = \text{cov}(X, X) = \frac{1}{N-1} \sum_i \bar{X}_i \cdot \bar{X}_i^T$

- Diagonalisierung der Kovarianzmatrix

$$C = Q \begin{pmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & & \\ \vdots & & \ddots & \vdots \\ 0 & \lambda_2 & \dots & \lambda_n \end{pmatrix} Q^T$$

wobei  $Q$  orthogonal ist und  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0$

- Die Spalten  $q_i$  von  $Q$  sind die Eigenvektoren. Diese bezeichnet man als Hauptachsen.

## 6.4 Aufgaben

### Aufgabe 6.1

Entscheiden Sie, ob die Division zweier Zahlen  $a$  und  $b$  ( $a/b$ ) gut oder schlecht konditioniert ist und begründen Sie Ihre Aussage.

### Aufgabe 6.2

Gegeben seien zwei Geraden durch die Gleichungen

$$G = \{(x, y) \mid 7x - 6y = 1\} \quad \text{und} \quad H = \{(x, y) \mid 6x - 5y = 1\},$$

welche sich im Punkt  $P = (1, 1)$  schneiden.

#### Teilaufgabe 6.2a

Ist das Schnittproblem bei Störung der Koeffizienten der Geradengleichung gut oder schlecht konditioniert?

#### Teilaufgabe 6.2b

Sei  $A$  die Matrix zur Bestimmung des Schnittpunktes  $G \cap H$ . Bestimmen Sie dann die Konditionszahlen dieser  $2 \times 2$ -Matrix bzgl. der Zeilensummennorm.

#### Teilaufgabe 6.2c

Sei  $D = \text{diag}(\lambda_1, \dots, \lambda_n)$  eine gegebene Diagonalmatrix mit  $|\lambda_{i-1}| \geq |\lambda_i| > 0$  für alle  $i$ . Bestimmen Sie die Konditionszahl der Matrix für eine von Ihnen gewählte Matrixnorm und geben Sie dabei vor allem an, welche Sie verwendet haben.

### Aufgabe 6.3

Betrachten Sie die beiden folgenden Funktionen

$$f, g: \mathbb{R}^+ \rightarrow \mathbb{R}^+, f(x) = x^\pi \quad g(x) = x^{1/\pi},$$

welche beide mittels eines nicht näher spezifizierten numerischen Verfahrens in der Nähe von 1 ausgewertet werden. Für welche Funktion ist dieses Problem besser konditioniert? Begründen Sie Ihre Antwort.

### Aufgabe 6.4

Gegeben sei die Singulärwertzerlegung (SVD) der Matrix  $A = U\Sigma V^T$ :

$$A = \frac{1}{25} \cdot \begin{pmatrix} 36 & 27 & -20 \\ 30 & -40 & 0 \\ -48 & -36 & -15 \end{pmatrix} = \frac{1}{5} \cdot \underbrace{\begin{pmatrix} 3 & 0 & -4 \\ 0 & 5 & 0 \\ -4 & 0 & -3 \end{pmatrix}}_{=: U} \cdot \underbrace{\begin{pmatrix} 3 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 1 \end{pmatrix}}_{=: \Sigma} \cdot \frac{1}{5} \cdot \underbrace{\begin{pmatrix} 4 & 3 & 0 \\ 3 & -2 & 0 \\ 0 & 0 & 5 \end{pmatrix}^T}_{=: V^T}$$

#### Teilaufgabe 6.4a

Bestimmen Sie für die Matrix  $A$  die Singulärwerte, den Rang, das Bild, den Kern, sowie die Konditionszahl bzgl. der Spektralnorm ( $\|\cdot\|_2$ )

#### Teilaufgabe 6.4b

Geben Sie die Definitionen der Spektralnorm, der Zeilensummennorm, der Spaltensummennorm und der Frobeniusnorm an und bestimmen Sie  $\|A\|_2$ ,  $\|A\|_\infty$ ,  $\|A\|_1$  und  $\|A\|_F$ .

#### Teilaufgabe 6.4c

Bestimmen Sie  $x$  unter Verwendung der SVD, so dass gilt

$$A \cdot x = \begin{pmatrix} 9 \\ 0 \\ -12 \end{pmatrix}.$$

**Teilaufgabe 6.4d**

Bestimmen Sie die Matrix vom Rang 1, die  $A$  im Sinne der Frobeniusnorm am besten approximiert.



# ITERATIVE VERFAHREN

## 7.1 Allgemeine Theorie von iterativen Verfahren

### Definition 7.1 (Fixpunkt)

Sei  $M \neq \emptyset$  und  $\Phi \in \text{Abb}(M, M)$  eine Abbildung. Ein  $x \in M$ , das den Zusammenhang

$$\Phi(x) = x \quad (7.1)$$

erfülle heie **Fixpunkt** von  $\Phi$ .

Wir betrachten also eine Abbildung  $\Phi : M \rightarrow M$ . Sei  $M \subseteq V$ , wobei  $V$  ein normierter Vektorraum ist. Wir betrachten fur einen Startwert  $x_0 \in M$  eine rekursiv definierte Folge

$$x_{n+1} := \Phi(x_n), \quad (*)$$

und stellen uns die Frage, was der Grenzwert der Folge sei. Aussage daruber verschafft der folgende Satz:

### Satz 7.1 (Fixpunkt)

Falls die oben definierte Folge  $(x_n)_{n \in \mathbb{N}}$  konvergiert mit  $\lim_n x_n := x_*$  und  $\Phi$  stetig sei, so ist der Grenzwert  $x_*$  **immer ein Fixpunkt** von  $\Phi$ .

$$\Phi(x_*) = x_*$$

### Definition 7.2 (Fixpunktiteration)

Sei  $\Phi \in \text{Abb}(M, M)$  mit  $M \subseteq V$ , wobei  $V$  ein normierter Vektorraum ist und  $x_0 \in M$ . Die Berechnung

$$x_{n+1} := \Phi(x_n) \quad (*)$$

heie **Fixpunktiteration**.

Die Fixpunktiteration hat ihren Namen daher, dass wenn  $\Phi$  stetig ist und die Fixpunktiteration konvergiert, so der Grenzwert immer ein Fixpunkt von  $\Phi$  ist. Wir wollen uns jetzt mit der Frage beschftigen, wann Fixpunktiterationen uberhaupt konvergieren. Ein wichtiges Kriterium dafur ist der Banach'sche Fixpunktsatz, auf den dieses Kapitel hinzielen mochte. Eine Voraussetzung dafur sind sogenannte Banachraume. Diese sollten aus C1 bekannt sein, wir wiederholen sie an dieser Stelle aber noch einmal.

### Definition 7.3 (Vollstandigkeit, Banachraum)

Ein normierter  $\mathbb{K}$ -VR  $(V, \|\cdot\|)$  heie **vollstandig** genau dann, wenn jede Cauchyfolge in  $V$  gegen ein Element in  $V$  konvergiert.

Ein Vektorraum heie **Banachraum** genau dann, wenn er normiert und vollstandig ist.

Ein Banachraum, dessen Norm durch ein Skalarprodukt erzeugt wird, heie **Hilbertraum**.

Nun zurück zur Frage, wann und unter welchen Bedingungen Fixpunktiterationen konvergieren. Man kann sich *graphisch* überlegen, dass die „Steilheit“ von  $\Phi$  dabei eine Rolle spielt. Wir definieren:

#### Definition 7.4 (Kontraktion)

Sei  $(V, \|\cdot\|)$  ein  $\mathbb{R}$ -Vektorraum,  $M \subseteq V$  und  $\Phi \in \text{Abb}(M, V)$ . Wir sagen  $\Phi$  heie **Kontraktion** genau dann, wenn eine Konstante  $0 \leq k < 1$  existiert, so dass fur alle  $x, y \in M$  gelte, dass

$$\|\Phi(x) - \Phi(y)\| \leq k \|x - y\|. \quad (\text{K})$$

Eine Konstante  $k$ , welche diesen Zusammenhang erfullt, heit auch **Kontraktionskonstante** von  $\Phi$ .

**i** Anschaulich bedeutet es, wenn  $\Phi$  eine Kontraktion ist, dass die Bilder von  $\Phi$  naher beieinander liegen als die Urbilder von  $\Phi$ .

#### Korrolar 7.4.1 (Stetigkeit)

Eine jede Kontraktion ist stetig.

Das in Definition 7.4 definierte Kriterium fur Kontraktionen (K) ist nur sehr umstandlich zu uberprufen. Folgender Satz vereinfacht dieses Kriterium, schrnkt aber gleichzeitig auch unsere Auswahlmenge an Funktionen ein:

#### Satz 7.2 (Kontraktionskriterium)

Im Fall  $V = \mathbb{R}$  ist fur die Kontraktionseigenschaft hinreichend, dass  $f$  differenzierbar ist mit  $k_* := \sup_{x \in M} |f'| < 1$ .  $k_*$  ist dann **Kontraktionskonstante** von  $\Phi$ .

**i** Es ist sicherlich auch moglich Satz 7.2 auch auf  $V = \mathbb{R}^n$  anzupassen, dann heie die Bedingung  $\sup_{x \in M} \|\mathcal{J}f(x)\| < 1$ . In diesem Fall ist es aber notwendig zu spezifizieren, welche Matrixnorm man verwendet.

Wir kommen nun zu einer Art „Hauptsatz“ in unserem Kapitel uber Fixpunktiterationen, dem Fixpunktsatz von Banach. Mit ihm wollen wir alles, was wir bislang uber Fixpunkte und ihre Iterationen gelernt haben zusammenfassen um eine moglichst allgemeine Aussage uber die Konvergenz von Fixpunktiterationen zu treffen.

#### Satz 7.3 (Fixpunktsatz von BANACH)

Sei  $(V, \|\cdot\|)$  ein **Banachraum** und sei  $\emptyset \neq M \subseteq V$  eine **abgeschlossene Teilmenge** von  $V$ . Sei  $\Phi \in \text{Abb}(M, V)$  **selbstabbildend**, sprich  $\Phi(M) \subseteq M$ , sowie eine **Kontraktion**. Dann hat  $\Phi$  **GENAU** einen Fixpunkt  $x_* \in M$  und  $x_*$  ist Grenzwert der Folge  $(*)$ , wobei  $x_0 \in M$  **beliebig** ist.

Sei  $k$  zudem noch eine Kontraktionskonstante von  $\Phi$ , so fallt der Approximationsfehler in jedem Iterationsschritt um mindestens den Faktor  $k$ , also

$$\|x_{n+1} - x_*\| \leq k \|x_n - x_*\| \quad \text{damit auch} \quad \|x_n - x_*\| \leq k^n \|x_0 - x_*\|$$

Satz 7.3 ist „**konstruktiv**“, damit erklart er insbesondere auch den Weg zum Fixpunkt und nicht nur Existenz oder Eindeutigkeit.

**i** Falls  $M = V$ , so ist sowohl die Eigenschaft der Abgeschlossenheit, als auch die Selbstabbildungseigenschaft **trivialerweise** erfullt.

### 7.1.1 Newtonverfahren

#### Verfahren 7.1 (*Das Newton-Verfahren als Fixpunktiteration*)

Das Newton-Verfahren

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$$

zur Bestimmung einer Nullstelle  $x_*$  von  $f$  ist eine Fixpunktiteration, und zwar für die Funktion

$$\Phi(x) := x - \frac{f(x)}{f'(x)}.$$

Die asymptotische Kontraktionskonstante ist hierbei 0.

**i** Das Newton-Verfahren hat eine Fehlerreduktion der Form  $\|x_{n+1} - x_*\| \approx c \cdot \|x_n - x_*\|^2$ , sollte man  $f \in \mathcal{C}^2$  und  $f'(x_*) \neq 0$  voraussetzen.

Auch diese Erkenntnis haben wir im zweiten Semester aus dem Newton-Verfahren gezogen. Wir definieren deswegen ebenso:

#### Definition 7.5 (*Lineare und quadratische Konvergenz*)

Sei  $(x_n) \subset (V, \|\cdot\|)$  eine Folge mit Grenzwert  $x_* \in V$  und  $V$  sei normiert. Das Verfahren heiÙe ...

... **linear konvergent** genau dann, wenn für ein  $k < 1$  der Zusammenhang  $|x_n - x_*| \leq k \cdot |x_{n-1} - x_*|$  erfüllt ist.

... **quadratisch konvergent** genau dann, wenn für ein  $c > 0$  der Zusammenhang  $\|x_n - x_*\| \approx c \cdot \|x_{n-1} - x_*\|^2$  erfüllt ist.

An dieser Stelle folgt nun ein abschließender Satz dieses Kapitels, hier an dieser Stelle allerdings ohne Beweis, eine Beweisidee sei mit dem Satz von Taylor genannt:

#### Satz 7.4 (*Hinreichendes Kriterium für quadratisches Konvergenzverhalten*)

Eine Fixpunktiteration sei mindestens **quadratisch konvergent** genau dann, wenn die asymptotische Kontraktionskonstante  $k_*$  gleich null ist.

**i** Fixpunktiterationen sind also im Allgemeinen – sofern die Voraussetzungen des Satzes 7.3 erfüllt sind – **linear konvergent**, das heißt der Fehler reduziert sich pro Iterationsschritt mindestens um die Kontraktionskonstante  $k$ . Als gute Näherung für die Fehlerabschätzung kann – bei hinreichender Nähe zu  $x_*$  – auch die asymptotische Kontraktionskonstante  $k_*$  verwendet werden.

#### Verfahren 7.2 (*Newtonverfahren im $\mathbb{R}^n$* )

Sei  $f \in \text{Abb}(\mathbb{R}^n, \mathbb{R}^n) \in \mathcal{C}^2(D)$  und es sei die Jacobimatrix  $\mathcal{J}f(x)$  invertierbar für alle  $x \in \mathbb{R}^n$ . So lautet das **Newtonverfahren** zur Bestimmung einer Nullstelle von  $f$ :

$$x_{n+1} := x_n - [\mathcal{J}f(x_n)]^{-1} \cdot f(x_n)$$

Das Verfahren ist – wie im skalaren Fall – **lokal quadratisch konvergent**.

#### Definition 7.6 (*Lokale quadratische Konvergenz*)

Es existiere eine Umgebung  $U_\varepsilon(x_*)$  um  $x_*$  derart, dass wenn der Startwert  $x_0 \in U_\varepsilon(x_*)$  aus ebendieser Umgebung gewählt sei, das Verfahren dann **quadratisch konvergent** heiÙe. Lokal ist in diesem Fall dann tatsächlich als Umgebung zu verstehen.

**Verfahren 7.3 („Effizienteres“ Newtonverfahren im  $\mathbb{R}^n$ )**

- ① Löse das lineare Gleichungssystem  $(Jf)(x_n) \cdot \Delta x = -f(x_n)$
- ② Setze  $x_{n+1} := x_n + \Delta x$

**7.1.2 Sekantenverfahren**

Im Gegensatz zum Newtonverfahren ist hierbei keine Kenntnis der Ableitungen notwendig.

**Verfahren 7.4 (Sekantenverfahren)**

Sei  $f \in \text{Abb}(\mathbb{R}, \mathbb{R})$  und zwei Startwerte  $x_0$  und  $x_1$ . Bestimme dann den Schnittpunkt der Sekante durch  $(x_0, f(x_0))$  und  $(x_1, f(x_1))$  mit der  $x$ -Achse:

$$x_{n+1} := \frac{x_{i-1} \cdot f(x_i) - x_i \cdot f(x_{i-1})}{f(x_i) - f(x_{i-1})}$$

Das Verfahren hat die Konvergenzordnung  $p = \frac{\sqrt{5}+1}{2}$ .

**7.1.3 Bisektionsverfahren****Verfahren 7.5 (Bisektionsverfahren)**

Sei  $f \in \text{Abb}(\mathbb{R}, \mathbb{R})$  und zwei Startwerte  $x_0$  und  $x_1$  mit  $f(x_0) \cdot f(x_1) < 0$  gegeben, so gibt es im Intervall  $[x_0, x_1]$  mindestens eine Nullstelle.

Bestimme damit den Mittelpunkt  $x_{i+1} = \frac{1}{2} \cdot (x_i + x_{i-1})$  und betrachte danach das Intervall  $[x_{i-1}, x_{i+1}]$  falls  $f(x_{i-1}) \cdot f(x_{i+1}) < 0$  oder  $[x_{i+1}, x_i]$  falls  $f(x_{i+1}) \cdot f(x_i) < 0$ .

Das Verfahren hat gerade noch so eine lineare Konvergenzordnung.

**7.1.4 Regula falsi**

Im Prinzip eine Mischung aus Bisektions- und Sekantenverfahren, bei dem man wie beim Bisektionsverfahren mit zwei Punkten startet, danach aber *wie beim Sekantenverfahren* den Schnittpunkt der Sekante mit der  $x$ -Achse bestimmt um anschließend wieder das neue Intervall zu bestimmen. Konvergenzordnung ähnlich des Bisektionsverfahrens.

**7.2 Jacobi-Verfahren****Verfahren 7.6 (Jacobi-Verfahren / Gesamtschrittverfahren)**

Gegeben sei ein lineares Gleichungssystem  $Ax = b$  mit  $A = (a_{ij}) \in \mathbb{R}^n \times \mathbb{R}^n$  mit  $n \in \mathbb{N}$ . Seien alle Diagonaleinträge  $a_{ii}$  ungleich 0. Wir konstruieren uns gemäß dem Beispiel oben (??) ein  $\Phi(x) := -D^{-1}Rx + D^{-1}b$ . Die Fixpunktiteration  $x_{m+1} := \Phi(x_m)$  mit dem so konstruierten  $\Phi$  heie dann **Jacobi-** oder auch **Gesamtschrittverfahren**. Ein Iterationsschritt sehe dabei wie folgt fr alle  $i = 1, \dots, n$  aus:

$$(x_{m+1})_i := \frac{1}{a_{ii}} \left( b_i - \sum_{j \in \{1, \dots, n\} \setminus \{i\}} a_{ij}(x_m)_j \right)$$

**Satz 7.5 (Hinreichendes Konvergenzkriterium)**

Sei  $\Phi$  wie in Definition 7.1, dann ist  $\Phi$  eine Kontraktion bezüglich der  $\|\cdot\|_\infty$ -Norm auf ganz  $\mathbb{R}^n$ , falls die Systemmatrix  $A$  **diagonaldominant** ist. Damit konvergiert Verfahren 7.6.  
 Sei  $A$  unzerlegbar und schwach diagonaldominant, so konvergieren ebenfalls alle Verfahren. Eine Matrix  $A$  heißt unzerlegbar, wenn ihr Graph **stark zusammenhängend** ist.

Sei  $A$  in eine strikte untere Dreiecksmatrix  $L$ , eine Diagonalmatrix  $D$  und eine strikte obere Dreiecksmatrix  $R$  aufgeteilt, so ist die Iterationsmatrix von  $A$  bestimmt durch

$$V_J = -D^{-1}(L + R).$$

! Hierbei werden neuere, bessere Werte ignoriert, das Verfahren ist dafür allerdings **besser** parallelisierbar.

**7.3 Gauss-Seidel-Verfahren****Verfahren 7.7 (Gauß-Seidel-Verfahren)**

Gegeben sei ein lineares Gleichungssystem  $Ax = b$  mit  $A = (a_{ij}) \in \mathbb{R}^n \times \mathbb{R}^n$  mit  $n \in \mathbb{N}$ . Seien des Weiteren alle Diagonaleinträge  $a_{ii}$  ungleich 0. Ein Verfahren heiße **Gauß-Seidel-** oder auch **Einzelschrittverfahren** genau dann, wenn es zur Berechnung der  $i$ -ten Komponente des neuen Vektors die Iterationsvorschrift

$$x_{m+1,i} := \frac{1}{a_{ii}} \left( b_i - \sum_{j=1}^{i-1} a_{ij}x_{m+1,j} - \sum_{j=i+1}^n a_{ij}x_{m,j} \right),$$

bei der der Unterschied zum Jacobiverfahren (7.6) **blau** hervorgehoben wurde, verwendet.

Sei  $A$  wieder in eine strikte untere Dreiecksmatrix  $L$ , eine Diagonalmatrix  $D$  und eine strikte obere Dreiecksmatrix  $R$  aufgeteilt, so ist die Iterationsmatrix von  $A$  bestimmt durch

$$V_{GS} = -(L + D)^{-1}R.$$

**Anmerkungen**

i Das Gauß-Seidel-Verfahren konvergiert im Allgemeinen (bei tridiagonalen Matrizen) um den Faktor **zwei schneller** als das Jacobiverfahren. Einen Beweis dieses Satzes findet sich in Plato, Numerische Mathematik kompakt mit Korollar 10.42.

Das Gauß-Seidel-Verfahren lässt sich platzsparender implementieren, bei der Berechnung von  $x_{m+1,i}$  ist  $x_{m,i}$  überschreibbar, der Wert wird nicht mehr gebraucht.

**7.4 SOR-Verfahren**

Das Gauß-Seidel-Verfahren kann man relaxiert auffassen, man nennt dies dann sukzessiv überrelaxiert „*successive overrelaxion*“ (SOR). Die Idee hierbei ist den Zielwert entweder **bewusst** zu „überschießen“ oder **bewusst** das „Überschießen“ abzuschwächen. Auch hier gilt dann  $\omega \in (0, 2)$ , sonst divergiert das Verfahren. Man erhält dann:

$$x_n = H_\omega x_{n-1} + \omega(D + \omega L)^{-1}b,$$

wobei  $H_\omega$  definiert ist als

$$H_\omega := (D + \omega L)^{-1} [(1 - \omega)D - \omega R] = E_n - \omega(D + \omega L)^{-1} A.$$

### Verfahren 7.8 (SOR-Verfahren)

Gegeben sei ein lineares Gleichungssystem  $Ax = b$  mit  $A = (a_{ij}) \in \mathbb{R}^n \times \mathbb{R}^n$  mit  $n \in \mathbb{N}$ . Seien des Weiteren alle Diagonaleinträge  $a_{ii}$  ungleich 0. Ein Verfahren heie **SOR-Verfahren** genau dann, wenn es zur Berechnung der  $i$ -ten Komponente des neuen Vektors die Iterationsvorschrift

$$x_i^{m+1} := (1 - \omega) \cdot x_i^m + \omega \cdot \left( \frac{1}{a_{ii}} \left( b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{m+1} - \sum_{j=i+1}^n a_{ij} x_j^m \right) \right)$$

verwendet.

Sei  $A$  wieder in eine strikte untere Dreiecksmatrix  $L$ , eine Diagonalmatrix  $D$  und eine strikte obere Dreiecksmatrix  $R$  aufgeteilt, so ist die Iterationsmatrix von  $A$  bestimmt durch

$$V_{\text{SOR}} = - \left( \frac{1}{\omega} D + L \right)^{-1} \left( \frac{\omega - 1}{\omega} D + R \right).$$

Im Vergleich zu der Komplexitt von  $\mathcal{O}(n)$  bei Jacobi und Gau-Seidel kommt es bei SOR-Verfahren mit einem optimalen Parameter zu einer  $\mathcal{O}(\sqrt{n})$  Komplexitt.

# NICHTLINEARE OPTIMIERUNG

## Satz 8.1 (*Lagrange-Multiplikatoren bei mehreren Nebenbedingungen*)

Sei  $f : D \subseteq \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $g : D \rightarrow \mathbb{R}^m$  mit  $m \leq n$ . Seien  $f, g \in \mathcal{C}^1(D)$ . Ferner besitze  $f$  in  $\xi$  unter der Nebenbedingung  $g(\xi) = 0$  eine Extremstelle. Dann ist eine **notwendige** Bedingung für die Existenz dieser Extremstelle, dass ...

$$(I) \quad \nabla f(x_0) = \sum_{i=1}^m \lambda_i \cdot \nabla g_i(\xi)$$

$$(II) \quad g(\xi) = 0$$

$\lambda \in \mathbb{R}^m$  beschreibt wieder den Lagrangemultiplikator, das dazugehörige Funktional ist  $\mathcal{L}(\xi, \lambda) := f(\xi) - \langle \lambda, g(\xi) \rangle$ .

## 8.1 Optimierung mittels Newtonverfahren

### Verfahren 8.1 (*Newtonverfahren zur Optimierung*)

Sei  $f \in \text{Abb}(\mathbb{R}^n, \mathbb{R})$  und es sei die Hessematrix  $\mathcal{H}f(x)$  invertierbar für alle  $x \in \mathbb{R}^n$ . So lautet das **Newtonverfahren** zur Bestimmung eines Optimums von  $f$ :

$$x_{n+1} := x_n - [\mathcal{H}f(x_n)]^{-1} \cdot \nabla f(x_n)$$

Das Problem ist der hohe Aufwand beim Berechnen der  $n^2/2$  vielen zweiten Ableitungen.

## 8.2 Abstiegsverfahren

### Verfahren 8.2 (*Abstiegsverfahren*)

Sei ein Startwert  $x_0 \in \mathbb{R}^n$  sowie eine Startsuchrichtung  $s_0 \in \mathbb{R}^n$  gegeben. Dann heißt die Iterationsvorschrift

$$x_{n+1} = x_n + \alpha_n \cdot s_n,$$

mit  $\alpha_n > 0$  als Suchweite und  $s_0$  als Suchrichtung **Abstiegsverfahren**, falls  $s_n \circ \nabla f(x_n) < 0$  gilt.

Die optimale Schrittweite kann man durch das näherungsweise Lösen des eindimensionalen Minimierungsproblems

$$t_n = \arg \min_{t > 0} \left\{ f(x_n + t \cdot s_n) \right\}$$

bestimmen, bei einem Newtonverfahren ist  $t_n = 1$  zu wählen, will man es besser dämpfen so auch  $t_n < 1$ . Ansonsten ist das nachfolgende Verfahren gut für Schrittweitenwahlen geeignet.

### 8.2.1 Abstiegsverfahren des goldenen Schnittes

#### Verfahren 8.3 (Verfahren des goldenen Schnittes)

Sei  $f \in \text{Abb}(\mathbb{R}, \mathbb{R})$  und drei Startwerte  $x_0, x_1$  und  $x_2$  mit  $f(x_0) > f(x_1) < f(x_2)$  gegeben. Wähle dann einen neuen Wert  $x_n$  mit  $x_{n-3} < x_n < x_{n-1}$ . Gilt nun  $f(x_n) < f(x_2)$ , so bildet man das neue Tripel  $(x_{n-2}, x_n, x_{n-1})$ , andernfalls mit  $(x_n, x_{n-2}, x_{n-1})$ .

$d$  sollte dabei im **goldenen Schnitt** gewählt worden sein.  $(\frac{\sqrt{5}+1}{2})$

### 8.2.2 Auffinden der Suchrichtung

**Newtonverfahren** Das Newtonverfahren wird zum Abstiegsverfahren, wenn die Hessematrix positiv definit ist.

#### Gradientenverfahren

#### Verfahren 8.4 (Gradienten-Verfahren)

Sei  $f$  „glatt“ und  $x_n$  eine hinreichende Näherung an die Minimalstelle von  $f$ . Dann ist

$$x_{n+1} := x_n - \alpha_n \nabla f(x_n) \text{ mit bspw. } \alpha_n := \frac{\langle \nabla f(x_n), \nabla f(x_n) \rangle}{\langle A \nabla f(x_n), \nabla f(x_n) \rangle}$$

im Allgemeinen eine bessere Näherung. Es gilt dann also  $s_n = -\nabla f(x_n)$ .

**Verfahren der konjugierten Gradienten (cg-Verfahren)** Ein **quadratisches Funktional** ist von der Form

$$F(x) = x^T A x + 2b^T x + c,$$

wobei  $A$  eine **positiv definite** quadratische Matrix,  $b$  und  $x$  Vektoren und  $c$  ein Skalar ist. Man nennt dann zwei Vektoren  $u$  und  $v$   **$A$ -konjugiert**, wenn

$$u^T A v = 0.$$

#### Satz 8.2 (Schrittzahl beim cg-Verfahren)

Wählt man als Suchrichtungen  $(s_i)$  eine Folge an paarweisen konjugierter Richtungen und dazu optimale Schrittweiten  $t_i$ , so ist man nach (spätestens)  $n$  Schritten im Minimum.

#### Verfahren 8.5 (Verfahren der konjugierten Gradienten)

① Wähle  $x^1$  beliebig und berechne

$$g^1 = b + A x^1, \quad s^1 = -g^1 \quad \text{und} \quad \alpha^1 = \frac{g^1 \circ g^1}{s^1 \circ A s^1}.$$

② Für  $k \in \{1, 2, \dots, n\}$  tue

2.1 Falls  $\|g^{k+1}\| < \varepsilon$  halte an.

2.2 Bestimme von  $x^k$  in Richtung  $s^k$  das Minimum  $x^{k+1} = x^k + \alpha^k s^k$ .



- 2.3 Aktualisiere dann den Gradienten durch  $g^{k+1} = g^k + \alpha^k A s^k$
- 2.4 Aktualisiere dann die Suchrichtung, so dass die neue Suchrichtung A-konjugiert zu allen bisherigen ist:

$$s^{k+1} = -g^{k+1} + \beta^k \cdot s^k \quad \text{mit} \quad \beta^k = \frac{g^{k+1} \circ g^{k+1}}{g^k \circ g^k}$$

- 2.5 Bestimme die neue optimale Schrittweite  $\alpha^{k+1}$  mit

$$\alpha^{k+1} = \frac{g^{k+1} \circ g^{k+1}}{s^{k+1} \circ A s^{k+1}}$$

Dieses Verfahren kann auch zur Lösung linearer Gleichungssysteme mit **positiv definiten Koeffizientenmatrix** verwendet werden.

### Verfahren 8.6 (Verfahren der konjugierten Gradienten zum Lösen allgemeiner LGS)

- 1 Wähle  $x^1$  beliebig und berechne

$$g^1 = \nabla f(x^0) \quad , \quad s^1 = -g^1 \quad \text{und} \quad \alpha^1 = \arg \min_t \left\{ f(x^0 + t s^0) \right\}.$$

- 2 Für  $k \in \{1, 2, \dots, n\}$  tue

- 2.1 Falls  $\|g^{k+1}\| < \varepsilon$  halte an.

- 2.2 Bestimme von  $x^k$  in Richtung  $s^k$  das Minimum  $x^{k+1} = x^k + \alpha^k s^k$ .

- 2.3 Aktualisiere dann den Gradienten durch  $g^{k+1} = \nabla f(x^{k+1})$

- 2.4 Aktualisiere dann die Suchrichtung, so dass die neue Suchrichtung A-konjugiert zu allen bisherigen ist:  $s^{k+1} = -g^{k+1} + \beta^k \cdot s^k$  mit

$$\beta^k = \frac{g^{k+1} \circ g^{k+1}}{g^k \circ g^k} \text{ (Fletcher/Reeves)} \quad \text{oder} \quad \beta^k = \frac{g^{k+1} \circ (g^{k+1} - g^k)}{g^k \circ g^k} \text{ (Pollack/Ribière)}$$

- 2.5 Bestimme die neue optimale Schrittweite  $\alpha^{k+1}$  mit

$$\alpha^{k+1} = \arg \min_t \left\{ f(x^{k+1} + t s^{k+1}) \right\}$$

**Quasinewtonverfahren** Wir wählen an der Hessematrix statt eine Folge positiv definiten Matrizen  $(H_k)_{k \in \mathbb{N}}$ , mit einer jeder die Bedingung

$$x_k - x_{k-1} = H_k (\nabla f(x_k) - \nabla f(x_{k-1}))$$

erfüllt ist. Diese Bedingung heißt dann **Quasinewtonbedingung**. Der Vorteil liegt in der erhaltenen superlinearen Konvergenz des Verfahrens.

# NUMERISCHE INTEGRATION

## 9.1 Monte-Carlo

## 9.2 Newton-Cotes-Formeln

Das Kapitel über die Newton-Cotes-Formeln befasst sich mit dem Sachverhalt, dass sich Polynome in Abhängigkeit dessen Grades exakt integrieren lassen.

Die grundsätzliche Idee dahinter sehe wie folgt aus:

Man bestimme zu einer gegebenen Funktion  $f(x)$  das Polynom  $p(x)$  durch die folglich vorgestellten Regeln, welches die Funktionswerte an den äquidistanten Stellen interpoliert und nehme somit als Näherungswert

$$\int_a^b p(x) dx \text{ anstelle vom ursprünglichen } \int_a^b f(x) dx$$

Hierbei gilt es zu beachten, dass zu einer Funktion  $f(x)$  mit  $n+1$  gegebenen Punkten

$$(x_0, f(x_0)), (x_1, f(x_1)), \dots, (x_n, f(x_n))$$

genau ein interpolierendes Polynom vom Grad  $n$  existiert.

### 9.2.1 Mittelpunktsregel

Es sei ein Punkt  $(x_0, f(x_0))$  mit  $x_0 = \frac{a+b}{2}$  gegeben.

$$\leftrightarrow \text{\#Punkte} = 1 \Leftrightarrow \text{Polynom Grad} = 0$$

Somit gilt für unser zu interpolierende Polynom:

$$p_0(x) = f(x_0) = \text{const}$$

$$\Leftrightarrow \int_a^b f(x) dx \approx \int_a^b p_0(x) dx = (b-a)f(x_0)$$

Dies ist auch als Rechtecksregel bekannt, da wir im Prinzip das Integral durch ein Rechteck approximieren.

Zusammengefasst:

- $x_0 = a$  : linke Rechtecksregel
- $x_0 = \frac{a+b}{2}$  : Mittelpunktsregel
- $x_0 = b$  : rechte Rechtecksregel

### 9.2.2 Trapezregel

Es seien zwei Punkte  $(x_0, f(x_0)), (x_1, f(x_1))$  mit  $x_0 = a, x_1 = b$  gegeben.

$$\Leftrightarrow \text{\#Punkte} = 2 \Leftrightarrow \text{Polynom Grad} = 1$$

Somit gilt für unser zu interpolierende Polynom:

$$p_1(x) = \frac{b-x}{b-a}f(x_0) + \frac{x-a}{b-a}f(x_1)$$

$$\Leftrightarrow \int_a^b f(x) dx \approx \int_a^b p_1(x) dx = (b-a) \frac{f(a)+f(b)}{2}$$

### 9.2.3 Simpsonregel

Es seien drei Punkte  $(x_0, f(x_0)), (x_1, f(x_1)), (x_2, f(x_2))$  mit  $x_0 = a, x_1 = \frac{a+b}{2}, x_2 = b$  gegeben.

$$\Leftrightarrow \text{\#Punkte} = 3 \Leftrightarrow \text{Polynom Grad} = 2$$

Somit gilt für unser zu interpolierende Polynom (analog zu Lagrange):

$$p_2(x) = \frac{(x-x_1)(x-x_2)}{(x_0-x_1)(x_0-x_2)}f(x_0) + \frac{(x-x_0)(x-x_2)}{(x_1-x_0)(x_1-x_2)}f(x_1) + \frac{(x-x_0)(x-x_1)}{(x_2-x_0)(x_2-x_1)}f(x_2)$$

$$\Leftrightarrow \int_a^b f(x) dx \approx \int_a^b p_2(x) dx = (b-a) \left( \frac{1}{6}f(x_0) + \frac{4}{6}f(x_1) + \frac{1}{6}f(x_2) \right)$$

# TABELLE DER ALGORITHMENAUFWÄNDE UND FEHLERABSCHÄTZUNGEN

Im Folgenden sei dann nur noch eine kleine Übersicht über die verschiedenen Aufwände und Fehlerabschätzungen einiger Verfahren und Algorithmen gegeben.

## A.1 Aufwände

	Algorithmus/Verfahren	Komplexität
Allgemeine Lineare Algebra	Matrix-Vektor-Multiplikation $Ab$ mit vollbesetzter Matrix $A$	$\mathcal{O}(n^2)$
	Matrix-Vektor-Multiplikation $Ab$ mit $m$ -diagonaler Matrix $A$	$\mathcal{O}(m \cdot n)$
	Matrix-Vektor-Multiplikation $Ab$ mit Rang-1-Matrix und gegebenem $A = uv^T$	$\mathcal{O}(n)$
	Berechnung der euklidischen Vektor-Norm	$\mathcal{O}(n)$
	Matrix-Matrix-Multiplikation $AB$ mit vollbesetzten Matrizen $A$ und $B$	$\mathcal{O}(n^3)$
	Matrix-Matrix-Multiplikation $AB$ mit vollbesetzter Matrix $A$ und $m$ -diagonaler Matrix $B$	$\mathcal{O}(m \cdot n^2)$
	Berechnung der Inversen einer invertierbaren $(n \times n)$ -Diagonalmatrix.	$\mathcal{O}(n)$
	Rang einer invertierbaren $n \times n$ -Matrix	$n$
Direkte Verfahren	Bestimmung der LR- bzw. QR-Zerlegung	$\mathcal{O}(n^3)$
	Bestimmung LR-Zerlegung von $m$ -diagonaler Matrix	$\mathcal{O}(m^2 \cdot n)$
	Unterscheidungsfaktor zwischen LR-Zerlegung und Householderspiegelungen	$q_{HR} \approx 2 \cdot q_{LR}$
	Unterscheidungsfaktor zwischen LR-Zerlegung und Givensrotationen	$q_{GR} \approx 4 \cdot q_{LR}$
	Unterscheidungsfaktor zwischen LR-Zerlegung und Cholesky	$q_{Ch} \approx \frac{1}{2} \cdot q_{LR}$
	Lösen des Gleichungssystems mit geg. LR/QR/SVD-Zerlegung	$\mathcal{O}(n^2)$
	Lösen des Gleichungssystems mit geg. LR/QR-Zerlegung bei $m$ -diagonaler Matrix	$\mathcal{O}(n \cdot m)$
	Vorwärts- und Rückwärtssubstitutionen	$\mathcal{O}(n^2)$
	Bestimmung der Determinante bei geg. LR-Zerlegung	$\mathcal{O}(n)$
	Bestimmung der Betragsdeterminante bei geg. PLR-Zerlegung	$\mathcal{O}(n)$
	Bestimmung der Determinante bei geg. PLR-Zerlegung	$\mathcal{O}(n^2)$
	Bestimmung der Betragsdeterminante bei geg. QR-Zerlegung	$\mathcal{O}(n)$

(Fortgesetzt auf nächster Seite  $\leftrightarrow$ )

		(Aufwandstabelle von vorheriger Seite fortgesetzt)	
		Algorithmus/Verfahren	Komplexität
Faltung		Anzahl der arithmetischen Operationen bei einem nicht separierbaren $M \times M$ Filter bei einem $N \times N$ Bildausschnitt	$2 \cdot N^2 \cdot M^2$
		Anzahl der arithmetischen Operationen bei einem separierbaren $M \times M$ Filter bei einem $N \times N$ Bildausschnitt	$4 \cdot N^2 \cdot M$
		Berechnung der Fourier-Transformation <i>ohne</i> FFT	$\mathcal{O}(n^2)$
		Berechnung der Fourier-Transformation <i>mit</i> FFT	$\mathcal{O}(n \cdot \log n)$
Bézier		Auswerten eines Punktes auf einer Bézierkurve mittels DE CASTELJAU	$\mathcal{O}(n^2)$
		Anzahl der Kontrollkurve einer Bézierkurve mit Grad $n$	$n + 1$
		Grad einer Bézierkurve mit $n$ Kontrollpunkten	$n - 1$
Optimierung		Aufwand eines Iterationsschrittes bei JACOBI, GAUSS-SEIDEL oder SOR	$\mathcal{O}(n^2)$
		Anzahl der Iterationsschritte bei JACOBI oder GAUSS-SEIDEL	$\mathcal{O}(n)$
		Anzahl der Iterationsschritte bei SOR	$\mathcal{O}(\sqrt{n})$
		Konvergenzordnung bei JACOBI, GAUSS-SEIDEL oder SOR	$\rho = 1$
		Bestimmung der Koeffizienten bei Interpolation mit Newtonbasis	$\mathcal{O}(n^2)$
		Sei $A$ eine dünnbesetzte $(n \times n)$ -Matrix mit höchstens $m$ von Null verschiedenen Einträgen pro Zeile. Welche Komplexität hat die Berechnung von $k$ Iterationsschritten des Jacobi/Gauß-Seidel/SOR-Verfahrens?	$\mathcal{O}(k \cdot (n \cdot m))$

## A.2 Fehlerabschätzungen

Verfahren	Fehlerabschätzung
Stückweise konstante Interpolation (NEAREST NEIGBOR)	$\mathcal{O}(h)$
Stückweise lineare Interpolation	$\mathcal{O}(h^2)$
Stückweise kubische Interpolation (CATMULL-ROM)	$\mathcal{O}(h^3)$
Midpoint Subdivision	$\mathcal{O}\left(\frac{1}{(\text{levelzahl})^4}\right)$
(Iterierte) TRAPEZregel	$\mathcal{O}(h^2)$
(Iterierte) SIMPSONregel	$\mathcal{O}(h^4)$

# PYTHON 3 — ZUSAMMENFASSUNG UND KODESCHNIPSEL

## B.1 Zusammenfassung der wichtigsten Befehle

Nachfolgend kommt eine kleiner Zusammenschnitt der elementarsten und wichtigsten Befehlskombination, welche Python 3 zu bieten hat.

Arrays	Types	Creating Arrays	Manipulation	NumPy															
■ 1D Array <table border="1" style="display: inline-table; vertical-align: middle;"> <tr><td>0</td><td>1</td><td>2</td></tr> <tr><td>91</td><td>1</td><td>102</td></tr> </table> ■ 2D Array <table border="1" style="display: inline-table; vertical-align: middle;"> <tr><td>0</td><td>1</td><td>2</td></tr> <tr><td>91</td><td>1</td><td>102</td></tr> <tr><td>5</td><td>64</td><td>123</td></tr> </table>	0	1	2	91	1	102	0	1	2	91	1	102	5	64	123	<b>Aggregate</b> np.int64 np.float32 np.bool np.string_  a.sum() a.cumsum([axis=]) a.min([axis=]) a.mean() a.corrcoef()	np.eye(n) np.zeros(shape, dtype=type) np.full(shape, const) <i>shape is a tuple</i>	<b>Aggregate</b> a.max([axis=]) a.median() - correlation coeff.	<b>Aggregate</b> Transposing Array <pre>&gt;&gt;&gt; i = np.transpose(b) &gt;&gt;&gt; i.T</pre> Changing Array Shape <pre>&gt;&gt;&gt; b.ravel() &gt;&gt;&gt; g.reshape(3,-2)</pre> Adding/Removing Elements <pre>&gt;&gt;&gt; h.resize((2,6)) &gt;&gt;&gt; np.append(h,g) &gt;&gt;&gt; np.insert(a, 1, 5) &gt;&gt;&gt; np.delete(a, [1])</pre>
0	1	2																	
91	1	102																	
0	1	2																	
91	1	102																	
5	64	123																	
<b>Sorting</b> a.sort([axis=])	<b>Base Types</b> integer, float, boolean, string, bytes <pre>int 783 0 -192 0b010 0o642 0xF3 float 9.23 0.0 -1.7e-6 bool True False str "One\nTwo" bytes b"toto\xfe\775"</pre>	<b>Container Types</b> ■ ordered sequences, fast index access, repeatable values <pre>list [1, 5, 9] tuple (1, 5, 9) str bytes (ordered sequences of chars / bytes)</pre> ■ key containers, no a priori order, fast key access, each key is unique <pre>dictionary dict {"key": "value"} collection set {"key1", "key2"}</pre>	<b>Conversions</b> type(expression) can specify integer number base in 2 <sup>nd</sup> parameter truncate decimal part rounding to 1 decimal (0 decimal → integer number) False for null x, empty container x, None or False x; True for other x representation string of x for display (cf. formatting on the back) code ↔ char literal representation string of x separator str and sequence of str → assembled str str splitted on whitespaces → list of str str splitted on separator str → list of str sequence of one type → list of another type (via list comprehension)																
<b>Identifiers</b> for variables, functions, modules, classes... names a...zA...Z_ followed by a...zA...Z_0...9 □ diacritics allowed but should be avoided □ language keywords forbidden □ lower/UPPER case discrimination © a toto x7 y_max BigOne © &#x2191; and for	<b>Variables assignment</b> = assignment ↔ binding of a name with a value 1) evaluation of right side expression value 2) assignment in order with left side names x=1.2+8+sin(y) a=b=c=0 assignment to same value y,z,r=9.2,-7.6,0 multiple assignments a,b=b,a values swap a,*b=seq } unpacking of sequence in *a,b=seq } item and list x+=3 increment ↔ x=x+3 and *= x-=2 decrement ↔ x=x-2 /= x=None « undefined » constant value %= del x remove name x ...	<b>Base Types</b> integer, float, boolean, string, bytes <pre>int 783 0 -192 0b010 0o642 0xF3 float 9.23 0.0 -1.7e-6 bool True False str "One\nTwo" bytes b"toto\xfe\775"</pre>	<b>Container Types</b> ■ ordered sequences, fast index access, repeatable values <pre>list [1, 5, 9] tuple (1, 5, 9) str bytes (ordered sequences of chars / bytes)</pre> ■ key containers, no a priori order, fast key access, each key is unique <pre>dictionary dict {"key": "value"} collection set {"key1", "key2"}</pre>																

### Sequence Containers Indexing

for lists, tuples, strings, bytes...

negative index	-5	-4	-3	-2	-1
positive index	0	1	2	3	4
<b>lst</b> = [10, 20, 30, 40, 50]					
positive slice	0	1	2	3	4
negative slice	-5	-4	-3	-2	-1

**Items count**  
`len(lst) → 5`  
 index from 0 (here from 0 to 4)

Individual access to **items** via `lst[index]`  
`lst[0] → 10` ⇒ first one      `lst[1] → 20`  
`lst[-1] → 50` ⇒ last one      `lst[-2] → 40`

On mutable sequences (**list**), remove with `del lst[3]` and modify with assignment `lst[4]=25`

Access to **sub-sequences** via `lst[start slice : end slice : step]`  
`lst[: -1] → [10, 20, 30, 40]`    `lst[ : : -1] → [50, 40, 30, 20, 10]`    `lst[1:3] → [20, 30]`    `lst[:3] → [10, 20, 30]`  
`lst[1: -1] → [20, 30, 40]`    `lst[ : : -2] → [50, 30, 10]`    `lst[-3: -1] → [30, 40]`    `lst[3: ] → [40, 50]`  
`lst[: : 2] → [10, 30, 50]`    `lst[ : ] → [10, 20, 30, 40, 50]` shallow copy of sequence

Missing slice indication → from start / up to end.  
 On mutable sequences (**list**), remove with `del lst[3:5]` and modify with assignment `lst[1:4]=[15,25]`

### Boolean Logic

Comparisons : < > <= >= == != (boolean results)    ≤ ≥ = ≠  
**a and b** logical and both simultaneously  
**a or b** logical or one or other or both  
 pitfall : **and** and **or** return value of **a** or **b** (under shortcut evaluation).  
 ⇒ ensure that **a** and **b** are booleans.  
**not a** logical not  
**True** } True and False constants  
**False** }

### Statements Blocks

```

parent statement:
┌ statement block 1...
│   ...
│   statement block 2...
│   ...
└ next statement after block 1
    
```

indentation !  
 configure editor to insert 4 spaces in place of an indentation tab.

### Modules/Names Imports

```

module truc ⇔ file truc.py
from monmod import nom1, nom2 as fct
    → direct access to names, renaming with as
import monmod → access via monmod.nom1 ...
    modules and packages searched in python path (cf sys.path)
    
```

### Conditional Statement

statement block executed only if a condition is true

```

if logical condition:
    → statements block
    
```

Can go with several **elif**, **elif...** and only one final **else**. Only the block of first true condition is executed.

```

if age <= 18:
    state = "Kid"
elif age > 65:
    state = "Retired"
else:
    state = "Active"
    
```

with a var **x**:  
`if bool(x) == True:` ⇔ `if x:`  
`if bool(x) == False:` ⇔ `if not x:`

### Maths

floating numbers... approximated values  
 Operators: + - \* / // % \*\*  
 Priority (...): × ÷ ↑ ↑ a<sup>b</sup>  
 integer ÷ ÷ remainder  
 @ → matrix × python3.5+ **numpy**  
`(1+5.3) * 2 → 12.6`  
`abs(-3.2) → 3.2`  
`round(3.57, 1) → 3.6`  
`pow(4, 3) → 64.0`  
 usual order of operations

### Maths

angles in radians  
`from math import sin, pi...`  
`sin(pi/4) → 0.707...`  
`cos(2*pi/3) → -0.4999...`  
`sqrt(81) → 9.0` ✓  
`log(e**2) → 2.0`  
`ceil(12.5) → 13`  
`floor(12.5) → 12`  
 modules **math, statistics, random,**  
**decimal, fractions, numpy, etc.** (cf. doc)

### Exceptions on Errors

Signaling an error:  
`raise ExcClass(...)`

Errors processing:  

```

try:
    → normal processing block
except Exception as e:
    → error processing block
    
```

 finally block for final processing in all cases.

### Conditional Loop Statement

statements block executed as long as condition is true

```

while logical condition:
    → statements block
    
```

**beware of infinite loops!**

```

s = 0 } initializations before the loop
i = 1 } condition with a least one variable value (here i)
while i <= 100:
    s = s + i**2
    i = i + 1
print("sum:", s)
    
```

Algo:  $s = \sum_{i=1}^{100} i^2$

### Iterative Loop Statement

statements block executed for each item of a container or iterator

```

for var in sequence:
    → statements block
    
```

Go over sequence's values  
`s = "Some text"` } initializations before the loop  
`cnt = 0`  
 loop variable, assignment managed by **for** statement

```

for c in s:
    if c == "e":
        cnt = cnt + 1
    print("found", cnt, "e")
    
```

Algo: count number of **e** in the string.

### Display

```

print("v=" 3, "cm :", x, ", ", y+4)
    
```

items to display : literal values, variables, expressions

print options:  
 □ `sep=" "` items separator, default space  
 □ `end="\n"` end of print, default new line  
 □ `file=sys.stdout` print to file, default standard output

### Input

```

s = input("Instructions: ")
    
```

**input** always returns a **string**, convert it to required type (cf. boxed Conversions on the other side).

### Generic Operations on Containers

`len(c)` → items count  
`min(c)` `max(c)` `sum(c)`  
`sorted(c)` → list sorted copy  
`val in c` → boolean, membership operator **in** (absence **not in**)  
`enumerate(c)` → iterator on (index, value)  
`zip(c1, c2...)` → iterator on tuples containing **c<sub>i</sub>** items at same index  
`all(c)` → **True** if **all** **c** items evaluated to true, else **False**  
`any(c)` → **True** if **at least one** item of **c** evaluated true, else **False**

Specific to **ordered sequences containers** (lists, tuples, strings, bytes...)  
`reversed(c)` → inverted iterator    `c*5` → duplicate    `c+c2` → concatenate  
`c.index(val)` → position    `c.count(val)` → events count

loop on dict/set ⇔ loop on keys sequences  
 use **slices** to loop on a subset of a sequence

Go over sequence's **index**  
 □ modify item at index  
 □ access items around index (before / after)

```

lst = [11, 18, 9, 12, 23, 4, 17]
lost = []
for idx in range(len(lst)):
    val = lst[idx]
    if val > 15:
        lost.append(val)
        lst[idx] = 15
print("modif:", lst, "-lost:", lost)
    
```

Algo: limit values greater than 15, memorizing of lost values.

Go simultaneously over sequence's **index** and **values**:  
`for idx, val in enumerate(lst):`

### Integer Sequences

```

range([start,] end [,step])
    
```

**start** default 0, **end** not included in sequence, **step** signed, default 1

```

range(5) → 0 1 2 3 4      range(2, 12, 3) → 2 5 8 11
range(3, 8) → 3 4 5 6 7    range(20, 5, -5) → 20 15 10
range(len(seq)) → sequence of index of values in seq
    
```

**range** provides an immutable sequence of int constructed as needed

good habit : don't modify loop variable

### Operations on Lists

☞ modify original list

```

lst.append(val)      add item at end
lst.extend(seq)     add sequence of items at end
lst.insert(idx, val) insert item at index
lst.remove(val)     remove first item with value val
lst.pop([idx]) → value remove & return item at index idx (default last)
lst.sort() lst.reverse() sort / reverse liste in place
    
```

### Operations on Dictionaries

```

d[key]=value      d.clear()
d[key] → value    del d[key]
d.update(d2)      { update/add
                   { associations
d.keys()          → iterable views on
d.values()        keys/values/associations
d.items()
d.pop(key[,default]) → value
d.popitem() → (key, value)
d.get(key[,default]) → value
d.setdefault(key[,default]) → value
    
```

### Operations on Sets

Operators:

- | → union (vertical bar char)
- & → intersection
- ^ → difference/symmetric diff.
- < <= > >= → inclusion relations

Operators also exist as methods.

```

s.update(s2) s.copy()
s.add(key) s.remove(key)
s.discard(key) s.clear()
s.pop()
    
```

### Function Definition

```

function name (identifier)
named parameters
def fct(x, y, z):
    """documentation"""
    # statements block, res computation, etc.
    return res ← result value of the call, if no computed
                  result to return: return None
    
```

☞ parameters and all variables of this block exist only *in* the block and *during* the function call (think of a "black box")

Advanced: **def fct**(x, y, z, \*args, a=3, b=5, \*\*kwargs):

- \*args variable positional arguments (→ tuple), default values,
- \*\*kwargs variable named arguments (→ dict)

### Function Call

```

r = fct(3, i+2, 2*i)
storage/use of returned value    one argument per parameter
    
```

☞ this is the use of function name with parentheses which does the call

Advanced: \*sequence \*\*dict

### Files

storing data on disk, and reading it back

```

f = open("file.txt", "w", encoding="utf8")
    
```

file variable for operations	name of file on disk (+path...)	opening mode	encoding of chars for text files:
		□ 'r' read	utf8 ascii
		□ 'w' write	latin1 ...
		□ 'a' append	
		□ ... '+' 'x' 'b' 't'	

cf. modules **os**, **os.path** and **pathlib**

writing	reading
<b>f.write</b> ("coucou")	<b>f.read</b> ([n]) → next chars
<b>f.writelines</b> (list of lines)	if n not specified, read up to end!
	<b>f.readlines</b> ([n]) → list of next lines
	<b>f.readline</b> () → next line

☞ text mode **t** by default (read/write **str**), possible binary mode **b** (read/write **bytes**). Convert from/to required type!

**f.close**() ☞ dont forget to close the file after use!

**f.flush**() write cache      **f.truncate**([size]) resize

reading/writing progress sequentially in the file, modifiable with:

**f.tell**() → position      **f.seek**(position[,origin])

Very common: opening with a guarded block (automatic closing) and reading loop on lines of a text file:

```

with open(...) as f:
    for line in f:
        # processing of line
    
```

### Operations on Strings

```

s.startswith(prefix[,start[,end]])
s.endswith(suffix[,start[,end]]) s.strip([chars])
s.count(sub[,start[,end]]) s.partition(sep) → (before, sep, after)
s.index(sub[,start[,end]]) s.find(sub[,start[,end]])
s.is...() tests on chars categories (ex. s.isalpha())
s.upper() s.lower() s.title() s.swapcase()
s.casefold() s.capitalize() s.center([width,fill])
s.ljust([width,fill]) s.rjust([width,fill]) s.zfill([width])
s.encode(encoding) s.split([sep]) s.join(seq)
    
```

### Formatting

formatting directives      values to format

```

"modele{} {} {}".format(x, y, r) → str
"{selection:formatting!conversion}"
    
```

☞ Selection:

```

2
nom
0.nom
4[key]
0[2]
    
```

Examples:

```

"{:+2.3f}".format(45.72793)
→ '+45.728'
"{1:>10s}".format(8, "toto")
→ '          toto'
"{x!r}".format(x="I'm")
→ "'I\m'"
    
```

☞ Formatting:

*fill char* *alignment* *sign* *mini width* . *precision-maxwidth* *type*

```

<> ^ = + - space 0 at start for filling with 0
integer: b binary, c char, d decimal (default), o octal, x or X hexa...
float: e or E exponential, f or F fixed point, g or G appropriate (default),
string: s ... % percent
☞ Conversion : s (readable text) or r (literal representation)
    
```