

# Mathematische Einführung in Data Science

11. Oktober 2023

Sven-Ake Wegner<sup>1</sup>

---

<sup>1</sup>Fachbereich Mathematik, Universität Hamburg, Bundesstraße 55, 20146 Hamburg,  
Kommentare an [datasciencebuch@gmail.com](mailto:datasciencebuch@gmail.com) sind sehr willkommen!

# Vorwort

Kenntnisse in den Bereichen Data Science und Machine Learning werden von Absolventinnen und Absolventen eines Mathematikstudiums immer häufiger erwartet und von Studentinnen und Studenten der Mathematik dementsprechend nachgefragt. Die Idee hinter dem vorliegenden Text ist es, kanonische Themen aus den vorgenannten Bereichen passgenau für die vorgenannte Zielgruppe aufzubereiten. Hierbei wird ein grundlegendes und sorgfältiges Verständnis der behandelten Methoden, die Frage warum diese zum Ziel führen, und was deren Grenzen sind, prioritär behandelt.

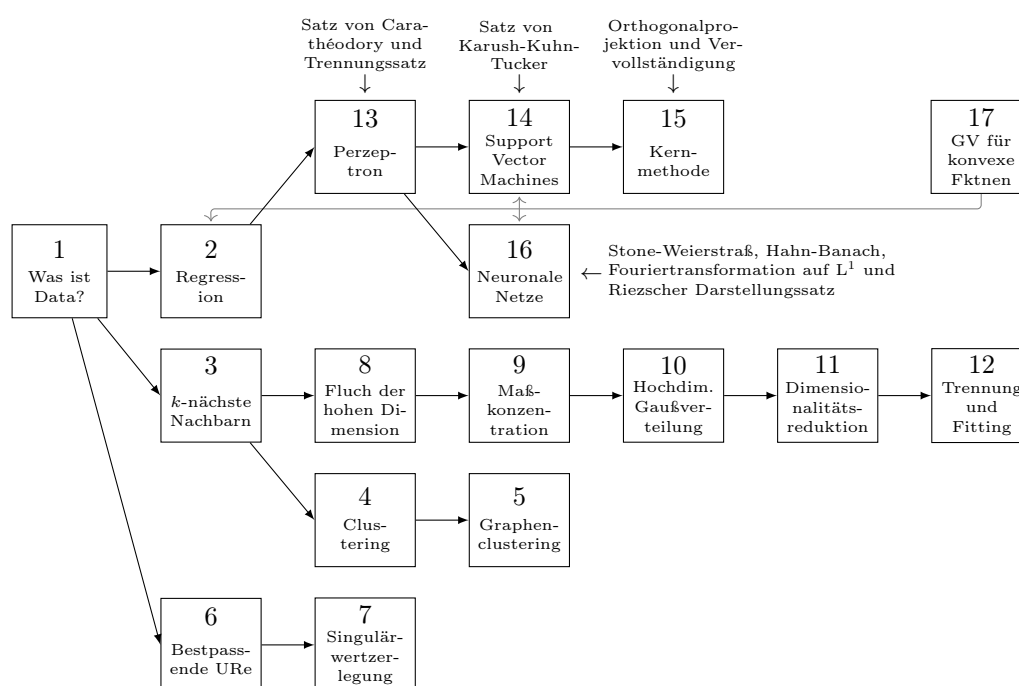
Das vorliegende Manuskript basiert in der Tat auf mehreren Vorlesungen, die der Autor in den letzten Jahren für Studierende der Mathematik, sowohl im Bachelor als auch für das gymnasiale Lehramt, gehalten hat. Es ist geeignet als Grundlage für eine 4+2 Vorlesung ab dem dritten Studienjahr und setzt die Inhalte der Grundvorlesungen in Analysis, Maßtheorie, Linearer Algebra und Wahrscheinlichkeitslehre voraus. Einige Kapitel erfordern darüber hinaus Kenntnisse in Optimierung und Funktionalanalysis. Benötigte Vorkenntnisse jenseits des Grundvorlesungen werden jeweils am Kapitelanfang ausgewiesen und können darüber hinaus im Diagramm auf der folgenden Seite abgelesen werden. Vorkenntnisse in Informatik oder Numerik sind natürlich hilfreich, aber nicht unbedingt erforderlich.

Wir folgen weitgehend dem in der mathematischen Literatur üblichen Satz-Beweis-Stil, ergänzt durch ausführliche Erläuterungen in Prosa. Dies wird komplementiert durch 121 unterrichtserprobte Aufgaben. Darunter sind sowohl theoretische Aufgaben wie auch Aufgaben bei denen implementiert werden muss. Zur besseren Lesbarkeit verwenden wir im folgenden das generische Maskulinum. Die in dieser Arbeit verwendeten Personenbezeichnungen beziehen sich aber stets auf alle Geschlechter.

Es folgt ein kurzer Abriss der im Text behandelten Themen. Wir beginnen mit einem einleitenden Kapitel 1, in welchem wir erläutern was wir in diesem Text unter *Daten* verstehen, und mit wie gearteten Methoden wir aus diesen Erkenntnisse welcher Art zu gewinnen suchen. Das erste richtige Kapitel 2 behandelt dann zunächst klassische Regressionsmethoden; der Leser wird aber hier schon viele Ideen kennenlernen, die später immer wieder auftauchen werden. In Kapitel 3 behandeln wir den sehr einfachen und anschaulichen  $k$ -NN Algorithmus, diskutieren mehrere Preprocessing-Methoden und wenden beides auf Beispiele aus den Bereichen Textmining und Produktbewertungen an. In den Kapiteln 4 und 5 behandeln wir Clusteringmethoden für Datenmengen in metrischen Räumen und dann auf Graphen.

Es folgen die Kapitel 6 und 7 in welchen bestpassende Unterräume, deren Zusammenhang mit der Singulärwertzerlegung von Matrizen, und als Anwendung davon die Hauptkomponentenanalyse, Dimensionalitätsreduktion und kollaboratives Filtern diskutiert werden. In den Kapiteln 8–12 wenden wir uns dann hochdimensionalen Datenmengen zu. Als erstes diskutieren wir die Eigenheiten der Gleichverteilung und der Gaußverteilung in hochdimensionalen Räumen und zeigen dann Strategien auf, mit denen Daten, die von mehreren unabhängigen Gaußverteilungen stammen, getrennt und deren Parameter geschätzt werden können. In Kapitel 13 behandeln wir den Perzeptronalgorithmus, gefolgt von Support Vector Machines in Kapitel 14 und der Kernmethode für SVMs in Kapitel 15. Zum Schluss kommen wir zu neuronalen Netzen, wobei wir in Kapitel 16.1 erst die Expressivität behandeln und dann in Kapitel 16.2 die Rückwärtspropagation, mit der die Parameter eines neuronalen Netzes an gegebene Daten angepasst werden können. Im finalen Kapitel 17 diskutieren wir das Gradientenverfahren im Kontext konvexer Funktionen.

Das folgende Diagramm zeigt die Abhängigkeiten zwischen den Kapiteln auf.



Wie oben angedeutet, gehen Resultate aus dem letzten Kapitel 17 in mehrere frühere Kapitel ein. Mancher Leser wird es daher bevorzugen, mit diesem letzten Kapitel zu beginnen. Andererseits verwenden wir in den Kapiteln 2 und 14 lediglich zwei im Verlauf von Kapitel 17 behandelte Resultate über die Existenz und Eindeutigkeit von Minimierern, welche man leicht on-the-spot nachlesen kann. In Kapitel 16 verweisen wir auf eine genaue Diskussion des Gradientenverfahrens unter gutartigen Bedingungen in Kapitel 17.

Ernster zu nehmen sind die folgenden, ebenfalls im Diagramm angegebenen, Voraussetzungen bei den Kapiteln 13–16, die wir an entsprechender Stelle jeweils ohne

Beweis notieren werden: In Kapitel 13 handelt es sich um Charathéodorys Charakterisierung der konvexen Hülle und den Trennungssatz für kompakte konvexe Mengen. Kapitel 14 setzt massiv die Theorie von Karush-Kuhn-Tucker zur Optimierung konvexer Funktionen unter Ungleichungsnebenbedingungen ein. Schließlich benötigt Kapitel 15 zwei Resultate aus der Hilbertraumtheorie und Kapitel 16 dann gleich mehrere tiefgehende Resultate aus der Funktionalanalysis, wobei diese aber nur in die zweite Hälfte des Unterkapitels 16.1 zur Expressivität eingehen.

Die ohne Beweis verwendeten Resultate in den Kapiteln 13 und 15 sind intuitiv zugänglich. Im Gegensatz dazu muss gesagt sein — ohne abschrecken zu wollen — dass insbesondere der Beweis des Satzes 14.11 zur Berechnung der SVM durch ein quadratisches Optimierungsproblem, sowie die Beweise der Expressivitätsresultate 16.20, 16.21, 16.22 und 16.25 für neuronale Netze am besten einer Optimierungs- bzw. Funktionalanalysisvorlesung nachgeschaltet sein sollten, wenn man diese vollständig verstehen will.

Ich bedanke mich herzlich bei allen Teilnehmern meiner Vorlesungen, sowie bei vielen meiner Kollegen für zahlreiche interessante und hilfreiche Diskussionen. Mein Dank gilt außerdem der gesamten Data Science und Machine Learning Community, deren Bücher, Lecture Notes, wissenschaftliche Artikel, Blogs, Videos und Beiträge in Foren mich für das Themengebiet begeistert haben, und von denen ich alles, was in diesem Manuskript behandelt wird, gelernt habe. Ganz besonderer Dank gilt dabei den Autoren der Bücher [BHK20, LRU12, Ver18, SSBD14] auf deren Vorarbeit das vorliegende Manuskript in großem Maße aufbaut. Jedes Kapitel enthält am Ende einen kurzen Abschnitt mit genaueren Referenzen.

Hamburg, im Oktober 2023

Sven-Ake Wegner

# Inhaltsverzeichnis

Vorwort	2
1 Was ist Data (Science)?	6
2 Affin-lineare, polynomiale und logistische Regression	12
3 Der $k$ -NN Algorithmus	38
4 Clustering	53
5 Graphenclustering	62
6 Bestpassende Unterräume	79
7 Singulärwertzerlegung	86
8 Fluch und Segen der hohen Dimension	110
9 Maßkonzentration	120
10 Gaußsche Zufallsvektoren in hohen Dimensionen	130
11 Dimensionalitätsreduktion à la Johnson-Lindenstrauss	140
12 Trennung von Gaußianen und Parameteranpassung	146
13 Das Perzeptron	164
14 Support Vector Machines	173
15 Die Kernmethode	192
16 Neuronale Netze	207
17 Gradientenverfahren für konvexe Funktionen	240
A Ausgewählte Resultate der Wahrscheinlichkeitstheorie	258
Literaturverzeichnis	265

# Kapitel 1

## Was ist Data (Science)?

Wir erläutern zuerst, was wir in diesem Text unter dem Begriff *Daten* verstehen werden, und sprechen gleichzeitig die Warnung aus, dass sich die genaue Definition von Kapitel zu Kapitel mitunter leicht ändert.

**Definition 1.1.** Seien  $X$  und  $Y$  Mengen.

- (i) Eine endliche Teilmenge  $D \subseteq X$  heißt *ungelabelte Datenmenge* in  $X$ , ihre Elemente heißen *Datenpunkte*.
- (ii) Eine endliche Menge  $D \subseteq X \times Y$  heißt *gelabelte Datenmenge*. Ist  $(x, y) \in D$ , so heißt  $x$  der *Featureteil* des Datenpunktes  $(x, y)$  und  $y$  sein *Label*.
- (iii) Eine gelabelte Datenmenge  $D \subseteq X \times Y$  heißt *kategorisch gelabelt*, wenn  $Y$  endlich ist und *kontinuierlich gelabelt*, falls  $Y$  ein Kontinuum ist.<sup>1</sup>
- (iv) Ist  $X = X_1 \times \cdots \times X_d$  und  $x = (x_1, \dots, x_d)$ , so nennen wir die  $x_i$ 's die *Features* von  $x$ . Ist  $Y = Y_1 \times \cdots \times Y_m$ , so sprechen wir von *mehrdimensionalen Labels*.

Die Menge  $X$  in der obigen Definition 1.1 ist häufig gleich dem Raum  $\mathbb{R}^d$  und dann stets mit der euklidischen Norm und dem Standardskalarprodukt ausgestattet. Ist  $X \subseteq \mathbb{R}^d$ , so können wir  $X$  immerhin mit der euklidischen Metrik versehen. In manchen Kapiteln ist  $(X, \rho)$  aber auch ein abstrakter metrischer Raum und in einigen Fällen braucht  $\rho$  nicht mal eine Metrik zu sein, sondern nur ein sogenanntes Abstandsmaß. Die Menge  $Y$  kann im kategoriellen Fall ohne Einschränkung als gleich  $\{1, \dots, m\}$  angenommen werden.

**Bemerkung 1.2.** In manchen Situationen ist es geboten, Dopplungen von Datenpunkten zu erlauben. Dies kann man erreichen, indem man  $X$  durch  $X \times \mathbb{N}$  ersetzt. In einer Menge  $D \subseteq X \times \mathbb{N}$  kann dann  $x \in X$  z.B. in der Form von  $(x, 1)$  und  $(x, 2)$  zweimal enthalten sein. Um die Dinge nicht unnötig kompliziert zu machen, schreiben wir im folgenden  $x_i$  oder  $x^{(i)} \in D$  statt  $(x, i) \in D \times \mathbb{N}$ . Sprechen wir von einer Datenmenge  $D = \{x_1, \dots, x_n\} \subseteq X$  so sei immer stillschweigend vorausgesetzt, dass derselbe Punkt, mit unterschiedlichem Index, mehrfach vorkommen kann.

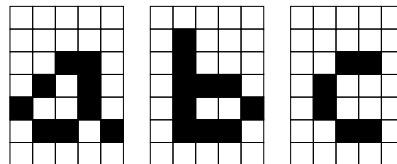
<sup>1</sup>Man kann sich hier abgeschlossene, offene, halboffene, beschränkte oder unbeschränkte Intervalle in  $\mathbb{R}$  vorstellen, die mindestens zwei Punkte haben, aber wir wollen eventuell auch mal etwas anderes zulassen.

In konkreten Anwendungen ist es eine Frage der Modellierung, ob man eine oder mehrere Koordinaten eines Datenvektors  $(x_1, \dots, x_d)$  als Label auszeichnet und wenn ja, welche dies sein sollen. Wir betrachten die folgenden Beispiele.

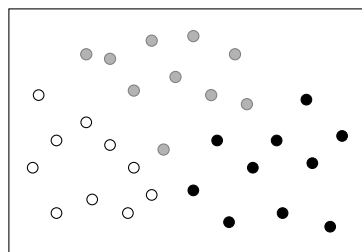
**Beispiel 1.3.** (i) Gegeben seien 10 Studenten, die eine Klausur zur Vorlesung ‘Mathematische Einführung in Data Science’ schreiben. Für jeden Studenten erfassen wir in der Woche vor der Klausur die Vorbereitungszeit auf die Klausur in Stunden, die auf sozialen Medien verbrachte Zeit, ebenfalls in Stunden, und schließlich das Klausurergebnis in Prozent. Wir können die folgende Tabelle als eine ungelabelte Datenmenge  $D \subseteq \mathbb{R}^3$  auffassen, oder z.B. als gelabelte Datenmenge  $D \subseteq \mathbb{R}^2 \times [0, 100]$ , also die Vorbereitungszeit und Zeit auf sozialen Medien als Features und das Klausurergebnis als kontinuierliches Label.

Student	Vorbereitung in h	Soziale Medien in h	Klausurergebnis in %
1	0.0	20.0	0.0
2	1.5	8.5	2.0
3	2.0	6.0	7.0
4	2.0	6.0	10.5
5	8.0	10.0	29.5
6	8.5	3.0	49.0
7	9.5	0.0	59.5
8	12.0	2.0	63.5
9	18.0	4.0	85.0
10	19.0	0.5	98.0

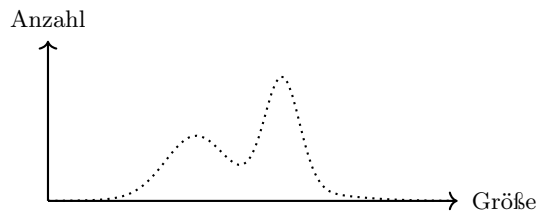
(ii) Wir betrachten handgeschriebene Buchstaben, die entsprechend der folgenden Abbildung als  $(7 \times 5)$ -Matrizen aus Einsen und Nullen geschrieben werden können. Ist uns dann jeweils noch bekannt, welcher Buchstabe hier geschrieben wurde, so erhalten wir eine gelabelte Datenmenge  $D \subseteq \mathbb{R}^{7 \times 5} \times \{a, b, c, \dots\}$ .



(iii) Sei  $X \subseteq \mathbb{R}^2$  ein Quader und  $Y = \{1, 2, 3\}$ . Das folgende Bild stellt eine kategoriell gelabelte Datenmenge dar, bei welcher jeder Datenpunkt zwei Features hat. Im Bild entspricht Label 1 einem weißen, Label 2 einem grauen, und Label 3 einem schwarzen Punkt.



(iv) Wir betrachten eine Vorlesung ‘Mathematische Einführung in Data Science’ und notieren für jeden Teilnehmer dessen Körpergröße. Auf diese Weise erhalten wir eine ungelabelte Datenmenge  $D \subseteq (0, \infty)$ , beachte hierbei Bemerkung 1.2. Das Bild zeigt die Verteilung der Körpergrößen.



Ist eine gelabelte Datenmenge  $D \subseteq X \times Y$  gegeben, so besteht eine zentrale Aufgabe darin, der Datenmenge eine Funktion  $f: X \rightarrow Y$  zuzuordnen. Hierbei kann man drei Sichtweisen einnehmen, die allerdings nicht scharf voneinander getrennt sind, sondern ineinander übergehen:

1. Es gibt eine uns unbekannte ‘echte’ Funktion  $f_0: X \rightarrow Y$  und die Datenmenge ist von der Form  $D = \{(x_1, f_0(x_1)), \dots, (x_n, f_0(x_n))\}$  mit  $x_i \in X$ . Wir wollen  $f_0$  durch  $f$  approximieren, also  $f(x) \approx f_0(x)$  für alle  $x \in X$  erreichen. Man nennt dies die *Approximationsperspektive*.

2. Die Datenmenge entsteht durch die zufällige Störung einer Funktion  $f_0: X \rightarrow Y$ , z.B. indem  $y_i = f_0(x_i) + \varepsilon_i$  gilt wobei  $\varepsilon_i \in \mathbb{R}$  Realisierungen einer normalverteilten Zufallsvariable sind. Wir suchen dasjenige  $f$ , von dem die Daten ‘am wahrscheinlichsten’ stammen. Man nennt dies die *Wahrscheinlichkeitsperspektive*.

3. Wir gehen nicht davon aus, dass es eine echte Funktion gibt, sondern passen durch die Minimierung oder Maximierung einer Zielfunktion  $\phi = \phi(f, D)$   $f$  an die Daten an. Man nennt dies die *Optimierungsperspektive*.

Die Funktion  $f$  nennen wir, je nach Kontext und Anwendung, *Regressor*, *Klassifizierer*, *Prediktor* oder auch *Approximant*. Die zur Bestimmung von  $f$  benutzte Datenmenge nennen wir dann die *Trainingsdaten* und den Prozess der Bestimmung von  $f$  bezeichnen wir als *überwachtes Lernen*.

Führt man in Beispiel 1.3(i) z.B. eine affin-lineare Regression durch, so kann man sich auf den Standpunkt stellen, dass tatsächlich ein ‘kausaler Zusammenhang’ der Form

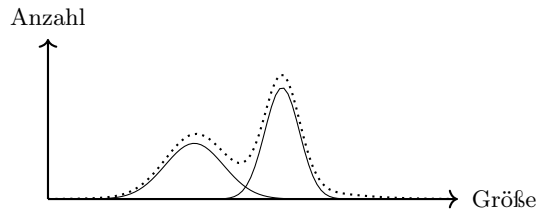
$$\text{Klausurergebnis} = a_1 \cdot \text{Vorbereitungszeit} + a_2 \cdot \text{Fernsehzeit} + b$$

besteht und dass man die reellen Konstanten  $a_1$ ,  $a_2$  und  $b$  mithilfe der Daten bestimmt. In Beispiel 1.3(ii) fällt es eher schwer mit einer Idee für eine zugrundeliegende Funktion aufzuwarten; stattdessen kann man hier eine geeignet definierte Fehlerfunktion über eine geeignete Funktionenklasse minimieren und dann darauf bauen, dass ein neu von Hand geschriebener Buchstabe richtig erkannt wird. Stammen schließlich kontinuierliche Label z.B. von einer physikalischen Messung, so ist die Vorstellung naheliegend, dass ein normalverteilter Messfehler vorliegt.

Liegen ungelabelte Datenmengen  $D \subseteq X$  vor, wie etwa in Beispiel 1.3(iv), so ist es auch hier eine natürliche Aufgabe, den Daten eine Funktion zuzuordnen. Im vorgenannten Beispiel liegt es z.B. nahe zu vermuten, dass die skizzierte Verteilung dadurch zustande kommt, dass zwei Normalverteilungen (Körpergröße der männli-



chen und Körpergröße der weiblichen Teilnehmer) superponiert werden:

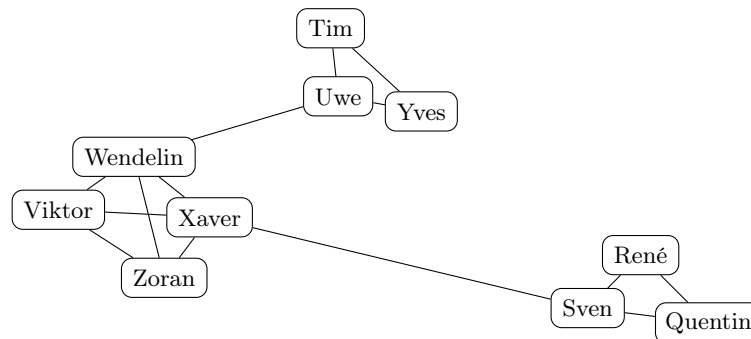


Gesucht ist dann eine Funktion  $f: D \rightarrow \{m, w\}$  die einer Körpergröße ein Geschlecht zuordnet, also die zwei ‘Cluster’ trennt. Ist dies geschehen, so besteht eine weitere Aufgabe darin, die Parameter der zwei einzelnen Verteilungen zu schätzen.

Prozesse, bei denen ungelabelte Datenmengen gegeben sind, und Muster innerhalb der Daten erkannt, oder Vorhersagen anhand der ungelabelten Daten gemacht werden, bezeichnen wir als *unüberwachtes Lernen*.

Wir geben weitere Beispiele.

**Beispiel 1.4.** (i) Ein soziales Netzwerk wird illustriert durch das folgende Bild, bei dem eine Verbindungslinie zwischen den Nutzern angibt, dass diese befreundet sind. Eine natürliche Aufgabe besteht darin die ‘Cluster’ zu finden. Im Bild unten sind letztere einfach erkennbar, aber bei einem größeren sozialen Netzwerk natürlich nicht.



(ii) Die nachfolgende Tabelle (entnommen aus [LRU12]) enthält Bewertungen von 0–5 für sieben Filme durch fünf Bewerberinnen. Natürliche Aufgaben sind die Erkennung von Mustern, wie z.B. welche Bewerberinnen ähnlichen Filmgeschmack haben, und die Vorhersage von Bewertungen, wenn neue Bewerberinnen oder neue Filme ins Spiel kommen.

	Alien	Casablanca	Star Wars	Titanic	Matrix
Abbey	0	2	0	2	1
Bailey	1	0	1	0	1
Caitlin	5	0	5	0	5
Daisy	0	4	0	4	2
Edith	3	0	3	0	3
Fae	0	5	0	5	0
Gail	4	0	4	0	4

(iii) Wir betrachten das folgende Graustufenbild welches durch eine  $(320 \times 240)$ -Matrix aus Zahlen zwischen Null und Eins dargestellt werden kann. Eine natürliche

Aufgabe ist hier die Datenkompression, bei der wir eine Methode suchen, mit der das Bild aus weniger als  $(320 \cdot 240)$ -vielen reellen Zahlen rekonstruiert bzw. approximiert werden kann.



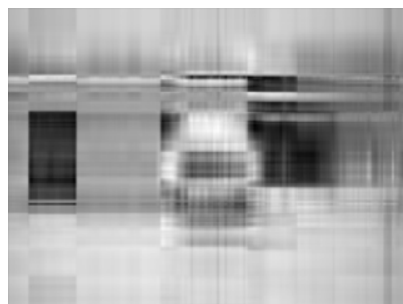
Wir notieren, dass a priori keines der Beispiele 1.4(i)–(iii) von unserer anfänglichen Definition 1.1 einer Datenmenge erfasst wird. Im Fall des sozialen Netzwerkes können wir allerdings die Nutzer mit  $1, 2, 3, \dots, d$  durchnummerieren und dann den  $i$ -ten Nutzer durch einen Vektor  $a_i := (a_{i1}, \dots, a_{id}) \in \mathbb{R}^d$  darstellen, wobei  $a_{ij}$  gleich eins ist, falls der  $i$ -te mit dem  $j$ -ten Nutzer befreundet ist. Zweckmäßiger ist es allerdings, statt der Datenmenge  $D = \{a_1, \dots, a_d\}$ , die in diesem Fall symmetrische *Datenmatrix*

$$A = \begin{bmatrix} a_{11} & \cdots & a_{1d} \\ \vdots & & \vdots \\ a_{d1} & \cdots & a_{dd} \end{bmatrix}$$

zu verwenden.

Bei der Bewertungstabelle in Beispiel 1.4(ii) erhalten wir eine Matrix frei Haus, könnten aber auch hier die Zeilen oder Spalten als Datenpunkte modellieren. Beachte hierbei, dass dies auf zwei verschiedene Datenmengen führt: Entweder auf eine Menge von Bewerterinnen, dargestellt durch Punkte in  $\mathbb{R}^5$ , oder auf eine Menge von Filmen, die dann durch Punkte in  $\mathbb{R}^7$  gegeben sind.

Im Fall des Graustufenbilds in Beispiel 1.4(iii) kommen die Daten ebenfalls in Matrixform, nämlich als Element von  $\mathbb{R}^{320 \times 240}$ , und es scheint eher unnatürlich, z.B. die Zeilen oder Spalten als eine Menge von einzelnen Datenpunkten zu interpretieren. Behandelt man die Matrix als ganzes, so kann die angesprochene Kompression z.B. per Approximation durch eine Matrix niedrigeren Ranges erreicht werden. Die folgenden Bilder zeigen zwei solche Approximationen.



rk = 5



rk = 15

Fasst man jetzt doch die Zeilen des Bildes in Beispiel 1.4(iii) als Datenpunkte in  $\mathbb{R}^{320}$  auf, so entsprechen die letzteren Bilder der Projektion dieser 240 Datenpunkte auf einen 5- bzw. 15-dimensionalen Unterraum des  $\mathbb{R}^{320}$  — wobei dieser Raum natürlich in einer geschickten Weise auszuwählen ist.

Anwendungen, die in diesem Sinne die *Dimensionalität* einer Datenmenge reduzieren sind insbesondere deswegen von hoher Bedeutung, als dass viele Datenmengen in natürlicher Weise in einem Raum mit sehr hoher Dimension leben. In manchen Fällen ist hierbei die Dimension des Raumes sogar deutlich größer als die Anzahl der gegebenen Datenpunkte: Betrachte z.B. wieder die Teilnehmer einer Vorlesung ‘Mathematische Einführung in Data Science’, dargestellt durch alle ihre Social Media Posts, d.h. jeder Teilnehmer wird durch einen Featurevektor gegeben, der alle Texte, Bilder und Videos enthält, die von den Teilnehmern jemals gepostet wurden.

Wir schließen dieses Kapitel mit der folgenden Bemerkung.

**Bemerkung 1.5.** Ordnet man einer gelabelten oder ungelabelten Datenmenge  $D$  eine Funktion  $f: X \rightarrow Y$  via eines Algorithmus zu, so wird dies oft als *maschinelles Lernen* bezeichnet und man sagt, dass die Funktion anhand der Daten *gelernt* wird. Wir raten dem Leser bei der Benutzung dieser Bezeichnungen zur Vorsicht. In der Tat haben wir im Kontext von Beispiel 1.3(ii) angedeutet, dass das ‘Erlernen’ einer Handschrift der numerischen Minimierung einer Zielfunktion entspricht, und wir werden noch sehen, dass das ‘Erlernen’ des Filmgeschmacks von Abbey, Bailey, Caitlin, Daisy, Edith, Fae und Gail in Beispiel 1.4(ii) durch die Bestimmung von Eigenwerten erreicht werden wird.

## Kapitel 2

# Affin-lineare, polynomiale und logistische Regression

Wir beginnen nun mit klassischen Regressionsmethoden. Dies mag auf den ersten Blick unspektakulär wirken, in der Tat kommen hier aber bereits viele Ideen, Konzepte und technische Tricks vor, auf die wir im späteren Verlauf zum Teil noch mehrfach zurückkommen werden. Beispielsweise wird uns die Minimierung einer konvexen Kostenfunktion im Kontext der Support Vector Machines wieder begegnen und die Maximum-Likelihood-Methode werden wir z.B. im Kapitel zu hochdimensionalen Räumen einsetzen. Im Unterkapitel über polynomiale Regression werden wir, wie später auch bei der Kernmethode, ein nichtlineares Problem auf ein lineares zurückführen. Schließlich steht das Perzeptron, und damit der zentrale Baustein neuronaler Netze, in enger Verbindung zur Methode der logistischen Regression, welche wir am Ende des aktuellen Kapitels behandeln.

Wir werden im Folgenden immer wieder auf konkrete Zusammenhänge mit anderen Kapiteln hinweisen und entsprechend referenzieren. An einigen Stellen werden wir auch Vorwärtsreferenzen auf spätere Kapitel machen, benötigen hier aber keine sehr tiefliegenden Resultate, sondern nur einige Aussagen über konvexe Funktionen sowie Fakten zur Normalverteilung.

Im Verlauf dieses Kapitels verweisen wir an einigen Stellen auf Kapitel 17, benutzen von dort aber nur einfache Resultate über Existenz und Eindeutigkeit von Minimierern konvexer Funktionen.

### 2.1 Affin-lineare Regression in einer Dimension

Gegeben sei eine Datenmenge

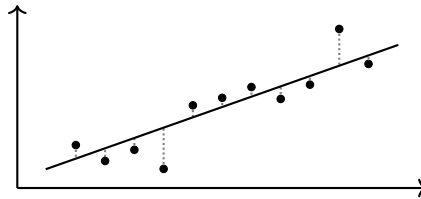
$$D := \{(x_i, y_i) \in \mathbb{R}^2 \mid i = 1, \dots, n\}, \quad (2.1)$$

wobei jeder Datenpunkt aus einem reellen Feature  $x_i$  und einem reellen Label  $y_i$  besteht. Bei der linearen Regression ist es unser Ziel, eine affin-lineare Funktion

$f: \mathbb{R} \rightarrow \mathbb{R}$  zu finden, sodass die Summe

$$\sum_{i=1}^n (f(x_i) - y_i)^2$$

der quadratischen Abstände von Funktionswert und Label über alle Features minimal ist. Anschaulich bedeutet dies, dass im folgenden Bild diejenige Gerade gesucht wird, für die die Summe der Quadrate der Längen der gepunkteten Strecken den kleinstmöglichen Wert annimmt.



A priori ist natürlich nicht klar, ob eine solche Gerade, bzw. affin-lineare Funktion, überhaupt existiert und, wenn dem so ist, ob sie eindeutig bestimmt ist. Schließlich stellt sich die Frage, ob  $f$  explizit aus den Daten berechnet werden kann. Alle drei Punkte beantwortet der folgende Satz.

**Satz 2.1.** Sei  $D := \{(x_i, y_i) \mid i = 1, \dots, n\} \subseteq \mathbb{R}^2$  eine Datenmenge, bei der nicht alle  $x_i$  gleich sind. Dann gibt es genau eine affin-lineare Funktion

$$f^* = \operatorname{argmin} \left\{ \sum_{i=1}^n (f(x_i) - y_i)^2 \mid f: \mathbb{R} \rightarrow \mathbb{R} \text{ affin-linear} \right\}$$

und es gilt  $f^*(x) = a^*x + b^*$  für  $x \in \mathbb{R}$  mit

$$a^* = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\left(\sum_{i=1}^n x_i y_i\right) - n\bar{x}\bar{y}}{\left(\sum_{i=1}^n x_i^2\right) - n\bar{x}^2} \quad \text{und} \quad b^* = \bar{y} - a^*\bar{x}.$$

Hierbei sind  $\bar{x} = (x_1 + \dots + x_n)/n$  und  $\bar{y} = (y_1 + \dots + y_n)/n$  die Mittelwerte der Features und der Label.

*Beweis.* ① Wir zeigen zunächst, dass die zwei Ausdrücke für  $a^*$  übereinstimmen. Unter Benutzung von  $x_1 + \dots + x_n = n\bar{x}$ , und Entsprechendem für  $\bar{y}$ , erhalten wir für die Zähler

$$\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \sum_{i=1}^n x_i y_i - \bar{x} \sum_{i=1}^n x_i - \bar{x} \sum_{i=1}^n y_i + n\bar{x}\bar{y} = \left(\sum_{i=1}^n x_i y_i\right) - n\bar{x}\bar{y}.$$

Die Gleichheit der Nenner folgt als Spezialfall  $y_i = x_i$ . Insbesondere sind beide Nenner ungleich Null, da per Voraussetzung nicht alle  $x_i$  gleich sind. Gezeigt werden

muss als nächstes, dass die Funktion

$$\phi: \{f: \mathbb{R} \rightarrow \mathbb{R} \mid f \text{ affin-linear}\} \rightarrow \mathbb{R}, \quad \phi(f) := \sum_{i=1}^n (f(x_i) - y_i)^2$$

einen Minimierer besitzt. Da affin-lineare Funktionen eindeutig durch zwei reelle Parameter  $a$  und  $b$  beschrieben sind, können wir die oben angegebene Funktion  $\phi$  als Funktion dieser Parameter lesen, also

$$\phi: \mathbb{R}^2 \rightarrow \mathbb{R}, \quad \phi(a, b) = \sum_{i=1}^n (ax_i + b - y_i)^2$$

minimieren. Per Konstruktion ist  $\phi \geq 0$  und eine stetige Funktion. Um die Existenz eines Minimierers zu garantieren, genügt es zu zeigen, dass der Limes Inferior von  $\phi(a, b)$  für  $\|(a, b)\| \rightarrow \infty$  unendlich ist. Angenommen, letzteres gilt nicht, so gibt es eine Folge  $(a^{(k)}, b^{(k)})_{k \in \mathbb{N}}$  mit

$$\|(a^{(k)}, b^{(k)})\| \rightarrow \infty \quad \text{und} \quad \phi(a^{(k)}, b^{(k)}) = \sum_{i=1}^n (a^{(k)}x_i + b^{(k)} - y_i)^2 \rightarrow c \in \mathbb{R}.$$

Wir wählen  $i$  und  $j$  derart, dass  $x_i \neq x_j$  gilt. Da in der Summe rechts nur nichtnegative Terme addiert werden, müssen alle Summanden beschränkt sein. Daraus folgt, dass sowohl  $(a^{(k)}x_i + b^{(k)} - y_i)_{k \in \mathbb{N}}$  als auch  $(a^{(k)}x_j + b^{(k)} - y_j)_{k \in \mathbb{N}}$  beschränkt sind, und damit ist auch die Differenz  $(a^{(k)}(x_i - x_j) - y_i + y_j)_{k \in \mathbb{N}}$  beschränkt. Da  $x_i - x_j \neq 0$  ist, bedeutet letzteres, dass  $(a^{(k)})_{k \in \mathbb{N}}$  beschränkt ist und folglich muss  $|b^{(k)}| \rightarrow \infty$  gelten. Aus beidem zusammen folgt  $\phi(a^{(k)}, b^{(k)}) \rightarrow \infty$  im Widerspruch zur Annahme. Der Limes Inferior von  $\phi(a, b)$  für  $\|(a, b)\| \rightarrow \infty$  ist also gleich  $\infty$  und wir erhalten die Existenz eines Minimierers.

② Als nächstes betrachten wir

$$\nabla \phi(a, b) = \begin{bmatrix} \frac{\partial \phi}{\partial a} \\ \frac{\partial \phi}{\partial b} \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^n 2(ax_i + b - y_i)x_i \\ \sum_{i=1}^n 2(ax_i + b - y_i)1 \end{bmatrix} \stackrel{!}{=} 0.$$

Die zweite Gleichung ist äquivalent zu

$$0 = \sum_{i=1}^n (ax_i + b - y_i) = a \sum_{i=1}^n x_i + nb - \sum_{i=1}^n y_i = an\bar{x} + n\bar{b} - n\bar{y},$$

woraus sich  $b = \bar{y} - a\bar{x}$  ergibt. Einsetzen in die erste Gleichung liefert

$$\begin{aligned} 0 &= \sum_{i=1}^n (ax_i + \bar{y} - a\bar{x} - y_i)x_i = a \sum_{i=1}^n x_i^2 + \bar{y} \sum_{i=1}^n x_i - a\bar{x} \sum_{i=1}^n x_i - \sum_{i=1}^n y_i x_i \\ &= a \sum_{i=1}^n x_i^2 + \bar{y}n\bar{x} - a\bar{x}n\bar{x} - \sum_{i=1}^n x_i y_i \end{aligned}$$

$$= a \left( \sum_{i=1}^n x_i^2 - n\bar{x}^2 \right) - \left( \sum_{i=1}^n x_i y_i - n\bar{x}\bar{y} \right),$$

was sich nach  $a$  auflösen lässt. Die Ableitung der überall differenzierbaren Funktion  $\phi$  verschwindet also in genau einem Punkt, und zwar in  $(a^*, b^*)$ , wobei  $a^*$  und  $b^*$  durch die im Satz vermerkten Formeln gegeben sind. Zusammen mit dem ersten Teil des Beweises folgt, dass dies der eindeutig bestimmte Minimierer von  $\phi$  sein muss.  $\square$

**Definition 2.2.** Sei  $D$  eine Datenmenge wie in Satz 2.1. Dann heißt die Funktion  $f^*: \mathbb{R} \rightarrow \mathbb{R}$ ,  $f^*(x) = a^*x + b^*$  aus Satz 2.1 *affin-linearer Regressor* für  $D$  und ihr Graph die *Regressionsgerade* zu  $D$ .

**Bemerkung 2.3.** (i) Unter den Voraussetzungen von Satz 2.1 gilt stets  $f^*(\bar{x}) = \bar{y}$ , d.h. die Regressionsgerade geht immer durch den Punkt  $(\bar{x}, \bar{y})$ .

(ii) Die im Beweis verwendete Funktion  $\phi$  nennt man *Fehlerfunktion*, *Kostenfunktion* oder *Zielfunktion*. Das darin untergebrachte Quadrat hat zur Folge, dass  $\phi$  differenzierbar ist. Verwendet man stattdessen z.B.

$$\psi(f) := \sum_{i=1}^n |f(x_i) - y_i|$$

so kann man zwar immer noch zeigen, dass ein Minimierer existiert, aber dieser ist im Allgemeinen nicht mehr eindeutig, siehe Aufgabe 2.2, und kann vom affin-linearen Regressor verschieden sein.

Der obige Ansatz zur Definition eines Regressors ist ein Beispiel für die in Kapitel 1 erläuterte Optimierungssichtweise, bei der eine von uns ausgesuchte Fehlerfunktion minimiert wird über eine von uns ausgesuchte Klasse von Funktionen.

Verwendet man die Fehlerfunktion  $\phi$  wie in Satz 2.1, so spricht man von der *Methode der kleinsten Quadrate*. Neben dem genannten Effekt, dass die Quadrate  $\phi$  differenzierbar machen, mag ihre Einführung auf den ersten Blick etwas willkürlich wirken. Denn während  $\psi$  in Bemerkung 2.3(ii) alle Abweichungen gleich behandelt, ‘bestraft’  $\phi$  Abweichungen echt größer eins härter als Abweichungen echt kleiner als eins. Auf den zweiten Blick wird sich jetzt aber zeigen, dass die Quadrate sehr natürlich sind.

Wir nehmen dazu die auch in Kapitel 1 bereits erwähnte Approximationsperspektive ein, d.h. wir nehmen an, dass es tatsächlich eine affin-lineare Funktion  $f: \mathbb{R} \rightarrow \mathbb{R}$ ,  $f(x) = ax + b$ , mit uns allerdings unbekannten Parametern  $a, b \in \mathbb{R}$  gibt. Die Daten  $D$  entstehen so, dass für ein Feature  $x_i$  das Label durch  $y_i := f(x_i) + \varepsilon_i$  gegeben ist, wobei die  $\varepsilon_i$  voneinander unabhängige normalverteilte Störungen sind. Basierend auf einer so erzeugten Datenmenge  $D$  wollen wir nun eine affin-lineare Funktion  $f^*: \mathbb{R} \rightarrow \mathbb{R}$ ,  $f^*(x) = a^*x + b^*$  derart bestimmen, dass die Wahrscheinlichkeit dafür, dass die Daten der Stichprobe von  $f^*$  stammen, maximal ist im Vergleich zu allen anderen affin-linearen Funktionen, von denen die Daten auch hätten stammen

können. Um dies präziser zu fassen, betrachten wir die *Likelihood-Funktion*

$$L: \{f: \mathbb{R} \rightarrow \mathbb{R} \mid f \text{ affin-linear}\} \rightarrow \mathbb{R}, \quad L(f) := \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(f(x_i)-y_i)^2}{2\sigma^2}} \quad (2.2)$$

für eine Datenmenge  $D$  und  $\sigma > 0$ . Ist dann  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum,  $\mathcal{E}_1, \dots, \mathcal{E}_n: \Omega \rightarrow \mathbb{R}$  unabhängige Zufallsvariablen mit  $\mathcal{E}_i \sim \mathcal{N}(0, \sigma^2)$ , und  $Y_i(f) := f(x_i) + \mathcal{E}_i$ , bei gegebenem affin-linearen  $f$ , so haben wir  $Y_i(f) \sim \mathcal{N}(f(x_i), \sigma^2)$  nach A.19 und folglich für feste Realisierungen  $y_i$  der  $Y_i(f)$  und für  $\delta > 0$

$$P[|Y_i(f) - y_i| < \delta] = \int_{x_i-\delta}^{x_i+\delta} \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(f(t)-y_i)^2}{2\sigma^2}} dt.$$

Per Unabhängigkeit ergibt sich

$$P[|Y_i(f) - y_i| < \delta \text{ für alle } i] = \int_{Q_\delta(x)} \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(f(t_i)-y_i)^2}{2\sigma^2}} d(t_1, \dots, t_n)$$

wenn wir mit  $Q_\delta(x) = [x_1 - \delta, x_1 + \delta] \times \dots \times [x_n - \delta, x_n + \delta]$  den Würfel um  $x = (x_1, \dots, x_n) \in \mathbb{R}^n$  mit Seitenlänge  $2\delta$  abkürzen. Wir stellen uns nun vor, dass  $\delta > 0$  eine sehr kleine, aber feste Konstante ist. Dann besagt das obige, dass die Wahrscheinlichkeit dafür, dass  $f(x_i) + \mathcal{E}_i \approx y_i$  für alle  $i$  gilt, groß ausfällt, wenn  $L(f)$  groß ist. Aus diesem Grund nennen wir einen Maximierer  $f^*$  von  $L$  einen *Maximum-Likelihood-Schätzer*. Der folgende Satz zeigt, dass ein solcher existiert und stellt den Zusammenhang zur Methode der kleinsten Quadrate her.

**Satz 2.4.** *Sei  $D := \{(x_i, y_i) \mid i = 1, \dots, n\} \subseteq \mathbb{R}^2$  eine Datenmenge, bei der nicht alle  $x_i$  gleich sind. Dann gibt es für jedes  $\sigma > 0$  genau einen Maximierer*

$$f^* = \operatorname{argmax}\{L(f) \mid f \text{ affin-linear}\}$$

*der Likelihood-Funktion aus (2.2) und  $f^*: \mathbb{R} \rightarrow \mathbb{R}$ ,  $f^*(x) = a^*x + b^*$  stimmt mit dem affin-linearen Regressor aus Satz 2.1 überein.*

*Beweis.* Anstelle der Likelihood-Funktion  $L$  maximieren wir die sogenannte *Log-Likelihood-Funktion*  $\ell := \log \circ L$ , für die sich ergibt

$$\begin{aligned} \ell(f) &= \log \left( \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(f(x_i)-y_i)^2}{2\sigma^2}} \right) \\ &= n \log \left( \frac{1}{\sqrt{2\pi}\sigma} \right) - \frac{1}{2\sigma^2} \sum_{i=1}^n (f(x_i) - y_i)^2 \\ &= c_1 + c_2 \phi(f) \end{aligned} \quad (2.3)$$

wobei  $c_1, c_2 \in \mathbb{R}$  von  $f$  unabhängige Konstanten sind und  $\phi$  die Kostenfunktion aus Satz 2.1 ist. Da  $c_2 < 0$  ist, folgen sofort alle Behauptungen aus Satz 2.1.  $\square$



**Definition 2.5.** Wir bezeichnen mit

$$\begin{aligned}\text{cov}: \mathbb{R}^n \times \mathbb{R}^n &\rightarrow \mathbb{R}, \quad \text{cov}(x, y) := \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \\ \text{var}: \mathbb{R}^n &\rightarrow \mathbb{R}, \quad \text{var}(x) := \text{cov}(x, x)\end{aligned}$$

*Kovarianz und Varianz.* Der Beweis von Satz 2.1 hat gezeigt, dass  $\text{cov}(x, y) = \overline{xy} - \bar{x}\bar{y}$  und  $\text{var}(x) = \overline{x^2} - \bar{x}^2$  gelten. Wird obiges in den Datenpunkten der Stichprobe einer Zufallsvariable ausgewertet, so spricht man auch von der *Stichprobenkovarianz* oder *Stichprobenvarianz*.

Mithilfe von Varianz und Kovarianz können die Parameter, durch welche die affin-lineare Funktion  $f^*$  in den Sätzen 2.1 und 2.4 gegeben ist, wie folgt geschrieben werden

$$a^* = \frac{\text{cov}(x, y)}{\text{var}(x)} \quad \text{und} \quad b^* = \bar{y} - a^* \bar{x}.$$

Wir weisen darauf hin, dass in der Definition der (Ko-)varianz häufig der Vorfaktor  $\frac{1}{n-1}$  statt  $\frac{1}{n}$  verwendet wird. Bei obiger Formel für  $a^*$  macht das keinen Unterschied und bei großem  $n$  spielt es auch für Varianz und Kovarianz einzeln keine Rolle. Es gibt aber gute Gründe den Faktor  $\frac{1}{n-1}$  zu verwenden, vergleiche Aufgabe 2.14.

Die Annahme, dass hinter den Daten eine unbekannte affin-lineare Funktion mit einer normalverteilten Störung steht, führt zu der Frage, ob für große Stichproben die geschätzten Parameter  $a^*$  und  $b^*$  nahe bei den Parametern  $a$  und  $b$  der ‘echten’ Funktion liegen. Der folgende Satz beantwortet diese Frage asymptotisch.

**Satz 2.6.** Sei  $(x_i)_{i \in \mathbb{N}} \subseteq \mathbb{R}$  eine beschränkte Folge reeller Zahlen, sei  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum, seien  $\mathcal{E}_i: \Omega \rightarrow \mathbb{R}$  unabhängige Zufallsvariablen mit  $\mathcal{E}_i \sim N(0, \sigma^2)$  für  $i \in \mathbb{N}$  und sei  $Y_i := ax_i + b + \mathcal{E}_i$ , wobei  $a, b \in \mathbb{R}$  und  $\sigma > 0$ . Für  $n \in \mathbb{N}$  setzen wir voraus, dass  $x^{(n)} := (x_1, \dots, x_n)$  nicht konstant ist, dass die Folge der Mittelwerte  $(\overline{x^{(n)}})_{n \in \mathbb{N}}$  konvergiert und die Folge der Varianzen  $(\text{var } x^{(n)})_{n \in \mathbb{N}}$  gegen  $v \neq 0$  konvergiert. Wir kürzen ab  $Y^{(n)} := (Y_1, \dots, Y_n)$ . Dann konvergieren die Folgen von Zufallsvariablen

$$A^{(n)} := \frac{\text{cov}(x^{(n)}, Y^{(n)})}{\text{var}(x^{(n)})} \xrightarrow{n \rightarrow \infty} a \quad \text{und} \quad B^{(n)} := \overline{Y^{(n)}} - A^{(n)} \overline{x^{(n)}} \xrightarrow{n \rightarrow \infty} b.$$

in Wahrscheinlichkeit. Explizit heißt letzteres, dass

$$\lim_{n \rightarrow \infty} P[|A^{(n)} - a| > \varepsilon] = \lim_{n \rightarrow \infty} P[|B^{(n)} - b| > \varepsilon] = 0$$

für jedes  $\varepsilon > 0$  gilt.

*Beweis.* Mithilfe der in Definition 2.5 angegebenen Formeln, und mit der Abkürzung  $\mathcal{E}^{(n)} = (\mathcal{E}_1, \dots, \mathcal{E}_n)$ , ergibt sich

$$\text{cov}(x^{(n)}, Y^{(n)}) = \overline{x^{(n)} Y^{(n)}} - \overline{x^{(n)}} \overline{Y^{(n)}}$$

$$\begin{aligned}
&= \frac{1}{n} \sum_{i=1}^n x_i (ax_i + b + \mathcal{E}_i) - \overline{x^{(n)}} \left( \frac{1}{n} \sum_{i=1}^n ax_i + b + \mathcal{E}_i \right) \\
&= a \overline{x^{(n)^2}} + \overline{x^{(n)}} b + \overline{(x\mathcal{E})^{(n)}} - a \overline{x^{(n)^2}} - \overline{x^{(n)}} b - \overline{x^{(n)}} \overline{\mathcal{E}^{(n)}} \\
&= a \operatorname{var}(x^{(n)}) + \overline{(x\mathcal{E})^{(n)}} - \overline{x^{(n)}} \overline{\mathcal{E}^{(n)}}
\end{aligned}$$

Der erste Summand konvergiert per Voraussetzung sogar punktweise gegen die konstante Zufallsvariable  $v \neq 0$ . Da wir  $\mathcal{E}_i \sim \mathcal{N}(0, \sigma^2)$  haben, liefert das Gesetz der großen Zahl, siehe Satz A.14, dass  $\overline{\mathcal{E}^{(n)}} \rightarrow 0$  in Wahrscheinlichkeit gilt. Weil die Folge der Mittelwerte  $(\overline{x^{(n)}})_{n \in \mathbb{N}}$  per Voraussetzung beschränkt ist, geht der letzte Summand in Wahrscheinlichkeit gegen Null. Es bleibt der mittlere Summand

$$\overline{(x\mathcal{E})^{(n)}} = \frac{1}{n} \sum_{i=1}^n x_i \mathcal{E}_i$$

bei dem  $x_i \mathcal{E}_i \sim \mathcal{N}(0, x_i \sigma^2)$  gilt. Da der Erwartungswert  $\mathbb{E}(\overline{(x\mathcal{E})^{(n)}}) = 0$  ist, erhalten wir für die Varianz

$$\begin{aligned}
\operatorname{Var} \overline{(x\mathcal{E})^{(n)}} &= \frac{1}{n^2} \mathbb{E}((x_1 \mathcal{E}_1 + \dots + x_n \mathcal{E}_n)^2) \\
&= \frac{1}{n^2} \sum_{i=1}^n x_i^2 \mathbb{E}(\mathcal{E}_i^2) + 2 \sum_{i < j} x_i x_j \mathbb{E}(\mathcal{E}_i) \mathbb{E}(\mathcal{E}_j) \\
&= \frac{1}{n} \sigma^2 \overline{x^{(n)^2}}.
\end{aligned}$$

Die Tschebyscheff-Ungleichung liefert für  $\varepsilon > 0$

$$\mathbb{P}[|\overline{(x\mathcal{E})^{(n)}} - 0| \geq \varepsilon] \leq \frac{\operatorname{Var} \overline{(x\mathcal{E})^{(n)}}}{\varepsilon^2} = \frac{\sigma^2 \overline{x^{(n)^2}}}{n \varepsilon^2} \xrightarrow{n \rightarrow \infty} 0$$

da  $\overline{x^{(n)^2}} = \operatorname{var} \overline{x^{(n)}} + \overline{x^{(n)}}^2$  als Summe zweier konvergenter Folgen selbst konvergent ist. Es folgt nun sofort  $A^{(n)} \rightarrow av/v = a$  in Wahrscheinlichkeit. Daraus, wegen  $\overline{\mathcal{E}^{(n)}} \rightarrow 0$  und weil  $(\overline{x^{(n)}})_{n \in \mathbb{N}}$  per Voraussetzung konvergiert, erhalten wir auch

$$B^{(n)} = a \overline{x^{(n)}} + b + \overline{\mathcal{E}^{(n)}} - A^{(n)} \overline{x^{(n)}} \rightarrow b$$

in Wahrscheinlichkeit wie gewünscht.  $\square$

Satz 2.6 besagt anschaulich, dass für eine Datenmenge  $D$ , bei der die  $y_i$  aus den  $x_i$  durch eine normalverteilte additive Störung einer affin-linearen Funktion  $f: \mathbb{R} \rightarrow \mathbb{R}$ ,  $f(x) = ax + b$  entstehen, die Wahrscheinlichkeit dafür, dass die geschätzten Parameter  $a^*$  und  $b^*$  beliebig nah an den ‘echten’ Parametern  $a$  und  $b$  liegen, wiederum beliebig nah an eins gehalten werden kann, wenn der Umfang der Datenmenge nur geeignet groß gewählt wird. Wir verweisen auf Aufgabe 2.9 für eine experimentelle Behandlung.

Nach diesem probabilistischen Abstecher nehmen wir nun wieder die Optimierungsperspektive ein. Wie in Satz 2.1 bereits vorgeführt, können wir beliebige Daten aus der Klasse der affin-linearen Funktionen heraus approximieren. Allerdings kann

es ohne ein Resultat wie das in Satz 2.6 passieren, dass für eine Datenmenge keine gute Approximation durch eine affin-lineare Funktion möglich ist. Der folgende Begriff quantifiziert, ‘wie gut’ die Daten durch eine affin-lineare Funktion beschrieben werden können.

**Definition 2.7.** Sei  $D := \{(x_i, y_i) \mid i = 1, \dots, n\} \subseteq \mathbb{R}^2$  eine Datenmenge, bei der sowohl nicht alle  $x_i$  als auch nicht alle  $y_i$  gleich sind. Dann heißt

$$r_{xy} = \frac{\text{cov}(x, y)}{\text{var}(x)^{1/2} \text{var}(y)^{1/2}} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\left(\sum_{i=1}^n (x_i - \bar{x})^2\right)^{1/2} \left(\sum_{i=1}^n (y_i - \bar{y})^2\right)^{1/2}}$$

der (affin-lineare) Regressionskoeffizient von  $D$ .

**Satz 2.8.** Sei  $D := \{(x_i, y_i) \mid i = 1, \dots, n\} \subseteq \mathbb{R}^2$  eine Datenmenge, bei der sowohl nicht alle  $x_i$  als auch nicht alle  $y_i$  gleich sind. Dann ist  $r_{xy} \in [-1, 1]$  und es gilt genau dann  $r_{xy} = \pm 1$  wenn alle  $(x_i, y_i)$  auf einer Geraden mit  $\pm$ ver Steigung liegen.

*Beweis.* ① Mit  $u := (x_i - \bar{x})_{i=1, \dots, n}$ ,  $v := (y_i - \bar{y})_{i=1, \dots, n} \in \mathbb{R}^n$  folgt die Behauptung aus der Cauchy-Schwarz-Bunjakowski-Ungleichung  $|\langle u, v \rangle| \leq \|u\| \|v\|$  in  $\mathbb{R}^n$ , denn die zweite Formel in Definition 2.7 zeigt gerade

$$r_{xy} = \frac{\langle u, v \rangle}{\|u\| \|v\|}.$$

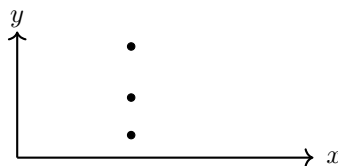
② In der Cauchy-Schwarz-Bunjakowski-Ungleichung gilt die Gleichheit  $|\langle u, v \rangle| = \|u\| \|v\|$  bekanntlich genau dann, wenn  $u$  und  $v$  linear abhängig sind. Wegen  $u \neq 0 \neq v$  ist dies äquivalent dazu, dass ein  $\lambda \in \mathbb{R}$  existiert mit  $v = \lambda u$ . Ausführlich heißt letzteres

$$\forall i = 1, \dots, n: y_i - \bar{y} = \lambda(x_i - \bar{x})$$

was bedeutet, dass alle  $(x_i, y_i)$  auf einer Gerade mit Steigung  $\lambda$  liegen, wenn  $r_{xy} \in \{+1, -1\}$  gilt. Ist letzteres der Fall, so ist  $\text{sign}(r_{xy}) = \text{sign}(\langle u, v \rangle) = \text{sign}(\lambda)$ , woraus die Aussagen über die Steigung folgen.  $\square$

Der Satz suggeriert, dass sich Daten um so besser durch eine affin-lineare Funktion beschreiben lassen, je näher  $r_{x,y}$  bei  $+1$  oder  $-1$  liegt. Diese Heuristik bestätigen wir experimentell in Aufgabe 2.9.

Jetzt wollen wir noch denjenigen Fall untersuchen, den wir in Satz 2.1 ausgeschlossen haben, nämlich dass alle  $x_i$  in unserer Datenmenge  $D := \{(x_i, y_i) \mid i = 1, \dots, n\}$  gleich sind. In diesem Fall existiert in Satz 2.1 kein eindeutig bestimmter Minimierer  $f^*: \mathbb{R}_x \rightarrow \mathbb{R}_y$ ,  $x \mapsto f^*(x)$ .



Andererseits liegen jetzt per Voraussetzung alle Datenpunkte auf einer Geraden, und diese kann auch durch eine affin-lineare Funktion dargestellt werden, wenn man nur die Achsen vertauscht und statt  $D$  die gespiegelte Datenmenge  $\tilde{D} := \{(y_i, x_i) \mid i = 1, \dots, n\}$  verwendet. Hierfür benötigt man dann, dass nicht alle  $y_i$  gleich sind und folglich kann man in diesem Sinne mit Satz 2.1 alle Datenmengen behandeln die mindestens zwei verschiedene Punkte enthalten.

Setzen wir andererseits voraus, dass sowohl nicht alle  $x_i$  als auch nicht alle  $y_i$  gleich sind, so kann Satz 2.1 auf  $D$  und auf  $\tilde{D}$  angewendet werden. Wir nennen dann die zu  $D$  gehörende Gerade die *erste Regressionsgerade* und die zu  $\tilde{D}$  gehörende die *zweite Regressionsgerade*. Letztere wird durch  $g: \mathbb{R}_y \rightarrow \mathbb{R}_x$ ,  $g(y) = \tilde{a}y + \tilde{b}$  mit

$$\tilde{a} = \frac{\text{cov}(x,y)}{\text{var}(y)} \quad \text{und} \quad \tilde{b} = \bar{x} - \tilde{a}\bar{y}$$

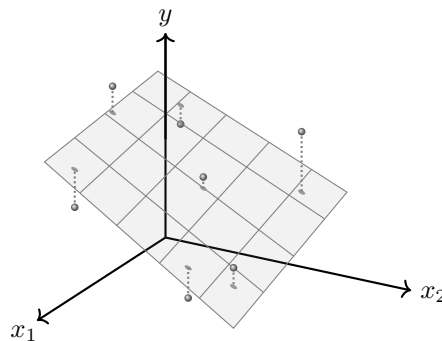
beschrieben und wir sehen, dass diese, wie auch die erste Regressionsgerade, durch den Punkt  $(\bar{x}, \bar{y})$  verläuft, wenn wir beide Geraden in das  $xy$ -Koordinatensystem einzeichnen. Im letzteren ist die Steigung der zweiten Regressionsgerade gleich der Steigung von  $g^{-1}$ , also  $1/\tilde{a}$ , falls  $\tilde{a} \neq 0$  ist. Im Allgemeinen sind beide Geraden nicht gleich, siehe Aufgabe 2.1. Sie fallen genau dann zusammen, wenn  $r_{xy} = \pm 1$  gilt, siehe Aufgabe 2.4.

## 2.2 Mehrdimensionale affin-lineare Regression

Das bisher diskutierte Verfahren nennt man häufig *einfache lineare Regression* in Abgrenzung zu *multivariabler* und *multivariater linearer Regression* bei der mehrdimensionale Features bzw. mehrdimensionale Label behandelt werden. Wir beginnen mit dem multivariablen Fall, d.h. es ist eine Datenmenge

$$D = \{(x_i, y_i) \in \mathbb{R}^d \times \mathbb{R} \mid i = 1, \dots, n\}$$

gegeben, die durch eine affin-lineare Funktion  $f: \mathbb{R}^d \rightarrow \mathbb{R}$ ,  $f(x) = \langle a, x \rangle + b$  mit  $a \in \mathbb{R}^d$  und  $b \in \mathbb{R}$ , approximiert werden soll. Für  $d = 2$  ergibt sich das folgende Bild in welchem die Datenpunkte durch eine Ebene approximiert werden.



Für beliebiges  $d$  erhalten wir das folgende, Satz 2.1 erweiternde, Resultat.

**Satz 2.9.** Sei  $D = \{(x_i, y_i) \in \mathbb{R}^d \times \mathbb{R} \mid i = 1, \dots, n\}$  eine Datenmenge mit  $x_i = (x_{i1}, x_{i2}, \dots, x_{id})$  für  $i = 1, \dots, n$ . Wir setzen voraus, dass die Matrix

$$X := \begin{bmatrix} 1 & x_{11} & \cdots & x_{1d} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n1} & \cdots & x_{nd} \end{bmatrix} \in \mathbb{R}^{n \times (d+1)}$$

Rang  $d + 1$  hat. Dann gibt es genau eine affin-lineare Funktion  $f^*: \mathbb{R}^d \rightarrow \mathbb{R}$ , d.h.  $f(x) = \langle a^*, x \rangle + b^*$  mit eindeutigen  $a^* = (a_1^*, \dots, a_d^*) \in \mathbb{R}^d$  und  $b^* \in \mathbb{R}$ , sodass

$$f^* = \operatorname{argmin} \left\{ \phi(f) := \sum_{i=1}^n (f(x_i) - y_i)^2 \mid f: \mathbb{R}^d \rightarrow \mathbb{R} \text{ affin-linear} \right\}$$

gilt. Dies wird erreicht für

$$\begin{bmatrix} b^* \\ a_1^* \\ \vdots \\ a_d^* \end{bmatrix} = (X^\top X)^{-1} X^\top \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}.$$

*Beweis.* Wie die obige Formel für die Parameter der affin-linearen Funktion suggeriert, ist es zweckmäßig diese als Vektor in  $\mathbb{R}^{d+1}$  aufzufassen. In diesem Sinne schreiben wir eine affin-lineare Funktion  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  per

$$f(z) = \langle a, z \rangle + b = \sum_{i=1}^n a_i z_i + b \cdot 1 = \left\langle \begin{bmatrix} b \\ a_1 \\ \vdots \\ a_d \end{bmatrix}, \begin{bmatrix} 1 \\ z_1 \\ \vdots \\ z_d \end{bmatrix} \right\rangle$$

für  $z = (1, z_1, \dots, z_d) \in \mathbb{R}^d$ . Damit können wir über Funktionen der Form

$$f: \mathbb{R}^{d+1} \rightarrow \mathbb{R}, \quad f = \langle \tilde{a}, \cdot \rangle,$$

bzw. den sie eindeutig bestimmenden Vektor  $\tilde{a} = (b, a_1, \dots, a_d) \in \mathbb{R}^{d+1}$  optimieren, wenn wir bei den Daten jeweils eine 1 vor dem ersten Eintrag ergänzen. Für die Zielfunktion ergibt sich dann

$$\begin{aligned} \phi(\tilde{a}) &= \sum_{i=1}^n \left( \left\langle \begin{bmatrix} b \\ a_1 \\ \vdots \\ a_d \end{bmatrix}, \begin{bmatrix} 1 \\ x_{i1} \\ \vdots \\ x_{id} \end{bmatrix} \right\rangle - y_i \right)^2 = \sum_{i=1}^n ((X\tilde{a})_i - y_i)^2 \\ &= \langle X\tilde{a} - y, X\tilde{a} - y \rangle = \langle X\tilde{a}, X\tilde{a} \rangle - \langle X\tilde{a}, y \rangle - \langle y, X\tilde{a} \rangle + \langle y, y \rangle \\ &= \langle \tilde{a}, X^\top X \tilde{a} \rangle - 2\langle \tilde{a}, X^\top y \rangle + \|y\|^2 \end{aligned}$$

wobei  $X$  die im Satz definierte Matrix bezeichnet. Wir behaupten jetzt

$$\phi'(\tilde{a})h = 2\langle h, X^\top X \tilde{a} \rangle - 2\langle h, X^\top y \rangle.$$

Der zweite Term ist dabei klar, da  $\tilde{a} \mapsto 2\langle \tilde{a}, X^\top y \rangle$  linear ist. Für den ersten Term

berechnen wir

$$\begin{aligned} \frac{|\phi(\tilde{a}+h) - \phi(\tilde{a}) - 2\langle h, X^\top X \tilde{a} \rangle|}{\|h\|} &= \frac{1}{\|h\|} |\langle \tilde{a}, X^\top X \tilde{a} \rangle + \langle \tilde{a}, X^\top X h \rangle + \langle h, X^\top X \tilde{a} \rangle \\ &\quad + \langle h, X^\top X h \rangle - \langle \tilde{a}, X^\top X \tilde{a} \rangle - 2\langle h, X^\top X \tilde{a} \rangle| \\ &= \frac{|\langle Xh, Xh \rangle|}{\|h\|} \leq \|X\|_{\text{op}}^2 \|h\| \xrightarrow{h \rightarrow 0} 0, \end{aligned}$$

wobei  $\|X\|_{\text{op}}$  die Operatornorm von  $X$  bezeichnet, und die Behauptung gezeigt ist. Mithilfe der Formel sehen wir, dass  $\phi'(\tilde{a}) = 0$  genau dann gilt, wenn  $2\langle h, X^\top X \tilde{a} \rangle - 2\langle h, X^\top y \rangle = 0$  ist für alle  $h \in \mathbb{R}^{d+1}$ . Dies ist äquivalent zu

$$X^\top X \tilde{a} - X^\top y = 0.$$

Um diese letzte Gleichung nach  $\tilde{a}$  aufzulösen, zeigen wir, dass  $X^\top X \in \mathbb{R}^{(d+1) \times (d+1)}$  invertierbar ist: Gelte hierzu  $X^\top X v = 0$  für  $v \in \mathbb{R}^{d+1}$ . Daraus folgt  $\langle v, X^\top X v \rangle = \langle Xv, Xv \rangle = \|Xv\|^2 = 0$ , also  $Xv = 0$  was wegen  $\text{rk } X = d + 1$  impliziert, dass  $v = 0$  sein muss, und  $X^\top X$  ist invertierbar.

Aus dem bisherigen folgt, dass  $\tilde{a} = (X^\top X)^{-1} X^\top y$  der einzige kritische Punkt von  $\phi: \mathbb{R}^{d+1} \rightarrow \mathbb{R}$  ist. Im Gegensatz zum eindimensionalen Beweis ist es hier nicht leicht zu sehen, was  $\liminf_{\|\tilde{a}\| \rightarrow \infty} \phi(\tilde{a})$  ist. Stattdessen greifen wir nun auf die Theorie der konvexen Funktionen vor, die wir in Kapitel 17 entwickeln werden. Zunächst bemerken wir, dass für unsere Zielfunktion

$$\phi(\tilde{a}) = \sum_{i=1}^n ((X\tilde{a})_i - y_i)^2 = \|X\tilde{a} - y\|^2$$

gilt, wobei rechts die 2-Norm auf  $\mathbb{R}^n$  gemeint ist. Wir haben es also bei  $\phi$  mit der Verkettung einer affin-linearen Funktion  $\mathbb{R}^{d+1} \rightarrow \mathbb{R}^n$ ,  $\tilde{a} \mapsto X\tilde{a} + y$ , mit dem Quadrat der euklidischen Norm  $\|\cdot\|^2: \mathbb{R}^n \rightarrow \mathbb{R}$  zu tun. Wir werden in Beispiel 17.15 und Lemma 17.16 zeigen werden, ist  $\phi: \mathbb{R}^{d+1} \rightarrow \mathbb{R}$  daher konvex. Für konvexe Funktionen gilt aber immer, dass die kritischen Punkte bereits Minimierer sind, siehe Folgerung 17.5, und wir sind mit dem Beweis durch.  $\square$

Spezialisieren wir  $d = 1$  in Satz 2.9, so vereinfacht sich die Matrix  $X$  zu

$$X := \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \in \mathbb{R}^{n \times 2}$$

und wir sehen, dass die Rangbedingung  $\text{rk } X = 2$  aus Satz 2.9 genau dann erfüllt ist, wenn die Bedingung aus Satz 2.1 gilt, nämlich nicht alle  $x_1, \dots, x_n$  gleich sind. Aus den jeweiligen Eindeutigkeitsaussagen der Sätze folgt nun sofort, dass beide die Parameter des gleichen affin-linearen Regressors liefern. Für den Fall, dass der Leser Satz 2.1 ganz und gar überspringen möchte, lässt sich aber auch direkt verifizieren, dass die Formel

$$\begin{bmatrix} b^* \\ a^* \end{bmatrix} = (X^\top X)^{-1} X^\top \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

aus Satz 2.9 auf die in Satz 2.1 gezeigten Formeln für  $a^*$  und  $b^*$  führt, siehe Aufgabe 2.5. Der Spezialfall  $d = 2$  von Satz 2.9 liefert eine Approximation von Punkten in  $\mathbb{R}^3$  durch eine Ebene, wie im Bild auf Seite 20 dargestellt. Auch hier können die drei Parameter  $a_1^*$ ,  $a_2^*$  und  $b^*$  explizit mithilfe von Stichprobenvarianz und -kovarianz angegeben werden, siehe Aufgabe 2.6.

Als nächstes untersuchen wir die Frage, wie im multivariablen Fall die Qualität der Approximation der Daten durch einen affin-linearen Regressor gemessen werden kann. Dazu verallgemeinern wir den Regressionskoeffizienten wie folgt.

**Definition 2.10.** Sei  $D = \{(x_i, y_i) \in \mathbb{R}^d \times \mathbb{R} \mid i = 1, \dots, n\}$  eine Datenmenge mit  $x_i = (x_{i1}, x_{i2}, \dots, x_{id})$  für  $i = 1, \dots, n$ . Wir setzen voraus, dass die Matrix

$$X := \begin{bmatrix} 1 & x_{11} & \cdots & x_{1d} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n1} & \cdots & x_{nd} \end{bmatrix} \in \mathbb{R}^{n \times (d+1)}$$

Rang  $d + 1$  hat und derart, dass nicht alle  $y_i$  gleich sind. Sei  $f^*: \mathbb{R}^d \rightarrow \mathbb{R}$  der affin-lineare Regressor für  $D$  aus Satz 2.9. Dann heißt

$$R^2 := \frac{\|f^*(x) - \bar{y}\|^2}{\|y - \bar{y}\|^2}$$

das *Bestimmtheitsmaß* von  $D$ , wobei  $\bar{y} = (\bar{y}, \dots, \bar{y})$  und  $f^*(x) = (f^*(x_1), \dots, f^*(x_n))$  als Vektoren in  $\mathbb{R}^n$  zu lesen sind und in Zähler und Nenner die euklidische Norm auf  $\mathbb{R}^n$  genommen wird.

Wir notieren zunächst, dass im Fall  $d = 1$  das Bestimmtheitsmaß mit dem Quadrat des Regressionskoeffizienten übereinstimmt:  $R^2 = r_{xy}^2$ , siehe Aufgabe 2.8. Im Gegensatz zum eindimensionalen Fall definieren wir für  $d > 1$  nur  $r^2$  und nicht  $R$  selbst. Im multivariablen Fall können wir insbesondere keine Aussage analog zum letzten Teil von Satz 2.8 erwarten. Wir notieren allerdings, dass  $R^2$  per Definition die Abweichung der Approximation vom Mittelwert mit der Abweichung der Daten vom Mittelwert vergleicht und daher zu erwarten ist, dass  $R^2$  genau dann nah bei 1 liegt, wenn die Daten gut durch den affin-linearen Regressor approximiert werden. Dies bestätigt der folgende Satz.

**Satz 2.11.** Seien  $D$ ,  $X$  und  $f^*$  wie in Definition 2.10.

- (i) Es gilt stets  $0 \leq R^2 = 1 - \frac{\|y - f^*(x)\|^2}{\|y - \bar{y}\|^2} \leq 1$ .
- (ii) Es gilt  $R^2 = 1$  genau dann, wenn  $f^*(x_i) = y_i$  für alle  $i = 1, \dots, n$  gilt, also genau dann wenn alle  $(x_i, y_i)$  auf einer Hyperebene liegen.

*Beweis.* (i) Es genügt die Formel für  $R^2$  zu zeigen. Die zwei Abschätzungen folgen dann unmittelbar. Wir verwenden hierfür die Abkürzungen

$$\hat{y} := f^*(x) = \underset{\substack{\uparrow \\ \text{Satz} \\ 2.9}}{X(X^\top X)^{-1}X^\top} y =: Py \quad \text{und} \quad Q := I - P.$$

Direktes Nachrechnen zeigt, dass  $P^2 = P$  und  $Q^\top = Q$  gelten, also  $P$  idempotent und  $Q$  symmetrisch ist. Es folgt dann sofort, dass  $QP = 0$  gilt. Unter Ausnutzung des vorherigen zeigen wir nun zunächst mehrere Identitäten.

- ① Durch Einsetzen verifiziert man sofort  $\langle \bar{y}, \bar{y} \rangle = n\bar{y}^2$  und  $\langle y, \bar{y} \rangle = n\bar{y}^2$ .
- ② Es gilt  $\langle y, \hat{y} \rangle = \langle y - \hat{y} + \hat{y}, \hat{y} \rangle = \langle Qy, Py \rangle + \langle \hat{y}, \hat{y} \rangle = \langle y, QPy \rangle + \langle \hat{y}, \hat{y} \rangle = \langle \hat{y}, \hat{y} \rangle$ .  
 $\uparrow_{Q=Q^\top} \quad \uparrow_{QP=0}$
- ③ Schließlich behaupten wir  $\bar{y} = \bar{\hat{y}}$ . Erstmal gilt

$$X^\top Qy = X^\top (y - Py) = X^\top y - X^\top X(X^\top X)^{-1} X^\top y = X^\top y - X^\top y = 0$$

weswegen das Matrixprodukt jeder Zeile von  $X^\top$  mit  $Qy$  gleich Null ist. Das bedeutet aber, dass das Skalarprodukt jeder Spalte von  $X$  mit  $Qy$  gleich Null ist. Dies stimmt insbesondere für die erste Spalte von  $X$  und es folgt  $\langle \mathbf{1}, Qy \rangle = 0$  mit der Abkürzung  $\mathbf{1} = (1, \dots, 1)$ . Damit folgt dann

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{n} \langle \mathbf{1}, y \rangle = \frac{1}{n} \langle \mathbf{1}, y - \hat{y} + \hat{y} \rangle = \frac{1}{n} \langle \mathbf{1}, Qy \rangle + \frac{1}{n} \langle \mathbf{1}, \hat{y} \rangle = \bar{\hat{y}}.$$

Nun zeigen wir die behauptete Gleichung. Wir beginnen mit

$$1 - \frac{\|y - f^*(x)\|^2}{\|y - \bar{y}\|^2} = \frac{\|y - \bar{y}\|^2 - \|y - \hat{y}\|^2}{\|y - \bar{y}\|^2} = \frac{-2\langle y, \bar{y} \rangle + \langle \bar{y}, \bar{y} \rangle + 2\langle y, \hat{y} \rangle - \langle \hat{y}, \hat{y} \rangle}{\|y - \bar{y}\|^2} =: (\circ)$$

und rechnen den Zähler weiter aus. Hier gelten

$$2\langle y, \bar{y} \rangle \stackrel{①}{=} 2\langle \bar{y}, \bar{y} \rangle \stackrel{③}{=} 2\langle \hat{y}, \hat{y} \rangle \stackrel{①}{=} 2\langle \hat{y}, \bar{y} \rangle \stackrel{③}{=} 2\langle \hat{y}, \bar{y} \rangle \quad \text{sowie} \quad 2\langle y, \hat{y} \rangle - \langle \hat{y}, \hat{y} \rangle \stackrel{②}{=} \langle \hat{y}, \hat{y} \rangle$$

und wir erhalten

$$(\circ) = \frac{\langle \hat{y}, \hat{y} \rangle - 2\langle \hat{y}, \bar{y} \rangle + \langle \bar{y}, \bar{y} \rangle}{\|y - \bar{y}\|^2} = \frac{\|\hat{y} - \bar{y}\|^2}{\|y - \bar{y}\|^2} = \frac{\|f^*(x) - \bar{y}\|^2}{\|y - \bar{y}\|^2} = R^2.$$

(ii) Nach (i) gilt  $R^2 = 1$  genau dann, wenn  $y - f^*(x) = 0$  ist.  $\square$

**Bemerkung 2.12.** Ist eine Datenmenge mit sowohl mehrdimensionalen Features als auch mehrdimensionalen Labels

$$D = \{(x_i, y_i) \in \mathbb{R}^d \times \mathbb{R}^k \mid i = 1, \dots, n\}$$

gegeben, so kann auch diese durch eine affin-lineare Funktion  $f^*: \mathbb{R}^d \rightarrow \mathbb{R}^k$ ,  $f^*(x) = A^*x + b^*$  mit  $A^* \in \mathbb{R}^{d \times k}$ ,  $b^* \in \mathbb{R}^k$  approximiert werden, indem die Fehlerfunktion

$$\phi(f) = \sum_{i=1}^n \|f(x_i) - y_i\|^2$$

minimiert wird. Hierfür muss die gleiche Voraussetzung wie in Satz 2.9 gemacht werden, nämlich  $\text{rk } X = d + 1$ . Aufgrund der Bauart von  $\phi$  kann dann Satz 2.9



summandenweise angewandt werden um jeweils die Koordinatenfunktionen von  $f^*$  zu erhalten. Dies liefert die Zeilen der Matrix  $A^*$  und die Einträge von  $b^*$ .

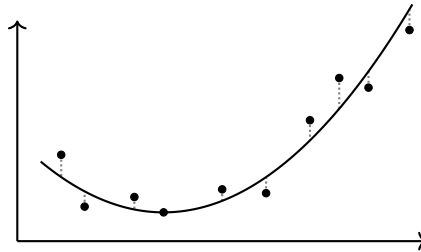
Bevor wir zum nächsten Thema übergehen, wollen wir bemerken, dass das Invertieren der Matrix  $X^\top X$  in hohen Dimensionen nicht effizient zu erledigen ist. In so fern garantieren die obigen Resultate zwar, dass es einen eindeutig bestimmten Regressor gibt, liefern aber keine praktikable Möglichkeit, diesen zu berechnen. Andererseits zeigt die von uns benutzte Beweismethode der Minimierung einer Kostenfunktion auf, wie  $f^*$  berechnet werden kann, nämlich (in der Notation von Satz 2.9) durch numerische Approximation von

$$\operatorname{argmin}_{\tilde{a} \in \mathbb{R}^{d+1}} \|X\tilde{a} - y\|^2.$$

Ein hierfür sehr populäres und anschaulich leicht nachzuvollziehendes Verfahren ist die sogenannte *Gradientenmethode*, die wir in Kapitel 17 behandeln werden.

## 2.3 Polynomiale Regression

Als nächstes behandeln wir Datenmengen, die nicht gut durch eine affin-lineare Funktion beschrieben werden können, aber stattdessen durch ein Polynom.



Wir bleiben bei der Methode der kleinsten Quadrate und erhalten den folgenden Satz, in welchem wir uns auf den eindimensionalen Fall beschränken.

**Satz 2.13.** Sei  $D = \{(x_i, y_i) \in \mathbb{R} \times \mathbb{R} \mid i = 1, \dots, n\}$  eine Datenmenge. Wir setzen voraus, dass die Matrix

$$X := \begin{bmatrix} 1 & x_1^1 & \dots & x_1^d \\ \vdots & \vdots & & \vdots \\ 1 & x_n^1 & \dots & x_n^d \end{bmatrix} \in \mathbb{R}^{n \times (d+1)}$$

Rang  $d + 1$  hat. Dann gibt es genau ein Polynom  $P^* \in \mathbb{R}[X]$  mit  $\operatorname{rk} P^* \leq d$ , d.h.  $P^* = a_0^* + \dots + a_1^* X^1 + a_d^* X^d$  mit einzigartigem  $a^* = (a_0^*, \dots, a_d^*) \in \mathbb{R}^{d+1}$ , sodass

$$P^* = \operatorname{argmin} \left\{ \phi(P) := \sum_{i=1}^n (P(x_i) - y_i)^2 \mid P \in \mathbb{R}[X] \text{ mit } \operatorname{rk} P \leq d \right\}$$

gilt. Das Polynom  $P^*$  heißt der polynomiale Regressor für die Datenmenge  $D$ . Seine

Koeffizienten sind gegeben durch

$$\begin{bmatrix} a_0^* \\ \vdots \\ a_d^* \end{bmatrix} = (X^\top X)^{-1} X^\top \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}.$$

*Beweis.* Wir definieren die Abbildung  $\psi: \mathbb{R} \rightarrow \mathbb{R}^{d+1}$ ,  $x \mapsto (1, x, x^2, \dots, x^d)$ , und betrachten die Datenmenge

$$\hat{D} := \{(\psi(x_i), y_i) \mid i = 1, \dots, n\} \subseteq \mathbb{R}^{d+1} \times \mathbb{R}.$$

Wir sehen dann, dass die zu  $\hat{D}$  gehörende Datenmatrix  $X$  aus Satz 2.9 mit der in Satz 2.13 angegebene übereinstimmt. Da wir  $\text{rk } X = d + 1$  vorausgesetzt haben, können wir Satz 2.9 anwenden und erhalten, in der Notation des aktuellen Satzes, den eindeutig bestimmten Minimierer  $f^* = \langle (a_1^*, \dots, a_d^*), \cdot \rangle + a_0^*$  der Zielfunktion

$$\begin{aligned} \phi(f) &= \sum_{i=1}^n (f(\psi(x_i)) - y_i)^2 \\ &= \sum_{i=1}^n \left( \left\langle \begin{bmatrix} a_1 \\ \vdots \\ a_d \end{bmatrix}, \begin{bmatrix} x_i^1 \\ \vdots \\ x_i^d \end{bmatrix} \right\rangle + a_0 - y_i \right)^2 \\ &= \sum_{i=1}^n ([a_0 + a_1 x_i^1 + \dots + a_d x_i^d] - y_i)^2 = \phi(P) \end{aligned} \tag{2.4}$$

wobei wir zuerst die affin-lineare Funktion  $f: \mathbb{R}^d \rightarrow \mathbb{R}$ ,  $f = \langle (a_1, \dots, a_n), \cdot \rangle + a_0$  mit den Parametern  $a_0, \dots, a_d$  identifiziert haben und dann eben diese Parameter wiederum mit dem Polynom  $P = a_0 + a_1 X + a_2 X^2 + \dots + a_d X^d$  identifizieren.  $\square$

Wir merken an, dass die Behandlung von nicht-linearen Problemen durch ‘Einbettung’ der Daten in einen höherdimensionalen Raum, in welchem diese dann aufgrund der Bauart der Einbettung eine (affin-)lineare Lösung erlauben, uns noch mehrfach begegnen wird, siehe z.B. Kapitel 15. Wir sehen außerdem, dass die Methode des obigen Beweises ohne Probleme auf Polynome in mehreren Variablen verallgemeinert werden kann: Hierfür muss lediglich die ‘Einbettung’  $\psi$  so gewählt werden, dass rechts Basisvektoren des entsprechenden Raumes von Polynomen stehen. Wir notieren schließlich, dass im Fall von  $n \leq d + 1$  der Minimierer gleich dem Interpolationspolynom ist.

## 2.4 Logistische Regression

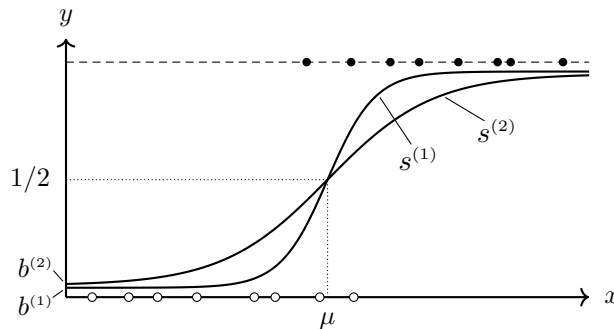
Als letztes wollen wir in diesem Kapitel die Methode der logistischen Regression behandeln. Wir beginnen mit einer Datenmenge

$$D = \{(x_i, y_i) \in \mathbb{R}^d \times \{0, 1\} \mid i = 1, \dots, n\}$$

bei der im Gegensatz zu allem bisherigen nur die Label Null und Eins vorkommen. Ist eine solche Datenmenge gegeben, so ist man daran interessiert, einen *Klassifizierer* zu bestimmen, hier also eine Abbildung  $f: \mathbb{R} \rightarrow \{0, 1\}$ , sodass  $f(x_i) = y_i$  für möglichst viele  $i \in \{1, \dots, n\}$  gilt. In den Kapiteln 3 und 13–15 werden wir uns noch genauer mit Klassifizierern beschäftigen. An dieser Stelle zeigen wir, wie ein solcher mithilfe einer Regressionsmethode gewonnen werden kann. Wie bei der affin-linearen und polynomialen Regression benötigen wir zuerst eine Klasse von Funktionen aus der heraus wir approximieren wollen. Hier wählen wir die Familie der *logistischen Funktionen*, d.h. Funktionen der Form

$$f: \mathbb{R}^d \rightarrow (0, 1), \quad f(z) = \frac{1}{1 + e^{-(w_1 z_1 + \dots + w_d z_d + b)}}$$

bei denen wir  $w_1, \dots, w_d$  und  $b \in \mathbb{R}$  variieren. Im Fall  $d = 1$  kann man  $s := 1/w_1$  und  $\mu := -b/w_1$  substituieren und erhält  $f(z) = 1/(1 + e^{-(x-\mu)/s})$  mit *Verschiebungsparameter*  $\mu$  und *Skalierungsparameter*  $s$ . Das folgende Bild zeigt zwei logistische Funktionen mit Parametern  $(b^{(1)}, w^{(1)})$  und  $(b^{(2)}, w^{(2)})$  bzw.  $(\mu, s^{(1)})$  und  $(\mu, s^{(2)})$ , wobei  $0 < s^{(1)} < s^{(2)}$  ist.



Das Bild suggeriert, dass wir die Funktion  $f$  durch Verschieben und Skalieren an die Datenpunkte anpassen können. Da die im Bild dargestellten Datenpunkten ‘überlappen’ kann man sich insbesondere vorstellen, dass sich ein Gleichgewicht einstellen wird, wenn man an die Minimierung der Summe der (quadratischen) Abstände denkt. Das ist zwar nicht die Strategie die wir unten verfolgen werden um einen *logistischen Regressor*  $f^*$  für die Daten zu finden, liefert aber dennoch die richtige Intuition. Hat man einmal  $f^*$ , so kann man durch Rundung der Werte von  $f^*$  daraus leicht einen  $\{0, 1\}$ -wertigen Klassifizierer machen.

Wie es sich bereits bei der affin-linearen Regression bewährt hat, ergänzen wir unsere Datenpunkte mit einer Eins — diesmal aber an der letzten Stelle<sup>1</sup>. Für  $x = (x_1, \dots, x_d) \in \mathbb{R}^d$  schreiben wir  $\hat{x} := (x_1, \dots, x_d, 1)$  für den (*um Eins*) *erweiterten Datenpunkt* und  $(w, b) = (w_1, \dots, w_d, b)$  für den *zusammengefassten Gewichtsvektor*.

<sup>1</sup>Wir ändern hier die Konvention im Vergleich zu Kapitel 2.2, erreichen so allerdings Konsistenz mit Kapitel 13 zum Perzeptron und Kapitel 16 zu neuronalen Netzen, vergleiche insbesondere Definition 13.9.

Damit erhalten wir

$$\langle w, z \rangle + b = \sum_{i=1}^d w_i z_i + b \cdot 1 = \left\langle \begin{bmatrix} w_1 \\ \vdots \\ w_d \\ b \end{bmatrix}, \begin{bmatrix} x_1 \\ \vdots \\ x_d \\ 1 \end{bmatrix} \right\rangle = \langle w, \hat{x} \rangle$$

wobei wir (etwas missbräuchlich!) links mit  $w$  den normalen Gewichtsvektor und rechts, ebenfalls mit  $w$ , den zusammengefassten Gewichtsvektor bezeichnen. Benutzen wir jetzt noch die *Sigmoidfunktion*, gewissermaßen den Prototyp einer logistischen Funktion mit Verschiebungsparameter Null und Skalierungsparameter Eins,

$$\text{sig}: \mathbb{R} \rightarrow (0, 1), \text{sig}(t) := \frac{1}{1 + e^{-t}},$$

so kann eine beliebige logistische Funktion  $f: \mathbb{R}^d \rightarrow (0, 1)$  via  $f(z) = \text{sig}(\langle w, \cdot \rangle)$  durch  $w \in \mathbb{R}^{d+1}$  beschrieben werden.

Da man die Werte der logistischen Funktion durch Rundung häufig als Wahrscheinlichkeiten auffasst, liegt es nahe, die Maximum-Likelihood-Methode anzuwenden, um zu präzisieren, wie die Parameter  $w \in \mathbb{R}^{d+1}$  an die Datenmenge angepasst werden sollen. Wir betrachten dazu auch hier erst einmal den Fall, dass unsere Datenpunkte in folgender Weise tatsächlich von einer logistischen Funktion  $f$  stammen: Ist ein Feature  $x_i$  gegeben, so nehmen wir an, dass das Label  $y_i$  gleich Eins ist mit Wahrscheinlichkeit  $f(x_i)$  und dementsprechend gleich Null mit Wahrscheinlichkeit  $1 - f(x_i)$ . Formal betrachten wir wieder einen Wahrscheinlichkeitsraum  $(\Omega, \Sigma, P)$  und unabhängige Zufallsvariablen  $Y_1, \dots, Y_n: \Omega \rightarrow \mathbb{R}$ , welche jeweils Bernoulli-verteilt sind, d.h.  $Y_i \sim \mathcal{B}(f(x_i))$  bei festen  $x_1, \dots, x_n$ . Dann gilt

$$P[Y_i(f) = y_i \text{ für alle } i] = \prod_{i=1}^n f(x_i)^{y_i} (1 - f(x_i))^{1-y_i}$$

und entsprechend definieren wir die Likelihood-Funktion als

$$\begin{aligned} L: \{f: \mathbb{R} \rightarrow (0, 1) \mid f \text{ logistische Funktion}\} &\rightarrow \mathbb{R} \\ L(f) &:= \prod_{i=1}^n f(x_i)^{y_i} (1 - f(x_i))^{1-y_i}. \end{aligned} \tag{2.5}$$

Ein Maximierer  $f^*$  von  $L$  maximiert dann die Wahrscheinlichkeit, dass die Daten in  $D$  in der oben beschriebenen Weise von der logistischen Funktion  $f^*$  stammen. Im Gegensatz zur affin-linearen Regression können wir, wie das folgende Beispiel zeigt, mit diesem Ansatz im Allgemeinen nicht erwarten, dass stets ein Maximierer existiert.

**Beispiel 2.14.** Wir betrachten die Datenmenge  $D = \{(-1, 0), (1, 1)\} \subseteq \mathbb{R} \times \{0, 1\}$  und die logistischen Funktionen  $f: \mathbb{R} \rightarrow (0, 1)$ ,  $f(x) := (1 + e^{-(wx+b)})^{-1}$  mit  $w, b \in \mathbb{R}$ . Dann gilt

$$L(f) = f(-1)^0 (1 - f(-1))^{1-(-1)} f(1)^1 (1 - f(1))^{1-1}$$

$$= \left(1 - \frac{1}{1 + e^{-(-w+b)}}\right) \left(\frac{1}{1 + e^{-(w+b)}}\right) < 1$$

für alle  $f$ . Für festes  $b$  und  $w \rightarrow \infty$  konvergiert  $L(f)$  gegen 1. Folglich kann  $L$  keinen Maximierer besitzen.

Skizziert man die Datenmenge aus Beispiel 2.14, so sieht man, dass gerade die Tatsache, dass die Daten nicht wie im Bild auf Seite 27 überlappen, für den obigen Effekt verantwortlich ist. In der Tat zeigt Aufgabe 2.11, dass für die Datenmenge  $D := \{(-1, 1), (0, 0), (1, 1)\}$  die Likelihood-Funktion einen Maximierer besitzt und dass dieser mithilfe analytischer Methoden explizit berechnet werden kann. Im Gegensatz zur affin-linearen und polynomialen Regression, kann im Allgemeinen aber keine explizite Formel für die Parameter von  $f^*$  angegeben werden. Wohl aber garantiert eine oben bereits angedeutete ‘Überlappingsbedingung’ die Existenz und, wie im Fall der mehrdimensionalen affin-linearen Regression, eine Rangbedingung die Eindeutigkeit.

Um dies beides zu zeigen, fassen wir ab jetzt die Likelihood-Funktion aus (2.5) als Funktion  $L = L(w)$  mit  $w \in \mathbb{R}^{d+1}$  auf und betrachten dann die *negative Log-Likelihood-Funktion*, d.h. wir definieren

$$\ell: \mathbb{R}^{d+1} \rightarrow \mathbb{R}, \quad \ell(w) := -\log \prod_{i=1}^n \text{sig}(\langle w, \hat{x}_i \rangle)^{y_i} (1 - \text{sig}(\langle w, \hat{x}_i \rangle))^{1-y_i} \quad (2.6)$$

und notieren zunächst die folgende Formel für  $\ell$ .

**Lemma 2.15.** *Sei  $D = \{(x_i, y_i) \in \mathbb{R}^d \times \{0, 1\} \mid i = 1, \dots, n\}$  eine Datenmenge und  $\ell$  die zugehörige Log-Likelihood-Funktion wie in (2.6). Dann gilt für  $w \in \mathbb{R}^{d+1}$*

$$\ell(w) = \sum_{i=1}^n -y_i \langle w, \hat{x}_i \rangle + \log(1 + e^{\langle w, \hat{x}_i \rangle}).$$

*Beweis.* Direktes Ausrechnen zeigt

$$\begin{aligned} \ell(w) &= -\sum_{i=1}^n y_i \log\left(\frac{1}{1 + e^{-\langle w, \hat{x}_i \rangle}}\right) + (1 - y_i) \log\left(1 - \frac{1}{1 + e^{-\langle w, \hat{x}_i \rangle}}\right) \\ &= -\sum_{i=1}^n y_i \log\left(\frac{1}{1 + e^{-\langle w, \hat{x}_i \rangle}}\right) + \log\left(\frac{e^{-\langle w, \hat{x}_i \rangle}}{1 + e^{-\langle w, \hat{x}_i \rangle}}\right) - y_i \log\left(\frac{e^{-\langle w, \hat{x}_i \rangle}}{1 + e^{-\langle w, \hat{x}_i \rangle}}\right) \\ &= -\sum_{i=1}^n y_i \log\left(\frac{1}{1 + e^{-\langle w, \hat{x}_i \rangle}} \cdot \frac{1 + e^{-\langle w, \hat{x}_i \rangle}}{e^{-\langle w, \hat{x}_i \rangle}}\right) + \log\left(\frac{1}{1 + e^{-\langle w, \hat{x}_i \rangle}}\right) \\ &= -\sum_{i=1}^n y_i \langle w, \hat{x}_i \rangle - \log(1 + e^{\langle w, \hat{x}_i \rangle}) \end{aligned}$$

wie behauptet. □

Mit der oben angegebenen Darstellung der Funktion  $\ell$  ist es nicht schwer, deren Ableitung zu bestimmen. Wir überlassen dies den Lesenden als Aufgabe 2.11 und

notieren hier nur das Ergebnis, nämlich

$$\nabla \ell(w) = \sum_{i=1}^n (\text{sig}(\langle w, \hat{x}_i \rangle) - y_i) \hat{x}_i.$$

Angenommen, wir haben eine Lösung des Gleichungssystems  $\nabla \ell(w) = 0$  gefunden. Um dann zu schließen, dass es sich dabei um einen Minimierer handelt, zeigen wir als nächstes, dass  $\ell$  konvex ist. Beachte, dass die Datenmatrix  $X$  unten anders definiert ist als in Satz 2.9, dass aber die dortige Rangbedingung zu der folgenden äquivalent ist.

**Lemma 2.16.** *Sei  $D = \{(x_i, y_i) \in \mathbb{R}^d \times \{0, 1\} \mid i = 1, \dots, n\}$  eine Datenmenge. Dann ist die zugehörige negative Log-Likelihood-Funktion  $\ell$  wie in (2.6) konvex. Hat zusätzlich die Datenmatrix*

$$X := \begin{bmatrix} x_{11} & \cdots & x_{1d} & 1 \\ \vdots & & \vdots & \vdots \\ x_{n1} & \cdots & x_{nd} & 1 \end{bmatrix} \in \mathbb{R}^{n \times (d+1)}$$

*Rang  $d + 1$ , so ist  $\ell$  sogar strikt konvex.*

*Beweis.* Wir betrachten  $g: \mathbb{R} \rightarrow \mathbb{R}$ ,  $g(t) = -\log \text{sig}(t) = \log(1 + e^{-t})$  und sehen, dass  $g'(t) = \text{sig}(t) - 1$  strikt monoton wächst, woraus folgt, dass  $g$  strikt konvex ist. Weiter betrachten wir  $h: \mathbb{R} \rightarrow \mathbb{R}$ ,  $h(t) = -\log(1 - \text{sig}(t)) = g(t) + t$ . Hier gilt  $h'(t) = \text{sig}(t) + 1$  und es folgt, dass auch  $h$  strikt konvex ist. Jetzt fassen wir im Beweis von Lemma 2.15 anders zusammen und erhalten

$$\begin{aligned} \ell(w) &= - \sum_{i=1}^n y_i \log\left(\frac{1}{1 + e^{-\langle w, \hat{x}_i \rangle}}\right) + (1 - y_i) \log\left(1 - \frac{1}{1 + e^{-\langle w, \hat{x}_i \rangle}}\right) \\ &= \sum_{i=1}^n y_i (-\log -\langle w, \hat{x}_i \rangle) + (1 - y_i) (-\log(1 - \text{sig}(\langle w, \hat{x}_i \rangle))) \\ &= \sum_{i=1}^n y_i g(\langle w, \hat{x}_i \rangle) + (1 - y_i) h(\langle w, \hat{x}_i \rangle). \end{aligned}$$

Weil  $y_i \in \{0, 1\}$  ist, bleibt in der Formel für  $\ell$  für jedes  $i$  jeweils entweder nur der Term mit  $g$  oder nur der Term mit  $h$  übrig. Wir können daher die Funktion  $\ell$  schreiben als

$$\ell(w) = \sum_{i=1}^n k_i(\langle w, \hat{x}_i \rangle)$$

mit  $k_i \in \{g, h\}$  für  $i = 1, \dots, n$ . Um die (strikte) Konvexität einzusehen, seien  $\lambda \in (0, 1)$  und  $v \neq w \in \mathbb{R}^{d+1}$ . Dann gilt

$$\begin{aligned} \ell(\lambda w + (1 - \lambda)v) &= \sum_{i=1}^n k_i(\lambda \langle w, \hat{x}_i \rangle + (1 - \lambda) \langle v, \hat{x}_i \rangle) \\ &\leq \sum_{i=1}^n \lambda k_i(\langle w, \hat{x}_i \rangle) + (1 - \lambda) k_i(\langle v, \hat{x}_i \rangle) \\ &= \lambda \ell(w) + (1 - \lambda) \ell(v) \end{aligned}$$

Da die  $k_i$  alle strikt konvex sind, erhalten wir in der obigen Rechnung eine strikte Abschätzung, wenn nur  $\langle w, \hat{x}_i \rangle \neq \langle v, \hat{x}_i \rangle$  für mindestens ein  $i \in \{1, \dots, n\}$  gilt. Wäre das nicht so, dann würde  $\langle \hat{x}_i, u \rangle = 0$  für alle  $i$  gelten mit  $u := v - w \in \mathbb{R}^{d+1} \setminus \{0\}$  und daher

$$Xu = \begin{bmatrix} x_{11} & \cdots & x_{1d} & 1 \\ \vdots & & \vdots & \vdots \\ x_{n1} & \cdots & x_{nd} & 1 \end{bmatrix} \begin{bmatrix} u_1 \\ \vdots \\ u_{d+1} \end{bmatrix} = \begin{bmatrix} \langle \hat{x}_1, u \rangle \\ \vdots \\ \langle \hat{x}_n, u \rangle \end{bmatrix} = 0$$

im Widerspruch zu  $\text{rk } X = d + 1$ .  $\square$

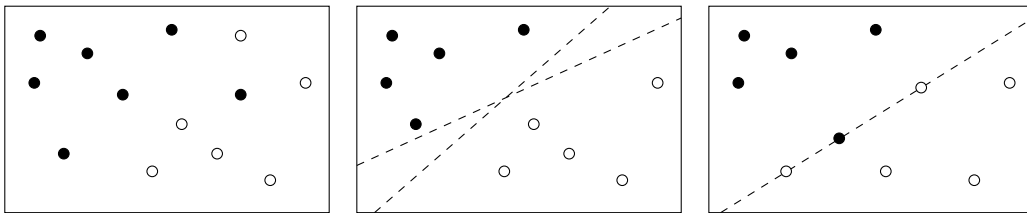
Haben wir eine Lösung  $w^*$  von  $\nabla \ell(w) = 0$  und ist die Rangbedingung in Lemma 2.15 erfüllt, so folgt, dass  $w^*$  der eindeutige Minimierer von  $\ell$  ist und somit den Maximum Likelihood Schätzer liefert. Für größere Datenmengen ist es allerdings nicht möglich  $\nabla \ell(w) = 0$  analytisch zu lösen und es muss stattdessen ein numerisches Verfahren angewandt werden. In Kapitel 17 werden wir die sogenannte Gradientenmethode behandeln, und zeigen, dass diese besonders gut funktioniert, wenn eine konvexe Funktion zu minimieren ist. Dies ist, neben der Eindeutigkeitsfrage, ein weiterer Grund dafür, dass man in diesem Kontext die negative Log-Likelihood-Funktion einführt. Unter der Rangbedingung ist dann sogar garantiert, dass es höchstens einen Minimierer geben kann, wenn man, im Vorgriff auf Kapitel 17 das Resultat in Proposition 17.14 benutzt. Im letzten Satz dieses Kapitels wollen wir jetzt noch die Existenzfrage behandeln. Die hierfür nötige Voraussetzung ist die folgende.

**Definition 2.17.** Sei  $D = \{(x_i, y_i) \in \mathbb{R}^d \times \{0, 1\} \mid i = 1, \dots, n\}$  eine Datenmenge. Wir sagen, dass die Daten *überlappen*, wenn für jedes  $w \in \mathbb{R}^{d+1}$  ein  $k$  existiert sodass  $y_k = 1$  und  $\langle w, \hat{x}_k \rangle < 0$  oder  $y_k = 0$  und  $\langle w, \hat{x}_k \rangle > 0$  gelten.

Betrachtet man für  $w \in \mathbb{R}^{d+1}$  den Klassifizierer

$$\mathbb{R}^d \rightarrow \{0, 1\}, x \mapsto \begin{cases} 1 & \text{falls } -(\langle w, x \rangle) \geq 1/2, \\ 0 & \text{sonst,} \end{cases}$$

der durch Rundung aus dem Regressor hervorgeht, so impliziert Definition 2.17, dass für jedes  $w$  mindestens ein Datenpunkt aus  $D$  durch diesen falsch klassifiziert wird. Die Bedingung ist aber in der Tat echt stärker, betrachte z.B. die eindimensionale Datenmenge  $D = \{(0, 0), (0, 1)\}$ . In zwei Dimensionen illustrieren die folgenden drei Bilder den Sachverhalt.



Links überlappen die Daten, in der Mitte und rechts nicht. Dabei gibt es im mittleren Bild allerdings sogar unendlich viele  $w$ 's die zu einem korrekten Klassifizierer führen würden, während es im rechten Bild kein solches  $w$  gibt.

Wir werden nun zeigen, dass im Fall überlappender Daten, also wie im Bild links, mindestens ein Maximum-Likelihood-Schätzer im Sinne eines Maximierers von  $L$  aus (2.5) existiert. *Linear trennbare Datenmengen* wie im mittleren Bild werden wir in späteren Kapiteln 13–14 mit anderen Methoden behandeln und dort auch skizzieren, auf welche Weise man deren im rechten Bild dargestellten Grenzfall behandeln kann.

**Satz 2.18.** *Sei  $D = \{(x_i, y_i) \in \mathbb{R}^d \times \{0, 1\} \mid i = 1, \dots, n\}$  eine überlappende Datenmenge und sei  $\ell$  die zugehörige negative Log-Likelihood-Funktion wie in (2.6). Dann existiert mindestens ein Minimierer  $w^* \in \arg\min_{w \in \mathbb{R}^{d+1}} \ell(w)$ .*

*Beweis.* Da die Funktion  $\ell$  stetig und auf ganz  $\mathbb{R}^{d+1}$  definiert ist, genügt es, ihr Verhalten im Unendlichen zu untersuchen. Sei dazu  $(w_k)_{k \in \mathbb{N}} \subseteq \mathbb{R}^{d+1}$  eine Folge mit  $\|w_k\| \rightarrow \infty$ . Wir behaupten, dass dann auch  $\ell(w_k) \rightarrow \infty$  für  $k \rightarrow \infty$  gilt. Hierfür benötigen wir mehrere Vorbereitungen.

① Unser erstes Ziel ist es, zu zeigen, dass die Funktion

$$\ell^\infty: \mathbb{R}^{d+1} \setminus \{0\} \rightarrow (0, \infty), \quad \ell^\infty(w) := \lim_{t \rightarrow \infty} \frac{\ell(tw)}{t}$$

wohldefiniert ist. Mithilfe der Formel für  $\ell$  aus Lemma 2.15 sehen wir, dass sich im ersten Term die  $t$ 's kürzen und dass der zweite Term gegen Null geht, wenn  $\langle w, \hat{x}_i \rangle \leq 0$  ist, und andernfalls gerade gegen letzteres Skalarprodukt konvergiert. Der Grenzwert

$$\begin{aligned} \lim_{t \rightarrow \infty} \frac{\ell(tw)}{t} &= \lim_{t \rightarrow \infty} \sum_{i=1}^n -y_i \frac{1}{t} \langle tw, \hat{x}_i \rangle + \frac{1}{t} \log(1 + e^{\langle tw, \hat{x}_i \rangle}) \\ &= \sum_{\substack{i=1 \\ y_i=1}}^n -\langle w, \hat{x}_i \rangle + \sum_{\substack{i=1 \\ \langle w, \hat{x}_i \rangle > 0}}^n \langle w, \hat{x}_i \rangle \end{aligned}$$

existiert also schonmal. Sei nun zunächst  $k \in \{1, \dots, n\}$  derart, dass  $y_k = 1$  und  $\langle w, \hat{x}_k \rangle < 0$  gilt. Dann folgt aus der Abschätzung

$$\begin{aligned} \ell(tz) &= -y_k t \langle w, \hat{x}_k \rangle + \underbrace{\log(1 + e^{t \langle w, \hat{x}_k \rangle})}_{\geq 0} + \sum_{i \neq k} -y_i t \langle w, \hat{x}_i \rangle + \log(1 + e^{t \langle w, \hat{x}_i \rangle}) \\ &\geq -t \langle w, \hat{x}_k \rangle + \sum_{\substack{i \neq k \\ y_i=0}} \underbrace{\log(1 + e^{t \langle w, \hat{x}_i \rangle})}_{\geq 0} + \sum_{\substack{i \neq k \\ y_i=1}} -t \langle w, \hat{x}_i \rangle + \underbrace{\log(1 + e^{t \langle w, \hat{x}_i \rangle})}_{\geq t \langle w, \hat{x}_i \rangle} \\ &\geq -t \langle w, \hat{x}_k \rangle \end{aligned}$$

dass der Grenzwert für jedes  $w \neq 0$  echt positiv ist. Ist  $k \in \{1, \dots, n\}$  derart, dass  $y_k = 0$  und  $\langle w, \hat{x}_i \rangle > 0$ , so erhalten wir

$$\ell(tz) = 0 + \underbrace{\log(1 + e^{t \langle w, \hat{x}_k \rangle})}_{\geq t \langle w, \hat{x}_k \rangle} + \sum_{i \neq k} -y_i t \langle w, \hat{x}_i \rangle + \log(1 + e^{t \langle w, \hat{x}_i \rangle}) \geq t \langle w, \hat{x}_k \rangle$$

und sehen wie oben, dass die Summe über  $i \neq k$  größer gleich Null ist.



② Für  $v, w \in \mathbb{R}^{d+1}$  schätzen wir als nächstes

$$\begin{aligned}
 |\ell(v) - \ell(w)| &\leq \sum_{i=1}^n |-y_i \langle v, \hat{x}_i \rangle + \log(1 + e^{\langle v, \hat{x}_i \rangle}) + y_i \langle w, \hat{x}_i \rangle - \log(1 + e^{\langle w, \hat{x}_i \rangle})| \\
 &\leq \sum_{i=1}^n |\langle v - w, \hat{x}_i \rangle| + \left| \frac{d}{ds} \log(1 + e^s) \right| \cdot |\langle v, \hat{x}_i \rangle - \langle w, \hat{x}_i \rangle| \\
 &\leq \sum_{i=1}^n 2 \|v - w\| \|\hat{x}_i\| \\
 &\leq L \|v - w\|
 \end{aligned}$$

ab und sehen so, dass  $\ell$  Lipschitz-stetig mit Konstante  $L := 2 \max_{i=1, \dots, n} \|\hat{x}_i\|$  ist.

③ Sei nun  $(w_k)_{k \in \mathbb{N}}$  wie am Anfang des Beweises gegeben. Durch Übergang zu einer Teilfolge können wir ohne Einschränkung annehmen, dass alle  $w_k$  ungleich Null sind. Dann hat die Folge  $(v_k)_{k \in \mathbb{N}} \subseteq \partial B_1(0)$  mit  $v_k := w_k / \|w_k\|$  eine konvergente Teilfolge mit Grenzwert  $v$  in der Einheitssphäre. Ohne Einschränkung nehmen wir an, dass bereits  $v_k \rightarrow v$  für  $k \rightarrow \infty$  gilt und behaupten

$$\lim_{k \rightarrow \infty} \frac{\ell(w_k)}{\|w_k\|} = \ell^\infty(v) \stackrel{\textcircled{1}}{>} 0$$

woraus insbesondere  $\ell(w_k) \rightarrow \infty$  folgt. Um die Gleichheit oben zu zeigen, setzen wir  $t_k := \|w_k\|$  und erhalten  $\ell(w_k) / \|w_k\| = \ell(t_k v_k) / t_k$ . Für  $\varepsilon > 0$  wählen wir  $k_0 \in \mathbb{N}$  derart dass für  $k \geq k_0$  sowohl  $|\ell(t_k v) / t_k - \ell^\infty(v)| < \varepsilon/2$  als auch  $|v_k - v| < \varepsilon/2$  ausfallen. Für  $k \geq k_0$  ergibt sich

$$\begin{aligned}
 \left| \frac{\ell(w_k)}{\|w_k\|} - \ell^\infty(v) \right| &\leq \left| \frac{\ell(t_k v_k)}{t_k} - \frac{\ell(t_k v)}{t_k} \right| + \left| \frac{\ell(t_k v)}{t_k} - \ell^\infty(v) \right| \\
 &\stackrel{\textcircled{2}}{\leq} \frac{L |t_k v_k - t_k v|}{t_k} + \frac{\varepsilon}{2} < \varepsilon
 \end{aligned}$$

und wir sind fertig. □

Wir verzichten an dieser Stelle darauf, in einem konkreten Beispiel die Parameterbestimmung mit der Gradientenmethode zu illustrieren. Erstens werden wir letztere erst in Kapitel 17 im Detail erklären und zweitens gibt es in vielen Programmiersprachen fertige Pakete, welche numerisch die Parameter  $w_1^*, \dots, w_d^*$  und  $b^*$  bestimmen, vgl. Aufgabe 2.12. Hat man  $w^*$  wie in Satz 2.16 gefunden, so nennen wir

$$f^*: \mathbb{R}^d \rightarrow (0, 1), \quad f^*(z) = \text{sig}(\langle w^*, \hat{z} \rangle).$$

einen *logistischen Regressor* für die Datenmenge  $D$ . Wir haben oben bereits erklärt, dass man daraus durch Runden einen  $\{0, 1\}$ -wertigen *Klassifizierer* gewinnen kann.

**Bemerkung 2.19.** Bevor wir das Kapitel beenden, wollen wir noch auf folgendes hinweisen: Durch Umstellen und Logarithmieren liefert die Definition einer logisti-

schen Funktion  $f = \text{sig}(\langle w, \cdot \rangle)$  mit  $w \equiv (w, b) \in \mathbb{R}^{d+1}$  die Gleichung

$$\log \frac{f}{1-f} = w_1 z_1 + \cdots + w_d z_d + b$$

für  $z = (z_1, \dots, z_d) \in \mathbb{R}^d$ . Man kann diese Gleichung so interpretieren, als dass man den Logarithmus der ‘Chance’  $\frac{f}{1-f}$  ( $f$ , dass das Label von  $z$  gleich Eins ist, zu  $1-f$ , dass das Label von  $z$  gleich Null ist) durch eine affin-lineare Funktion der Features ausdrückt. Der Versuch allerdings, die Datenmenge via der linken Funktion zu transformieren und dann ‘gewöhnliche’ affin-lineare Regression anzuwenden, schlägt hier — im Gegensatz zur polynomialen Regression in Kapitel 2.3 — aber fehl. Hierzu müsste nämlich die offenbar nicht wohldefinierte Datenmenge

$$\hat{D} := \left\{ (x_i, \log \frac{y_i}{1-y_i}) \mid i = 1, \dots, n \right\}$$

benutzt werden. Hat man Daten, die zwar kategoriell interpretiert werden sollen, bei denen die Label aber in  $(0, 1)$  liegen, so kann obige Methode jedoch durchaus angewandt werden, vergleiche Aufgabe 2.13.

## Referenzen

Die am Anfang dieses Kapitels vorgestellten Standardresultate zur linearen Regression sind überall zu finden. Den Begriff der ‘zweiten Regressionsgerade’ haben wir von aus dem Vorlesungsmanuskript [Rin08] übernommen. Die Konsistenzaussage in Satz 2.6 folgt dem Beweis [Sha15, Lecture 3]. Der Abschnitt zu polynomialer Regression orientiert sich an [SSBD14, Section 9.2.2]. Eine gute Referenz zu Existenz- und Eindeutigkeitsaussagen des logistischen Regressors ist [AA84]. Dort werden drei Fälle eingeführt: overlap, separation und quasi-separation. Unser Beweis in 2.18 basiert auf dem Konzept der ‘recession function’, siehe z.B. [Giu03]. Der Autor bedankt sich bei T. Schmidt für mehrere hilfreiche Diskussionen in diesem Kontext.

Im Zusammenhang mit den Maximum-Likelihood-Funktionen (2.2) und (2.5) und als Voraussetzung in Satz 2.6 haben wir endlich bzw. abzählbar viele unabhängige Zufallsvariablen mit vorgeschriebenen Verteilungen auf einem festen Wahrscheinlichkeitsraum  $(\Omega, \Sigma, P)$  betrachtet. Das derartige existiert folgt aus dem sogenannten Klonsatz und einer Variante davon [Beh13, Sätze 4.5.1 und 4.5.2], die wir im Anhang als Satz A.13 aufführen.

## Aufgaben

**Aufgabe 2.1.** Bestimmen Sie für die Daten  $x = (1, 2, 2, 3, 4)$  und  $y = (2, 2, 3, 3, 3)$  beide Regressionsgeraden und die Regressionskoeffizienten  $r_{xy}$  und  $r_{yx}$ , und zwar von Hand (Taschenrechner oder CAS ist erlaubt, aber kein fertiges Python-Paket!). Skizzieren Sie die Datenpunkte und die Regressionsgeraden.

**Aufgabe 2.2.** In dieser Aufgabe bezeichnet  $\phi$  die Zielfunktion aus Satz 2.1 und  $\psi$  diejenige aus Bemerkung 2.3(ii).

(i) Finden Sie eine Datenmenge  $D$  sodass gilt

$$\operatorname{argmin}\{\phi(f) \mid f: \mathbb{R} \rightarrow \mathbb{R} \text{ affin-linear}\} \neq \operatorname{argmin}\{\psi(f) \mid f: \mathbb{R} \rightarrow \mathbb{R} \text{ affin-linear}\}.$$

- (ii) Finden Sie eine Datenmenge, sodass nicht alle  $x$ -Werte gleich sind und es trotzdem mehr als eine Funktion  $f^*$  gibt mit

$$\psi(f^*) = \min\{\psi(f) \mid f: \mathbb{R} \rightarrow \mathbb{R} \text{ affin-linear}\}.$$

**Aufgabe 2.3.** Führen Sie für das folgende Dataset die Methode der linearen Regression. Die Klausurvorbereitungszeit soll hierbei die  $x$ -Variable sein und das Klausurergebnis die  $y$ -Variable. Zeichnen Sie alle Datenpunkte und auch die Regressionsgerade.

Student	Klausurvorbereitung in h	Klausurergebnis in %
1	21.0	82
2	18.0	69
3	15.0	29
4	8.0	41
5	8.0	44
6	1.5	8
7	0.0	10

**Aufgabe 2.4.** Sei  $D = \{(x_i, y_i) \in \mathbb{R}^2 \mid i = 1, \dots, n\}$  eine Datenmenge bei der weder alle  $x_i$  gleich sind noch alle  $y_i$  gleich sind. Zeigen Sie, dass genau dann die erste und die zweite Regressionsgerade zusammenfallen, wenn  $r_{xy} = \pm 1$  gilt.

**Aufgabe 2.5.** Sei  $D$  ein Dataset wie in Satz 2.1 und  $X$  definiert wie in Satz 2.9. Zeigen Sie durch direktes Nachrechnen, dass

$$\begin{bmatrix} b^* \\ a^* \end{bmatrix} := (X^\top X)^{-1} X^\top \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

genau diejenigen Formeln für  $a^*$  und  $b^*$  liefert, die in Satz 2.1 angegeben wurden.

**Aufgabe 2.6.** Zeigen Sie, dass unter den Voraussetzungen von Satz 2.9, aber im Spezialfall  $N = 2$ , der lineare Regressor durch  $f: \mathbb{R}^2 \rightarrow \mathbb{R}$ ,  $f = \langle a, \cdot \rangle + b$  mit

$$\begin{aligned} b &= \bar{y} - a_1 \bar{x}_1 - a_2 \bar{x}_2 \\ a_1 &= \frac{\text{cov}(x_1, y) \text{cov}(x_2, x_2) - \text{cov}(x_2, y) \text{cov}(x_1, x_2)}{\text{cov}(x_1, x_1) \text{cov}(x_2, x_2) - \text{cov}(x_1, x_2)^2} \\ a_2 &= \frac{\text{cov}(x_2, y) \text{cov}(x_1, x_1) - \text{cov}(x_1, y) \text{cov}(x_1, x_2)}{\text{cov}(x_1, x_1) \text{cov}(x_2, x_2) - \text{cov}(x_1, x_2)^2} \end{aligned}$$

gegeben wird, wobei  $a = (a_1, a_2)$ .

*Hinweis:* Lösen Sie in der Notation von Satz 2.9 das LGS  $X^\top X \begin{bmatrix} b \\ a \end{bmatrix} = X^\top y$ . Lassen Sie dabei zur Schreiberleichterung bei auftretenden Summen den Summationsindex weg, wenn keine Missverständnisse auftreten können, z.B.  $\Sigma x_1 y := \sum_{i=1}^n x_{1i} y_i$ . Um das Ergebnis mithilfe der Kovarianz auszudrücken zeigen Sie dann zuerst, dass  $n \text{cov}(r, s) = (\Sigma r s) - n \bar{r} \bar{s}$  für  $r, s \in \mathbb{R}^n$  gilt.

**Aufgabe 2.7.** Implementieren Sie die Formeln aus Aufgabe 2.6 in Python und wenden Sie diese auf die folgende Datenmenge an, wobei Klausurvorbereitungszeit und Zeit auf sozialen Medien die  $x$ -Variablen und das Klausurergebnis die  $y$ -Variable ist. Erstellen Sie einen Plot der Datenpunkte und der Regressionsebene.

Student	Klausurvorbereitung in h	Soziale Medien in h	Klausurergebnis in %
1	0.0	20.0	0.0
2	1.5	8.5	2.0
3	2.0	6.0	7.0
4	2.0	6.0	10.5
5	8.0	10.0	29.5
6	8.5	3.0	49.0
7	9.5	0.0	59.5
8	12.0	2.0	63.5
9	18.0	4.0	85.0
10	19.0	0.5	98.0

*Hinweis:* Mit den in `numpy` verfügbaren Funktionen `average` und `size` können Sie sehr effizient eine eigene Funktion `cov` schreiben. Zur Kontrolle können Sie per `sklearn` auch nochmal multiple lineare Regression durchführen. Hierbei müssen die  $x$ -Variablen als  $10 \times 2$ -Matrix übergeben werden.

**Aufgabe 2.8.** Sei  $D = \{(x_i, y_i) \in \mathbb{R}^2 \mid i = 1, \dots, n\}$  eine Datenmenge bei der nicht alle  $x_i$  gleich sind, und sei  $f: \mathbb{R} \rightarrow \mathbb{R}$  ihr linearer Regressor. Zeigen Sie, dass gilt

$$r_{xy}^2 = 1 - \frac{\sum_{i=1}^n (y_i - f(x_i))^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = \frac{\sum_{i=1}^n (f(x_i) - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}.$$

**Aufgabe 2.9.** Wir betrachten die Funktion  $g: [-100, 100] \rightarrow \mathbb{R}$ ,  $g(x) = 3x - 2$ . Erzeugen Sie, für variables  $n$ , in Python eine Datenmenge, deren Features  $n$ -viele gleichmäßig zufällig gewählte Punkte in  $[-10, 10]$  sind. Einem Feature  $x_i$  weisen Sie dann das Label  $y_i := g(x_i) + \varepsilon_i$  zu, wobei die  $\varepsilon_i$  unabhängige Samples der Normalverteilung  $\mathcal{N}(0, 1)$  sind. Führen Sie nun für  $n = 10, 20, 30, \dots$  für die so generierte Datenmenge  $D_n$  affin-lineare Regression durch, z.B. durch Verwendung des `sklearn`-Paketes.

**Aufgabe 2.10.** Finden Sie dasjenige Polynom  $P \in \mathbb{R}[X]$  mit  $\text{rk } P \leq 2$  welches die quadratischen Abstände zu den folgenden Daten minimiert.

Datenpunkt	1	2	3	4	5
$x$	0.0	1.0	-1.5	3.5	4.0
$y$	2.0	2.2	0.0	2.5	3.0

**Aufgabe 2.11.** Sei  $D = \{(x_i, y_i) \in \mathbb{R}^d \times \{0, 1\} \mid i = 1, \dots, n\}$  eine Datenmenge und  $\ell$  die zugehörige negative Log-Likelihood-Funktion aus (2.6). Zeigen Sie, dass

$$\nabla \ell(w) = \sum_{i=1}^n (-(\langle w, x_i \rangle) - y_i) \hat{x}_i$$

gilt. Zeigen Sie weiter, dass  $\nabla \ell$  für die Datenmenge  $D = \{(-1, 1), (0, 0), (1, 1)\}$  genau eine Nullstelle hat und berechnen Sie diese. Weil  $\ell$  konvex ist, muss diese Nullstelle  $w^* \in \mathbb{R}^2$  den logistischen Regressor liefern.

**Aufgabe 2.12.** Die folgende Datenmenge ist diejenige, die im Bild auf Seite 27 dargestellt ist. Bestimmen Sie mittels Python den logistischen Regressor und plotten Sie diesen.

$$D = \{(0.5, 0), (0.35, 0), (0.5, 0), (0.35, 0), (0.1, 0), (0.72, 0), (0.8, 0), (0.24, 0), (1.10, 0), (2.7, 0.97, 0), (1.5, 1), (1.9, 1), (1.65, 1), (1.35, 1), (1.7, 1), (1.24, 1), (1.09, 1), (0.92, 1)\}$$

**Aufgabe 2.13.** Um Personen dazu zu bewegen, an einer Umfrage teilzunehmen, werden diesen unterschiedliche Geldbeträge als Aufwandsentschädigung angeboten. In der folgenden Tabelle sind die angebotenen Beträge notiert und es ist jeweils vermerkt, wieviele Person das Angebot angenommen haben.

EUR	0.50	1	2	3	5	10	15	25	50
Annahmequote	1/43	2/50	4/48	5/32	30/37	15/32	55/100	49/50	19/20

- (i) Erstellen Sie anhand der Tabelle eine Datenmenge  $D$  mit kategoriellen Labels und ermitteln Sie für diese den logistischen Regressor.
- (ii) Betrachten Sie nun die Datenmenge

$$\hat{D} := \{(x_i, \log \frac{y_i}{1-y_i}) \mid i = 1, \dots, n\}$$

bei der die Label  $y_i$  die unveränderten Quotienten aus der Tabelle sind. Ermitteln Sie für  $\hat{D}$  den affin-linearen Regressor mit Parametern  $(a, b) \in \mathbb{R}^2$ . Vergleichen Sie jetzt die logistische Funktion

$$f: \mathbb{R} \rightarrow (0, 1), \quad f(x) = \frac{1}{1 + e^{-(ax+b)}}$$

mit dem logistischen Regressor aus (i).

**Aufgabe 2.14.** Sei  $(\Omega, \Sigma)$  ein Messraum und  $X: \Omega \rightarrow \mathbb{R}$  eine Zufallsvariable mit endlichem Erwartungswert  $\mu$  und endlicher Varianz  $\sigma$ . Seien  $X_1, \dots, X_n$  unabhängige Kopien von  $X$ . Dann definieren

$$\mu_s := \frac{1}{n} \sum_{j=1}^n X_j \quad \text{und} \quad \sigma_s^2 := \frac{1}{n} \sum_{j=1}^n (X_j - \mu_s)^2.$$

neue Zufallsvariablen. Zeigen Sie, dass für den Erwartungswert der letzteren  $E(\sigma_s^2) = \frac{n-1}{n} \sigma^2$  gilt und vergleichen Sie mit der Bemerkung nach Definition 2.5.

## Kapitel 3

# Der $k$ -NN Algorithmus

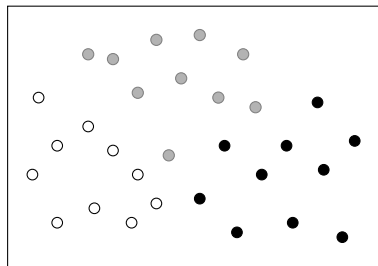
Am Ende des vorhergehenden Kapitels 2 haben wir mit der logistischen Regression eine Methode diskutiert, die einer binär gelabelten Datenmenge einen Klassifizierer zuordnet. Diese Methode hat *beweisbare* Eigenschaften, wie z.B. dass sie einer Maximum Likelihood Heuristik gehorcht oder dass sie bei überlappenden Daten genau einen Klassifizierer liefert. In späteren Kapiteln 13–14 werden wir noch sehen, dass für binär gelabelte und linear trennbare Datenmengen auch eine rigorose Theorie entwickelt werden kann.

In diesem Kapitel wollen wir praktisch gar keine Voraussetzungen an die Datenmengen machen und Methoden vorführen, mit denen Klassifizierer und Regressoren ohne viele Hilfsmittel und Annahmen gewonnen werden können. Dies hat den Vorteil der universellen Einsetzbarkeit und bietet die Möglichkeit, das hier Präsentierte später mit raffinierteren Methoden zu vergleichen. Der Nachteil besteht natürlich darin, dass wir keine tiefergehende beweisbare Theorie erwarten können.

### 3.1 $k$ -NN Klassifizierer

Wir beginnen mit einer gelabelten endlichen Datenmenge  $D \subseteq X \times Y$  wobei  $X$  ein metrischer Raum ist und  $Y$  eine endliche Menge. Ist  $(x, y) \in D$  so nennen wir wie üblich  $x$  das Feature und  $y$  das (nun kategorielle) Label des Datenpunktes.

**Beispiel 3.1.** Im folgenden Bild ist  $X \subseteq \mathbb{R}^2$  ein Quader und  $Y = \{1, 2, 3\}$ .



Wir statten  $X$  mit der euklidischen Metrik aus und verstehen weiße Punkte als mit Label 1, graue Punkte als mit Label 2 und schwarze Punkte als mit Label 3 versehen.

Unser Ziel ist es, eine Funktion  $f: X \rightarrow Y$  zu finden, die ‘möglichst gut zu den Daten passt’. Im einfachsten Sinne kann letzteres heißen, dass  $f(x) = y$  für alle, oder zumindest für möglichst viele, der Datenpunkte  $(x, y) \in D$  gelten soll. Hat man dann Punkte  $x \in X \setminus D$  gegeben, deren Label entweder unbekannt sind oder die bisher nicht mit einem Label versehen wurden, so kann man via  $y := f(x)$  ein Label *vorhersagen* oder *zuweisen*. Eine sehr naheliegende Idee, sich eine solche Funktion  $f$  zu verschaffen, ist die folgende.

**Definition 3.2.** Sei  $(X, \rho)$  ein metrischer Raum,  $Y$  eine endliche Menge,  $D \subseteq X \times Y$  eine Datenmenge und  $k \in \mathbb{N}$ . Wir definieren eine Funktion  $f: X \rightarrow Y$  wie folgt. Für gegebenes  $x \in X$  wählen wir zuerst  $x_1, \dots, x_k \in D$  aus mit

$$x_1 \in \operatorname{argmin}_{z \in D} \rho(x, z) \quad \text{sowie} \quad x_j \in \operatorname{argmin}_{z \in D \setminus \{x_1, \dots, x_{j-1}\}} \rho(x, z) \quad \text{für } j \geq 2$$

und nennen diese Punkte die  $k$ -*nächsten Nachbarn* von  $x$ . Seien  $y_1, \dots, y_k$  die Label der  $x_1, \dots, x_k$ . Wir betrachten  $N: Y \rightarrow \mathbb{N}$ ,  $N(y) := \#\{i \mid y_i = y\}$ . Schließlich wählen wir den Funktionswert wie folgt aus

$$f(x) \in \operatorname{argmax}_{i=1, \dots, n} N(y)$$

und nennen  $f: X \rightarrow Y$  den  $k$ -*NN Klassifizierer mit Mehrheitswahl*.

In Prosa lässt sich Definition 3.2 sehr viel weniger technisch beschreiben: Zu  $x \in X$  nimmt man den am nächsten an  $x$  gelegenen gelabelten Punkt, dann den zweit-nächsten usw. Dann schaut man, welches Label unter den  $k$ -nächsten Nachbarn am häufigsten vertreten ist und weist dieses dem Punkt  $x$  zu. Haben hierbei mehrere Punkte aus  $D$  den gleichen Abstand zu  $x$ , oder tritt am Ende unter den Nachbarn ein Gleichstand mehrerer Label auf, so trifft man eine beliebige Wahl.

Die Funktion  $f$  kann durch den folgenden Algorithmus bestimmt werden.

**Algorithmus 3.3.** Sei  $(X, \rho)$  ein metrischer Raum,  $Y$  eine endliche Menge,  $D \subseteq X \times Y$  eine Datenmenge und  $k \in \mathbb{N}$ . Der folgende Pseudocode stellt den  $k$ -NN Algorithmus mit Mehrheitswahl dar.

```

1: function K-NN KLASSIFIZIERER( $D, k, x$ )
2:    $D' \leftarrow D$ ,  $N \leftarrow \emptyset$ 
3:   for  $j \leftarrow 1$  to  $k$  do
4:      $z^* \leftarrow \operatorname{argmin}_{z \in D'} \rho(x, z)$ 
5:      $A \leftarrow A \cup \{z^*\}$ ,  $D' \leftarrow D' \setminus \{z^*\}$ 
6:   for  $y$  in  $Y$  do
7:      $N_y \leftarrow \#\{a \in A \mid \pi_2(a) = y\}$ 
8:    $\ell \leftarrow \operatorname{argmax}_{y \in Y} N_y$ 
9:   return  $\ell$ 

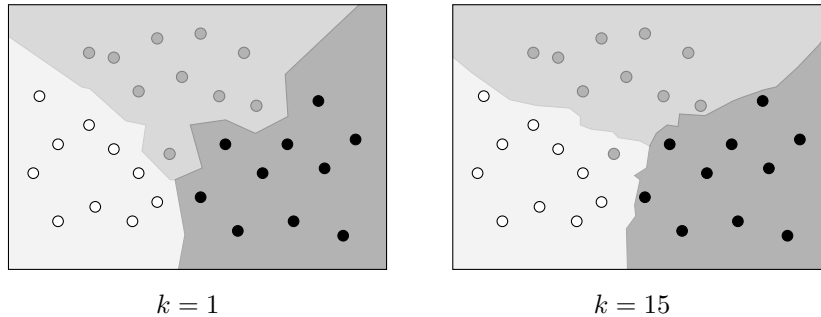
```

Hierbei bezeichnet  $\pi_2(x) = y$  die Projektion auf den zweiten Eintrag von  $(x, y) \in$

*D. Außerdem haben wir das explizite Durchnummerieren der Menge der  $k$ -nächsten Nachbarn durch die Verwendung der Menge  $A$  ersetzt.*

Unsere Definition 3.2 garantiert, dass  $f$  eine Abbildung ist. Implementiert man das obige, so erreicht man dies praktisch automatisch, wenn man z.B. durch  $D$  immer in einer vorgegebene Reihenfolge durchgeht und im Fall eines nicht einelementigen Argmins den als erstes gefundenen Punkt auswählt. Entsprechend kann man beim Argmax für die Mehrheitswahl vorgehen.

**Beispiel 3.4.** Wir kommen nun zur Datenmenge in Beispiel 3.1 zurück und visualisieren hier den  $k$ -NN Klassifizierer für zwei verschiedene  $k$ 's. Dabei geben die eingefärbten Bereiche jeweils den Wert der Funktion  $f$  auf diesen Bereichen an.



Beim Vergleich der Bilder fällt auf, dass die Linien, an denen die Farben umschlagen, für  $k = 15$  weniger zackig sind als für  $k = 1$ . Insbesondere sieht man bei  $k = 1$  eine Ausbuchtung um den grauen Punkt mit niedrigster Ordinate. Bei  $k = 15$  ist diese nicht mehr da, weil die weißen und schwarzen Punkte um den vorgenannten grauen Punkt herum, diesen bei der Mehrheitswahl überstimmen. Der Preis für die Entfernung der Ausbuchtungen ist allerdings, dass für den vorgenannten grauen Punkt  $(x, y)$  mit niedrigster Ordinate nun  $f(x) \neq y$  gilt.

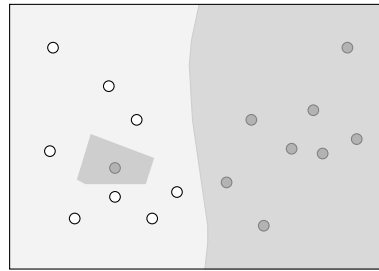
**Definition 3.5.** Sei  $(X, \rho)$  ein metrischer Raum,  $Y$  eine endliche Menge,  $D \subseteq X \times Y$  eine Datenmenge,  $k \in \mathbb{N}$  und  $f: X \rightarrow Y$  ein  $k$ -NN Klassifizierer.

- (i) Die Mengen  $X_y := \{x \in X \mid f(x) = y\}$  für  $y \in Y$  heißen *Entscheidungsbereiche*.
- (ii) Die Ränder  $\partial X_y$  heißen *Entscheidungsgrenzen*.

In Beispiel 3.4 hatten wir beobachtet, dass die Erhöhung von  $k$  einerseits die Entscheidungsgrenzen glättet und Ausbuchtungen entfernt, andererseits zu ‘Missklassifikationen’ führt. Auf den ersten Blick scheint die Missklassifikation das deutlich größere Übel zu sein und man ist gewillt den in Beispiel 3.4 links dargestellten Klassifizierer dem rechten vorzuziehen. Unser Ziel ist es zwar nicht unbedingt, den Leser vom Gegenteil zu überzeugen, aber doch eine Lanze für den rechten Klassifizierer zu brechen. Wir betrachten dazu ein weiteres Beispiel.

**Beispiel 3.6.** Sei wieder  $X \subseteq \mathbb{R}^2$  ein Quader und diesmal  $Y = \{1, 2\}$ . Wir betrachten eine binär gelabelte Menge wie im folgenden Bild. Dargestellt ist unten außerdem der  $k$ -NN Klassifizierer für  $k = 1$ . Wir überlassen es dem Leser sich davon zu überzeugen, dass für  $k \geq 3$  die graue Insel auf der linken Seite verschwindet.



 $k = 1$ 

Geht man davon aus, dass es eine echte unterliegende Funktion  $g: X \rightarrow Y$  gibt, die man durch Auswahl einer Stichprobe  $D$  via  $k$ -NN zu approximieren sucht, so kann es natürlich sein, dass  $g$  genau so aussieht wie im obigen Bild — es kann aber auch sein, dass bei der Bestimmung des ‘echten’ Labels des Punktes auf der Insel ein Fehler unterlaufen ist und dieser eigentlich weiß sein sollte! Ist letzteres der Fall, so passt der 1-NN Algorithmus die Funktion  $f$  zu *gut* an die Datenmenge an, anstatt von den Daten zu *verallgemeinern*.

Das folgende ist kein formal definierter mathematischer Begriff, aber von sehr großer Bedeutung im Kontext maschinellen Lernens. Aus diesem Grund gestehen wir ihm eine nummerierte Definition zu.

**Definition 3.7.** Sei  $(X, \rho)$  ein metrischer Raum,  $Y$  eine endliche Menge,  $D \subseteq X \times Y$  eine Datenmenge und  $k \in \mathbb{N}$ . Ist  $f: X \rightarrow Y$  ‘zu gut’ an die Datenmenge  $D$  angepasst, so spricht man von *Overfitting*.

In Beispielen kann man mit der folgenden Heuristik überprüfen, ob Overfitting vorliegt: Ist die Datenmenge  $D$  gegeben, so partitioniert man diese in *Trainingsdaten*  $D_1$  und *Testdaten*  $D_2$ . Dann bestimmt man einen Klassifizierer  $f: X \rightarrow Y$  anhand von  $D_1$ , und stellt fest, welcher Anteil der Punkte aus  $D_2$  durch  $f$  korrekt klassifiziert wird. Ist dieser Anteil eher klein, so kann dies auf Overfitting hindeuten.

Als nächstes diskutieren wir die Frage, wie aufwendig, oder ‘teuer’, die algorithmische Berechnung des  $k$ -NN Klassifizierers ist, wenn wir Daten mit Featurevektoren in  $\mathbb{R}^d$  gegeben haben und die euklidische Metrik benutzen. Dies drücken wir aus, indem wir die Anzahl Multiplikationen zählen, die bei einem Durchlauf des Algorithmus 3.3 anfallen.

**Satz 3.8.** Sei  $D \subseteq \mathbb{R}^d \times Y$  mit  $\#D = n$  und  $\#Y < \infty$ , sowie  $k \in \mathbb{N}$ . Sei  $\mathbb{R}^d$  mit der euklidischen Metrik ausgestattet und  $x \in \mathbb{R}^d$  fest. Die Berechnung der  $k$ -nächsten Nachbarn von  $x$  kann so implementiert werden, dass dabei höchstens  $(C \cdot n \cdot d \cdot k)$ -viele Multiplikationen anfallen, wobei  $C \in \mathbb{N}$  eine von  $d$ ,  $k$  und  $n$  unabhängige Konstante ist.

*Beweis.* Der  $k$ -NN Algorithmus berechnet  $n + (n - 1) + \dots + (n - k + 1)$ -mal eine Distanz von Punkten in  $\mathbb{R}^d$ . In der euklidischen Metrik müssen für eine solche Distanz

$d$ -viele Multiplikationen ausgeführt werden, wenn wir für die Berechnung des Argmin die Wurzel weglassen. Dies führt auf

$$d \cdot (n + (n - 1) + \cdots + (n - k + 1)) \leq C \cdot d \cdot k \cdot n$$

Multiplikationen mit einem geeigneten  $C \in \mathbb{N}$ . □

Man kann das obige Resultat noch verbessern, indem man die Abstände aller Datenpunkte  $z$  mit  $(z, y) \in D$  zum festen  $x$  nur einmal berechnet und abspeichert. Entscheidend ist aber, dass auch dann die Dimension  $d$  multiplikativ eingeht. D.h. bei hochdimensionalen Daten ist selbst bei nicht so großer Kardinalität von  $D$  die Laufzeit des  $k$ -NN Algorithmus' lang. Wir kommen in den Kapiteln 8–10 auf dieses Problem zurück.

**Bemerkung 3.9.** Einfache Mehrheitswahl ist nur eine Möglichkeit um  $f(x)$  anhand der  $k$ -nächsten Nachbarn von  $x \in X$  festzulegen. Man kann hier auch die tatsächlichen Abstände einfließen lassen im Sinne, dass ein weiter weg liegender Nachbar weniger Einfluss auf die Klassifizierung hat als ein näherer. Beispielsweise kann man in Definition 3.2 die Funktion  $N(y) := \#\{i \mid y_i = y\}$  durch

$$N(y) = \sum_{\substack{i=1 \\ y_i=y}}^k \frac{c}{1 + \rho(x_i, x)} \quad \text{oder} \quad N(y) = \sum_{\substack{i=1 \\ y_i=y}}^k e^{-c\rho(x_i, x)^2}$$

ersetzen, wobei man mit  $c > 0$  kontrolliert, wie schnell mit wachsendem Abstand der Beitrag kleiner wird. Letzteres nennt man *Gaußsche Gewichtung*.

## 3.2 $k$ -NN Regressoren

Sei  $D \subseteq X \times Y$  eine endliche Datenmenge, wobei  $(X, \rho)$  ein metrischer Raum und  $Y$  ein  $\mathbb{R}$ -Vektorraum ist. Für  $k \in \mathbb{N}$  und  $x \in X$  wählen wir  $k$ -nächste Nachbarn  $x_1, \dots, x_n$  von  $x$  und bezeichnen mit  $y_1, \dots, y_n$  deren Label. Setzt man dann

$$f: X \rightarrow Y, \quad f(x) = \frac{1}{k} \sum_{i=1}^k y_i,$$

so nennen wir  $f$  den  *$k$ -NN Regressor mit arithmetischem Mittel*.

Ähnlich zu Bemerkung 3.9 kann man auch hier wieder den Abstand einfließen lassen, und zum Beispiel

$$f: X \rightarrow Y, \quad f(x) = \frac{\sum_{i=1}^n w(x_i, x) \cdot y_i}{\sum_{i=1}^n w(x_i, x)}$$

verwenden mit  $w(z, x) = \frac{c}{1 + \rho(z, x)}$ ,  $w(z, x) = e^{-c\rho(z, x)^2}$  oder anderen Gewichtsfunktionen, die echt positiv sind und die mit fallendem Abstand wachsen.

### 3.3 Preprocessing

Wir betrachten im folgenden Datenmengen  $D \subseteq \mathbb{R}^d \times Y$  und konzentrieren uns auf die Menge der Featurevektoren, die wir für den Moment mit

$$F = \{x^{(1)}, \dots, x^{(n)}\} \subseteq \mathbb{R}^d$$

bezeichnen. Benutzen wir zur Berechnung eines  $k$ -NN Prediktors im Sinne der obigen Ausführungen die euklidische Metrik, oder allgemeiner z.B. die von  $\|\cdot\|_p$  mit  $p > 0$  induzierte Metrik, so werden alle Koordinaten gleich behandelt. Bei der Datenmenge ist es aber gut möglich, dass die einzelnen Features ganz unterschiedliche Größenordnungen haben. Kennt man hier ein bestimmtes Muster, so kann man ‘ad hoc’ gegensteuern, z.B. durch Verwendung einer gewichteten Metrik

$$d: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}, \quad \rho_{p,w}(x, y) := \sum_{i=1}^d w_i |x_i - y_i|^p$$

für  $p > 0$  und mit geeignet gewählten  $w_1, \dots, w_d > 0$ .

Hat man allerdings keine Kenntnis über die Bedeutung der einzelnen Features, so liegt es nahe zu versuchen, diese alle erst einmal möglichst gleich zu behandeln. Hierzu kann man es mit den folgenden Verfahren versuchen.

**Definition 3.10.** Sei  $F = \{x^{(1)}, \dots, x^{(n)}\} \subseteq \mathbb{R}^d$  gegeben, seien  $a < b$  reelle Zahlen und  $\|\cdot\|$  eine Norm auf  $\mathbb{R}^n$ . Wir bezeichnen mit  $\bar{x}_j^{(\cdot)} = \frac{1}{n}(x_j^{(1)} + \dots + x_j^{(n)})$  den Mittelwert der  $j$ -ten Einträge über alle Datenpunkte und  $\sigma_j^2 = \text{var}(x_j^{(\cdot)})$  die Varianz im selben Sinne. Dann heißt die Menge  $\tilde{F} := \{\tilde{x}^{(1)}, \dots, \tilde{x}^{(n)}\}$  mit

- (i)  $\tilde{x}^{(i)} = \left( \frac{(x_1^{(i)} - \min_{j=1, \dots, n} x_1^{(j)})(b-a)}{\max_{j=1, \dots, n} x_1^{(j)} - \min_{j=1, \dots, n} x_1^{(j)}}, \dots \right)$  die *Minmax-Normalisierung* auf  $[a, b]$ .
- (ii)  $\tilde{x}^{(i)} = \left( x_1^{(i)} - \frac{1}{n} \sum_{j=1}^n x_1^{(j)}, \dots \right)$  die *Zentrierung* von  $F$ .
- (iii)  $\tilde{x}^{(i)} = \left( \frac{x_1^{(i)} - \bar{x}_1}{\sigma_1}, \dots \right)$  die *Standardisierung* von  $F$ .
- (iv)  $\tilde{x}^{(i)} = \left( \frac{x_1^{(i)}}{\|(x_1^{(1)}, \dots, x_1^{(n)})\|}, \dots \right)$  die *Normierung* von  $F$  bezüglich  $\|\cdot\|$ .

Die obigen Verfahren können natürlich auch kombiniert werden. Zum Beispiel ist (iii) eine Kombination von (ii) und (iv) wenn wir die 2-Norm auf  $\mathbb{R}^n$  nehmen. Wir könnten aber auch erst zentrieren und dann mit  $\|\cdot\|_\infty$  normalisieren. Dann würde jeder Eintrag in  $[-1, 1]$  liegen und der koordinatenweise Mittelwert wäre Null. Beachte aber, dass dies etwas anderes liefert als Min-max-Normalisierung auf  $[-1, 1]$ .

Wir kommen in Aufgabe 3.5 insbesondere auf die Normalisierung nochmal zurück. Neben der ‘Homogenisierung’ der Features einer gegebenen Datenmenge, hat insbesondere das Zentrieren noch eine weitere wichtige Anwendung. Wir betrachten dazu Daten mit Featurevektoren aus  $\mathbb{R}^d$  und Labels aus  $\mathbb{R}$ . Im Unterschied zu vorher

nehmen wir jetzt aber an, dass uns kein Datenpunkt

$$(x, y) = (x_1, \dots, x_d, y)$$

vollständig bekannt ist, sondern stattdessen eine Datenmenge gegeben ist, bei der jeweils nur *ein Teil* der obigen  $(d+1)$ -vielen Koordinaten bekannt ist. Ein klassisches Beispiel ist das folgende.

**Beispiel 3.11.** Gegeben sei eine Tabelle mit Produktbewertungen, in der weiße Zellen bedeuten, dass uns an dieser Stelle die Bewertung unbekannt ist.

	Produkt 1	Produkt 2	Produkt 3	Produkt 4	Produkt 5	Produkt 6	Produkt 7
Kunde 1	★★★★☆	★★★★★	★★★★★	★☆☆☆☆		★★☆☆☆	
Kunde 2		★☆☆☆☆	★★☆☆☆	★★★★★	★★★★★		★☆☆☆☆
Kunde 3	★★☆☆☆	★★★☆☆	★★★☆☆	★★★☆☆		★★★☆☆	
Kunde 4	★★★★★			★★★☆☆	★★★☆☆	★★★★★	★★★★☆
Kunde 5	★★★★★	★★★★☆		★☆☆☆☆	★★☆☆☆		★★★★☆

Eine natürliche Aufgabe ist dann z.B. die Vorhersage der Bewertung für Produkt 7 durch Kunde 1. Dazu ist man gewillt, die Kunden als die Datenpunkte in  $\mathbb{R}^6$  zu betrachten, gegeben durch die ihre Bewertungen für die Produkte 1–6, und den Eintrag unter Produkt 7 als Label. Auf dieser Basis könnte man dann einen Prediktor  $f: \mathbb{R}^6 \rightarrow \{1, \dots, 5\}$ , z.B. via des  $k$ -NN Algorithmus 3.3, bestimmen — wenn alle Einträge außer dem bei Kunde 1 und Produkt 7 vorhanden wären.

Um in Situationen wie in Beispiel 3.11 trotz fehlender Koordinaten und Label Vorhersagen machen zu können, müssen die unbekannten Einträge zunächst irgendwie gefüllt werden. Eine Möglichkeit ist es dabei, die Featurevektoren zuerst zu zentrieren und dann die leeren Einträge mit Nullen zu füllen. Anstatt dies abstrakt zu notieren, führen wir es anhand der obigen Bewertungsmatrix vor.

**Beispiel 3.12.** (i) Wir betrachten die Tabelle aus Beispiel 3.11 aufgeteilt in Feature und Label wie dort beschrieben. Zentrieren der Features und Auffüllen mit Nullen führt auf die Datenmenge  $D = \{(x^{(i)}, y^{(i)}) \mid i = 2, 4, 5\}$  und ungelabelte Punkte  $x^{(1)}$ ,  $x^{(3)}$  wie in der folgenden Tabelle angegeben.

	Feature 1	Feature 2	Feature 3	Feature 4	Feature 5	Feature 6	Label
$(x^{(1)}, y^{(1)})$	0.60	1.60	1.60	−2.40	0.00	−1.40	
$(x^{(2)}, y^{(2)})$	0.00	−2.25	−1.25	1.75	1.75	0.00	1
$(x^{(3)}, y^{(3)})$	−0.80	0.20	0.20	0.20	0.00	0.20	
$(x^{(4)}, y^{(4)})$	1.25	0.00	0.00	−1.75	−0.75	1.25	4
$(x^{(5)}, y^{(5)})$	2.00	1.00	0.00	−2.00	−1.00	0.0	4

Verwendet man die euklidische Metrik  $\rho$  auf  $\mathbb{R}^6$  so ergibt sich

$$\rho(x^{(1)}, x^{(2)}) = 6.749, \quad \rho(x^{(1)}, x^{(4)}) = 3.361 \quad \text{und} \quad \rho(x^{(1)}, x^{(5)}) = 2.828,$$

woran sich die  $k$ -nächsten Nachbarn ablesen lassen. Wir fordern den Leser auf, nochmal die ursprüngliche Tabelle im Licht der berechneten Abstände anzuschauen! Für

die Bewertung von Produkt 7 durch Kunde 1 ergeben sich für jedes  $k = 1, 2, 3$  und bei Mehrheitswahl jeweils vier Sterne. Verwendet man das arithmetische Mittel, gefolgt von Rundung, so erhält man bei  $k = 1, 2$  wieder vier Sterne, aber bei  $k = 3$  jetzt drei Sterne.

(ii) Die obige Methode führt dazu, dass die Informationen über Kunde 3 vollkommen ungenutzt bleiben, da für diesen die Bewertung für Produkt 7 unbekannt ist. Will man diese Information auch noch einbauen, kann man über die gesamten Zeilen zentrieren, anstatt nur über deren erste sechs Einträge. Dies liefert dann  $D = \{(x^{(i)}, y^{(i)}) \mid i = 1, \dots, 5\}$  wie folgt.

	Feature 1	Feature 2	Feature 3	Feature 4	Feature 5	Feature 6	Label
$(x^{(1)}, y^{(1)})$	0.60	1.60	1.60	-2.40	0.00	-1.40	0.00
$(x^{(2)}, y^{(2)})$	0.00	-1.80	-0.80	2.20	2.20	0.00	-1.80
$(x^{(3)}, y^{(3)})$	-0.80	0.20	0.20	0.20	0.00	0.20	0.00
$(x^{(4)}, y^{(4)})$	1.20	0.00	0.00	-1.80	-0.80	1.20	0.20
$(x^{(5)}, y^{(5)})$	1.80	0.80	0.00	-2.20	-1.20	0.00	0.80

Jetzt berechnet man die Abstände der Zeilen bezüglich euklidischer Metrik auf  $\mathbb{R}^7$  und erhält

$$\begin{aligned} \rho\left(\begin{bmatrix} x^{(1)} \\ y^{(1)} \end{bmatrix}, \begin{bmatrix} x^{(2)} \\ y^{(2)} \end{bmatrix}\right) &= 6.991, & \rho\left(\begin{bmatrix} x^{(1)} \\ y^{(1)} \end{bmatrix}, \begin{bmatrix} x^{(3)} \\ y^{(3)} \end{bmatrix}\right) &= 3.899, \\ \rho\left(\begin{bmatrix} x^{(1)} \\ y^{(1)} \end{bmatrix}, \begin{bmatrix} x^{(4)} \\ y^{(4)} \end{bmatrix}\right) &= 3.644, & \rho\left(\begin{bmatrix} x^{(1)} \\ y^{(1)} \end{bmatrix}, \begin{bmatrix} x^{(5)} \\ y^{(5)} \end{bmatrix}\right) &= 2.953, \end{aligned}$$

und die  $k$ -nächsten Nachbarn können wieder abgelesen werden. Für die Mehrheitswahl und auch für das arithmetische Mittel verwendet man dann aber die Originalbewertungen im Fall von Kunde 2, 4 und 5 und bei Kunde 3 macht man die Zentrierung rückgängig, d.h. man setzt hier  $(2 + 3 + 3 + 3 + 3)/5 \approx 3$  Sterne an. Für die Bewertung von Produkt 7 durch Kunde 1 ergeben sich damit für jedes  $k = 1, 2, 3$  und bei Mehrheitswahl wieder jeweils vier Sterne. Mit dem arithmetischen Mittel kommt man nun für  $k = 1, 2, 3$  ebenfalls jeweils auf vier Sterne.

**Bemerkung 3.13.** (i) Als erstes notieren wir, dass mit der einmal zentrierten Tabelle in Beispiel 3.12(ii) jetzt auch alle anderen vormals leeren Tabelleneinträge ausgefüllt werden könnten. Wir erhalten somit eine Methode, die es erlaubt eine ‘löchrige’ Datenmatrix zu füllen. In Kapitel 7.3 werden wir noch deutlich raffiniertere Techniken hierzu kennenlernen.

(ii) Die obigen Methoden zur Vorhersage sind kollaborativ im Sinne, dass die Vorhersage der Bewertung von Kunde 1 für Produkt 7 nicht nur auf den anderen Bewertungen von Kunde 1 basiert, sondern auf dem Vergleich von Bewertungen von Kunde 1 mit denen aller anderen Kunden. Man spricht von *User-to-User Collaborative Filtering*.

(iii) Anstatt die Kunden als Punkte in  $\mathbb{R}^7$  zu lesen, hätten wir auch die Produkte als Punkte in  $\mathbb{R}^5$  lesen können. Macht man dies, so ist eine spaltenweise Zentrierung angebracht, vergleiche Aufgabe 3.6(ii–iii).

(iv) Anstatt der euklidischen Metrik kann man auch andere Metriken, oder sogar

sogenannte Abstandsmaße, die keine Metrik sind, verwenden. Das behandeln wir im folgenden Unterkapitel.

### 3.4 Kosinusähnlichkeit

In den Kapiteln 3.1–3.3 haben wir stets mit einem zugrundeliegenden metrischen Raum  $(X, \rho)$  gearbeitet. In diesem Kapitel wollen wir dies verallgemeinern. Als erstes bemerken wir, dass sowohl  $k$ -NN Klassifizierer als  $k$ -NN Regressoren im Sinne der vorhergehenden Kapitel auch dann noch definiert werden können, wenn  $\rho$  nicht definit ist und auch nicht die Dreiecksungleichung erfüllt.

**Definition 3.14.** Sei  $X$  eine nichtleere Menge. Eine Funktion  $\rho: X \times X \rightarrow \mathbb{R}$  heißt *Abstandsmaß*, falls für alle  $x, y \in X$  gilt

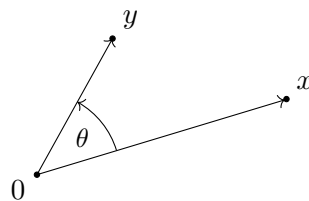
- (AM1)  $\rho(x, y) \geq 0$  und  $\rho(x, x) = 0$ , (Positivität)
- (AM2)  $\rho(x, x) = 0$ , (‘schwache’ Definitheit)
- (AM3)  $\rho(x, y) = \rho(y, x)$ . (Symmetrie)

Für Abstandsmaße gilt weiterhin die Heuristik, dass kleine Werte bedeuten, dass Punkte ähnliche Eigenschaften haben. Bei der unten definierten Kosinusähnlichkeit ist die Sache anders gelagert.

**Definition 3.15.** Sei  $\mathbb{R}^d$  ausgestattet mit dem Standardskalarprodukt und der davon induzierten euklidischen Norm. Für  $x, y \in \mathbb{R}^d \setminus \{0\}$  heißt dann

$$\text{cossim}(x, y) := \frac{\langle x, y \rangle}{\|x\| \|y\|} = \frac{\sum_{i=1}^d x_i y_i}{\left(\sum_{i=1}^d x_i^2\right)^{1/2} \left(\sum_{i=1}^d y_i^2\right)^{1/2}} \in [-1, 1]$$

die *Kosinusähnlichkeit* von  $x$  und  $y$ . Bezeichnet man mit  $\theta = \angle(x, y)$  den Winkel



zwischen  $x$  und  $y$ , so gilt  $\langle x, y \rangle = \|x\| \|y\| \cos(\theta)$  und es folgt  $\text{cossim}(x, y) = \cos(\theta)$ , was die Bezeichnung ‘Kosinusähnlichkeit’ begründet.

Im Vergleich zu einem Abstandsmaß ist die Heuristik bei der Kosinusähnlichkeit, dass  $x$  und  $y$  ähnlich sind, wenn  $\text{cossim}(x, y)$  möglichst groß ist. Ist  $\text{cossim}(x, y)$  nah bei Null, so sind  $x$  und  $y$  eher unähnlich. Der Fall negativer Kosinusähnlichkeit drückt etwas aus, was durch ein Abstandsmaß gar nicht beschrieben werden kann.

Schränken wir uns auf eine Teilmenge von  $\mathbb{R}^d \setminus \{0\}$  ein, so induziert die Kosinusähnlichkeit in folgender Weise ein Abstandsmaß.

**Lemma 3.16.** Sei  $\mathbb{R}_{\geq 0}^d := \{(x_1, \dots, x_d) \in \mathbb{R}^d \setminus \{0\} \mid x_i \geq 0 \text{ für alle } i = 1, \dots, d\}$ . Dann ist  $\text{cosdist}: \mathbb{R}_{\geq 0}^d \times \mathbb{R}_{\geq 0}^d \rightarrow \mathbb{R}$ ,  $\text{cosdist}(x, y) := 1 - \text{cossim}(x, y)$  ein Abstandsmaß, welches wir als Kosinusabstand bezeichnen.

*Beweis.* Alle drei zu prüfenden Eigenschaften folgen direkt aus Eigenschaften von (Standard-)Skalarprodukt und Norm.  $\square$

Wir überlassen es dem Leser als Aufgabe 3.7 zu verifizieren, dass  $\rho$  wie in Lemma 3.16 nicht positiv definit ist und auch nicht die Dreiecksungleichung erfüllt. Ein gutes Beispiel für das obige stammt aus dem Bereich des Textminings.

**Beispiel 3.17.** Wir betrachten die folgenden drei Texte.

Am 22. November 1963 wurde der amerikanische Präsident John Fitzgerald Kennedy in Dallas ermordet. Kennedy war zu diesem Zeitpunkt seit zwei Jahren im Amt. Der Präsident wurde auf einer Fahrt durch Dallas in einem offenen Wagen mit einem Gewehr erschossen. Seine Frau Jackie Kennedy befand sich ebenfalls im Wagen, blieb aber unverletzt. John Fitzgerald Kennedy wurde am 25. November 1963 in Arlington beigesetzt.

Am 22. November 1963 ermordete Lee Harvey Oswald den damaligen US Präsident John Fitzgerald Kennedy in Dallas. Oswald erschoss Kennedy mit einem Gewehr, während dieser zusammen mit seiner Frau Jackie Kennedy in einem offenen Wagen durch Dallas fuhr. Oswald wurde am 24. November 1963 bei seiner Überstellung ins Bezirksgefängnis von Dallas durch den Nachtclubbesitzer Jack Ruby in einer Tiefgarage erschossen.

Am John Fitzgerald Kennedy Airport in Dallas gibt es eine große Auswahl an Restaurants. Besonders zu empfehlen sind das Porterhouse Steak sowie das T-bone Steak jeweils serviert mit gegrillten Tomaten im Restaurant Sive Steak. Aber auch Vegetarier kommen im John Fitzgerald Kennedy Airport auf ihre Kosten: Sive Steak bietet zum Beispiel einen Black Been Burger an. Im Jikji Cafe gibt es vegetarische Gerichte aus koreanischer Küche.

Wir ordnen diesen nun auf die folgende Weise drei Vektoren in  $\mathbb{R}_{\geq 0}^d$  zu. Zunächst betrachten wir die Menge aller in allen drei Texten vorkommenden Worte

$$W := \{\text{Am, 22, November, 1963, wurde, der, } \dots\},$$

setzen  $d := \#W$  und wählen eine Bijektion  $\{1, \dots, d\} \rightarrow W$ . Via dieser Bijektion können wir Vektoren in  $\mathbb{R}^d$  als  $(x_w)_{w \in W}$  schreiben. Insbesondere können wir für die Texte 1, 2 und 3 von oben deren *Vektorisierungen*

$$x^{(i)} = (x_w^{(i)})_{w \in W} \text{ per } x_w^{(i)} := \text{Anzahl der Vorkommnisse des Wortes } w \text{ in Text } i$$

für  $i = 1, 2, 3$  definieren. Per Konstruktion gilt dann  $x^{(1)}, x^{(2)}, x^{(3)} \in \mathbb{R}_{\geq 0}^d$ . Im obigen Beispiel ist  $d = 115$  und man kann die paarweisen Kosinusähnlichkeiten

$$\text{cossim}(x^{(1)}, x^{(2)}) = 0.61, \text{ cossim}(x^{(1)}, x^{(3)}) = 0.28, \text{ cossim}(x^{(2)}, x^{(3)}) = 0.19$$

mit einem einfachen Programm berechnen, vergleiche Aufgabe 3.8. Wir sehen also, dass sich die Texte 1 und 2 deutlich ähnlicher sind als 1 und 2 bzw. 1 und 3. Verwendet man die Kosinusabstände

$$\text{cosdist}(x^{(1)}, x^{(2)}) = 0.39, \text{ cosdist}(x^{(1)}, x^{(3)}) = 0.72, \text{ cosdist}(x^{(2)}, x^{(3)}) = 0.81$$

und sind beispielsweise die von einem Leser vergebenen binären Label ‘👍’ und ‘👎’ für eine geeignete Anzahl Texte bekannt, so kann man mithilfe des  $k$ -NN Algorithmus 3.3 einen Prediktor bestimmen, der für ungelabelte Texte vorhersagt, ob dieser Leser sie als gut oder nicht gut bewerten wird.

**Bemerkung 3.18.** (i) Über den oben benutzten *raw count* hinausgehend, gibt es weitere Vektorisierungsmöglichkeiten, die z.B. kompensieren, dass der raw count aller Wörter mit der Textlänge wächst, dass Wörter wie ‘der’, ‘die’, ‘das’ usw. in jedem Text oft vorkommen, oder dass oben ‘ermordet’ und ‘ermordete’ als verschieden behandelt werden. Einige solcher Methoden diskutieren wir in Aufgabe 3.8.

(ii) Textmining ist ein gutes Beispiel dafür, dass man sehr natürlich und sehr schnell zu hochdimensionalen Daten kommt. Die drei Vektorisierungen unserer kurzen Beispieltexthe sind bereits Elemente von  $\mathbb{R}^{115}$ .

(iii) Anstelle des Kosinusabstandes könnte man beim Textvergleich auch die euklidische Metrik verwenden. Im Allgemeinen ist der Kosinusabstand hier aber besser geeignet. Dies sieht man ein, wenn man ein Wort betrachtet, dass in einem Text bereits oft vorkommt. Der euklidische Abstand zu einem anderen Text fällt dann minimal aus, wenn dieses Wort in ihm exakt genauso oft vorkommt wie im ersten Text. Der Kosinusabstand wird hingegen immer kleiner, je öfter das Wort im zweiten Text auftritt. Wir verweisen auch auf das künstliche, aber instruktive, Beispiel in Aufgabe 3.9.

Man kann die Kosinusähnlichkeit auch für Bewertungsvorhersagen einsetzen und erhält dann natürlich andere Ergebnisse als mit der vorher benutzten euklidischen Metrik.

**Beispiel 3.19.** Wir betrachten die Produktbewertungen aus Beispiel 3.11 und stellen unten die originale Bewertungstabelle und die zeilenweise zentrierte und dann mit Nullen aufgefüllte Tabelle gegenüber. Die Einträge links bezeichnen wir mit  $x_j^{(i)}$ , die rechts mit  $\tilde{x}_j^{(i)}$ , verzichten jetzt also darauf, eine Koordinate als Label auszuzeichnen.

	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$	$x_7$		$\tilde{x}_1$	$\tilde{x}_2$	$\tilde{x}_3$	$\tilde{x}_4$	$\tilde{x}_5$	$\tilde{x}_6$	$\tilde{x}_7$
$x^{(1)}$	4	5	5	1		2		$\tilde{x}^{(1)}$	0.6	1.6	1.6	-2.4	0.0	-1.4	0.0
$x^{(2)}$		1	2	5	5		1	$\tilde{x}^{(2)}$	0.0	-1.8	-0.8	2.2	2.2	0.0	-1.8
$x^{(3)}$	2	3	3	3		3		$\tilde{x}^{(3)}$	-0.8	0.2	0.2	0.2	0.0	0.2	0.0
$x^{(4)}$	5			2	3	5	4	$\tilde{x}^{(4)}$	1.2	0.0	0.0	-1.8	-0.8	1.2	0.2
$x^{(5)}$	5	4		1	2		4	$\tilde{x}^{(5)}$	1.8	0.8	0.0	-2.2	-1.2	0.0	0.8

Die Kosinusähnlichkeiten von  $\tilde{x}^{(1)}$  zu den restlichen zentrierten Datenpunkten lauten

$$\text{cossim}(\tilde{x}^{(1)}, \tilde{x}^{(2)}) = -0.634, \text{cossim}(\tilde{x}^{(1)}, \tilde{x}^{(3)}) = -0.185,$$

$$\text{cossim}(\tilde{x}^{(1)}, \tilde{x}^{(4)}) = 0.355, \text{cossim}(\tilde{x}^{(1)}, \tilde{x}^{(5)}) = 0.640$$

und wir fordern den Leser auf, die Ähnlichkeiten mit den Abständen in Beispiel 3.12(ii) zu vergleichen. Eine naheliegende Idee für Vorhersagen ist es nun, die echten



Bewertungen mit diesen Kosinusähnlichkeiten zu gewichten, sehr ähnlich zum  $k$ -NN Regressor in Unterkapitel 3.2. Als Beispiel nehmen wir wieder die Bewertung von Kunde 1 für Produkt 7, und verwenden die 2-kosinusähnlichsten Kunden, oder mit anderen Worten die 2-nächsten Kunden bzgl. Kosinusabstand. Dies führt auf

$$\hat{x}_7^{(1)} = \frac{\text{cossim}(\tilde{x}^{(1)}, \tilde{x}^{(5)}) \cdot x_7^{(5)} + \text{cossim}(\tilde{x}^{(1)}, \tilde{x}^{(4)}) \cdot x_7^{(4)}}{\text{cossim}(\tilde{x}^{(1)}, \tilde{x}^{(5)}) + \text{cossim}(\tilde{x}^{(1)}, \tilde{x}^{(4)})} = \frac{0.640 \cdot 4 + 0.355 \cdot 4}{0.640 + 0.355} = 4.$$

$\uparrow$   
 Vorher-  
 sage

Verwenden wir allerdings die 3-ähnlichsten Kunden bei denen eine Bewertung für Produkt 7 vorliegt, so erhalten wir

$$\hat{x}_7^{(1)} = \frac{0.640 \cdot 4 + 0.355 \cdot 4 - 0.634 \cdot 1}{0.640 + 0.355 - 0.634} \approx 9.24.$$

$\uparrow$   
 Vorher-  
 sage

Dass sich durch die negative Kosinusähnlichkeit von  $x^{(1)}$  und  $x^{(2)}$  und die niedrige Bewertung  $x_7^{(1)}$  eine *Erhöhung* der Vorhersage ergibt, scheint eventuell sinnvoll, andererseits schießt diese hier deutlich über das Ziel hinaus. Bevor dieses letzte Ergebnis den Leser zu sehr abschreckt, weisen wir darauf hin, dass bei einer großen Datenmenge und nicht zu groß gewähltem  $k$  zu erwarten ist, dass für jede Vorhersage stets nur Kunden mit jeweils großer Ähnlichkeit in die Formel eingehen, und wir daher eher selten in der letzten beschriebenen Situation sein werden.

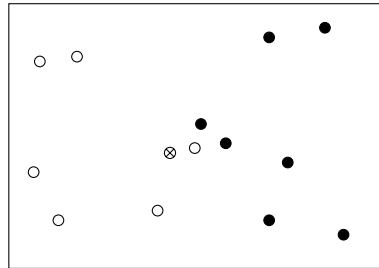
Das obige Beispiel zeigt, dass die möglicherweise negativen Werte der Kosinusähnlichkeit mit Vorsicht zu genießen sind. Andererseits bieten Sie auch in bestimmten Fällen die Möglichkeit nicht nur zwischen ‘ähnlich’ und ‘nicht ähnlich’ zu unterscheiden, sondern der Skala noch so etwas wie ‘entgegengesetzt’ hinzuzufügen. Wir erwähnen noch, dass manche Autoren in der Vorhersageformel des Beispiels 3.19 durch die Summe der Beträge der Kosinusähnlichkeiten teilen. Macht man das oben, so *verringert* sich allerdings die Vorhersage.

## Referenzen

Durch seine Einfachheit und universelle Einsetzbarkeit ist der  $k$ -NN Algorithmus sehr populär und wird in vielen Büchern und Vorlesungen über Data Science oder Machine Learning behandelt. Gleiches gilt für die behandelten Preprocessingmethoden. Unsere Hauptquelle in diesem Kapitel ist [LRU12]. Mehr Details zum kollaborativen Filtern findet man in [SKKR01] und [Agg16].

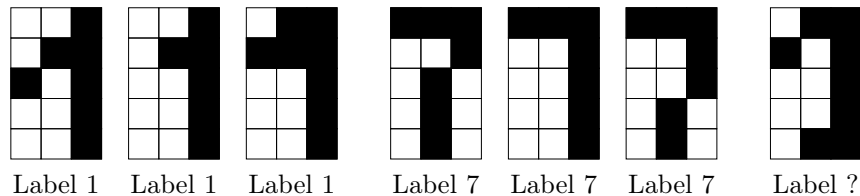
## Aufgaben

**Aufgabe 3.1.** Wir betrachten die folgende Datenmenge  $D \subseteq \mathbb{R}^2 \times \{1, 2\}$  wobei weiß gezeichnete Punkte Label 1 haben und schwarz gezeichnete Label 2. Das Label des mit einem Kreuz markierten Punktes ist uns unbekannt.

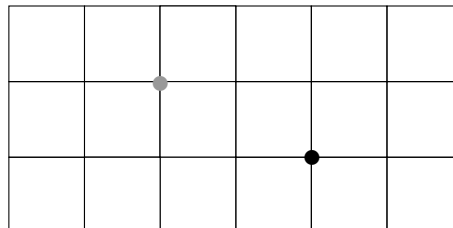


Bestimmen Sie (Augenmaß genügt!) für den mit einem Kreuz markierten Punkt das Label gemäß  $k$ -NN mit Mehrheitswahl für  $k = 1, \dots, 7$ .

**Aufgabe 3.2.** Wir betrachten eine Datenmenge aus sechs Bildern mit angegebenen Labels. Klassifizieren Sie das neue Bild mithilfe von 1-NN mit Mehrheitswahl. Benutzen Sie als Abstand die von der  $\ell^1$ -Norm auf  $\mathbb{R}^{15}$  induzierte Metrik.



**Aufgabe 3.3.** Skizzieren Sie die Entscheidungsbereiche und Entscheidungsgrenzen bzgl. 1-NN, wenn wir die von der 2-Norm, der 1-Norm bzw. der  $\infty$ -Norm induzierte Metrik benutzen.



**Aufgabe 3.4.** Betrachten Sie nochmal das Dataset aus Aufgabe 2.3. Benutzen Sie 2-NN mit Mittelwertbildung. Skizzieren Sie den so definierten Prediktor und die Daten. Vergleichen Sie das Ergebnis mit dem aus Aufgabe 2.3. Finden Sie 2-NN hier sinnvoll?

**Aufgabe 3.5.** Wir betrachten das folgende 2-dimensionale Dataset mit Labels in  $\{-1, +1\}$ .

Datenpunkt	Feature 1	Feature 2	Label
1	111	0.3	+1
2	82	0.4	+1
3	110	0.4	+1
4	148	0.7	+1
5	91	0.5	+1
6	133	0.9	+1
7	71	1.2	-1
8	111	1.1	-1
9	123	1.6	-1
10	151	1.9	-1
11	89	1.3	-1
12	99	2.0	-1

- (i) Implementieren Sie  $k$ -NN mit euklidischer Metrik und Mehrheitwahl derart dass man  $k \in 2\mathbb{N} - 1$  und einen Punkt  $x \in \mathbb{R}^2$  wählen kann und dann ausgegeben wird, ob dieser Label  $+1$  oder  $-1$  bekommen soll.
- (ii) Fällt Ihnen irgendein Problem auf? Falls nicht, dann wenden Sie  $k$ -NN auf  $x = (112, 1.9)$  an und erzeugen Sie einen Plot bei dem Ordinate und Abszisse gleich skaliert sind.
- (iii) Minmax-normalisieren Sie die Daten auf  $[0, 1]$  und den Punkt  $x$  entsprechend. Führen Sie dann  $k$ -NN mit Mehrheitwahl erneut aus. Erläutern Sie das (überraschende?) Ergebnis.

**Aufgabe 3.6.** Wir betrachten die folgende Tabelle von Produktbewertungen, bei der jeweils eine Punktzahl zwischen 0 und 10 vergeben werden darf und bei der uns mehrere Bewertungen unbekannt sind.

	Produkt 1	Produkt 2	Produkt 3	Produkt 4
Kundin 1	10	7	<input type="text"/>	4
Kundin 2	<input type="text"/>	<input type="text"/>	5	3
Kundin 3	3	0	3	6
Kundin 4	7	9	<input type="text"/>	5
Kundin 5	10	9	8	<input type="text"/>

- (i) Welche Bewertungen für Produkt 4 durch Kundin 5 ergeben sich mit den in Beispiel 3.12 und welche mit der in Beispiel 3.19 erläuterten Methode?
- (ii) In (i) haben wir User-to-User Collaborative Filtering benutzt. Was ergibt sich, wenn stattdessen Item-to-Item Collaborative Filtering angewandt wird?
- (iii) Was sind Vor- und Nachteile der User-to-User bzw. der Item-to-Item Methode?

**Aufgabe 3.7.** Zeigen Sie, dass Kosinusdistanz auf  $\mathbb{R}_{>0}^2$  weder positiv definit ist noch die Dreiecksungleichung erfüllt.

**Aufgabe 3.8.** Vektorisieren die Texte aus Beispiel 3.17 in einer Programmiersprache Ihrer Wahl unter Benutzung geeigneter Pakete (in Python z.B. `Re`) und verifizieren Sie die dort notierten Kosinusähnlichkeiten. Berechnen Sie zusätzlich auch die euklidischen Abstände der Texte.

**Aufgabe 3.9.** Die folgenden drei Listen sind Vektorisierungen von Texten via raw count:

$$\begin{aligned}
 T1 &= [ 12, 0, 4, 1, 0, 1, 12, 10 ] \\
 T2 &= [ 11, 1, 0, 2, 0, 2, 11, 10 ] \\
 T3 &= [ 24, 0, 8, 2, 0, 2, 24, 20 ]
 \end{aligned}$$

Was ist der 1-NN von T1 wenn man nach euklidischem Abstand geht? Welcher Text ist am ähnlichsten zu T1 im Sinne der Kosinusähnlichkeit?

**Aufgabe 3.10.** Seien mehrere Textdokumente  $D_1, \dots, D_n$  wie in Beispiel 3.17 gegeben und sei  $W$  die durchnummerierte Menge aller Wörter. Ist  $D$  einer der Texte und  $w \in W$ , so bezeichnen wir mit  $\text{rc}(w, D)$  die Anzahl der Vorkommnisse des Wortes  $w$  im Text  $T$ . Wir definieren weiter die *Termfrequenz* und die *inverse Dokumentenfrequenz* per

$$\text{tf}(w, D) := \frac{\text{rc}(w, D)}{\max\{\text{rc}(v, D) \mid v \in D\}} \quad \text{und} \quad \text{idf}(w, D_1, \dots, D_n) := \log\left(\frac{n}{\#\{i \mid w \in D_i\}}\right).$$

Das Produkt aus Termfrequenz und inverser Dokumentenfrequenz bezeichnen wir mit

$$\text{tf-idf}(w, D) := \text{tf}(w, D) \cdot \text{idf}(w, D_1, \dots, D_m)$$

wobei  $D \in \{D_1, \dots, D_m\}$ . Erläutern Sie die folgenden Heuristiken zum Vergleich der Vektorisierung von  $D$  durch  $(\text{rc}(w, D))_{w \in W}$ ,  $(\text{tf}(w, D))_{w \in W}$  oder  $(\text{tf-idf}(w, D))_{w \in W}$ :

- (i) Im Vergleich zum Raw Count kompensiert die Termfrequenz, dass in längeren Dokumenten tendenziell alle Wörter entsprechend häufiger vorkommen als in kürzeren.
- (ii) Durch die Multiplikation mit der inversen Dokumentenfrequenz kann man ferner darauf reagieren, dass für den Textvergleich irrelevante Wörter wie ‘der’, ‘die’, ‘das’, ‘und’ usw. in praktisch jedem sehr langen Dokument sehr häufig vorkommen.

# Kapitel 4

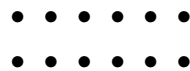
## Clustering

In den vorhergehenden Kapiteln haben wir uns mit Klassifizierern und Regressoren beschäftigt. Dabei haben wir uns stets eine gelabelten Datenmenge vorgegeben, auf deren Basis dann ein Prediktor für neue und ungelabelte Datenpunkte bestimmt wurde. In diesem Kapitel betrachten wir generell eine *ungelabelte Datenmenge*  $D \subseteq (X, \rho)$  bei der  $\rho: X \times X \rightarrow \mathbb{R}$  ein Abstandsmaß im Sinne von Definition 3.14 ist. Ziel ist es, eine Partition zu finden, d.h. Mengen  $C_1, \dots, C_n \subseteq D$  derart, dass gilt

$$D = C_1 \dot{\cup} \dots \dot{\cup} C_n$$

wobei idealerweise Datenpunkte  $x, y \in D$  für die  $\rho(x, y)$  klein ist, zur selben Menge  $C_i$  gehören sollen. Die  $C_i$  nennen wir dann die *Cluster*. Anstatt zu versuchen, das letztgesagte abstrakt zu formalisieren, schauen wir das folgende sehr illustative Beispiel (entnommen aus [SSBD14]) an.

**Beispiel 4.1.** Wir fassen die im nachfolgenden Bild dargestellte Menge  $D$  als Teilmenge von  $\mathbb{R}^2$  auf und statt  $\mathbb{R}^2$  mit der euklidischen Metrik  $\rho$  aus.



Die Menge  $D$  soll nun in zwei Cluster aufgeteilt werden, d.h.  $D = C_1 \dot{\cup} C_2$ , wobei Punkte, die nah beieinander liegen, zum selben Cluster gehören sollen. Man kommt hier schnell auf zwei naheliegende, aber doch ganz verschiedene Lösungen:



Da die vertikalen Abstände zwischen den Punkten etwas größer sind als die horizontalen, sieht man, dass im linken Bild für jeden Punkt gilt, dass jeweils sein 1-nächster Nachbar im selben Cluster liegt. Dies hat dann aber zur Folge, dass innerhalb des Clusters Abstände entstehen, die viel größer sind als die zwischen den Clustern. Rechts sind die Durchmesser der Cluster deutlich kleiner. Man bezahlt dies aber damit, dass es Punkte gibt, deren 1-nächster Nachbar nicht, aber der 2-nächste Nachbar sehr wohl, im selben Cluster liegt.

Beispiel 4.1 macht deutlich, dass die oben diskutierte Heuristik ‘Punkte die nah zusammen liegen, sollten zum selben Cluster gehören’ eingeschränkt werden muss. Wir stellen nun zwei verschiedene Möglichkeiten vor, dies zu tun, und die das obige Beispiel verallgemeinern.

## 4.1 Verknüpfungsbasiertes Clustering

Die erste Möglichkeit, welche dem linken Bild aus Beispiel 4.1 entspricht, ist das *verknüpfungsbasierte Clustering*. Hierfür erinnern wir zuerst an die folgende Definition.

**Definition 4.2.** Sei  $X$  eine Menge und  $\rho: X \times X \rightarrow \mathbb{R}$  ein Abstandsmaß. Für  $A, B \subseteq X$  definieren wir

$$\rho(A, B) = \min_{\substack{a \in A \\ b \in B}} \rho(a, b)$$

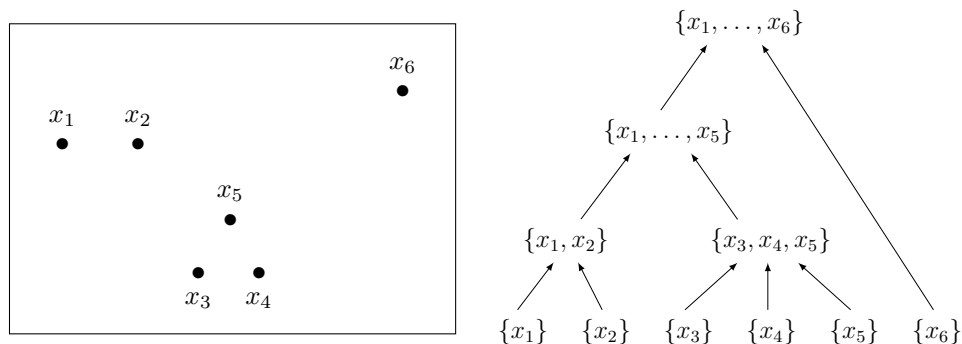
und nennen  $\rho(A, B)$  den *Abstand* von  $A$  und  $B$ .

Die Idee hinter verknüpfungsbasiertem Clustering ist nun sehr einfach. Ist in der Situation von Definition 4.2 eine Datenmenge

$$D = \{x_1, \dots, x_n\} \subseteq (X, \rho)$$

gegeben, so beginnen wir mit dem *diskreten Clustering*  $D = \{x_1\} \cup \dots \cup \{x_n\}$ . Von diesem ausgehend vereinigen wir schrittweise diejenigen Cluster mit minimalem Abstand. Obiges Verfahren liefert dann eine Folge von Clusterings, bis am Ende das *triviale Clustering*  $D = \{x_1, \dots, x_n\}$  herauskommt. Die Folge der Clusterings kann man durch ein *Dendrogramm* darstellen, was im folgenden Beispiel (sehr ähnlich zu [SSBD14]) illustriert wird.

**Beispiel 4.3.** Sei  $X \subseteq \mathbb{R}^d$  ein Quader,  $\rho$  die euklidische Metrik auf  $X$  und  $D$  die unten skizzierte Menge bestehend aus sechs Punkten.



Wir weisen darauf hin, dass im Dendrogramm auf der rechten Seite die Zeilen nicht den Runden des Algorithmus entsprechen; z.B. haben wir in der ersten Runde nur  $\{x_3, x_4, x_5\}$  zusammengelegt und  $\{x_1, x_2\}$  erst in der zweiten Runde.

Ist man nicht am gesamten Dendrogramm interessiert, so muss man den vorgenannten Prozess an geeigneter Stelle abbrechen. Hierfür kann man die Anzahl der ‘Runden’ im oben erläuterten Vorgehen begrenzen oder man kann die Anzahl der Cluster beschränken. Eine dritte Alternative ist die Betrachtung des Minimums der paarweisen Abstände zwischen den Clustern

$$\min_{i \neq j} \rho(C_i, C_j),$$

die sich in jeder Runde vergrößert, und für die man eine obere Schranke festlegen kann. Im folgenden Pseudocode geben wir diese Variante wieder und lassen es als Aufgabe 4.1 ein Abbruchkriterium über die Anzahl der Cluster in Pseudocode zu formulieren.

**Algorithmus 4.4.** Sei  $D \subseteq (X, \rho)$ . Der folgende Pseudocode stellt das verküpfungsbasierte Clustering dar, wobei abgebrochen wird, bevor der minimale Abstand zwischen den Clustern in der nächsten Runde erstmalig unter einen einzugebenden Wert  $\delta > 0$  fallen würde.

```

1: function VERKNÜPFUNGSBASIERTES CLUSTERING ( $X, \rho, D, \delta$ )
2:    $n \leftarrow \#D$ 
3:   for  $i \leftarrow 1$  to  $n$  do
4:      $C_i \leftarrow \{x_i\}$ 
5:   while  $\min_{i \neq j} \rho(C_i, C_j) \geq \delta$  do
6:      $m \leftarrow 0$ 
7:      $(i^*, j^*) \leftarrow \operatorname{argmin}_{i \neq j} \rho(C_i, C_j)$ 
8:     for  $k \leftarrow 1$  to  $n$  do
9:       if  $k = \min(i^*, j^*)$  then
10:         $C_k \leftarrow C_{i^*} \cup C_{j^*}$ 
11:       if  $k = \max(i^*, j^*)$  then
12:         $m \leftarrow 1$ 
13:       else
14:         $C_k \leftarrow C_{k+m}$ 
15:      $n \leftarrow n - 1$ 
16:   return  $C_1, \dots, C_k$ 

```

Ist das Argmin in Zeile 7 nicht einelementig, so wählen wir dort einen beliebigen Minimierer  $(i^*, j^*) \in \operatorname{argmin}_{i \neq j} \rho(C_i, C_j)$  aus.

**Bemerkung 4.5.** Statt des üblichen Abstands von Teilmengen wie in Definition 4.3 kann man die folgenden Funktionen benutzen, was dann natürlich zu anderen Clustering führt.

- (i)  $\rho_1(A, B) := \frac{1}{|A||B|} \sum_{\substack{a \in A \\ b \in B}} \rho(a, b)$  (Mittelwert aller paarweisen Abstände)
- (ii)  $\rho_2(A, B) := \max_{\substack{a \in A \\ b \in B}} \rho(a, b)$  (Maximum aller paarweisen Abstände)

Wir kommen hierauf in Aufgabe 4.1 zurück.

## 4.2 Kostenminimierendes Clustering

Im Gegensatz zum verknüpfungsbasierten Clustering betrachten wir nun eine *Kostenfunktion* oder auch *Zielfunktion* auf der Menge aller möglichen Clusterings. Dabei ordnen wir Zerlegungen  $X = C_1 \cup \dots \cup C_k$  mit ungewünschten Eigenschaften (wie z.B. zu großem Durchmesser) hohe Kosten zu und kommen dann durch einen Minimierungsprozess auf Clusterings, die die gewünschten Eigenschaften haben.

Für eine Datenmenge  $D$  und  $k \in \mathbb{N}$  bezeichnen wir im folgenden mit

$$\mathcal{C}_k := \{(C_1, \dots, C_k) \in \mathcal{P}(D)^k \mid D = C_1 \dot{\cup} \dots \dot{\cup} C_k\}$$

die Menge der  $k$ -Clusterings von  $D$ . Ist weiter  $D \subseteq X$  und  $\rho: X \times X \rightarrow \mathbb{R}$  eine Abstandsmaß und  $Z: \mathcal{C}_k \rightarrow \mathbb{R}$  eine vorgegebene *Kostenfunktion*, so nennen wir einen Minimierer  $(C_1, \dots, C_k) \in \operatorname{argmin}_{C \in \mathcal{C}_k} Z(C)$  ein bezüglich  $Z$  *kostenminimierendes Clustering* von  $D$ . Das populärste Beispiel einer solchen Kostenfunktion ist das folgende.

**Definition 4.6.** Sei  $X$  eine Menge und  $\rho$  ein Abstandsmaß auf  $X$  derart, dass für jede endliche Menge  $A \subseteq X$  ein *Schwerpunkt*

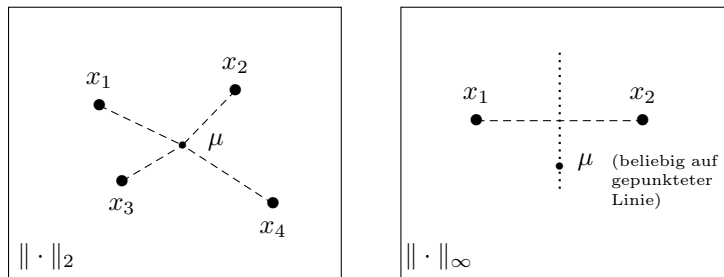
$$\mu(A) \in \operatorname{argmin}_{\mu \in X} \sum_{x \in A} \rho(x, \mu)^2$$

in  $X$  existiert. Sei  $D \subseteq X$  eine Datenmenge, sei  $k \in \mathbb{N}$  und sei  $\mathcal{C}_k$  die Menge aller  $k$ -Clusterings von  $D$ . Für  $k \in \mathbb{N}$  heißt

$$K: \mathcal{C} \rightarrow \mathbb{R}, \quad K(C_1, \dots, C_k) := \sum_{i=1}^k \sum_{x \in C_i} \rho(x, \mu(C_i))^2$$

die *k-means-Kostenfunktion* auf  $(X, \rho)$ .

**Bemerkung 4.7.** (i) Ist  $X = \mathbb{R}^d$  ausgestattet mit der euklidischen Metrik, so existiert  $\mu(A)$  für jedes endliche  $A \subseteq \mathbb{R}^d$ , vergleiche Bemerkung 4.10. Stattet man  $X = \mathbb{R}^d$  mit der von der  $\|\cdot\|_\infty$  induzierten Norm aus, so existieren ebenfalls stets Schwerpunkte, sie sind aber im Allgemeinen nicht eindeutig.



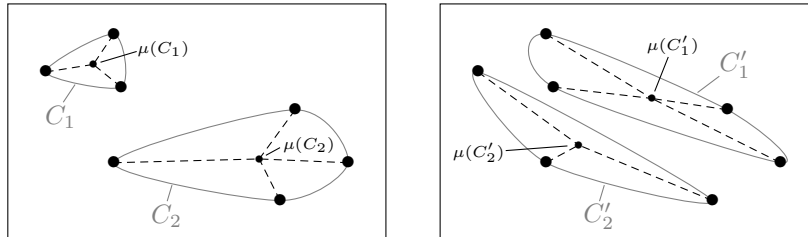
(ii) In der Situation von Definition 4.6 gilt stets

$$K(C_1, \dots, C_k) = \min_{\mu_1, \dots, \mu_k \in X} \sum_{i=1}^k \sum_{x \in C_i} \rho(x, \mu_i)^2.$$



Hierbei ist ' $\geq$ ' klar und für ' $\leq$ ' vertauscht man zuerst das Minimum und die Summe über  $i$  und sieht dann, dass links gerade ein Minimierer eingesetzt wurde.

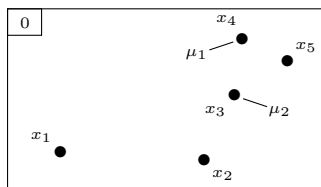
(iii) Im folgenden Bild sehen wir für sieben Punkte jeweils zwei verschiedene Clusterings  $(C_1, C_2)$  bzw.  $(C'_1, C'_2)$  und die jeweiligen Schwerpunkte bzgl. euklidischer Metrik.



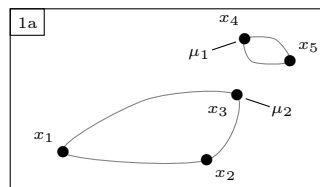
Die 2-means-Kostenfunktion  $K$  berechnet die Summe der quadratischen Abstände der Datenpunkte vom jeweiligen Schwerpunkt des Clusters. Im Beispiel sehen wir sofort, dass  $K(C_1, C_2) < K(C'_1, C'_2)$  gilt.

Um einen Minimierer für  $K$  zu finden, verwenden wir den *k-means-Algorithmus*. Wir wählen zu Beginn  $\mu_1, \dots, \mu_k \in D$  beliebig, aber paarweise verschieden, aus. Dann definieren wir die Cluster, indem wir einen Datenpunkt  $x \in D$  dem Cluster  $C_i$  zuordnen, wenn  $x$  am nächsten an  $\mu_i$  liegt; falls dies für mehrere  $i$  gilt, wählen wir unter diesen aus. Dann berechnen wir die Schwerpunkte der Cluster und wiederholen den Zuordnungsschritt.

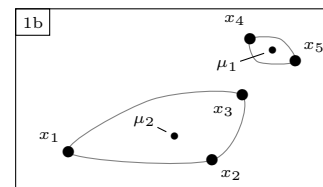
**Beispiel 4.8.** Wir beginnen mit dem folgenden Beispiel einer Datenmenge mit fünf Elementen im  $\mathbb{R}^2$  mit euklidischer Metrik und  $k = 2$ .



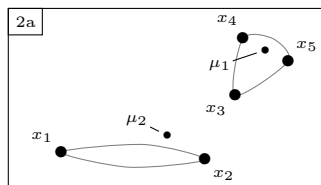
Zu Beginn werden  $\mu_1 \neq \mu_2$  beliebig aus  $\{x_1, \dots, x_5\}$  gewählt; im Beispiel oben  $\mu_1 = x_4$  und  $\mu_2 = x_3$ .



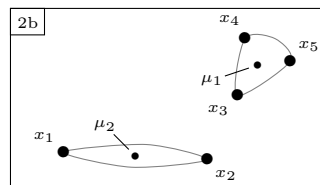
Dann werden diejenigen Punkte, die näher an  $\mu_1$  liegen als an  $\mu_2$ , dem Cluster  $C_1$  hinzugefügt und diejenigen, die näher an  $\mu_2$  liegen als an  $\mu_1$ , dem Cluster  $C_2$ . Es ergibt sich  $C_1 = \{x_4, x_5\}$  und  $C_2 = \{x_1, x_2, x_3\}$ .



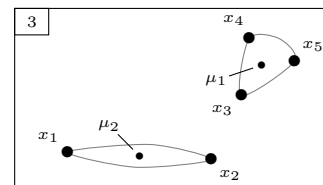
Als Nächstes werden die Schwerpunkte von  $C_1$  und  $C_2$  berechnet und die  $\mu_1, \mu_2$  entsprechend upgedatet.



Nun werden die Cluster neu eingeteilt, und zwar nach demselben Prinzip wie vorher, aber bezüglich der neuen Schwerpunkte. Damit ergibt sich als Update für die Cluster  $C_1 = \{x_3, x_4, x_5\}$  und  $C_2 = \{x_1, x_2\}$ .



Jetzt werden wieder die Schwerpunkte von  $C_1$  und  $C_2$  berechnet und  $\mu_1, \mu_2$  ein weiteres Mal upgedatet.



Schließlich wäre wieder die Neueinteilung der Cluster dran, aber man sieht, dass sich diese nicht mehr ändern und somit auch die Schwerpunkte ab hier konstant bleiben. Daher bricht man den Prozess hier ab und nimmt das letzte Clustering  $C_1 = \{x_3, x_4, x_5\}$  und  $C_2 = \{x_1, x_2\}$  als Ergebnis.

Wir formulieren nun den Algorithmus in Pseudocode.

**Algorithmus 4.9.** Sei  $X$  eine Menge und  $\rho$  ein Abstandsmaß auf  $X$  derart, dass für jede endliche Menge  $A \subseteq X$  ein Schwerpunkt  $\mu(A) \in X$  existiert. Sei  $D \subseteq X$  eine Datenmenge und sei  $k \in \mathbb{N}$ . Der folgende Pseudocode liefert einen Minimierer der  $k$ -means-Kostenfunktion.

```

1: function K-MEANS ( $\mathcal{X}, d, X, k$ )
2:    $\mu_1, \dots, \mu_k \leftarrow$  paarweise verschieden aus  $X$  wählen
3:   for  $i \leftarrow 1$  to  $k$  do
4:      $C_i \leftarrow \{x \in X \mid i \in \operatorname{argmin}_{j=1, \dots, k} \rho(x, \mu_j)\}$ 
5:    $U \leftarrow \text{True}$ 
6:   while  $U = \text{True}$  do
7:      $U \leftarrow \text{False}$ 
8:     for  $i \leftarrow 1$  to  $k$  do
9:        $C'_i \leftarrow \{x \in X \mid i \in \operatorname{argmin}_{j=1, \dots, k} \rho(x, \mu_j)\}$ 
10:       $\mu_i \leftarrow \mu(C_i)$ 
11:      if  $C'_i \neq C_i$  then
12:         $C_i \leftarrow C'_i$ 
13:         $U \leftarrow \text{True}$ 
14:   return  $C_1, \dots, C_k$ 

```

**Bemerkung 4.10.** (i) Ist  $X = \mathbb{R}^d$  und  $\rho$  die euklidische Metrik, dann existiert der Schwerpunkt  $\mu(A)$  jeder endlichen Menge  $A \subseteq \mathcal{X}$ , er eindeutig bestimmt und gleich dem Mittelwert

$$\mu(A) = \frac{1}{|A|} \sum_{a \in A} a,$$

siehe Aufgabe 4.5. Der  $k$ -means-Algorithmus berechnet also in diesem Fall  $k$ -viele Mittelwerte. Statt der Schwerpunkte kann man wahlweise auch den *Medoid* oder den *Geometrischer Median*

$$\mu(A) \in \operatorname{argmin}_{\mu \in X} \sum_{x \in A} \rho(x, \mu)^2 \quad \text{bzw.} \quad \mu(A) \in \operatorname{argmin}_{\mu \in \mathcal{X}} \sum_{x \in A} \rho(x, \mu)$$

verwenden. Die entsprechenden Algorithmen nennt man dann  $k$ -medoid- bzw.  $k$ -median-Algorithmus.

(ii) Kostenfunktionen müssen nicht auf Schwerpunkten, oder irgendeiner Version von Mittelwerten, basieren; z.B. kann man auch verlangen, dass die Summe aller paarweisen Abstände zwischen Punkten eines Clusters minimiert wird. Dies entspräche der Kostenfunktion

$$Z(C_1, \dots, C_k) = \sum_{i=1}^k \sum_{x, y \in C_i} \rho(x, y).$$

In der Formulierung von Algorithmus 4.9 ist a priori nicht klar, dass dieser terminiert, d.h. dass irgendwann kein Update der Cluster mehr erfolgt und die While-Schleife in Zeile 6–13 verlassen wird. Im allgemeinen kann dies auch nicht erwartet

werden, vergleiche Bemerkung 4.12. Es gilt aber immerhin das folgende.

**Satz 4.11.** *Sei  $X$  eine Menge und  $\rho$  ein Abstandsmaß auf  $X$  derart, dass für jede endliche Menge  $A \subseteq X$  ein Schwerpunkt  $\mu(A) \in X$  existiert. Sei  $D \subseteq X$  eine Datenmenge und sei  $k \in \mathbb{N}$ . Dann ist die Folge  $(K(C_1^{(j)}, \dots, C_k^{(j)})_{j \in \mathbb{N}}$  der Auswertungen der  $k$ -means Kostenfunktion auf den vom  $k$ -means-Algorithmus produzierten Clusterings monoton fallend.*

*Beweis.* Seien  $\mu_1^{(0)}, \dots, \mu_k^{(0)}$  die zu Beginn festgelegten Punkte. Für  $j \geq 1$  bezeichnen wir dann mit  $(C_1^{(j)}, \dots, C_k^{(j)})$  die Cluster und mit  $\mu_1^{(j)} = \mu(C_1^{(j)}), \dots, \mu_k^{(j)} = \mu(C_k^{(j)})$  deren Schwerpunkte in der  $j$ -ten Runde des Algorithmus. Für  $j \geq 2$  gilt dann

$$\begin{aligned} K(C_1^{(j)}, \dots, C_k^{(j)}) &= \min_{\mu_1, \dots, \mu_k \in X} \sum_{i=1}^k \sum_{x \in C_i^{(j)}} \rho(x, \mu_i)^2 \\ &\leq \sum_{i=1}^k \sum_{x \in C_i^{(j)}} \rho(x, \mu(C_i^{(j-1)}))^2 \\ &\leq \sum_{i=1}^k \sum_{x \in C_i^{(j-1)}} \rho(x, \mu(C_i^{(j-1)}))^2 \\ &= K(C_1^{(j-1)}, \dots, C_k^{(j-1)}) \end{aligned}$$

wobei wir für die Gleichungen die Definition 4.6 der  $k$ -means-Kostenfunktion und Bemerkung 4.7(ii) verwendet haben. Die erste Ungleichung ist lediglich eine Spezialisierung. Für die zweite beachten wir, dass in Zeile 9 von Algorithmus 4.9 in der  $j$ -ten Runde  $C_i^{(j)}$  gerade so definiert wird, dass  $\rho(x, \mu(C_i^{(j-1)}))$  für  $x \in C_i^{(j)}$  minimal ist. Damit kann aber die Summe über obige Abstände bei jeder anderen Clusterzuordnung höchstens größer werden.

Wir raten dem Leser, sich zur Veranschaulichung des letzten Arguments nochmal das Beispiel 4.8 für  $j = 2$  anzuschauen: Dort entspricht Bild 2b der vorletzten Zeile in der Abschätzung, und Bild 1b der darüber. Wir haben dann

$$\sum_{j=3}^5 \rho(x_j, \mu_1)^2 + \sum_{j=1}^2 \rho(x_j, \mu_2)^2 \leq \sum_{j=4}^5 \rho(x_j, \mu_1)^2 + \sum_{j=1}^3 \rho(x_j, \mu_2)^2,$$

weil wir ja gerade wegen  $\rho(x_3, \mu_1) < \rho(x_3, \mu_2)$  den Punkt  $x_3$  vom Cluster  $C_2$  in Bild 1b zu Cluster  $C_2$  in Bild 2a verschoben haben.  $\square$

**Bemerkung 4.12.** (i) Satz 4.11 sagt uns, dass wir durch längeres Laufenlassen von  $k$ -means nichts kaputt machen können, gibt aber keine Garantie, dass hierdurch die Folge der Clusterings irgendwann stationär wird, siehe Aufgabe 4.7. Um dies zu adressieren, kann man zusätzlich zur Abbruchbedingung in Algorithmus 4.9 eine maximale Anzahl an Iterationen für die while-Schleife festlegen.

(ii) Es kann passieren, dass die Folge der Clusterings  $C_1^{(j)}, \dots, C_k^{(j)}$  ab einem  $j_0$  stationär wird, ohne dass  $C_1^{(j_0)}, \dots, C_k^{(j_0)}$  ein Minimierer für die  $k$ -means-Zielfunktion

ist, vergleiche Aufgabe 4.6.

(iii) Bei ungünstiger Wahl der Startwerte kann der Fall eintreten, dass ein oder mehrere  $C_i$  leer werden. Daher empfiehlt es sich, bei einer gegebenen Datenmenge den  $k$ -means-Algorithmus, mit verschiedenen zufällig gewählten Startwerten, jeweils mehrfach laufen lassen und dann die Ergebnisse vergleichen.

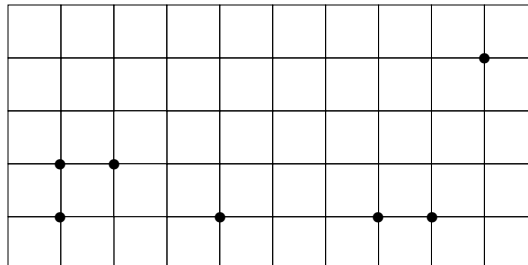
## Referenzen

Die in diesem Kapitel behandelten Clusteringalgorithmen sind Standard. Wir folgen relativ eng der hervorragenden Darstellung [SSBD14, Chapter 22.1–22.2], ergänzen aber einige Details und Beispiele. Insbesondere stimmt Aufgabe 4.6 mit [SSBD14, Exercise 22.2] überein.

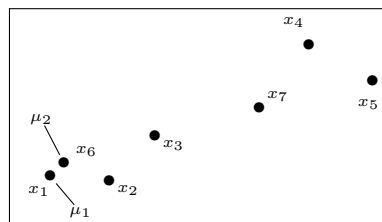
## Aufgaben

**Aufgabe 4.1.** Notieren Sie eine Version des verknüpfungsbasierten Clusterings in Pseudocode, bei dem man die Anzahl der Cluster vorwählen kann, und der dann ein Clustering mit genau der vorgegebenen Anzahl Cluster ausgibt. Diskutieren Sie hierbei möglicherweise eintretende Pathologien.

**Aufgabe 4.2.** Führen Sie für die folgende Datenmenge zweimal verknüpfungsbasiertes Clustering aus und zeichnen Sie die Dendrogramme. Nutzen Sie beide Male die euklidische Metrik  $\rho$  auf  $\mathbb{R}^2$  aber einmal den üblichen Abstand  $\rho(A, B) = \min_{a \in A, b \in B} \rho(a, b)$  für Teilmengen  $A, B$  der Datenmenge und einmal den alternativen Abstand  $\rho_2(A, B) = \max_{a \in A, b \in B} \rho(a, b)$ .



**Aufgabe 4.3.** Führen Sie auf der folgenden Datenmenge den 2-means-Algorithmus aus, wobei zu Beginn  $\mu_1 = x_1$  und  $\mu_2 = x_6$  gewählt sei und die euklidische Metrik auf  $\mathbb{R}^2$  zugrunde liege. Tragen Sie jeweils die Clusterings  $C_i$  und die sich ergebenden Schwerpunkte  $\mu_i$  (Augenmaß genügt!) ein.



**Aufgabe 4.4.** Implementieren Sie den  $k$ -means Algorithmus und clustern Sie damit die [hier](#) hinterlegte Datenmenge aus dem `sklearn`-Paket: Erstellen Sie in Ergänzung zur dortigen Liste `L` der echten Labels eine neue Liste `N` mit den per Algorithmus bestimmten Labels.

Vergleichen Sie den vorhandenen Plot mit einem neuen Plot in welchem die Farben der Datenpunkte durch  $N$  festgelegt werden. Probieren Sie dabei verschiedene  $k$ 's aus und achten Sie insbesondere auf den orangenen Punkt im lila Cluster.

**Aufgabe 4.5.** Sei  $A \subseteq \mathbb{R}^d$  und bezeichne  $\rho: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  die euklidische Metrik. Zeigen Sie, dass für jede endliche Menge  $A \subseteq \mathbb{R}^d$  gilt

$$\frac{1}{|A|} \sum_{a \in A} a = \operatorname{argmin}_{\mu \in \mathbb{R}^d} \sum_{a \in A} \rho(a, \mu)^2.$$

Insbesondere ist also hier der Schwerpunkt eindeutig und konsistent mit der üblicherweise in der Linearen Algebra gegebenen Definition.

**Aufgabe 4.6.** Gegeben sei die Datenmenge  $\{1, 2, 3, 4\} \subseteq \mathbb{R}$  auf das wir den 2-means-Algorithmus anwenden. Wir nehmen an, dass mit den Mittelpunkten  $\mu_1 = 2$  und  $\mu_2 = 4$  begonnen wird und dass, falls  $\operatorname{argmin}_{j=1,2} \|x - \mu_j\|$  nicht einelementig ist, stets  $i = 1$  als Minimierer gewählt wird. Zeigen Sie, dass die Folge der Clusterings  $(C_1^{(j)}, C_2^{(j)})_{j \in \mathbb{N}}$  stationär wird, ohne dass das stationäre Clustering die Zielfunktion minimiert.

**Aufgabe 4.7.** Sei

$$D = \{(-1, -1), (0, -1), (1, -1), (-1, 1), (0, 1), (1, 1), (0, -0.5), (0, 0.5)\} \subseteq \mathbb{R}^2.$$

Wir wenden den 2-means Algorithmus auf  $D$  an, aber diesmal mit zufälligen Startwerten  $\mu_1$  und  $\mu_2$  und zufälliger Auswahl von  $i \in \operatorname{argmin}_{j=1,2} \|x - \mu_j\|$  im Fall, dass letztere Menge zweielementig ist.

- (i) Zeigen Sie, dass nach endlich vielen Iterationen die Mittelwerte  $\mu_1 = (-0.5, 0)$  und  $\mu_2 = (0.5, 0)$  herauskommen *können*.
- (ii) Zeigen Sie, ausgehend von (i), dass es passieren *kann*, dass die Folge der Clusterings  $(C_1^{(j)}, C_2^{(j)})_{j \in \mathbb{N}}$  nicht stationär wird.
- (iii) Wie wahrscheinlich ist es, dass (ii) eintritt?

# Kapitel 5

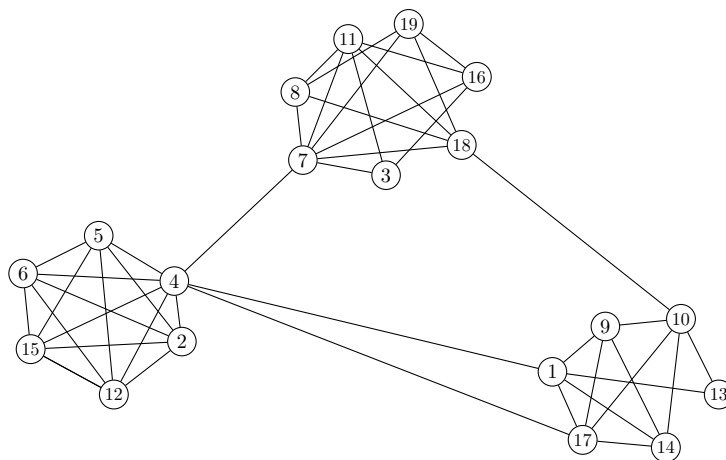
## Graphenclustering

In Kapitel 4 über Clustering waren uns Datenmengen  $D$  immer als Teilmengen eines Raumes  $(X, \rho)$  gegeben, wobei  $\rho$  Abstandsmaß war, und wir haben dabei stets angenommen, dass es möglich ist, den Abstand  $\rho(x, y)$  zweier Datenpunkten direkt, d.h., ohne Kenntnis der restlichen Datenpunkte zu berechnen. In diesem Kapitel beginnen wir völlig anders und zwar mit einer Datenmenge, die durch einem sogenannten *Graph* gegeben ist.

**Definition 5.1.** Ein *Graph*  $G = (V, E)$  ist ein Paar bestehend aus einer endlichen Menge  $V$ , genannt *Vertices* oder *Knoten*, und einer Teilmenge  $E \subseteq \{\{v, w\} \mid v, w \in V, v \neq w\}$  deren Elemente wir *Kanten* nennen.

Graphen wie in Definition 5.1 sind *ungerichtet*, haben keine *Schleifen* ( $\circlearrowleft$ ) und keine *mehrfachen Kanten* ( $\circ \! \! \! \circ$ ). Beachte, dass Graphen in der Literatur manchmal auch anders definiert werden.

**Beispiel 5.2.** Das folgende Beispiel zeigen einen Graphen mit neunzehn Vertices. Die Kanten sind durch Verbindungslinien angegeben.



Man sieht, dass von jedem Vertex in einer drei Teilmengen  $C_1 = \{2, 4, 5, 6, 12, 15\}$ ,  $C_2 = \{1, 9, 10, 13, 14, 17\}$  und  $C_3 = \{3, 7, 8, 11, 16, 18, 19\}$  weitaus weniger Kanten in

eine andere Teilmenge gehen als Kanten zu Punkten innerhalb derselben Teilmenge bleiben. Es ist daher sehr natürlich die  $C_1, C_2, C_3$  als *Cluster des Graphen* zu bezeichnen.

Natürlich auftretende Graphen sind zum Beispiel soziale Netzwerke, bei denen die Nutzer den Vertices entsprechen und die Kanten dem ‘befreundet sein’. Ein solcher sehr großer Graph kommt dann natürlich nicht als Zeichnung wie oben daher, sondern nur z.B. nur durch eine Liste, die alle Freundschaften angibt. Formal führen wir den folgenden Begriff ein.

**Definition 5.3.** Sei  $G = (V, E)$  ein Graph mit  $V = \{1, \dots, n\}$ . Dann heißt

$$A = (a_{ij})_{i,j \in \mathbb{N}} \text{ mit } a_{ij} = \begin{cases} 1 & \text{falls } \{i, j\} \in E, \\ 0 & \text{sonst} \end{cases}$$

die *Adjazenzmatrix* von  $G$  und

$$L = (\ell_{ij})_{i,j \in \mathbb{N}} \text{ mit } \ell_{ij} = \begin{cases} -1 & \text{falls } \{i, j\} \in E, \\ \deg(i) & \text{falls } i = j, \\ 0 & \text{sonst} \end{cases}$$

die *Laplacematrix* oder der *Laplacian* von  $G$ , wobei  $\deg(i) = \#\{e \in E \mid \exists j \in V: \{i, j\} = e\}$  der *Grad* von  $i \in V$  ist, d.h. die Anzahl der Kanten die *inzident* mit Vertex  $i$  sind (z.B. hat der Knoten  $\star$  Grad 5).

Wir notieren die folgenden einfachen Eigenschaften.

**Bemerkung 5.4.** (i) Die Matrizen  $A$  und  $L$  sind symmetrisch.

(ii) Die  $i$ -te Zeilensumme/Spaltensumme von  $A$  ist gleich  $\deg(i)$ .

(iii) Die  $i$ -te Zeilensumme/Spaltensumme von  $L$  ist gleich 0.

(iv) Es gilt  $L = D - A$  wobei  $D = \text{diag}(\deg(1), \dots, \deg(n))$ .

**Beispiel 5.5.** Als einfaches Beispiel betrachten wir den sogenannten vollständigen Graph mit drei Vertices

$$G = \begin{array}{c} \textcircled{1} \\ \diagdown \quad \diagup \\ \textcircled{2} \text{---} \textcircled{3} \end{array} \quad \text{mit} \quad A = \begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix} \quad \text{und} \quad L = \begin{bmatrix} 2 & -1 & -1 \\ -1 & 2 & -1 \\ -1 & -1 & 2 \end{bmatrix}.$$

Wir notieren die folgende Eigenschaft von Laplacematrizen.

**Proposition 5.6.** Für jeden Graph  $G$  ist  $\mathbf{1} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$  ein Eigenvektor des Laplacians mit Eigenwert Null.

*Beweis.*  $L\mathbf{1}$  hat als Einträge genau die Zeilensummen von  $L$ . Da diese alle Null sind, folgt die Gleichung  $L\mathbf{1} = 0\mathbf{1}$ .  $\square$

Aus der Linearen Algebra wissen wir, dass für eine symmetrische Matrix  $M \in \mathbb{R}^{n \times n}$  alle Eigenwerte reell sind und für jeden Eigenwert  $\lambda \in \sigma(M)$  die geometrische Vielfachheit ( $= \dim \ker(M - \lambda)$ ) gleich der algebraischen Vielfachheit ( $=$  maximale Potenz mit der man  $x - \lambda$  aus dem charakteristischen Polynom  $\det(M - x)$  herausdividieren kann). Mit Vielfachheiten haben wir Eigenwerte  $\lambda_1 \leq \dots \leq \lambda_n$  und können eine Orthonormalbasis aus Eigenvektoren  $v_1, \dots, v_n$  finden. Dies wenden wir jetzt auf die Laplacematrizen von Graphen an.

Wir beginnen mit zwei Beispielen.

**Beispiel 5.7.** Wir bleiben bei dem Graph aus Beispiel 5.5 und der dort notierten Laplacematrix. Mit Proposition 5.6 haben wir bereits den Eigenwert 0 mit zugehörigem Eigenvektor  $\mathbf{1}$ . Berechnung der restlichen Eigenwerte liefert mit Vielfachheiten  $\lambda_1 = 0$ ,  $\lambda_2 = 3$ ,  $\lambda_3 = 3$  und eine zugehörige Orthonormalbasis

$$v_1 = \frac{1}{\sqrt{3}} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}, \quad v_2 = \frac{1}{\sqrt{6}} \begin{bmatrix} 1 \\ 1 \\ -2 \end{bmatrix}, \quad v_3 = \frac{1}{\sqrt{2}} \begin{bmatrix} -1 \\ 1 \\ 0 \end{bmatrix}.$$

**Beispiel 5.8.** Als nächstes betrachten wir einen Graphen mit sechs Knoten, bei dem es zwischen  $\{1, 2, 3\}$  und  $\{4, 5, 6\}$  keine Kanten gibt, also

$$G = \begin{array}{cc} \textcircled{1} & \textcircled{4} \\ \diagdown \quad \diagup & \diagdown \quad \diagup \\ \textcircled{2} \text{---} \textcircled{3} & \textcircled{5} \text{---} \textcircled{6} \end{array} \quad \text{und} \quad L = \begin{bmatrix} 2 & -1 & -1 & 0 & 0 & 0 \\ -1 & 2 & -1 & 0 & 0 & 0 \\ -1 & -1 & 2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 2 & -1 & -1 \\ 0 & 0 & 0 & -1 & 2 & -1 \\ 0 & 0 & 0 & -1 & -1 & 2 \end{bmatrix}.$$

Mit Beispiel 5.7 ergeben sich die Eigenwerte  $\lambda_1 = 0$ ,  $\lambda_2 = 0$ ,  $\lambda_3 = 3$ ,  $\lambda_4 = 3$ ,  $\lambda_5 = 3$ ,  $\lambda_6 = 3$  und Eigenvektoren

$$v_1 = \frac{1}{\sqrt{3}} \begin{bmatrix} 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \quad v_2 = \frac{1}{\sqrt{3}} \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 1 \\ 1 \end{bmatrix}, \quad v_3 = \frac{1}{\sqrt{6}} \begin{bmatrix} 1 \\ 1 \\ -2 \\ 0 \\ 0 \\ 0 \end{bmatrix},$$

$$v_4 = \frac{1}{\sqrt{6}} \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 1 \\ -2 \end{bmatrix}, \quad v_5 = \frac{1}{\sqrt{2}} \begin{bmatrix} -1 \\ -1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \quad v_6 = \frac{1}{\sqrt{2}} \begin{bmatrix} 0 \\ 0 \\ 0 \\ -1 \\ -1 \\ 0 \end{bmatrix}.$$

Aufgrund dieser Beispiele vermuten wir, dass der Laplacian eines jeden Graphs nur nicht-negative Eigenwerte hat, und dass ferner der kleinste Eigenwert  $\lambda_1 = 0$  ist. Der zweitkleinste Eigenwert  $\lambda_2$  kann Null sein oder ungleich Null und dies scheint etwas über die Cluster auszusagen: In Beispiel 5.7 ist  $\lambda_2 \neq 0$  und es gibt keine Cluster, in Beispiel 5.8 ist  $\lambda_2 = 0$  und es gibt offenbar zwei Cluster.

Um dies zu beweisen, zerlegen wir den Laplacian in eine Summe von Laplacianen von sehr einfachen Graphen.



**Definition 5.9.** Ist  $G = (V, E)$  der Graph mit Vertices  $V = \{1, \dots, n\}$  und genau einer Kante  $E = \{e\}$  so bezeichnen wir mit  $L_e$  dessen Laplacematrix, beispielsweise:

$$L_{\{1,2\}} = \begin{bmatrix} 1 & -1 & 0 & \dots & 0 \\ -1 & 1 & 0 & \dots & 0 \\ 0 & \dots & \dots & \dots & 0 \\ \vdots & & & & \vdots \\ 0 & \dots & \dots & \dots & 0 \end{bmatrix}, \quad L_{\{1,3\}} = \begin{bmatrix} 1 & 0 & -1 & 0 & \dots & 0 \\ 0 & 0 & 0 & 0 & \dots & 0 \\ -1 & 0 & 1 & 0 & \dots & 0 \\ 0 & \dots & \dots & \dots & \dots & 0 \\ \vdots & & & & & \vdots \\ 0 & \dots & \dots & \dots & \dots & 0 \end{bmatrix}.$$

**Lemma 5.10.** Sei  $G$  ein Graph mit Laplacian  $L$ . Dann gilt

- (i)  $L = \sum_{e \in E} L_e$ ,
- (ii)  $\forall x \in \mathbb{R}^n: \langle x, Lx \rangle = \sum_{\{i,j\} \in E} (x_i - x_j)^2 \geq 0$ ,
- (iii) Alle Eigenwerte von  $L$  sind größer gleich Null.

*Beweis.* (i) Dies ist einer der seltenen Fälle, in denen es legitim ist, einen ‘Beweis durch Beispiel’ zu führen, da man an dem einen folgenden Beispiel in der Tat sofort sieht, dass dieselbe Rechnung auch für jeden anderen Graphen funktioniert:

$$\begin{aligned} L(\text{graph with edge (1,2)}) + L(\text{graph with edge (1,3)}) + L(\text{graph with edge (2,3)}) &= \begin{bmatrix} 1 & -1 & 0 \\ -1 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} + \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & -1 \\ 0 & -1 & 1 \end{bmatrix} + \begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{bmatrix} \\ &= \begin{bmatrix} 2 & -1 & -1 \\ -1 & 2 & -1 \\ -1 & -1 & 2 \end{bmatrix} = L(\text{triangle graph}). \end{aligned}$$

(ii) Wir nehmen ohne Einschränkung an, dass  $V = \{1, \dots, n\}$  gilt. Unter Benutzung von (i) und der Definition von  $L_{\{i,j\}}$  gilt

$$\begin{aligned} \langle x, Lx \rangle &= \langle x, (\sum_{e \in E} L_e)x \rangle = \sum_{e \in E} \langle x, L_e x \rangle \\ &= \sum_{\{i,j\} \in E} \left\langle \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} x_i \\ x_j \end{bmatrix} \right\rangle \\ &= \sum_{\{i,j\} \in E} \left\langle \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \begin{bmatrix} x_i - x_j \\ -x_i + x_j \end{bmatrix} \right\rangle \\ &= \sum_{\{i,j\} \in E} x_i(x_i - x_j) + x_j(-x_i + x_j) \\ &= \sum_{\{i,j\} \in E} (x_i - x_j)^2. \end{aligned}$$

(iii) Wegen (ii) ist  $L$  positiv semidefinit und es folgt, dass alle Eigenwerte größer gleich Null sind.  $\square$

Im Folgenden bezeichnen wir mit  $0 = \lambda_1 \leq \lambda_2 \leq \dots$  die Eigenwerte des Laplacians, wenn ein Graph  $G$  gegeben ist und nennen diese auch manchmal die *Eigenwerte des Graphen*.

**Definition 5.11.** Sei  $G = (V, E)$  ein Graph.

- (i) Wir sagen, dass  $G$  *zusammenhängend* ist, wenn es für beliebige  $i \neq j \in V$  Kanten  $\{i, k_1\}, \{k_1, k_2\}, \dots, \{k_\ell, j\}$  gibt, die  $i$  und  $j$  verbinden.
- (ii) Eine Teilmenge  $V' \subseteq V$  heißt *Zusammenhangskomponente*, falls  $G' = (V', E')$  mit  $E' = \{\{i, j\} \in E \mid i, j \in V'\}$  zusammenhängend und unter Inklusion maximal mit dieser Eigenschaft ist (d.h. für alle  $V' \subseteq V'' \subseteq V$  mit  $G''$  zusammenhängend folgt  $G' = G''$ ).

Offenbar ist die Zugehörigkeit zu einer Zusammenhangskomponente eine Äquivalenzrelation und die Zusammenhangskomponenten bilden eine Partition von  $V$ .

**Satz 5.12.** *Sei  $G$  ein Graph. Die Anzahl der Zusammenhangskomponenten stimmt mit der Vielfachheit des Eigenwerts Null ( $= \dim \ker L$ ) des Laplacians überein. Insbesondere ist  $G$  zusammenhängend genau dann wenn  $\lambda_2 > 0$ .*

*Beweis.* ‘ $\leq$ ’ Analog zu Beispiel 5.8 erhalten wir für jede Zusammenhangskomponente  $C_i$  den Eigenvektor  $\mathbb{1}_{C_i} = (\mathbb{1}_{C_i}(j))_{j=1, \dots, n}$  zum Eigenwert Null, und diese sind orthogonal.

‘ $\geq$ ’ Angenommen, es gilt ‘ $<$ ’, sagen wir z.B. es gibt  $k$  Zusammenhangskomponenten, aber es gilt  $\dim \ker L \geq k + 1$ . Dann können wir die linear unabhängigen Eigenvektoren  $v_1 = \mathbb{1}_{C_1}, \dots, v_k = \mathbb{1}_{C_k}$  von oben durch  $v_{k+1} \in \ker L$  zu einem linear unabhängigen System  $\{v_1, \dots, v_{k+1}\}$  erweitern. Es folgt  $Lv_{k+1} = 0$  und daher

$$0 = \langle v_{k+1}, Lv_{k+1} \rangle = \sum_{\substack{\uparrow \\ \text{Lem.} \\ 5.10}}_{\{i,j\} \in E} ((v_{k+1})_i - (v_{k+1})_j)^2.$$

Es muss also  $(v_{k+1})_i - (v_{k+1})_j = 0$  sein für alle Kanten  $\{i, j\}$ . Damit ist die Abbildung  $V \rightarrow \mathbb{R}, i \mapsto (v_{k+1})_i$  konstant auf Zusammenhangskomponenten  $C_1, \dots, C_k \subseteq V$ . Seien die Werte dort  $\alpha_1, \dots, \alpha_k$ . Dann ist

$$v_{k+1} = \alpha_1 \mathbb{1}_{C_1} + \dots + \alpha_k \mathbb{1}_{C_k} = \alpha_1 v_1 + \dots + \alpha_k v_k$$

im Widerspruch dazu, dass  $\{v_1, \dots, v_{k+1}\}$  linear unabhängig ist.  $\square$

Mit Satz 5.12 können wir Zusammenhangskomponenten, also die ‘besten’ aber auch die ‘trivialsten’ Cluster, ausfindig machen: Da wir deren Anzahl kennen, können wir irgendwo anfangen und jeweils alle Knoten, die zur Zusammenhangskomponente gehören, aufspüren indem wir von dort ausgehend entlang Kanten z.B. per Breitensuche immer weiter vom Ausgangsvertex weggehen.

Interessanter ist allerdings der Fall, dass  $G$  zusammenhängend ist, aber trotzdem Cluster hat, wie etwa in Beispiel 5.2. Hierfür brauchen wir das folgende Hilfsmittel aus der Linearen Algebra.

**Satz 5.13.** (Courant-Fischer-Formel) *Sei  $M \in \mathbb{R}^{n \times n}$  symmetrisch, d.h. wir haben mit Vielfachheiten reelle Eigenwerte  $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$  und können ein Orthogonalsys-*

stem  $\{v_1, \dots, v_n\}$  aus zugehörigen Eigenvektoren wählen. Dann gilt für  $k = 1, \dots, n$

$$\lambda_k = \min_{\substack{x \neq 0 \\ \langle x, v_i \rangle = 0 \\ \text{für } i=1, \dots, k-1}} \frac{\langle x, Mx \rangle}{\langle x, x \rangle} \quad \text{und} \quad v_k \in \operatorname{argmin}_{\substack{x \neq 0 \\ \langle x, v_i \rangle = 0 \\ \text{für } i=1, \dots, k-1}} \frac{\langle x, Mx \rangle}{\langle x, x \rangle}$$

wobei wir die Skalarproduktbedingung als leer verstehen, wenn  $k = 1$  ist.

*Beweis.* Da  $\langle x, v_i \rangle = 0$  genau dann gilt, wenn  $\langle x, v_i / \|v_i\| \rangle = 0$ , können wir ohne Einschränkung  $\|v_i\| = 1$  annehmen und davon ausgehen, dass  $\{v_1, \dots, v_n\}$  eine Orthonormalbasis ist. Sei  $0 \neq x \in \mathbb{R}^n$  beliebig. Dann können wir  $x$  nach der Orthonormalbasis entwickeln, also  $x = \alpha_1 v_1 + \dots + \alpha_n v_n$  mit  $\langle x, v_i \rangle = \langle \alpha_1 v_1 + \dots + \alpha_n v_n, v_i \rangle = \alpha_1 \langle v_1, v_i \rangle + \dots + \alpha_n \langle v_n, v_i \rangle = \alpha_i$  für  $i = 1, \dots, n$  schreiben. Sei  $k$  fest. Für  $x \neq 0$  mit  $\langle x, v_i \rangle = 0$  für  $i = 1, \dots, k-1$  wie im Satz gilt

$$\begin{aligned} \langle x, Mx \rangle &= \left\langle \sum_{i=1}^n \alpha_i v_i, M \left( \sum_{j=1}^n \alpha_j v_j \right) \right\rangle = \sum_{i,j=1}^n \alpha_i \alpha_j \langle v_i, M v_j \rangle \\ &= \sum_{i,j=1}^n \alpha_i \alpha_j \langle v_i, \lambda_j v_j \rangle = \sum_{i,j=1}^n \alpha_i \alpha_j \lambda_j \langle v_i, v_j \rangle \\ &= \sum_{i=1}^n \alpha_i^2 \lambda_i \stackrel{(*)}{=} \sum_{i=k}^n \lambda_i \alpha_i^2 \geq \lambda_k \sum_{i=k}^n \alpha_i^2 \stackrel{(*)}{=} \lambda_k \sum_{i=1}^n \alpha_i^2. \end{aligned}$$

wobei wir in  $(*)$  benutzt haben, dass  $\alpha_i = \langle x, v_i \rangle = 0$  für  $i = 1, \dots, k-1$  gilt und in der Abschätzung dazwischen, dass die  $\lambda_i$  wachsend sortiert sind. Völlig analog sieht man

$$\langle x, x \rangle = \sum_{i=1}^n \alpha_i^2$$

und damit folgt  $\frac{\langle x, Mx \rangle}{\langle x, x \rangle} \geq \lambda_k$  für alle  $x$ , die wir oben betrachtet haben. D.h. wir haben

$$\min_{\substack{x \neq 0 \\ \langle x, v_i \rangle = 0 \\ \text{für } i=1, \dots, k-1}} \frac{\langle x, Mx \rangle}{\langle x, x \rangle} \geq \lambda_k.$$

Andererseits ist  $\frac{\langle v_k, M v_k \rangle}{\langle v_k, v_k \rangle} = \frac{\langle v_k, \lambda_k v_k \rangle}{\langle v_k, v_k \rangle} = \lambda_k$ . D.h. oben gilt die Gleichheit und  $v_k$  ist ein Minimierer.  $\square$

**Bemerkung 5.14.** (i) Im obigen Kontext wird  $\frac{\langle x, Mx \rangle}{\langle x, x \rangle}$  oft als *Rayleigh-Quotient* bezeichnet.

(ii) Satz 5.13 kann wie folgt verallgemeinert werden: Für symmetrisches  $M \in \mathbb{R}^{n \times n}$  mit Eigenwerten  $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$  gilt

$$\lambda_k = \max_{\substack{U \subseteq \mathbb{R}^n \\ \dim U = n-k+1}} \min_{\substack{x \in U \\ x \neq 0}} \frac{\langle x, Mx \rangle}{\langle x, x \rangle} = \min_{\substack{U \subseteq \mathbb{R}^n \\ \dim U = k}} \max_{\substack{x \in U \\ x \neq 0}} \frac{\langle x, Mx \rangle}{\langle x, x \rangle},$$

wobei das Maximum bzw. Minimum über alle Unterräume  $U \subseteq \mathbb{R}^n$  mit der an-

gegebenen Dimension genommen wird. In der Tat zeigt Satz 5.13 die Abschätzung “ $\leq$ ” in der ersten angegebenen Gleichung. Für die andere Richtung betrachten wir für  $U$  mit  $\dim U = n - k + 1$  den Schnitt  $U \cap \text{span}\{v_k, \dots, v_n\}$  und wählen  $x = \alpha_1 v_1 + \dots + \alpha_n v_n \neq 0$  in letzterem, wobei die  $v_i$  eine Basis aus Eigenvektoren zu den  $\lambda_i$  bilden. Dann folgt

$$\frac{\langle x, Mx \rangle}{\langle x, x \rangle} = \frac{\sum_{i=k}^n \lambda_i \alpha_i^2}{\sum_{i=k}^n \alpha_i^2} \geq \frac{\lambda_k \sum_{i=k}^n \alpha_i^2}{\sum_{i=k}^n \alpha_i^2} = \lambda_k$$

weil die  $\lambda_i$  wachsend sind. Da  $x \in U$  beliebig war, erhalten wir  $\max_{0 \neq x \in U} \frac{\langle x, Mx \rangle}{\langle x, x \rangle} \geq \lambda_k$  und da dieses nun für beliebige  $U$  mit der angegebenen Dimension gilt, muss auch

$$\max_{\substack{U \subseteq \mathbb{R}^n \\ \dim U = n-k+1}} \min_{\substack{x \in U \\ x \neq 0}} \frac{\langle x, Mx \rangle}{\langle x, x \rangle} \geq \lambda_k$$

gelten. Für die zweite Gleichung wendet man die erste auf  $-M$  an.

(iii) Anstatt oben jeweils  $\frac{\langle x, Mx \rangle}{\langle x, x \rangle}$  über alle  $x \neq 0$  zu minimieren/maximieren, kann man auch nur den Zähler  $\langle x, Mx \rangle$  aber über  $\|x\| = 1$ , minimieren/maximieren.

Angewandt auf den Laplacian eines Graphs erhalten wir nun aus der Courant-Fischer-Formel das folgende.

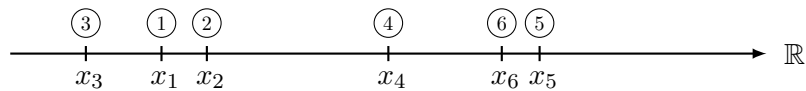
**Korollar 5.15.** *Sei  $G$  ein Graph,  $L$  sein Laplacian, und  $0 = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$  dessen Eigenwerte mit Vielfachheit. Wir wissen, dass  $v_1 = \mathbb{1}$  gewählt werden kann. Dann gelten*

$$\lambda_2 = \underset{\substack{\uparrow \\ \text{Satz} \\ 5.13}}{\min_{\substack{x \neq 0 \\ \langle x, \mathbb{1} \rangle = 0}}} \frac{\langle x, Lx \rangle}{\langle x, x \rangle} = \min_{\substack{\|x\|=1 \\ \langle x, \mathbb{1} \rangle = 0}} \langle x, Lx \rangle = \underset{\substack{\uparrow \\ \text{Lem.} \\ 5.10}}{\min_{\substack{\|x\|=1 \\ \langle x, \mathbb{1} \rangle = 0}}} \sum_{\{i,j\} \in E} (x_i - x_j)^2$$

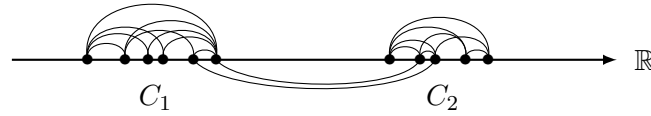
und

$$v_2 = \underset{\substack{\|x\|=1 \\ \langle x, \mathbb{1} \rangle = 0}}{\operatorname{argmin}} \langle x, Lx \rangle = \underset{\substack{\|x\|=1 \\ \langle x, \mathbb{1} \rangle = 0}}{\operatorname{argmin}} \sum_{\{i,j\} \in E} (x_i - x_j)^2. \quad \square$$

Korollar 5.15 führt auf die folgende Heuristik: Sei  $G$  ein Graph und es sei  $v_2 = (x_1, \dots, x_n)$  wie oben ein Eigenvektor mit Norm 1 zum zweitkleinsten Eigenwert des Laplacians von  $G$ . Wir notieren jetzt die Vertices als Punkte in  $\mathbb{R}$  und verwenden dabei  $x_1$  als Koordinate von Vertex 1,  $x_2$  als Koordinate von Vertex 2, usw.:



Dann bedeuten die Gleichungen in Korollar 5.15, dass sich dort, wo viele  $x_i$ 's nah zusammenliegen, ein Cluster befindet:



Denn würde man einen Punkt von  $C_1$  nach  $C_2$  verschieben, so hätte man viel mehr ‘lange Kanten’, die den Wert  $\sum_{\{i,j\} \in E} (x_i - x_j)^2$  vergrößern würden. Andererseits können  $C_1$  und  $C_2$  nicht näher zusammen geschoben werden, da die Bedingung  $\|v_2\| = 1$  gilt. Ist  $\lambda_2 \approx 0$ , so bedeutet dies, dass eine Anordnung möglich ist, bei der wenig Verbindungen von  $C_1$  nach  $C_2$  bestehen. Man spricht dann von einem *Flaschenhals*. Jetzt machen wir obiges formal.

**Definition 5.16.** Sei  $G = (V, E)$  ein Graph mit  $\deg(e) > 0$  für alle  $e \in E$ .

(i) Für  $S \subseteq V$  definieren wir das *Volumen* und den *Rand* per

$$\text{vol}(S) = \sum_{v \in S} \deg(v) \quad \text{und} \quad \partial S = \{\{i, j\} \in E \mid i \in S, j \in S^c\}.$$

(ii) Für  $\emptyset \neq S \subset V$  definieren wir die *Leitfähigkeit* per

$$\phi(S) = \frac{\#\partial S}{\min(\text{vol } S, \text{vol } S^c)}.$$

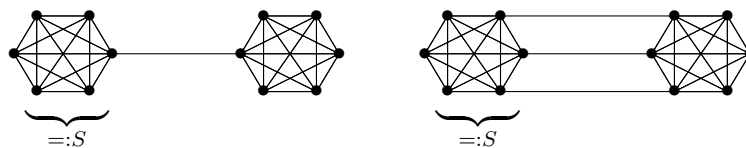
(iii) Schließlich definieren wir die *Cheegerkonstante* durch

$$C_G := \min_{\emptyset \neq S \subset V} \phi(S).$$

Sei  $G = (V, E)$  ein Graph. Dann ist  $C_G$  klein, falls  $S \subseteq V$  existiert sodass einerseits  $\#\partial S$  klein ist, also bei einer Zerlegung von  $V = S \cup S^c$  in zwei Cluster nicht zu viele Kanten zwischen den Clustern ‘zerschnitten’ werden. Andererseits muss  $\min(\text{vol } S, \text{vol } S^c)$  groß sein; wählt man also  $S$  sehr klein, um  $\#\partial S$  klein zu halten, so macht einem der Nenner von  $\phi(S)$  den Bruch wieder groß.

Etwas ungenau, aber sehr anschaulich, kann man sich die Minimierung von  $\phi$  so vorstellen, als dass man nach dem *besten Schnitt durch  $G$*  sucht, wobei (a) möglichst wenig Kanten zerschnitten werden sollen und *gleichzeitig* (b) beide entstehenden Teile möglichst gleichviele Kanten enthalten sollen.

Als instruktives Beispiel betrachten wir einen Graph der aus zwei *Cliquen* besteht. Diese sind links durch eine Kante und rechts durch drei Kanten verbunden.



Durch Abzählen erkennt man, dass im linken Bild

$$\phi(S) = \frac{1}{5 \cdot 5 + 1 \cdot 6} \approx 0.03$$

gilt und dass jede andere Wahl von  $S$  auf einen größeren Bruch führt. Der optimale Schnitt ist also der, der die beiden Cliques trennt und die Cheegerkonstante des linken Graphen ist also 0.03. Im rechten Graphen ergibt sich analog

$$C_G = \phi(S) = \frac{3}{3 \cdot 5 + 3 \cdot 6} \approx 0.09.$$

Die Cheegerkonstante ist links also deutlich kleiner. Dies kann man so interpretieren, als dass der ‘Flaschenhals’ zwischen den Cliques links ausgeprägter ist als rechts.

Wollen wir nun die Cluster eines durch seine Adjazenzmatrix gegebenen Graphen finden, so müssen wir sukzessive möglichst gute Schnitte  $S \cup S^c$  durchführen, solange die Cheegerkonstante klein ist — und aufhören sobald diese groß wird. Auf die Frage, was hierbei mit *groß* und *klein* gemeint ist, kommen wir noch zurück.

Wir notieren die folgenden einfachen Konsequenzen von Definition 5.16.

**Bemerkung 5.17.** Sei  $G = (V, E)$  ein Graph.

- (i) Für  $S \subseteq V$  gilt  $\phi(S) = \phi(S^c)$ ,  $0 \leq \phi(S) \leq 1$ , also insbesondere  $0 \leq C_G \leq 1$ .
- (ii) Es existiert  $S \subseteq V$  mit  $C_G = \phi(S) = \frac{\#\partial S}{\text{vol } S}$  denn wegen  $\partial S = \partial(S^c)$  ist mit  $S$  auch immer  $S^c$  ein Minimierer.
- (iii) Es gilt  $\text{vol}(V) = \text{vol}(S) + \text{vol}(S^c)$  für alle  $\emptyset \neq S \subset V$ .

Für das Hauptresultat dieses Kapitels, eine Abschätzung der Cheegerkonstanten nach oben und unten durch Eigenwerte, müssen wir die Laplacematrix, die uns ja bei den Zusammenhangskomponenten schon gute Dienste geleistet hat, normalisieren.

**Definition 5.18.** Sei  $G$  ein Graph mit Knoten  $\{1, \dots, n\}$  und  $\deg(i) > 0$  für alle  $i = 1, \dots, n$ . Es sei  $L$  dessen Laplacian und  $D = \text{diag}(\deg(1), \dots, \deg(n)) =: \text{diag}(d_1, \dots, d_n)$ . Die Matrix  $\mathcal{L} := D^{-1/2} L D^{-1/2}$  mit  $D^{-1/2} = \text{diag}(d_1^{-1/2}, \dots, d_n^{-1/2})$  heißt *normalisierter Laplacian*.

Wir notieren zunächst ein Analogon zu Korollar 5.15.

**Lemma 5.19.** Sei  $G$  ein Graph mit Knoten  $\{1, \dots, n\}$  und  $\deg(i) > 0$  für alle  $i = 1, \dots, n$ . Wir nummerieren die Eigenwerte des Laplacians  $\mathcal{L}$  aufsteigend durch. Dann gilt  $0 = \lambda_1(\mathcal{L}) \leq \lambda_2(\mathcal{L}) \leq \dots$  und

$$\lambda_2(\mathcal{L}) = \min_{\substack{x \neq 0 \\ D x = 0}} \frac{\sum_{\{i,j\} \in E} (x_i - x_j)^2}{\sum_{i=1}^n x_i^2 d_i}$$

*Beweis.* Für  $x \in \mathbb{R}^n$  gilt

$$\langle x, \mathcal{L} x \rangle = \langle x, D^{-1/2} L D^{-1/2} x \rangle = \langle D^{-1/2} x, L D^{-1/2} x \rangle = \langle y, L y \rangle \geq 0$$

$\uparrow$   
 Lemma 5.10

mit  $y := D^{-1/2}x$ , d.h.  $\mathcal{L}$  ist positiv semidefinit und alle Eigenwerte sind größer gleich Null. Ferner ist  $\lambda_1(\mathcal{L}) = 0$  mit Eigenvektor  $D^{1/2}\mathbf{1}$  denn

$$\mathcal{L}D^{1/2}\mathbf{1} = D^{-1/2}LD^{-1/2}D^{1/2}\mathbf{1} = D^{-1/2}L\mathbf{1} \underset{\substack{\uparrow \\ \text{Lemma} \\ 5.6}}{=} D^{-1/2}0\mathbf{1} = 0D^{1/2}\mathbf{1}.$$

Für den zweitkleinsten Eigenwert gilt schließlich

$$\begin{aligned} \lambda_2(\mathcal{L}) &= \min_{\substack{\uparrow \\ \text{Lemma} \\ 5.13} \atop x \neq 0, \langle x, D^{1/2}\mathbf{1} \rangle = 0} \frac{\langle x, \mathcal{L}x \rangle}{\langle x, x \rangle} \stackrel{(*)}{=} \min_{\substack{y \neq 0 \\ \langle D^{1/2}y, D^{1/2}\mathbf{1} \rangle = 0}} \frac{\langle D^{1/2}y, \mathcal{L}D^{1/2}y \rangle}{\langle D^{1/2}y, D^{1/2}y \rangle} \\ &= \min_{\substack{y \neq 0 \\ \langle y, D\mathbf{1} \rangle = 0}} \frac{\langle y, Ly \rangle}{\langle y, Dy \rangle} = \min_{\substack{y \neq 0 \\ Dy=0}} \frac{\sum_{\{i,j\} \in E} (y_i - y_j)^2}{\sum_{i=1}^n y_i^2 d_i} \end{aligned}$$

wobei wir in  $(*)$  die Substitution  $y := D^{-1/2}x$  durchgeführt haben und dann die Identität  $L = D^{1/2}\mathcal{L}D^{1/2}$  benutzt haben.  $\square$

Wir bemerken, dass im Allgemeinen natürlich  $\lambda_2(L) \neq \lambda_2(\mathcal{L})$  gilt. Wir werden aber noch sehen, dass  $\lambda_2(\mathcal{L}) > 0$  den Zusammenhang von  $G$  charakterisiert. Ist  $d$  regulär, so gilt  $\lambda_2(\mathcal{L}) = \frac{1}{d}\lambda_2(L)$ , vergleiche Aufgabe 5.4.

Für den Beweis der von Satz 5.21 benötigen wir das folgende technische Lemma.

**Lemma 5.20.** *Sei  $G = (V, E)$  ein Graph mit  $V = \{1, \dots, n\}$ . Für  $i = 1, \dots, n$  definieren wir  $S_i := \{1, \dots, i\}$  und  $S_0 := \emptyset$ . Sei*

$$r := \max\{i \in \{1, \dots, n\} \mid 2 \text{ vol } S_i \leq \text{vol } V\}.$$

- (i) Für  $0 \leq k \leq r$  gilt  $\min(\text{vol } S_k, \text{vol } S_k^c) = \text{vol } S_k$  und  $\text{vol } S_k - \text{vol } S_{k+1} = -d_{k+1}$ .
- (ii) Für  $r \leq k \leq n$  gilt  $\min(\text{vol } S_k, \text{vol } S_k^c) = \text{vol } S_k^c$  und  $\text{vol } S_k^c - \text{vol } S_{k+1}^c = d_{k+1}$ .

*Beweis.* (i) Sei  $0 \leq k \leq r$ . Dann gilt  $2 \text{ vol } S_k \leq \text{vol } V = \text{vol } S_k + \text{vol } S_k^c$  per Definition von  $r$ . Es folgt  $\text{vol } S_k \leq \text{vol } S_k^c$ . Weiter gilt

$$\text{vol } S_k - \text{vol } S_{k+1} = \sum_{i=1}^k d_i - \sum_{i=1}^{k+1} d_i = -d_{k+1}.$$

(ii) Sei  $r \leq k \leq n$ . Dann gilt  $2 \text{ vol } S_k \geq \text{vol } V = \text{vol } S_k + \text{vol } S_k^c$ , also  $\text{vol } S_k \geq \text{vol } S_k^c$ . Ferner erhalten wir

$$\text{vol } S_k^c - \text{vol } S_{k+1}^c = \sum_{i=k+1}^n d_i - \sum_{i=k+2}^n d_i = d_{k+1}$$

wie behauptet.  $\square$

Nun kommen wir zum Hauptergebnis dieses Kapitels.

**Satz 5.21.** (Cheeger-Ungleichung) Sei  $G = (V, E)$  ein Graph mit  $\deg(i) > 0$  für alle  $i \in V$  und seien  $0 = \lambda_1 \leq \lambda_2 \leq \dots$  die Eigenwerte des normalisierter Laplacians von  $G$ . Dann gilt

$$\frac{\lambda_2}{2} \leq C_G \leq \sqrt{2\lambda_2}.$$

Es folgt, dass  $\lambda_2 \approx 0$  genau dann gilt wenn  $C_G \approx 0$  und formalisiert damit die oben propagierte Heuristik, nach der es genau im Fall  $\lambda_2 \approx 0$  möglich ist, den Graph derart in zwei Teile zu schneiden, sodass im Vergleich zu der Anzahl der Kanten in beiden Teilen wenig Kanten zwischen ihnen zerschnitten werden.

*Beweis.* Wie in früheren Beweisen nehmen wir ohne Einschränkung  $V = \{1, \dots, n\}$  an. Wir bezeichnen mit  $\mathbf{1}_S = (\mathbf{1}_{S,i})_{i=1,\dots,n} \in \mathbb{R}^n$  den Vektor mit Einträgen  $\mathbf{1}_{S,i} = 1$  falls  $i \in S$  und  $\mathbf{1}_{S,i} = 0$  für  $i \notin S$ . Die zwei Abschätzungen im Satz zeigen wir nun einzeln. Wir beginnen mit der linken.

① Erste Cheegerungleichung:  $\lambda_2 \leq 2 C_G$ .

Sei  $S \subseteq V$  ein Minimierer für die Cheegerkonstante. Nach Bemerkung 5.17(ii) können wir ohne Einschränkung  $C_G = \frac{\#\partial S}{\text{vol } S}$  annehmen und das heißt insbesondere  $\min(\text{vol } S, \text{vol } S^c) = \text{vol } S$ . Wir definieren den Vektor

$$x := \mathbf{1}_S - \frac{\text{vol } S}{\text{vol } V} \mathbf{1}_V \in \mathbb{R}^n.$$

Dann gilt  $x \neq 0$  sowie

$$Dx = \sum_{i=1}^n d_i x_i = \sum_{i=1}^n d_i \mathbf{1}_{S,i} - \frac{\text{vol } S}{\text{vol } V} \sum_{i=1}^n d_i \mathbf{1}_{V,i} = \text{vol } S - \text{vol } S = 0,$$

und es folgt mit Lemma 5.19

$$\lambda_2 \leq \frac{\sum_{\{i,j\} \in E} (x_i - x_j)^2}{\sum_{i=1}^n x_i^2 d_i} =: \frac{Z}{N}.$$

Wir behandeln zuerst den Zähler. Da per Definition von  $x$  für  $i = 1, \dots, n$

$$x_i = \begin{cases} 1 - \frac{\text{vol } S}{\text{vol } V} & \text{falls } i \in S \\ -\frac{\text{vol } S}{\text{vol } V} & \text{falls } i \notin S \end{cases}$$

gilt, haben wir für  $i \neq j$

$$|x_i - x_j| = \begin{cases} 0 & \text{falls } i, j \in S \text{ oder } i, j \notin S \\ 1 & \text{sonst} \end{cases}$$

und daher ist

$$Z = \sum_{\{i,j\} \in E} (x_i - x_j)^2 = \#\{\{i,j\} \in E \mid i \in S \text{ and } j \notin S\} = \#\partial S.$$



Als nächstes schätzen wir den Nenner ab. Hier gilt

$$\begin{aligned}
 N &= \sum_{i=1}^n x_i^2 d_i = \sum_{i=1}^n \left( \mathbb{1}_{S,i} - \frac{\text{vol } S}{\text{vol } V} \mathbb{1}_{V,i} \right)^2 d_i \\
 &= \sum_{i=1}^n \mathbb{1}_{S,i}^2 d_i - 2 \frac{\text{vol } S}{\text{vol } V} \sum_{i=1}^n \mathbb{1}_{S,i} \mathbb{1}_{V,i} d_i + \left( \frac{\text{vol } S}{\text{vol } V} \right)^2 \sum_{i=1}^n \mathbb{1}_{V,i}^2 d_i \\
 &= \text{vol } S - 2 \frac{(\text{vol } S)^2}{\text{vol } V} + \frac{(\text{vol } S)^2}{\text{vol } V} \\
 &= \text{vol } S - \text{vol } V \cdot \frac{\text{vol } S}{\text{vol } S + \text{vol } S^c} \\
 &\geq \text{vol } S - \text{vol } V \cdot \frac{\text{vol } S}{2 \text{vol } S} \\
 &= \frac{\text{vol } S}{2}
 \end{aligned}$$

wobei wir  $\text{vol } S + \text{vol } S^c \geq 2 \min(\text{vol } S, \text{vol } S^c) = 2 \text{vol } S$  abgeschätzt haben. Zusammen erhalten wir

$$\lambda_2 \leq \frac{Z}{N} \leq \frac{|\partial S|}{(\text{vol } S)/2} = 2 \frac{|\partial S|}{\text{vol } S} = 2 C_G.$$

② Zweite Cheegerungleichung:  $\lambda_2 \geq \frac{C_G^2}{2}$ .

Für  $i = 1, \dots, n$  sei  $S_i := \{1, \dots, i\}$  und  $S_0 := \emptyset$ . Weiter definieren wir

$$\alpha := \min_{i=1, \dots, n} \phi(S_i) \geq \min_{\emptyset \neq S \subset V} \phi(S) = C_G \quad (5.1)$$

und unser Ziel im Folgenden wird es sein, zu zeigen, dass  $\lambda_2 \geq \alpha^2$  ist, denn damit folgt dann die fehlende Ungleichung. Um  $\lambda_2$  ins Spiel zu bringen, verwenden wir Lemma 5.19 und wählen dort einen Minimierer  $x \in \mathbb{R}^n$ . D.h. es gilt

$$\lambda_2 = \frac{\sum_{\{i,j\} \in E} (x_i - x_j)^2}{\sum_{i=1}^n x_i^2 d_i} \quad \text{mit} \quad x \neq 0 \text{ und } Dx = \sum_{i=1}^n d_i x_i = 0 \quad (5.2)$$

und wir schätzen nun als erstes den Nenner nach oben ab. Zunächst können wir ohne Einschränkung annehmen, dass  $x_1 \geq x_2 \geq \dots \geq x_n$  gilt. Dann definieren wir

$$r := \max\{i \in \{1, \dots, n\} \mid 2 \text{vol } S_i \leq \text{vol } V\}$$

und betrachten den Vektor  $(x_i - x_r)_{i=1, \dots, n}$ . Hier sind nun per obigem die Einträge fallend und zwar erst größer gleich Null, bei  $i = r$  gleich Null und danach kleiner gleich Null. Wir teilen  $(x_i - x_r)_{i=1, \dots, n}$  auf in Positivteil und Negativteil:

$$\begin{bmatrix} x_1 - x_r \\ \vdots \\ x_{r-1} - x_r \\ 0 \\ x_{r+1} - x_r \\ \vdots \\ x_n - x_r \end{bmatrix} = \begin{bmatrix} x_1 - x_r \\ \vdots \\ x_{r-1} - x_r \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} - \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 0 \\ x_r - x_{r+1} \\ \vdots \\ x_r - x_n \end{bmatrix} =: p + n.$$

Die Einträge von  $p$  sind nichtnegativ und fallend, die von  $n$  sind ebenfalls nichtnegativ aber wachsend. Insbesondere ist  $n_i \cdot p_i = 0$  für jedes  $i = 1, \dots, n$ . Jetzt schätzen wir ab, wobei in der Ungleichung der erste hinzugefügte Summand nach (5.2) Null und der zweite offenbar positiv ist:

$$\begin{aligned} \sum_{i=1}^n x_i^2 d_i &\leq \sum_{i=1}^n x_i^2 d_i - 2x_r \sum_{i=1}^n x_i d_i + x_r^2 \sum_{i=1}^n d_i = \sum_{i=1}^n (x_i - x_r)^2 d_i \\ &= \sum_{i=1}^n (p_i - n_i)^2 d_i = \sum_{i=1}^n (p_i^2 - 2p_i n_i + n_i^2) d_i = \sum_{i=1}^n (p_i^2 + n_i^2) d_i. \end{aligned} \quad (5.3)$$

Um den Zähler des Quotienten für  $\lambda_2$  in (5.2) abzuschätzen, betrachten wir die Summanden einzeln. Hier gilt

$$(x_i - x_j)^2 = [(p_i - n_i) - (p_j - n_j)]^2 \stackrel{(*)}{\geq} (p_i - p_j)^2 + (n_i - n_j)^2 \quad (5.4)$$

wobei die erste Gleichung durch Einsetzen sofort ersichtlich ist. Die mit  $(*)$  gekennzeichnete Ungleichung ist klar, falls  $i, j \geq r$  oder  $i, j \leq r$  gilt, denn dann ist  $p_i = p_j = 0$  bzw.  $n_i = n_j = 0$ . Gezeigt werden muss also ohne Einschränkung nur der Fall  $i \leq r \leq j$  wobei  $i \neq j$ . Hier haben wir  $p_i, n_j \geq 0$  und  $p_j = n_i = 0$  und es folgt

$$\begin{aligned} (x_i - x_j)^2 &= ((x_i - x_r) - (x_j - x_r))^2 \\ &= (x_i - x_r)^2 - 2(x_i - x_r)(x_j - x_r) + (x_j - x_r)^2 \\ &= p_i^2 - 2p_i(-n_j) + (-n_j)^2 \geq p_i^2 + (-n_j)^2 = (p_i - p_j)^2 + (n_i - n_j)^2. \end{aligned}$$

Für das Folgende notieren wir noch die beiden Ungleichungen

$$\begin{aligned} \forall a, c \geq 0, b, d > 0: \frac{a+c}{b+d} &\geq \min\left(\frac{a}{b}, \frac{c}{d}\right) \\ \forall u, v \in \mathbb{R}: (u+v)^2 &\leq 2(u^2 + v^2) \end{aligned} \quad (5.5)$$

deren einfache Beweise wir dem Leser überlassen. Die Anwendung unserer bisherigen Abschätzungen auf die in (5.2) notierte Formel für  $\lambda_2$  führt auf eine Formel in der Positiv- und Negativteile unvermischt auftreten

$$\begin{aligned} \lambda_2 &= \frac{\sum_{\{i,j\} \in E} (x_i - x_j)^2}{\sum_{i=1}^n x_i^2 d_i} \stackrel{(5.3)}{\geq} \frac{\sum_{\{i,j\} \in E} (p_i - p_j)^2 + (n_i - n_j)^2}{\sum_{i=1}^n (p_i^2 + n_i^2) d_i} \\ &\stackrel{(5.4)}{=} \frac{\sum_{\{i,j\} \in E} (p_i - p_j)^2 + \sum_{\{i,j\} \in E} (n_i - n_j)^2}{\sum_{i=1}^n p_i^2 d_i + \sum_{\{i,j\} \in E} n_i^2 d_i} \\ &\stackrel{(5.5)}{\geq} \min \left( \frac{\sum_{\{i,j\} \in E} (p_i - p_j)^2}{\sum_{i=1}^n p_i^2 d_i}, \frac{\sum_{\{i,j\} \in E} (n_i - n_j)^2}{\sum_{\{i,j\} \in E} n_i^2 d_i} \right) \\ &= \min \left( \frac{\sum_{\{i,j\} \in E} (p_i - p_j)^2}{\sum_{i=1}^n p_i^2 d_i}, \frac{\sum_{\{i,j\} \in E} (p_i + p_j)^2}{\sum_{\{i,j\} \in E} (p_i + p_j)^2}, \dots \right) \end{aligned}$$

$$=: \min \left( \frac{Z}{N}, \dots \right)$$

sodass wir im Folgenden nur noch den ausschließlich Positivteile enthaltenden Ausdruck  $Z/N$  abschätzen können und dabei beobachten werden, dass die Abschätzung des die Negativteile enthaltenden Ausdrucks analog funktioniert. Wir starten mit

$$\begin{aligned} N &= \sum_{i=1}^n p_i^2 d_i \cdot \sum_{\{i,j\} \in E} (p_i + p_j)^2 \stackrel{(5.5)}{\leq} \sum_{i=1}^n p_i d_i \cdot \sum_{\{i,j\} \in E} 2(p_i^2 + p_j^2) \\ &= 2 \left( \sum_{i=1}^n p_i^2 d_i \right) \cdot \left( \sum_{\{i,j\} \in E} p_i^2 + \sum_{\{i,j\} \in E} p_j^2 \right) = 2 \left( \sum_{i=1}^n p_i^2 d_i \right) \cdot 2 \left( \sum_{\{i,j\} \in E} p_i^2 \right) \\ &\stackrel{(\circ)}{=} 2 \left( \sum_{i=1}^n p_i^2 d_i \right) \cdot \left( \sum_{i=1}^n p_i^2 d_i \right) = 2 \left( \sum_{i=1}^n p_i^2 d_i \right)^2 \end{aligned}$$

wobei  $(\circ)$  so zustande kommt: Anstatt über alle Kanten zu summieren, summiert man über alle Knoten und multipliziert jeweils mit dem Grad — da dabei dann jede Kante zweimal vorkommt, korrigiert man dies durch Streichen des Vorfaktors. Ganz formal kann man wie folgt argumentieren

$$\begin{aligned} 2 \left( \sum_{\{i,j\} \in E} p_i^2 \right) &= 2 \left( \sum_{\substack{\{i,j\} \in E \\ i < j}} p_i^2 \right) = 2 \left( \sum_{\substack{i,j=1 \\ i < j}}^n \mathbb{1}_E(i,j) p_i^2 \right) = \sum_{i,j=1}^n \mathbb{1}_E(i,j) p_i^2 \\ &= \sum_{i=1}^n \left( p_i^2 \sum_{j=1}^n \mathbb{1}_E(i,j) \right) = \sum_{i=1}^n p_i^2 \cdot |\{ \{i,j\} \in E \}| = \sum_{i=1}^n p_i^2 d_i \end{aligned}$$

wobei wir die Abkürzung

$$\mathbb{1}_E: V \times V \rightarrow \mathbb{R}, \quad \mathbb{1}_E(i,j) = \begin{cases} 1 & \text{falls } \{i,j\} \in E, \\ 0 & \text{sonst} \end{cases} \quad (5.6)$$

benutzt haben. Als nächstes schätzen wir den Zähler ab. Wir empfehlen dem Leser für die mit  $(\Delta)$  gekennzeichneten Gleichungen jeweils den (dreieckigen) Bereich zu skizzieren über den summiert wird.

$$\begin{aligned} Z &= \sum_{\{i,j\} \in E} (p_i - p_j)^2 \cdot \sum_{\{i,j\} \in E} (p_i + p_j)^2 \\ &= \left\| ((p_i - p_j)^2)_{\{i,j\} \in E} \right\|_{\mathbb{R}^{|E|}}^2 \cdot \left\| ((p_i + p_j)^2)_{\{i,j\} \in E} \right\|_{\mathbb{R}^{|E|}}^2 \\ &\geq \left\langle ((p_i - p_j))_{\{i,j\} \in E}, ((p_i + p_j))_{\{i,j\} \in E} \right\rangle_{\mathbb{R}^{|E|}}^2 \\ &\stackrel{\substack{\uparrow \\ \text{CSB-} \\ \text{Ungl.}}}{=} \left( \sum_{\substack{i,j=1 \\ i < j}}^n \mathbb{1}_E(i,j) \cdot (p_i - p_j)(p_i + p_j) \right)^2 \\ &\stackrel{\substack{\uparrow \\ (5.6)}}{=} \left( \sum_{i=1}^{n-1} \sum_{j=i+1}^n \mathbb{1}_E(i,j) \sum_{k=i}^{j-1} p_k^2 - p_{k+1}^2 \right)^2 \\ &\stackrel{\substack{\uparrow \\ \text{Teleskop-} \\ \text{summe}}}{=} \end{aligned}$$

$$\begin{aligned}
&= \left( \sum_{i=1}^{n-1} \sum_{k=i}^{n-1} \sum_{j=k+1}^n \mathbb{1}_E(i, j) \cdot (p_k^2 - p_{k+1}^2) \right)^2 \\
&\stackrel{(\Delta)}{\uparrow} = \left( \sum_{k=1}^{n-1} \sum_{i=1}^k \sum_{j=k+1}^n \mathbb{1}_E(i, j) \cdot (p_k^2 - p_{k+1}^2) \right)^2 \\
&\stackrel{(\Delta)}{\uparrow} = \left( \sum_{k=1}^{n-1} (p_k^2 - p_{k+1}^2) \sum_{\substack{i \leq k < j \\ \{i, j\} \in E}} 1 \right)^2 \\
&= \left( \sum_{k=1}^{n-1} (p_k^2 - p_{k+1}^2) |\partial S_k| \right)^2 \\
&\geq \left( \sum_{k=1}^{n-1} (p_k^2 - p_{k+1}^2) \alpha \min(\text{vol } S_k, \text{vol } S_k^c) \right)^2 \\
&\stackrel{\substack{\uparrow \\ \text{Dfn } \alpha \\ (5.1)}}{=} \left( \sum_{k=1}^r (p_k^2 - p_{k+1}^2) \alpha \text{vol } S_k \right)^2 \\
&\stackrel{\substack{\uparrow \\ p_k=0 \text{ für } \\ k \geq r \text{ u. 5.20}}}{=} \alpha^2 \left( \sum_{k=1}^r p_k^2 \text{vol } S_k - \sum_{k=1}^r p_{k+1}^2 \text{vol } S_k \right)^2 \\
&= \alpha^2 \left( \sum_{k=0}^{r-1} p_{k+1}^2 \text{vol } S_{k+1} - \sum_{k=0}^{r-1} p_{k+1}^2 \text{vol } S_k \right)^2 \\
&\stackrel{\substack{\uparrow \\ \text{Indexshift} \\ \text{vol } S_0=0 \\ \text{u. } p_r=0}}{=} \alpha^2 \left( \sum_{k=0}^{r-1} p_{k+1}^2 d_{k+1} \right)^2 \\
&\stackrel{\substack{\uparrow \\ \text{Lemma} \\ 5.20}}{=} \alpha^2 \left( \sum_{k=1}^n p_k^2 d_k \right)^2 \\
&\stackrel{\substack{\uparrow \\ \text{Indexshift} \\ \text{u. } p_k=0 \\ \text{für } k \geq r}}{=} \alpha^2 \left( \sum_{k=1}^n p_k^2 d_k \right)^2
\end{aligned}$$

Zusammensetzen der Abschätzungen für Zähler und Nenner liefert

$$\frac{Z}{N} \geq \frac{\alpha^2 (\sum_{k=1}^n p_k^2 d_k)^2}{2 (\sum_{i=1}^n p_i^2 d_i)^2} = \frac{\alpha^2}{2}.$$

Für den Term mit den Negativteilen erhält man genau die gleiche Abschätzung durch analoge Argumentation und Benutzung der ‘dualen’ Aussagen in Lemma 5.20, bzw. mit  $n_k = 0$  für  $k \leq r$ . Wir überlassen diese Rechnung dem Leser. Wie bereits nach (5.1) erläutert, folgt  $\lambda_2 \geq \alpha^2/2 \geq C_G^2/2$ .  $\square$

Wir bemerken, dass die Cheeger-Ungleichung nach Jeff Cheeger benannt ist, der 1970 eine kontinuierliche Version, d.h. für Mannigfaltigkeiten statt Graphen, bewies. Daher stammen die Begriffe ‘Volumen’, ‘Rand’ usw. wie wir sie in diesem Kapitel benutzt haben. Der obige etwas vereinfachte Beweis stammt von Fan Chung aus dem Jahr 2007.

Als letztes wollen wir noch notieren, dass Satz 5.21 implizit eine obere Schranke für den Eigenwert  $\lambda_2(\mathcal{L})$  liefert.

**Korollar 5.22.** Sei  $G = (V, E)$  ein Graph mit  $\deg(i) > 0$  für alle  $i \in V$ . Dann gilt für den zweitkleinsten Eigenwert  $\lambda_2$  der normalisierten Laplacematrix von  $G$  die Abschätzung  $\lambda_2 \leq 8$  und für die Cheegerkonstante  $C_G \leq 4$ .  $\square$

## Referenzen

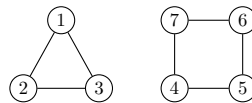
Die in diesem Kapitel benutzten Konzepte der spektralen Graphentheorie sind allesamt Standard und praktisch überall zu finden. Der Beweis von Satz 5.21 inklusive der vorhergehenden Lemmas ist [Chu07] entnommen. Der Autor bedankt sich bei T. Sauerwald für viele interessante Diskussionen über den Inhalt dieses Kapitels.

## Aufgaben

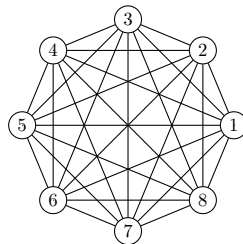
**Aufgabe 5.1.** Sei  $G$  ein  $d$ -regulärer Graph mit Adjazenzmatrix  $A$  und Laplacian  $L$ . Zeigen Sie, dass  $\sigma(L) = d - \sigma(A)$  gilt, genauer  $Av = \mu v$  gilt genau dann wenn  $Lv = (d - \mu)v$  gilt.

*Bemerkung:* Manchmal ist es rechnerisch einfacher, Eigenwerte und Eigenvektoren von  $A$  zu bestimmen und dann obiges zu benutzen um die Eigenwerte/Eigenvektoren des Laplacians zu bekommen.

**Aufgabe 5.2.** Bestimmen Sie jeweils den Laplacian  $L$ , dessen Eigenwerte und ein zugehöriges Orthonormalsystem für den folgenden unzusammenhängenden Graph  $G$ :



Machen Sie das gleiche für den *vollständigen* Graphen  $C_8$ , indem Sie zuerst zwei Eigenwerte erraten, dann deren Eigenräume bestimmen und dann weitermachen:



**Aufgabe 5.3.** Betrachten Sie den Graph  $G$  gegeben durch die folgende Adjazenzmatrix.

$$A = \begin{bmatrix} 0 & 0 & 1 & 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 & 1 & 0 \\ 1 & 0 & 0 & 1 & 0 & 0 & 0 & 1 \\ 1 & 0 & 1 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 & 1 & 0 & 1 \\ 0 & 1 & 0 & 0 & 1 & 0 & 0 & 1 \\ 1 & 1 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 & 1 & 0 & 0 \end{bmatrix}$$

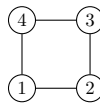
- (i) Zeichnen Sie  $G$  von Hand, indem Sie die Knoten so anordnen wie beim vollständigen Graph  $C_8$  in Aufgabe 5.2 und dann die Kanten eintragen.

- (ii) Berechnen Sie den Laplacian  $L$  und bestimmen Sie den zweitkleinsten Eigenwert  $\lambda_2$  und einen zugehörigen Eigenvektor  $v_2$ .
- (iii) Zeichnen Sie nun den Graph erneut, indem Sie die Knoten als Punkte in  $\mathbb{R}^2$  darstellen, sodass Knoten  $i$  als  $x$ -Koordinate den  $i$ -ten Eintrag von  $v_2$  hat. Die  $y$ -Koordinaten wählen Sie so, dass das Bild gut aussieht. Können Sie nun die Cluster sehen?
- (iv) Machen Sie das gleiche wie in (iii) für die zwei Graphen aus Aufgabe 5.2 und für den Graph  $G$  der durch die [hier](#) hinterlegte Adjazenzmatrix gegeben ist.
- (v) In der Vorlesung hatten wir die Heuristik propagiert, dass Punkte die nah zusammen liegen mit höherer Wahrscheinlichkeit durch eine Kante miteinander verbunden sind als Punkte die weit auseinander liegen und dass deshalb die Anordnung der Punkte die Cluster erkennen lässt. Gehen Sie nochmal durch die Beispiele und notieren Sie Effekte die zeigen, dass man vorsichtig sein muss.

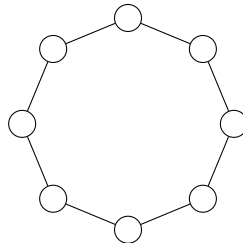
*Hinweis:* Z.B. in Pythin gibt es fertige Funktionen die anhand der Adjazenzmatrix den Laplacian bestimmen, bzw. Eigenwerte nach Größe sortiert und mit zugehörigen Eigenvektoren ausrechnen.

**Aufgabe 5.4.** Sei  $G$  ein  $d$ -regulärer Graph,  $L$  der (normale) Laplacian und  $\mathcal{L}$  der normalisierte Laplacian. Zeigen Sie  $\sigma(\mathcal{L}) = \frac{1}{d} \cdot \sigma(L)$ , genauer:  $Lv = \lambda v$  gilt genau dann wenn  $\mathcal{L}v = \frac{\lambda}{d}v$  gilt.

**Aufgabe 5.5.** Berechnen Sie die Cheeger-Konstante  $C_G$  des Graphen  $G$ :



Stellen Sie eine Vermutung auf, was die Cheeger-Konstante des *Kreises*  $K_n$



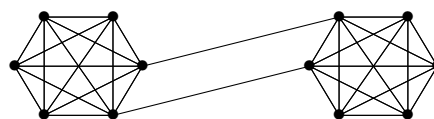
mit  $n \geq 2$  Knoten sein könnte und beweisen Sie diese.

**Aufgabe 5.6.** (i) Ein Graph  $G$  habe die Laplacematrix  $L$  und es gelte

$$\det(L - \lambda) = \lambda^3(\lambda - 2)(\lambda - 3)^4(\lambda - 4)^2$$

Wieviele Zusammenhangskomponenten hat  $G$ ?

- (ii) Bestimmen Sie die Cheeger-Konstante des folgenden Graphen:



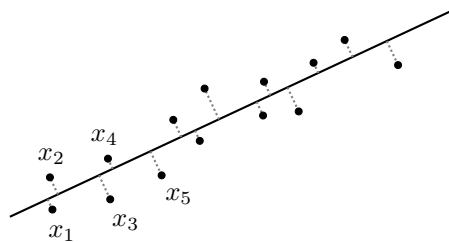
## Kapitel 6

# Bestpassende Unterräume

Im Kapitel 2 zur Regression haben wir eine affin-lineare Funktion  $f: \mathbb{R}^d \rightarrow \mathbb{R}$ , bzw. den affin-linearen Unterraum  $\text{gr}(f) \subseteq \mathbb{R}^{d+1}$ , an eine gelabelte Datenmenge  $\{(x_i, y_i) \in \mathbb{R}^d \times \mathbb{R} \mid i = 1, \dots, n\}$  angepasst. Die dort verwandte Methode der kleinsten Quadrate minimiert die Summe über die *Abstände in Labelrichtung* zwischen den Datenpunkten und dem affin-linearen Unterraum, vergleiche die Bilder vor den Sätzen 2.1 und 2.9.

Die in diesem Kapitel vorgestellte Theorie ist eng verbunden mit derjenigen der Singulärwertzerlegung, die im nachfolgenden Kapitel 7 entwickelt wird. Dort gehen wir dann auch auf die genauen Zusammenhänge ein

Im Gegensatz zur affin-linearen Regression betrachten wir in diesem Kapitel eine ungelabelte Datenmenge  $\{x_1, \dots, x_n\} \subseteq \mathbb{R}^d$  und approximieren diese durch einen Unterraum, indem wir die Summe der Quadrate der euklidischen Abstände zwischen Punkten und Unterraum minimieren. Im folgenden Bild soll also die Summe der Quadrate der Längen der gepunkteten Linien minimiert werden.



Wir beschränken uns im folgenden auf lineare Unterräume (anstelle affin-linearer) und notieren, dass man bei Datenmengen, für die dies keine guten Ergebnisse liefert, die gesamte Datenmenge zuerst zentrieren kann, vergleiche Kapitel 3.3.

**Definition 6.1.** Seien  $D := \{x_1, \dots, x_n\} \subseteq \mathbb{R}^d$  und  $1 \leq k \leq d$ . Ein Unterraum  $V \subseteq \mathbb{R}^d$  heißt *k-bestpassender Unterraum* für  $D$  falls  $\dim V = k$  ist und für jeden  $k$ -dimensionalen Unterraum  $W \subseteq \mathbb{R}^d$  die Abschätzung

$$\sum_{i=1}^n \text{dist}(x_i, V)^2 \leq \sum_{i=1}^n \text{dist}(x_i, W)^2$$

gilt, wobei  $\text{dist}(x, V) = \min_{v \in V} \|x - v\|$  der Abstand von  $x \in \mathbb{R}^d$  zu  $V$  ist.

Es ist einfach zu sehen, dass  $k$ -bestpassende Unterräume für eine gegebene Datenmenge  $D$  im Allgemeinen nicht eindeutig sind, vergleiche Aufgabe 6.3.

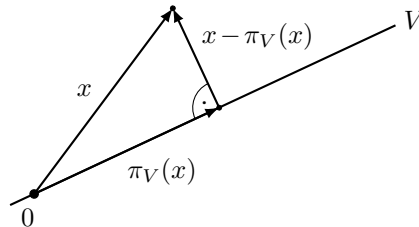
Im Folgenden sei der Raum  $\mathbb{R}^d$  stets ausgestattet mit dem Standardskalarprodukt  $\langle \cdot, \cdot \rangle$  und der davon induzierten euklidischen Norm  $\|\cdot\|$ . Wir schreiben  $x \perp y$  für orthogonale Vektoren,  $x \perp V$  wenn ein Vektor  $x$  orthogonal auf einem Unterraum  $V \subseteq \mathbb{R}^d$  steht und  $V^\perp$  für das orthogonale Komplement des Unterraums  $V$ . Wir bezeichnen mit

$$\pi_V: \mathbb{R}^d \rightarrow \mathbb{R}^d, u \mapsto \pi_V(x)$$

die Orthogonalprojektion und vermerken, dass wir diese per

$$\pi_V(x) = \sum_{j=1}^k \langle x, b_j \rangle b_j$$

berechnen können, wenn  $\{b_1, \dots, b_k\}$  eine beliebige Orthonormalbasis von  $V$  ist. Wir erinnern den Leser daran, dass in der Linearen Algebra gezeigt wird, dass stets  $x - \pi_V(x) \perp V$  und  $\pi_V(x) = \text{argmin}_{v \in V} \|x - v\|$  gelten. Betrachten wir nun die Projektion eines Datenpunktes  $x \in D$  auf einen potentiellen bestpassenden Unterraum  $V$ , so erhalten wir das folgende Bild.



Da wir bei festen Daten über  $V \subseteq \mathbb{R}^d$  optimieren, suggeriert das Bild, dass wir, anstatt das Abstandsquadrat von  $x$  zu  $U$  zu minimieren, auch das Normquadrat der Projektion maximieren können. In der Tat gilt das Folgende.

**Lemma 6.2.** *Sei eine Datenmenge  $D = \{x_1, \dots, x_n\} \in \mathbb{R}^d$  gegeben mit Koordinaten  $x_i = (x_{i1}, \dots, x_{id})$  für  $i = 1, \dots, n$ . Wir betrachten die Datenmatrix*

$$X = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} x_{11} & \cdots & x_{1d} \\ \vdots & & \vdots \\ x_{n1} & \cdots & x_{nd} \end{bmatrix} \in \mathbb{R}^{n \times d}$$

deren  $i$ -te Zeile gerade die Koordinaten des Punktes  $x_i$  enthält. Sei weiter  $1 \leq k \leq d$  und  $V \subseteq \mathbb{R}^d$  sei ein Unterraum mit Orthonormalbasis  $\{v_1, \dots, v_k\}$ . Dann ist  $V$  genau dann ein  $k$ -bestpassender Unterraum für  $D$  wenn für jedes Orthonormalsystem  $\{w_1, \dots, w_k\} \subseteq \mathbb{R}^d$  die Ungleichung

$$\sum_{j=1}^k \|Xv_j\|^2 \geq \sum_{j=1}^k \|Xw_j\|^2$$



*gilt.*

*Beweis.* Sei zuerst  $B \subseteq \mathbb{R}^d$  ein beliebiger Unterraum mit einer Orthonormalbasis  $\{b_1, \dots, b_k\}$  und sei  $\pi_B: \mathbb{R}^d \rightarrow \mathbb{R}^d$  die Orthogonalprojektion auf  $B$ . Dann gilt

$$\begin{aligned} \sum_{i=1}^n \|x_i\|^2 &= \sum_{i=1}^n \|x_i - \pi_B(x_i)\|^2 + \sum_{i=1}^n \|\pi_B(x_i)\|^2 \\ &\stackrel{\text{Pythagoras}}{=} \sum_{i=1}^n \text{dist}(x_i, B)^2 + \sum_{i=1}^n \sum_{j=1}^k |\langle x_i, b_j \rangle|^2 \\ &= \sum_{i=1}^n \text{dist}(x_i, B)^2 + \sum_{j=1}^k \|Xb_j\|^2. \end{aligned}$$

Ist nun  $V = \text{span}\{v_1, \dots, v_k\}$  wie im Lemma und  $W \subseteq \mathbb{R}^d$  ein beliebiger Unterraum mit Orthonormalbasis  $\{w_1, \dots, w_k\}$ , so folgt mit dem obigen die Äquivalenz

$$\sum_{i=1}^n \text{dist}(x_i, V)^2 \leq \sum_{i=1}^n \text{dist}(x_i, W)^2 \iff \sum_{j=1}^k \|Xv_j\|^2 \geq \sum_{j=1}^k \|Xw_j\|^2.$$

Per Definition 6.1 ist  $V$  ein  $k$ -bestpassender Unterraum genau dann wenn die linke Ungleichung für alle  $k$ -dimensionalen  $W \subseteq \mathbb{R}^d$  gilt. Dies ist folglich dazu äquivalent, dass die rechte Ungleichung für alle Orthonormalsysteme  $\{w_1, \dots, w_k\}$  gilt.  $\square$

Ist eine Datenmenge  $D = \{x_1, \dots, x_n\} \subseteq \mathbb{R}^d$  gegeben, so können wir nach Lemma 6.2 die zugehörige Datenmatrix  $X$  bilden und für diese die Optimierungsaufgabe

$$\{v_1, \dots, v_k\} \in \underset{\substack{\{\tilde{v}_1, \dots, \tilde{v}_k\} \subseteq \mathbb{R}^d \\ \text{Orthonormal-} \\ \text{basis}}}{\text{argmax}} \sum_{j=1}^n \|X\tilde{v}_j\|^2$$

bearbeiten. Ist ein entsprechender Maximierer gefunden, so erhalten wir per  $V := \text{span}\{v_1, \dots, v_k\}$  einen  $k$ -bestpassenden Unterraum. Auf den ersten Blick ist oben nicht klar, ob ein Maximierer existiert. Betrachtet man allerdings den Spezialfall  $k = 1$ , dann reduziert sich das Optimierungsproblem auf die Maximierung einer stetigen Funktion über die kompakte Einheitssphäre und die Existenz von

$$v_1 \in \underset{\|v\|=1}{\text{argmax}} \|Xv\|^2 = \underset{\|v\|=1}{\text{argmax}} \|Xv\|$$

ist gesichert. Entsprechend erhalten wir einen 1-bestpassenden Unterraum  $V_1 = \text{span}\{v_1\}$ . Die entscheidende Idee ist es nun, für  $k = 2$  nicht wieder von vorne anzufangen, sondern (etwas kühn!) darauf zu spekulieren, dass ein 2-bestpassender Unterraum  $V_2$  existiert, der  $V_1$  als Teilraum enthält. Dies führt auf das Optimierungsproblem

$$v_2 \in \underset{\substack{\|v\|=1 \\ v \perp v_1}}{\text{argmax}} \|Xv\|$$

bei welchem wieder eine stetige Funktion über eine kompakte Teilmenge, nämlich  $\partial B_1(0) \cap V_1^\perp$ , maximiert wird. Die Existenz von  $v_2$  ist also wieder gesichert. Was hingegen gezeigt werden muss, ist dass  $V_2 := \text{span}\{v_1, v_2\}$  tatsächlich ein 2-bestpassender Unterraum ist! Im folgenden Satz erledigen wir dies, und zwar nicht für  $k = 2$  sondern per mit einem Induktionsargument direkt für beliebiges  $1 \leq k \leq d$ .

**Satz 6.3.** *Sei  $D = \{x_1, \dots, x_n\} \in \mathbb{R}^d$  gegeben und  $X = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n \times d}$  die zugehörige Datenmatrix. Für  $1 \leq k \leq d$  können dann*

$$v_j \in \underset{\substack{\|v\|=1 \\ v \perp v_1, \dots, v_{k-1}}}{\operatorname{argmax}} \|Xv\|$$

*sukzessive für  $j = 1, \dots, n$  wobei die Orthogonalitätsbedingung im Fall  $j = 1$  als leer zu lesen ist. Mit diesen ist  $V_k := \text{span}\{v_1, \dots, v_k\}$  ein  $k$ -bestpassender Unterraum für  $D$ .*

*Beweis.* Da für jedes  $j$  jeweils eine stetige Funktion über eine kompakte Menge maximiert wird, ist die Existenz der  $v_1, \dots, v_k$  kein Problem. Um zu zeigen, dass  $V_k$  ein  $k$ -bestpassender Unterraum ist, definieren wir

$$K := \{k \in \{1, \dots, d\} \mid V_k \text{ ist } k\text{-bestpassender Unterraum}\}.$$

Dann gilt  $1 \in K$  nach Lemma 6.2 bzw. nach den Erläuterungen vor direkt vor Satz 6.3. Wir behaupten nun, dass die Implikation

$$\forall k \in \{1, \dots, d-1\}: k \in K \implies k+1 \in K$$

gilt. Sei hierzu  $1 \leq k < d$  mit  $k \in K$ . Wir wählen die  $v_1, \dots, v_{k+1}$  wie oben und definieren  $V_k$  und  $V_{k+1}$  als die entsprechenden Spane. Um zu zeigen, dass  $k+1 \in K$  gilt, fixieren wir einen beliebigen  $(k+1)$ -dimensionaler Unterraum  $W \subseteq \mathbb{R}^d$  und definieren

$$U := \text{span}\{\pi_W(v_1), \dots, \pi_W(v_k)\} \subseteq W.$$

Da  $\dim U \leq k$  gilt, können wir einen (suggestiv benannten!) Vektor  $w_{k+1} \in W$  mit  $w_{k+1} \perp U$  und  $\|w_{k+1}\| = 1$  wählen. Für  $i = 1, \dots, k$  gilt dann

$$\langle w, v_i \rangle = \langle w, \underbrace{v_i - \pi_W(v_i)}_{\perp W} \rangle + \langle w, \underbrace{\pi_W(v_i)}_{\in U} \rangle = 0 + 0 = 0$$

und damit  $w_{k+1} \perp v_1, \dots, v_k$ . Zusammen mit  $\|w_{k+1}\| = 1$  impliziert dies

$$\|Xw_{k+1}\|^2 \leq \max_{\substack{\|v\|=1 \\ v \perp v_1, \dots, v_k}} \|Xv\|^2 \underset{\substack{\uparrow \\ \text{Dfn} \\ v_{k+1}}}{=} \|Xv_{k+1}\|^2.$$

Andererseits können wir  $\{w_{k+1}\}$  zu einer Orthonormalbasis  $\{w_1, \dots, w_k, w_{k+1}\}$  von  $W$  ergänzen. Für das Orthonormalsystem  $\{w_1, \dots, w_k\} \subseteq \mathbb{R}^d$  gilt dann nach Lemma

6.2 die Ungleichung

$$\sum_{j=1}^k \|Aw_j\|^2 \leq \sum_{j=1}^k \|Av_j\|^2$$

weil  $V_k = \text{span}\{v_1, \dots, v_k\}$  per Voraussetzung ein  $k$ -bestpassender Unterraum ist. Addition der beiden abgesetzten Ungleichungen liefert

$$\sum_{j=1}^{k+1} \|Av_j\|^2 \geq \sum_{j=1}^{k+1} \|Aw_j\|^2,$$

woraus durch eine erneute Anwendung von Lemma 6.2 folgt, dass  $V_{k+1}$  ein  $(k+1)$ -bestpassender Unterraum ist.  $\square$

Der durch Satz 6.3 gegebenen Algorithmus zur Berechnung von  $V_k$  ist ein sogenannter *gieriger Algorithmus*: Bei der sukzessiven Wahl der  $v_j$  hätte es zum Beispiel passieren können, dass

$$\min_{\substack{V \subseteq \mathbb{R}^d \\ \dim V=2}} \sum_{i=1}^n \text{dist}(x_i, V)^2 < \min_{\substack{v_1 \in V \subseteq \mathbb{R}^d \\ \dim V=2}} \sum_{i=1}^n \text{dist}(x_i, V)^2$$

eintritt, also die ‘gierige’ Festlegung auf den ersten Basisvektor  $v_1$ , die zum bestpassenden Unterraum der Dimension 1 führt, zur Folge hat, dass das Argmin über beliebige 2-dimensionale Unterräume gar nicht mehr erreicht werden kann. Satz 6.3 zeigt gerade, dass dies *nicht* passiert und damit in diesem Fall die gierige Strategie Erfolg hat.

Ist eine Datenmenge gegeben, die man mit einem  $k$ -bestpassenden Unterraum approximieren möchte, so stellt sich in natürlicher Weise die Frage, welches  $k$  man wählen soll. Wir beginnen mit dem Folgenden.

**Lemma 6.4.** *Sei  $D = \{x_1, \dots, x_n\} \subseteq \mathbb{R}^d$  eine Datenmenge und sei  $r := \dim \text{span } D$ . Dann gilt  $V_r = \text{span } D$ .*

*Beweis.* Zunächst ist klar, dass  $V := \text{span}\{x_1, \dots, x_n\}$  ein  $r$ -bestpassender Unterraum ist, da  $\text{dist}(x_i, V) = 0$  für alle  $i$  gilt. Für den  $r$ -bestpassenden Unterraum  $V_r$  muss per Definition dann

$$\sum_{i=1}^n \text{dist}(x_i, V_r)^2 = \sum_{i=1}^n \text{dist}(x_i, V)^2 = 0$$

gelten. Da  $V_r \subseteq \mathbb{R}^d$  abgeschlossen ist, geht dies nur wenn  $x_i \in V_r$  für alle  $i$  gilt. Aus Dimensionsgründen folgt dann  $V_r = V$ .  $\square$

Sinnvollerweise ist also  $k < r$  zu wählen. Für  $k = r$  erhält man den Span der Datenmenge und bei einer Wahl von  $k > r$  sind die nach der Rekursionsformel in Satz 6.3 gewählten  $v_{r+1}, \dots, v_k$  nach obigem lediglich eine beliebige Ergänzung von  $v_1, \dots, v_k$  zu einem Orthonormalsystem.

Ist eine Datenmenge  $D$  gegeben, so ist es eventuell nicht-trivial, die Dimension von  $\text{span } D$  zu ermitteln, womit obiges nur bedingt nützlich ist. Der folgende Satz zeigt, wie man bei der rekursiven Berechnung der  $v_1, v_2, \dots$  erkennt, wann man bei  $v_r$  angekommen ist.

**Proposition 6.5.** *Sei  $D = \{x_1, \dots, x_n\} \in \mathbb{R}^d$  gegeben und  $X = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n \times d}$  die zugehörige Datenmatrix. Für  $1 \leq k \leq d$  sei*

$$\sigma_k = \|Xv_k\| = \max_{\substack{\|v\|=1 \\ v \perp v_1, \dots, v_{k-1}}} \|Xv\|$$

wobei wir wieder die Orthogonalitätsbedingung im Fall  $k = 1$  als leer auffassen. Dann gilt  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$  und  $\sigma_{r+1} = \dots = \sigma_d = 0$  mit  $r = \dim \text{span } D$ .

*Beweis.* Dass die  $\sigma_k$  fallend und nicht-negativ sind, ist per Konstruktion klar. Für  $k \geq r + 1$  folgt

$$\sigma_k^2 = \max_{\substack{\|v\|=1 \\ v \perp v_1, \dots, v_{k-1}}} \|Xv\|^2 = \min_{\substack{\|v\|=1 \\ v \perp V_{k-1}}} \sum_{i=1}^n \text{dist}(x_i, V_{k-1})^2 = 0$$

nach Lemma 6.4 und wegen  $V_{k-1} \supseteq V_r$ , also  $\sigma_{r+1}, \dots, \sigma_d = 0$ . Wir behaupten nun, dass andererseits  $\sigma_1, \dots, \sigma_r > 0$  gilt. Dies sieht man am besten per Widerspruch: Angenommen, es gibt ein  $1 \leq k \leq r$  mit  $\sigma_k = \|Xv_k\| = 0$ . Dann sind wegen  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_d \geq 0$  auch  $\sigma_{k+1}, \dots, \sigma_d$  gleich Null. Letzteres heißt dann

$$\|Xv_k\| = \|Xv_{k+1}\| = \dots = \|Xv_d\| = 0.$$

Damit haben wir  $v_k, \dots, v_d \in \ker X$  und folglich gilt  $\dim \ker X \geq d - (k - 1)$ . Weiter gilt  $\text{rk } X = \dim \text{span}\{x_1, \dots, x_n\} = r$  (Zeilenrang=Spaltenrang). Es folgt

$$\dim \ker X + \dim \text{ran } X \geq (d - k + 1) + r = d + 1 + (r - k) \geq d + 1$$

$\uparrow$   
 $k \leq r$

im Widerspruch zum Rangsatz. □

Die  $\sigma_j$  in Proposition 6.5 sind wichtige Kennzahlen der Matrix  $X$ , die wir im folgenden Kapitel zur *Singulärwertzerlegung* noch genauer untersuchen werden. Insbesondere werden wir in Korollar 7.19 die Approximationsqualität von  $V_k$ , d.h. die Summe  $\text{dist}(x_1, V_k)^2 + \dots + \text{dist}(x_n, V_k)^2$  via der  $\sigma_{k+1}, \dots, \sigma_r$  ausdrücken.

## Referenzen

Dieses Kapitel folgt hauptsächlich [BHK20, Chapter 3], wurde aber im Verlauf von mehreren durch den Autor gehaltenen Vorlesungen immer wieder verändert. Insbesondere werden im vorliegenden Text zunächst die bestpassenden Unterräume getrennt von der Singulärwertzerlegung eingeführt, die dann im folgenden Kapitel 7 behandelt wird. Bei Aufgabe 6.1 empfiehlt sich auch ein Blick auf [Ham].

## Aufgaben

**Aufgabe 6.1.** Berechnen Sie für  $(3, 2, 2), (2, 3, -2) \in \mathbb{R}^3$  die Vektoren  $v_1, v_2, v_3$  mithilfe der Rekursionsformel in Satz 6.3. Erstellen Sie einen Plot des 1-dimensionalen und 2-dimensionalen bestpassenden Unterraumes mit einem geeigneten Programm. Daher

*Hinweis:* Verwenden Sie Lagrange Multiplier und prüfen Sie Zwischenergebnisse per CAS.

**Aufgabe 6.2.** Bestimmen Sie, z.B. durch Benutzung geeigneter Pakete in Python, einen bestpassenden 1-dimensionalen Unterraum zu den Daten aus Aufgabe 2.3 und plotten Sie diesen. Vergleichen Sie mit der Regressionsgerade aus Aufgabe 2.3.

**Aufgabe 6.3.** Zeigen Sie dass  $n$ -viele Datenpunkte in  $\mathbb{R}^d$  einen  $k$ -dimensionalen bestpassenden Unterraum im Allgemeinen nicht eindeutig bestimmen. Dies gilt selbst dann, wenn der Rang der Datenmatrix  $X \in \mathbb{R}^{n \times d}$  echt größer als  $k$  ist.

*Hinweis:* Ein Beispiel in dem die Nichteindeutigkeit anschaulich klar ist, ist schnell gefunden. Für den Beweis empfiehlt es sich, Methode 1 oder Methode 1 $\frac{1}{2}$  aus dem nachfolgenden Kapitel 7 zu benutzen.

## Kapitel 7

# Singulärwertzerlegung

Ist  $A \in \mathbb{R}^{n \times n}$  eine symmetrische Matrix mit Spektrum  $\sigma(A) = \{\lambda_1, \dots, \lambda_n\}$ , so existiert nach dem aus der Linearen Algebra wohlbekannten Satz über die orthogonale Diagonalisierbarkeit stets eine orthogonale Matrix  $B \in \mathbb{R}^{n \times n}$ , sodass die Matrixgleichung

$$B^T A B = \text{diag}(\lambda_1, \dots, \lambda_n)$$

gilt. In diesem Kapitel verallgemeinern wir das obige auf beliebige nicht-quadratische Matrizen  $A \in \mathbb{R}^{n \times d}$ . Dazu müssen wir zunächst den Begriff des Eigenwertes verallgemeinern; beachte insbesondere dass die übliche Eigenwertgleichung  $Av = \lambda v$  für  $A \in \mathbb{R}^{n \times d}$  mit  $n \neq d$  nicht wohldefiniert ist.

**Definition 7.1.** Sei  $A \in \mathbb{R}^{n \times d}$ . Eine Zahl  $\sigma > 0$  heißt *Singulärwert* von  $A$ , falls  $v \in \mathbb{R}^d \setminus \{0\}$  und  $u \in \mathbb{R}^n \setminus \{0\}$  existieren mit  $Av = \sigma u$  und  $A^T u = \sigma v$ . Wir nennen  $v$  und  $u$  zu  $\sigma$  gehörende *Singulärvektoren*; genauer heißt  $v$  ein *rechter* und  $u$  ein *linker Singulärvektor*. Die Menge aller Singulärwerte bezeichnen wir mit  $s(A)$ .

Singulärwerte und Singulärvektoren einer Matrix  $A \in \mathbb{R}^{n \times d}$  sind eng verbunden mit den Eigenwerten und Eigenvektoren der Matrizen  $A^T A \in \mathbb{R}^{d \times d}$  und  $AA^T \in \mathbb{R}^{n \times n}$ . Für spätere Nutzung notieren wir die folgenden Eigenschaften.

**Lemma 7.2.** Sei  $A \in \mathbb{R}^{n \times d}$ . Dann gelten:

- (i)  $A^T A$  und  $AA^T$  sind symmetrisch und positiv semidefinit.
- (ii) Es gilt  $\sigma(A^T A) \setminus \{0\} = \sigma(AA^T) \setminus \{0\}$ .
- (iii) Für Eigenwerte  $\lambda \neq 0$  gilt  $\dim \ker(A^T A - \lambda) = \dim \ker(AA^T - \lambda)$ .
- (iv) Es gilt  $\text{rk } A = \text{rk } A^T A = \text{rk } AA^T = \text{rk } A^T$ .
- (v) Mit (geometrischer = algebraischer) Vielfachheit gezählt, gibt es genau  $\text{rk}(A)$ -viele von Null verschiedene Eigenwerte.

*Beweis.* (i) Die Symmetrie ist klar; um die positive Semidefinitheit zu sehen, berechnen wir für  $x \in \mathbb{R}^d$  und  $y \in \mathbb{R}^n$  jeweils

$$\langle x, A^T A x \rangle = \langle Ax, Ax \rangle = \|Ax\|^2 \geq 0 \quad \text{und} \quad \langle y, AA^T y \rangle = \langle A^T y, A^T y \rangle = \|A^T y\|^2 \geq 0.$$

(ii) Sei  $\lambda \in \sigma(A^\top A) \setminus \{0\}$ . Dann existiert  $v \in \mathbb{R}^d \setminus \{0\}$  mit  $A^\top A v = \lambda v$ . Multiplikation mit  $A$  von links liefert  $AA^\top A v = \lambda A v$ . Wir setzen  $u := A v \in \mathbb{R}^n$  und erhalten  $AA^\top u = \lambda u$  mit

$$\|u\|^2 = \|A v\|^2 \langle A v, A v \rangle = \langle v, A^\top A v \rangle = \langle v, \lambda v \rangle = \lambda \|v\|^2 \neq 0$$

also  $u \neq 0$  und daher  $\lambda \in \sigma(AA^\top) \setminus \{0\}$ . Die andere Inklusion zeigt man analog.

(iii) Seien  $v_1, \dots, v_k \in \ker(A^\top A - \lambda)$  linear unabhängig. Wir setzen  $u_j := A v_j$  und notieren  $u_j \in \ker(AA^\top - \lambda)$  für  $j = 1, \dots, k$  wegen (ii). Sei jetzt  $\alpha_1 u_1 + \dots + \alpha_k u_k = 0$ . Dann folgt

$$0 = \sum_{j=1}^k \alpha_j A v_j = \sum_{j=1}^k \alpha_j \lambda v_j = \lambda \sum_{j=1}^k \alpha_j v_j \xrightarrow{\lambda \neq 0} \sum_{j=1}^k \alpha_j v_j = 0$$

und damit  $\alpha_1 = \dots = \alpha_k = 0$ , was zeigt, dass  $u_1, \dots, u_k$  linear unabhängig sind. Analog sieht man, dass linear unabhängige  $u_1, \dots, u_k \in \ker(AA^\top - \lambda)$  per  $v_j := A^\top u_j$  auf  $k$ -viele linear unabhängige Vektoren in  $\ker(A^\top A - \lambda)$  führen.

(iv) Wir zeigen, dass  $\ker A = \ker A^\top A$  gilt; daraus folgt  $\text{rk } A = \text{rk } A^\top A$  nach dem Rangsatz. Die Inklusion ' $\subseteq$ ' ist klar. Für ' $\supseteq$ ' gelte  $A^\top A x = 0$ . Dann folgt

$$0 = \langle x, A^\top A x \rangle = \langle A x, A x \rangle = \|A x\|^2$$

und damit  $A x = 0$ . Die mittlere behauptete Gleichung ist klar und die letzte zeigt man analog zur ersten.

(iv) Nach dem Satz über orthogonale Diagonalisierbarkeit gibt es eine orthogonale Matrix  $B \in \mathbb{R}^{d \times d}$  mit

$$B^\top (A^\top A) B = \text{diag}(\lambda_1, \dots, \lambda_n)$$

wobei  $\lambda_1, \dots, \lambda_n$  die Eigenwerte von  $A^\top A$  mit Vielfachheiten sind. Die rechte Seite der Gleichung hat also als Rang die Anzahl der von Null verschiedenen Eigenwerte wenn diese mit Vielfachheiten gezählt werden. Der Rang der linken Seite in der Matrixgleichung ist gleich  $\text{rk } A^\top A$ .  $\square$

Jetzt stellen wir die Verbindung zwischen Singulärwerten und Singulärvektoren von  $A$  zu Eigenwerten und Eigenvektoren von  $AA^\top$  und  $A^\top A$  her.

**Lemma 7.3.** *Sei  $A \in \mathbb{R}^{n \times d}$ . Es gilt:*

- (i) *Sind  $u$  und  $v$  Singulärvektoren zum Singulärwert  $\sigma \in s(A)$ , so ist  $v$  Eigenvektor von  $A^\top A$  und  $u$  Eigenvektor von  $AA^\top$  jeweils zum Eigenwert  $\sigma^2$ .*

*Umgekehrt haben wir:*

- (ii) *Ist  $v$  Eigenvektor von  $A^\top A$  zum positiven Eigenwert  $\lambda > 0$ , dann sind  $v$  und  $u := \frac{1}{\sqrt{\lambda}} A v$  Singulärvektoren zum Singulärwert  $\sqrt{\lambda}$ .*
- (iii) *Ist  $u$  Eigenvektor von  $AA^\top$  zum positiven Eigenwert  $\lambda > 0$ , dann sind  $u$  und  $v := \frac{1}{\sqrt{\lambda}} A^\top u$  Singulärvektoren zum Singulärwert  $\sqrt{\lambda}$ .*

*Beweis.* (i) Es gelten

$$AA^\top v = A^\top \sigma u = \sigma \sigma v = \sigma^2 v \quad \text{und} \quad AA^\top u = A \sigma v = \sigma \sigma u = \sigma^2 u.$$

(ii) Hier haben wir  $Av = \sqrt{\lambda}u$  per Definition und

$$A^\top u = A^\top \left( \frac{1}{\sqrt{\lambda}} Av \right) = \frac{1}{\sqrt{\lambda}} A^\top Av = \frac{\lambda}{\sqrt{\lambda}} v = \sqrt{\lambda} v.$$

(iii) Wieder gilt  $A^\top u = \sqrt{\lambda}u$  per Definition und

$$Av = A \left( \frac{1}{\sqrt{\lambda}} A^\top u \right) = \frac{1}{\sqrt{\lambda}} AA^\top u = \frac{\lambda}{\sqrt{\lambda}} u = \sqrt{\lambda} u$$

was den Beweis beendet.  $\square$

Mit Lemma 7.3 folgt also schonmal, dass es für  $A \in \mathbb{R}^{n \times d}$  höchstens  $p := \min(n, d)$ -viele paarweise verschiedene Singulärwerte geben kann. Wie auch bei Eigenwerten, wollen wir Singulärwerte allerdings im folgenden ‘mit Vielfachheiten’ zählen. Dass dies möglich ist, folgt ebenfalls aus Lemma 7.3 in Kombination mit Lemma 7.2. Nach diesen gilt nämlich für  $\sigma \in s(A)$

$$\begin{aligned} \text{vfh}(\sigma) &:= \dim \left\{ v \in \mathbb{R}^d \mid \begin{array}{l} v = 0 \text{ oder } v \text{ rechter} \\ \text{SV von } A \text{ zum SW } \sigma \end{array} \right\} = \dim \left\{ u \in \mathbb{R}^n \mid \begin{array}{l} u = 0 \text{ oder } u \text{ linker} \\ \text{SV von } A \text{ zum SW } \sigma \end{array} \right\} \\ &= \text{Vielfachheit des Eigenwertes } \sigma^2 \text{ von } A^\top A, \text{ oder, äquivalent, von } AA^\top. \end{aligned}$$

Wir nennen die obige Zahl die *Vielfachheit des Singulärwertes*  $\sigma > 0$ . Zählen wir die Vielfachheiten mit, so gibt es nach Lemma 7.2 genau  $r := \text{rk}(A)$ -viele Singulärwerte.

**Lemma 7.4.** *Sei  $A \in \mathbb{R}^{n \times d}$ . Sei  $\{v_1, \dots, v_r\}$  ein Orthonormalsystem aus Eigenvektoren von  $A^\top A$  zu den Eigenwerten  $\sigma_1^2 \geq \dots \geq \sigma_r^2 > 0$ . Dann ist  $\{u_1, \dots, u_r\}$  mit  $u_i = \frac{1}{\sigma_i} Av_i$  ein Orthonormalsystem. Analoges gilt, wenn  $\{u_1, \dots, u_r\}$  gegeben ist, für  $\{v_1, \dots, v_r\}$  mit  $v_i = \frac{1}{\sigma_i} A^\top u_i$ .*

*Beweis.* Es gilt

$$\langle u_i, u_j \rangle = \left\langle \frac{1}{\sigma_i} Av_i, \frac{1}{\sigma_j} Av_j \right\rangle = \frac{1}{\sigma_i \sigma_j} \langle v_i, A^\top Av_j \rangle = \frac{\sigma_j}{\sigma_i} \langle v_i, v_j \rangle,$$

d.h. für  $i \neq j$  gilt  $\langle u_i, u_j \rangle = 0$  und für  $i = j$  gilt  $\|u_j\|^2 = 1$ , also  $\|u_j\| = 1$ . Den zweiten Teil erledigt man mit einer analogen Rechnung.  $\square$

Ergänzen wir die Orthonormalsysteme aus Lemma 7.4 zu Orthonormalbasen  $\mathcal{V} = \{v_1, \dots, v_d\}$  und  $\mathcal{U} = \{u_1, \dots, u_n\}$ , so folgt, dass die Darstellungsmatrix der linearen Abbildung  $\mathbb{R}^d \rightarrow \mathbb{R}^n$ ,  $x \mapsto Ax$ , gerade die mit Nullen geeignet aufgefüllte ‘Diagonalmatrix’  $\text{diag}(\sigma_1, \dots, \sigma_r) \in \mathbb{R}^{n \times d}$  ist, wenn wir  $\mathbb{R}^d$  mit der Basis  $\mathcal{V}$  und  $\mathbb{R}^n$  mit der Basis  $\mathcal{U}$  ausstatten. Organisieren wir die Basisvektoren in Matrizen, so erhalten wir eine Faktorisierung von  $A$  durch eine Diagonalmatrix.



**Definition 7.5.** Sei  $A \in \mathbb{R}^{n \times d}$  und  $r = \text{rk } A$ . Eine Faktorisierung

$$A = U \Sigma V^\top$$

mit orthogonalen Matrizen  $V \in \mathbb{R}^{d \times d}$ ,  $U \in \mathbb{R}^{n \times n}$  und einer Diagonalmatrix  $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_r) \in \mathbb{R}^{n \times d}$  mit  $\sigma_1 \geq \dots \geq \sigma_r > 0$  heißt *Singulärwertzerlegung*, oder kurz *SVD*, von  $A$ .

**Bemerkung 7.6.** (i) Da  $U$  und  $V$  orthogonal sind, können wir die Gleichung  $A = U \Sigma V^\top$  in Definition 7.5 zu

$$U^\top A V = \text{diag}(\sigma_1, \dots, \sigma_r) \in \mathbb{R}^{n \times d}$$

umstellen und erhalten das zu Beginn des Kapitels angekündigte Analogon des orthogonalen Diagonalisierungssatzes.

(ii) Setzt man  $\sigma_{r+1} := \dots := \sigma_p := 0$  für  $p = \min(n, d)$ , so gilt für die Diagonalmatrix in der Singulärwertzerlegung

$$\Sigma \in \left\{ \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_p \end{bmatrix}, \begin{bmatrix} \sigma_1 & & 0 & \dots & 0 \\ & \ddots & \vdots & & \vdots \\ & & \sigma_p & 0 & \dots & 0 \end{bmatrix}, \begin{bmatrix} \sigma_1 & & & \\ & \ddots & & \\ & & \sigma_p & \\ 0 & \dots & 0 & \\ \vdots & & \vdots & \\ 0 & \dots & 0 & \end{bmatrix} \right\}$$

je nachdem ob  $n = d$ ,  $n < d$  oder  $n > d$  ausfällt. Beachte aber, dass die  $\sigma_{r+1}, \dots, \sigma_p$  nach Definition 7.1 keine Singulärwerte sind.

Wir notieren die folgende einfache aber sehr hilfreiche Darstellung von Matrix-Matrix-Matrix-Produkten der Form  $U \Sigma V^\top$ .

**Lemma 7.7.** Seien  $U = (u_{ij}) \in \mathbb{R}^{n \times n}$  und  $V = (v_{ij}) \in \mathbb{R}^{d \times d}$  orthogonale Matrizen mit Spalten  $u_1, \dots, u_n \in \mathbb{R}^n$  bzw.  $v_1, \dots, v_d \in \mathbb{R}^d$ . Sei  $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_p) \in \mathbb{R}^{n \times d}$  eine Matrix mit  $\sigma_i \in \mathbb{R}$  und  $p = \min(n, d)$ . Dann gelten

$$U \Sigma V^\top = \sum_{i=1}^p u_i \sigma_i v_i^\top \quad \text{und} \quad (U \Sigma V^\top)_{ij} = \sum_{\ell=1}^p u_{i\ell} \sigma_\ell v_{j\ell}.$$

Ist  $1 \leq k \leq p$  und haben wir  $\sigma_{r+1} = \dots = \sigma_p = 0$ , so gilt die obige Gleichung mit  $k$  statt  $p$  auf der rechten Seite.

*Beweis.* Wir rechnen nach

$$\begin{aligned} \sum_{i=1}^p u_i \sigma_i v_i^\top &= \sum_{i=1}^p \begin{bmatrix} u_{1i} \\ \vdots \\ u_{ni} \end{bmatrix} [\sigma_i v_{1i} \dots \sigma_i v_{di}] \\ &= \sum_{i=1}^p \begin{bmatrix} u_{1i} \sigma_i v_{1i} & \dots & u_{1i} \sigma_i v_{di} \\ \vdots & & \vdots \\ u_{ni} \sigma_i v_{1i} & \dots & u_{ni} \sigma_i v_{di} \end{bmatrix} \\ &= \begin{bmatrix} \sum_{i=1}^p u_{1i} \sigma_i v_{1i} & \dots & \sum_{i=1}^p u_{1i} \sigma_i v_{di} \\ \vdots & & \vdots \\ \sum_{i=1}^p u_{ni} \sigma_i v_{1i} & \dots & \sum_{i=1}^p u_{ni} \sigma_i v_{di} \end{bmatrix} \end{aligned}$$

$$\begin{aligned}
&= \begin{bmatrix} u_{11} & \cdots & u_{1p} \\ \vdots & & \vdots \\ u_{n1} & \cdots & u_{np} \end{bmatrix} \cdot \begin{bmatrix} \sigma_1 v_{11} & \cdots & \sigma_1 v_{d1} \\ \vdots & & \vdots \\ \sigma_p v_{1p} & \cdots & \sigma_p v_{dp} \end{bmatrix} \\
&= \begin{bmatrix} u_{11} & \cdots & u_{1p} \\ \vdots & & \vdots \\ u_{n1} & \cdots & u_{np} \end{bmatrix} \cdot \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_p \end{bmatrix} \cdot \begin{bmatrix} v_{11} & \cdots & v_{d1} \\ \vdots & & \vdots \\ v_{1p} & \cdots & v_{dp} \end{bmatrix} = U \Sigma V^\top,
\end{aligned}$$

wobei die letzte Gleichung im Fall  $n = p$  klar ist. In den anderen beiden Fällen, vergleiche Bemerkung 7.6(ii), sieht man, dass das Abschneiden von Zeilen von  $V^\top$  bzw. von Spalten von  $U$  keinen Einfluß auf die Einträge von  $U \Sigma V^\top$  hat.  $\square$

Wir sind jetzt bereit für unseren ersten Satz zur Singulärwertzerlegung.

**Satz 7.8.** (Methode 1 zur Bestimmung einer SVD) *Sei  $A \in \mathbb{R}^{n \times d}$  und  $r = \text{rk } A$ . Die folgenden Schritte führen zu einer Singulärwertzerlegung von  $A$ :*

1. *Bilde die Matrix  $A^\top A$ , bestimme deren echt positive Eigenwerte und nummeriere diese absteigend  $\lambda_1 \geq \cdots \geq \lambda_r > 0$  (mit Vielfachheiten ergeben sich genau  $r$ -viele; insbesondere muss man  $r$  nicht a priori kennen!). Bestimme ein zugehöriges Orthonormalsystem aus Eigenvektoren  $v_1, \dots, v_r$ .*
2. *Setze  $\sigma_i = \sqrt{\lambda_i}$  für  $i = 1, \dots, r$ .*
3. *Setze  $u_i = \frac{1}{\sigma_i} A v_i$  für  $i = 1, \dots, r$ ; erhalte ein Orthonormalsystem  $\{u_1, \dots, u_r\}$ .*
4. *Ergänze zu Orthonormalbasen  $\mathcal{V} = \{v_1, \dots, v_d\}$  und  $\mathcal{U} = \{u_1, \dots, u_n\}$ .*
5. *Setze  $V = [v_1 \cdots v_d]$ ,  $U = [u_1 \cdots u_n]$  und  $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_r) \in \mathbb{R}^{n \times d}$ .*

*Hat  $AA^\top$  kleineres Format als  $A^\top A$ , so kann es einfacher sein, mit dieser Matrix zu beginnen und ansonsten analog zu verfahren. Dies bezeichnen wir als Methode 1 $\frac{1}{2}$ .*

*Beweis.* Dass die Schritte 1–5 auf zwei orthogonale Matrizen  $U$  und  $V$  führen, deren Spalten Singulärvektoren zu den Singulärwerten  $\sigma_i$  sind, folgt aus den Lemmas 7.2, 7.3 und 7.4. Gezeigt werden muss, dass  $A = U \Sigma V^\top$  gilt. Nach Lemma 7.7 genügt es hierfür zu prüfen, dass die beiden linearen Abbildungen

$$A: \mathbb{R}^d \rightarrow \mathbb{R}^n \quad \text{und} \quad B := \sum_{i=1}^r u_i \sigma_i v_i^\top: \mathbb{R}^d \rightarrow \mathbb{R}^n$$

auf der Basis  $\mathcal{V}$  übereinstimmen. Für  $1 \leq j \leq r$  gilt

$$\left( \sum_{i=1}^r u_i \sigma_i v_i^\top \right) v_j = \sum_{i=1}^r u_i \sigma_i v_i^\top v_j = \sum_{i=1}^r u_i \sigma_i \langle v_i, v_j \rangle = \sigma_j u_j = A v_j.$$

Für  $r < j \leq d$  haben wir  $B v_j = 0$  und es bleibt zu zeigen, dass  $A v_j = 0$  gilt. Wir wissen  $\sigma_1, \dots, \sigma_r \neq 0$  und wegen  $A^\top u_i = \sigma_i v_i$  müssen die  $v_1, \dots, v_r$  im Bild von  $A^\top$  liegen. Da  $\text{ran } A^\top$  Dimension  $r$  hat, sind also  $v_{r+1}, \dots, v_d$  orthogonal zu  $\text{ran } A^\top$ . D.h. für beliebiges  $u \in \mathbb{R}^d$  und  $r < j \leq d$  haben wir

$$\langle A v_j, u \rangle = \langle v_j, A^\top u \rangle = 0$$

und dies geht nur, wenn  $A v_j$  für die genannten  $j$  Null ist.  $\square$

Durch die multiplen Wahlmöglichkeiten ist klar, dass die Matrizen  $U$  und  $V$  im Allgemeinen nicht eindeutig durch  $A$  bestimmt sind. Die Matrix  $\Sigma$  ist es allerdings, und zwar aufgrund unserer Konvention die Singulärwerte  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$  immer absteigend zu sortieren.

**Proposition 7.9.** *Gelte  $A = U\Sigma V^\top$  mit orthogonalen Matrizen  $U \in \mathbb{R}^{n \times n}$ ,  $V = \mathbb{R}^{d \times d}$  und einer Diagonalmatrix  $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_p) \in \mathbb{R}^{n \times d}$  wobei  $p = \min(n, d)$  ist. Sei  $r = \text{rk } A$ . Dann sind genau  $r$ -viele der  $\sigma_i$  gleich Null. Die restlichen sind die Singulärwerte von  $A$  und kommen in der Diagonalmatrix entsprechend ihrer Vielfachheit oft vor.*

*Beweis.* Da  $U$  und  $V$  vollen Rang haben, folgt  $\text{rk } \Sigma = \text{rk } A = r$  und damit die erste Aussage. Für  $j \in \{1, \dots, p\}$  mit  $\sigma_j \neq 0$  gelten nach Lemma 7.7

$$Av_j = U\Sigma V^\top v_j = \sum_{i=1}^r u_i \sigma_i v_i^\top v_j = \sigma_j u_j$$

und

$$A^\top u_j = (U\Sigma V^\top)^\top u_j = \sum_{i=1}^r (u_i \sigma_i v_i^\top)^\top u_j = \sum_{i=1}^r v_i \sigma_i u_i^\top u_j = \sigma_j v_j.$$

Folglich ist  $\sigma_j$  Singulärwert mit der  $j$ -ten Spalte  $v_j$  von  $V$  und der  $j$ -ten Spalte  $u_j$  von  $U$  als zugehörigen Singulärvektoren.  $\square$

Als nächstes stellen wir den bereits am Ende von Kapitel 6 angekündigten Zusammenhang zwischen Singulärwertzerlegungen und bestpassenden Unterräumen her. Hierfür benötigen wir die Courant-Fischer-Formel, diesmal aber in leicht anderer Version als in Satz 5.13, nämlich für absteigend durchnummerierte Eigenwerte.

**Satz 7.10.** (Courant-Fischer-Formel für absteigend nummerierte Eigenwerte) *Sei  $M \in \mathbb{R}^{n \times n}$  eine symmetrische Matrix mit Eigenwerten  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$  und einer Orthonormalbasis aus Eigenvektoren  $\{v_1, \dots, v_n\}$ . Dann gilt für  $k = 1, \dots, n$*

$$\lambda_k = \max_{\substack{x \neq 0 \\ \langle x, v_i \rangle = 0 \\ \text{für } i=1, \dots, k-1}} \frac{\langle x, Mx \rangle}{\langle x, x \rangle} \quad \text{und} \quad v_k \in \operatorname{argmax}_{\substack{x \neq 0 \\ \langle x, v_i \rangle = 0 \\ \text{für } i=1, \dots, k-1}} \frac{\langle x, Mx \rangle}{\langle x, x \rangle}$$

wobei die Orthogonalitätsbedingung für  $k = 1$  als leer zu lesen ist.

*Beweis.* Seien  $M$  und  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$  wie angegeben. Die Matrix  $-M$  hat dann die Eigenwerte  $-\lambda_1 \leq -\lambda_2 \leq \dots \leq -\lambda_n$  und die  $v_1, \dots, v_n$  bilden (ohne Vorzeichenänderung!) eine Orthonormalbasis aus Eigenvektoren für  $-M$ . Nach Satz 5.13 gilt dann für  $k = 1, \dots, n$

$$-\lambda_k = \min_{\substack{x \neq 0 \\ \langle x, v_i \rangle = 0 \\ \text{für } i=1, \dots, k-1}} \frac{\langle x, -Mx \rangle}{\langle x, x \rangle} = - \max_{\substack{x \neq 0 \\ \langle x, v_i \rangle = 0 \\ \text{für } i=1, \dots, k-1}} \frac{\langle x, Mx \rangle}{\langle x, x \rangle},$$

sowie

$$v_k \in \underset{\substack{x \neq 0 \\ \langle x, v_i \rangle = 0 \\ \text{für } i=1, \dots, k-1}}{\operatorname{argmin}} \frac{\langle x, -Mx \rangle}{\langle x, x \rangle} = \underset{\substack{x \neq 0 \\ \langle x, v_i \rangle = 0 \\ \text{für } i=1, \dots, k-1}}{\operatorname{argmin}} \frac{\langle x, Mx \rangle}{\langle x, x \rangle}$$

was den Beweis beendet.  $\square$

Als erstes zeigen wir, dass die in Kapitel 6 vorgestellte algorithmische Berechnungsmethode für den  $k$ -bestpassenden Unterraum eine Singulärwertzerlegung liefert.

**Satz 7.11.** (Methode 2 zur Bestimmung einer SVD) *Sei  $A \in \mathbb{R}^{n \times d}$  und  $r = \operatorname{rk} A$ . Die folgenden Schritte führen zu einer Singulärwertzerlegung von  $A$ :*

1. Für  $k \geq 1$  bestimme iterativ

$$\sigma_k = \max_{\substack{\|v\|=1 \\ v \perp v_1, \dots, v_{k-1}}} \|Av\|, \quad v_k \in \underset{\substack{\|v\|=1 \\ v \perp v_1, \dots, v_{k-1}}}{\operatorname{argmax}} \|Av\|, \quad u_k := \frac{1}{\sigma_k} Av_k$$

*solange  $\sigma_k \neq 0$  gilt (dies ist genau für  $k = 1, \dots, r$  der Fall; man muss also  $r$  nicht a priori kennen!).*

2. Ergänze zu Orthonormalbasen  $\{v_1, \dots, v_d\}$  und  $\{u_1, \dots, u_n\}$ .

3. Setze  $V = [v_1 \cdots v_d]$ ,  $U = [u_1 \cdots u_d]$  und  $\Sigma = \operatorname{diag}(\sigma_1, \dots, \sigma_r) \in \mathbb{R}^{n \times d}$ .

*Startet man mit  $A^\top$  anstelle von  $A$  und geht ansonsten analog vor, so erhält man ebenfalls eine Singulärwertzerlegung. Dies bezeichnen wir als Methode 2 $\frac{1}{2}$ .*

*Beweis.* Seien  $\lambda_1 \geq \dots \geq \lambda_r > 0$  wie in Satz 7.8 die echt positiven Eigenwerte von  $A^\top A$  und  $\{v_1, \dots, v_r\}$  ein Orthonormalsystem aus Eigenvektoren. Die Courant-Fischer Formel 7.10 (mit Substitution  $v := x/\|x\|$ ) impliziert dann

$$\sigma_k = \lambda_k^{1/2} = \max_{\substack{\|v\|=1 \\ v \perp v_1, \dots, v_{k-1}}} (\langle v, A^\top A v \rangle)^{1/2} = \max_{\substack{\|v\|=1 \\ v \perp v_1, \dots, v_{k-1}}} (\langle Av, Av \rangle)^{1/2} = \max_{\substack{\|v\|=1 \\ v \perp v_1, \dots, v_{k-1}}} \|Av\|.$$

Per Konstruktion sind die  $v_1, \dots, v_r$  ein Orthonormalsystem. Gezeigt werden muss jetzt also, dass  $v_i$  Eigenvektor zum Eigenwert  $\lambda_i$  ist für  $i = 1, \dots, r$ . Da  $A^\top A$  symmetrisch ist, können wir eine orthogonale Matrix  $B \in \mathbb{R}^{d \times d}$  wählen mit  $B^\top A^\top A B = \operatorname{diag}(\lambda_1, \dots, \lambda_r, 0, \dots, 0)$ . Wir betrachten die isometrische Bijektion

$$B^\top: \mathbb{R}^d \rightarrow \mathbb{R}^d, \quad v \mapsto B^\top v$$

und vermerken, dass  $v \in \ker(A^\top A - \lambda_i)$  genau dann gilt wenn  $B^\top v \in \ker(\operatorname{diag}(\dots) - \lambda_i)$  ist. Letzteres gerade heißt gerade, dass  $(B^\top v)_j = 0$  ist für diejenigen  $j$  mit  $\lambda_j \neq \lambda_i$ . Weiter berechnen wir

$$\|Av\|^2 = \|ABB^\top v\|^2 = \|B^\top ABB^\top v\|^2 = \|\operatorname{diag}(\dots) B^\top v\|^2 = \sum_{j=1}^r \lambda_j^2 (B^\top v)_j^2.$$

Sei  $i_1$  die Vielfachheit von  $\lambda_1$  und  $1 \leq i \leq i_1$ . Wir behaupten

$$\operatorname{argmax}_{\|x\|=1} \sum_{j=1}^r \lambda_j^2 x_j^2 = \{x \in \partial B_1(0) \mid x_j = 0 \text{ für } j > i_1\} =: M_1.$$

In der Tat liefert Einsetzen eines jeden  $x \in M_1$  in die auf der linken Seite zu maximierende Funktion jeweils  $\lambda_1^2$ . Ist  $x = (x_1, \dots, x_d)$  gegeben mit  $\|x\| = 1$  und existiert  $j_0 > i_1$  mit  $x_{j_0} \neq 0$ , so haben wir

$$\sum_{j=1}^r \lambda_j^2 x_j^2 \leq \sum_{j=1}^r \lambda_1^2 x_j^2 \leq \lambda_1^2 \|x\|^2 = \lambda_1^2$$

wobei im Fall dass  $j_0 \leq r$  gilt, die erste Abschätzung strikt ist, weil dann der  $j_0$ -te Summand echt vergrößert wird. Ist  $j_0 > r$ , so ist die Summe über die  $x_i^2$  echt kleiner als  $\|x\|^2$  und daher die zweite Abschätzung strikt. Es folgt, dass  $M_1$  genau die angegebene Menge der Maximierer ist. Nach obigem ist  $B^\top v_i \in M_1$  und es folgt  $v_i \in \ker(A^\top A - \lambda_1)$  für  $1 \leq i \leq i_1$ . Weiterhin notieren wir

$$M_1 = \operatorname{span}\{B^\top v_1, \dots, B^\top v_{i_1}\} \cap \partial B_1(0).$$

Sei  $i_2$  die Vielfachheit von  $\lambda_{i_1+1}$  und ab jetzt  $i_1 < i \leq i_1 + i_2$ . Wir behaupten

$$\operatorname{argmax}_{\substack{\|x\|=1 \\ x \perp B^\top v_1, \dots, B^\top v_{i_1}}} \sum_{j=1}^r \lambda_j^2 x_j^2 = \{x \in \partial B_1(0) \mid x_j = 0 \text{ für } j \leq i_1 \text{ und für } j > i_2\} =: M_2.$$

Sei  $x \in M_2$ . Dann gilt  $\|x\| = 1$  und  $x \perp B^\top v_1, \dots, B^\top v_{i_1}$ . Einsetzen von  $x$  in die Summe auf der linken Seite liefert  $\lambda_{i_1+1}^2$ . Sei  $x$  gegeben mit  $\|x\| = 1$  und  $x \perp B^\top v_1, \dots, B^\top v_{i_1}$ , also  $x \perp M_1$  und damit  $x_j = 0$  für  $j \leq i_1$ . Falls es ein  $j_0 > i_2$  mit  $x_{j_0} \neq 0$  gibt, so sieht man

$$\sum_{j=1}^r \lambda_j^2 x_j^2 = \sum_{j=i_1+1}^r \lambda_j^2 x_j^2 \leq \sum_{j=i_1+1}^r \lambda_{i_1+1}^2 x_j^2 \leq \lambda_{i_1+1}^2 \|x\|^2 = \lambda_{i_1+1}^2$$

wobei wieder mindestens eine der Ungleichungen strikt ist. Damit haben wir die Gleichung für  $M_2$  und es folgt, dass  $v_{i_1+1}, \dots, v_{i_1+i_2}$  Eigenvektoren zu  $\lambda_{i_1+1}$  sind. Für alle weiteren Eigenwerte wiederholen wir die gleichen Argumente.  $\square$

Die Bedingung nach der in Satz 7.11 die  $v_1, \dots, v_r$  gewählt werden, stimmt mit derjenigen in Satz 6.3 überein. Man kann also sagen, dass die Berechnung der  $k$ -bestpassenden Unterräume für  $k = 1, \dots, r$  zu einer Datenmenge  $\{a_1, \dots, a_n\} \subseteq \mathbb{R}^d$  nebenbei eine Singulärwertzerlegung der Matrix  $A$  liefert deren Zeilen die  $a_i$  sind. Der nächste Satz zeigt, dass auch die Umkehrung wahr ist. Beachte, dass die Matrix  $V$  nicht eindeutig durch  $A$  bestimmt ist.

**Satz 7.12.** Sei  $D = \{a_1, \dots, a_n\} \subseteq \mathbb{R}^d$  eine Datenmenge bestehend aus Punkten

$a_i = (a_{i1}, \dots, a_{id})$  für  $i = 1, \dots, n$ . Wir betrachten die Datenmatrix

$$A = \begin{bmatrix} a_{11} & \cdots & a_{1d} \\ \vdots & & \vdots \\ a_{n1} & \cdots & a_{nd} \end{bmatrix} \in \mathbb{R}^{n \times d}$$

deren  $i$ -te Zeile gerade die Koordinaten des Punktes  $a_i$  enthält. Sei  $A = U\Sigma V^\top$  eine beliebige Singulärwertzerlegung von  $A$ . Dann spannen die  $k$ -ersten Spalten  $v_1, \dots, v_k$  von  $V$  einen  $k$ -bestpassenden Unterraum für  $D$  auf.

*Beweis.* Sei  $r = \text{rk } A$  und zuerst  $k \leq r$ . Dann sind nach Proposition 7.9 die Spalten  $v_1, \dots, v_k$  von  $V$  Singulärvektoren zu  $\sigma_1, \dots, \sigma_k > 0$ . Nach Lemma 7.3 ist daher  $v_i$  Eigenvektor von  $A^\top A$  zum Eigenwert  $\sigma_i$  für  $i = 1, \dots, k$ . Wie im Beweis des vorherigen Satzes 7.11 folgt

$$v_i \in \underset{\substack{\|v\|=1 \\ v \perp v_1, \dots, v_{i-1}}}{\text{argmax}} \|Av\|$$

für  $i = 1, \dots, k$ . Nach Satz 6.3 spannen die  $v_1, \dots, v_k$  also einen  $k$ -bestpassenden Unterraum für  $D$  auf. Ist  $k > r$ , so gilt nach Lemma 6.4

$$\text{span}\{a_1, \dots, a_n\} \subseteq \text{span}\{v_1, \dots, v_r\} \subseteq \text{span}\{v_1, \dots, v_k\}$$

und letzterer Raum ist trivialerweise  $k$ -bestpassend.  $\square$

**Bemerkung 7.13.** Man kann die Existenz einer Singulärwertzerlegung auch ohne die Theorie der orthogonalen Diagonalisierbarkeit beweisen sondern direkt mit dem Algorithmus in Satz 7.11 ('algorithmischer Beweis'): Seien hierzu  $U$ ,  $V$  und  $\Sigma$  so definiert wie in Satz 7.11 angegeben. Wir wissen also per Konstruktion, dass  $V^\top V = \text{id}_{\mathbb{R}^d}$  und  $U^\top U = \text{id}_{\mathbb{R}^n}$  gilt. Zu zeigen ist  $U^\top A V = \text{diag}(\sigma_1, \dots, \sigma_r)$ . Dazu rechnen wir das Matrixprodukt aus, wobei wir mit  $u_i^\top$  die Zeilen von  $U^\top$  bezeichnen:

$$U^\top A V = \begin{bmatrix} u_1^\top \\ \vdots \\ u_n^\top \end{bmatrix} A [v_1 \cdots v_n] = \begin{bmatrix} u_1^\top \\ \vdots \\ u_n^\top \end{bmatrix} [Av_1 \cdots Av_n] = \begin{bmatrix} u_1^\top Av_1 & \cdots & u_1^\top Av_d \\ \vdots & & \vdots \\ u_n^\top Av_1 & \cdots & u_n^\top Av_d \end{bmatrix}.$$

Für  $1 \leq i \leq r$  gilt dann  $u_i^\top Av_i = u_i^\top \sigma_i u_i = \sigma_i \langle u_i, u_i \rangle = \sigma_i$  wobei wir die Definition  $u_i = \frac{1}{\sigma_i} Av_i$  benutzt haben. Für  $1 \leq i < k \leq r$  gilt  $u_k^\top Av_i = u_k^\top \sigma_k u_i = \sigma_k \langle u_k, u_i \rangle = 0$ . Für  $k > r$  ist  $\|Av_k\| = 0$ , also sind alle dementsprechenden Einträge in der Produktmatrix Null:

$$U^\top A V = \begin{bmatrix} \sigma_1 & \boxed{u_1^\top Av_2} & \boxed{u_1^\top Av_3} & \cdots & \boxed{u_1^\top Av_d} \\ 0 & \sigma_2 & \boxed{u_2^\top Av_2} & \cdots & \boxed{u_2^\top Av_d} \\ \vdots & & & & \vdots \\ 0 & 0 & \cdots & \cdots & \boxed{u_n^\top Av_d} \end{bmatrix}.$$

Wir zeigen nun iterativ, dass die oben umrahmten Zeilen ebenfalls Null sind und nennen diese zu diesem Zweck von oben beginnend  $w_1^\top, w_2^\top, \dots$ . Weiter nennen wir die

Teilmatrix die genau unterhalb von  $w_i^\top$  liegt  $B_i$ . Für  $i = 1$  haben wir also

$$U^\top AV = \begin{bmatrix} \sigma_1 & w_1^\top \\ 0 & B_1 \end{bmatrix}$$

und behaupten, dass  $w_1 = 0$  ist. Dazu berechnen wir

$$\|U^\top AV \begin{bmatrix} \sigma_1 \\ w_1 \end{bmatrix}\|^2 = \left\| \begin{bmatrix} \sigma_1 & w_1^\top \\ 0 & B_1 \end{bmatrix} \begin{bmatrix} \sigma_1 \\ w_1 \end{bmatrix} \right\|^2 = \left\| \begin{bmatrix} \sigma_1^2 + \langle w_1, w_1 \rangle \\ B_1 w_1 \end{bmatrix} \right\|^2 \geq (\sigma_1^2 + \|w_1\|^2)^2$$

und

$$\left\| \begin{bmatrix} \sigma_1 \\ w_1 \end{bmatrix} \right\|^2 = \sigma_1^2 + \|w_1\|^2.$$

Weil  $U^\top$  und  $V$  Isometrien sind und mit  $v := \begin{bmatrix} \sigma_1 \\ w_1 \end{bmatrix} / \left\| \begin{bmatrix} \sigma_1 \\ w_1 \end{bmatrix} \right\|$  erhalten wir

$$\begin{aligned} \sigma_1 &= \max_{\|x\|=1} \|Ax\| = \|A\|_{\text{op}} = \|U^\top\|_{\text{op}} \|A\|_{\text{op}} \|V\|_{\text{op}} \\ &\geq \|U^\top AV\|_{\text{op}} = \max_{\|x\|=1} \|U^\top AVx\| \geq \frac{\|U^\top AV \begin{bmatrix} \sigma_1 \\ w_1 \end{bmatrix}\|}{\left\| \begin{bmatrix} \sigma_1 \\ w_1 \end{bmatrix} \right\|} \\ &\geq \frac{\sigma_1^2 + \|w_1\|^2}{\sqrt{\sigma_1^2 + \|w_1\|^2}} = \sqrt{\sigma_1^2 + \|w_1\|^2} \geq \sigma_1 \end{aligned}$$

was zeigt, dass alle Ungleichungen in der Tat Gleichungen waren. Das geht aber nur, wenn  $w_1 = 0$  ist. Sind nun bereits  $w_1, \dots, w_k = 0$ , so zeigt man  $w_{k+1} = 0$  indem man das letzte Argument wiederholt, aber mit

$$v = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ w_{k+1} \\ B_{k+1} \end{bmatrix} \Bigg/ \left\| \begin{bmatrix} 0 \\ \vdots \\ 0 \\ w_{k+1} \\ B_{k+1} \end{bmatrix} \right\|.$$

Bevor wir gleich zu weiteren Anwendungen, neben der Berechnung bestpassender Unterräume, der Singulärwertzerlegung kommen, wollen wir noch zeigen, dass die Singulärwertzerlegung den Satz über die orthogonale Diagonalisierbarkeit in der Tat verallgemeinert.

**Proposition 7.14.** *Sei  $A \in \mathbb{R}^{n \times n}$  eine symmetrische Matrix, sei  $r = \text{rk } A$  und seien  $\sigma_1 \geq \dots \geq \sigma_r > 0$  deren Singulärwerte mit rechten Singulärvektoren  $v_1, \dots, v_r$ , welche wir zu einer Orthonormalbasis  $\{v_1, \dots, v_n\}$  ergänzen. Dann sind  $\sigma_1, \dots, \sigma_r$  die echt positiven Eigenwerte von  $A$  und Null ist Eigenwert mit Vielfachheit  $n - r$ . Ferner gilt  $V^\top AV = \text{diag}(\sigma_1, \dots, \sigma_r, 0, \dots, 0)$  mit  $V = [v_1 \dots v_n]$ .*

*Beweis.* Sei  $\lambda$  ein Eigenwert ungleich Null. Dann gilt  $Av = \lambda v$  mit einem  $v \neq 0$ . Da  $A$  symmetrisch ist, folgt  $A^\top v = \lambda v$ , d.h.  $\lambda$  ist Singulärwert mit rechtem Singulärvektor  $v$  und linkem Singulärvektor ebenfalls  $v$ . Wir erhalten  $A = V\Sigma V^\top$  wenn wir für die  $\sigma_{r+1}, \dots, \sigma_n$  die Orthonormalsysteme  $\mathcal{U}$  und  $\mathcal{V}$  auf die gleiche Weise zu Orthonormalbasen erweitern.  $\square$

Multiplikation mit  $V^\top$  von links in Proposition 7.14 liefert also genau eine wie ganz am Anfang des Kapitels angegebene orthogonale Diagonalisierung  $V^\top AV =$

$\text{diag}(\sigma_1, \dots, \sigma_r, 0, \dots, 0)$ . Der Leser sei aber wie folgt gewarnt: Ist  $A$  zwar quadratisch aber nicht symmetrisch, so stimmen die  $\sigma_i$  im Allgemeinen *nicht* mit den Eigenwerten von  $A$  überein!

## 7.1 Dimensionalitätsreduktion

Berechnet man die Singulärwertzerlegung, z.B. von Bewertungsmatrizen  $A$ , siehe Beispiel 7.22 oder Bildmatrizen  $A$ , siehe Aufgabe 7.4, so fällt auf, dass die Singulärwerte  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$  häufig *sehr schnell fallen*, d.h. bereits für relativ kleines  $k$  haben wir

$$\sigma_{k+1} \approx \dots \approx \sigma_r \approx 0.$$

Die zentrale Idee ist es nun, in der Singulärwertzerlegung ohnehin kleine Singulärwerte auf exakt Null zu setzen und auf diese Weise eine Approximation

$$\check{A} := U\check{\Sigma}V^\top \quad \text{mit} \quad \check{\Sigma} := \text{diag}(\sigma_1, \dots, \sigma_k) \in \mathbb{R}^{n \times d} \quad (7.1)$$

der Ausgangsmatrix  $A$  zu definieren. Um  $A$  mit  $\check{A}$  zu vergleichen, verwenden wir die folgende Norm auf dem Raum  $\mathbb{R}^{n \times d}$ .

**Definition 7.15.** Sei  $A \in \mathbb{R}^{n \times d}$  gegeben mit  $A = (a_{ij})$ . Dann ist die *Frobeniusnorm* von  $A$  definiert durch

$$\|A\|_F := \left( \sum_{i=1}^n \sum_{j=1}^d |a_{ij}|^2 \right)^{1/2}.$$

Man sieht, dass die Frobeniusnorm nichts anderes ist als die euklidische Norm des Vektors  $(a_{11}, a_{12}, \dots, a_{nd}) \in \mathbb{R}^{n \cdot d}$  bei dem wir alle Matrixeinträge hintereinander schreiben. Daher ist es klar, dass es sich bei  $\|\cdot\|_F$  tatsächlich um eine Norm auf  $\mathbb{R}^{n \times d}$  handelt. Wir notieren den folgenden zentralen Zusammenhang zwischen Frobeniusnorm und Singulärwerten.

**Proposition 7.16.** Sei  $A \in \mathbb{R}^{n \times d}$ ,  $r = \text{rang}(A)$  und  $\sigma_1 \geq \dots \geq \sigma_r > 0$  seien die Singulärwerte von  $A$ . Dann gilt

$$\|A\|_F = \left( \sum_{j=1}^r \sigma_j^2 \right)^{1/2} = \left\| \begin{bmatrix} \sigma_1 \\ \vdots \\ \sigma_r \end{bmatrix} \right\|$$

wobei auf der rechten Seite  $\|\cdot\|$  für die euklidische Norm auf  $\mathbb{R}^r$  steht.

*Beweis.* Im folgenden bezeichnen wir mit  $\text{spur}(B)$  für eine quadratische Matrix  $B$  die Summe der Einträge auf der Hauptdiagonalen von  $B$ . Damit ergibt sich

$$\text{spur}(AA^\top) = \sum_{i=1}^n \sum_{j=1}^d a_{ij}a_{ij} = \sum_{j=1}^d \sum_{i=1}^n |a_{ij}|^2 = \|A\|_F^2.$$

Weiter sieht man durch Berechnung der Hauptdiagonalen von Matrix-Matrix-Produkten  $XY$  und  $YX$  für  $X \in \mathbb{R}^{m \times \ell}$  und  $Y \in \mathbb{R}^{\ell \times m}$ , dass  $\text{spur}(XY) = \text{spur}(YX)$



gilt. Jetzt berechnen wir

$$\begin{aligned}\|A\|_F^2 &= \text{spur}(AA^\top) = \text{spur}(U\Sigma V^\top(U\Sigma V^\top)^\top) = \text{spur}(U\Sigma V^\top V\Sigma^\top U^\top) \\ &= \text{spur}(U\Sigma\Sigma^\top U^\top) = \text{spur}(U^\top U\Sigma\Sigma^\top) = \text{spur}(\Sigma\Sigma^\top) = \sigma_1^2 + \dots + \sigma_r^2\end{aligned}$$

und sind fertig.  $\square$

Als nächstes quantifizieren wir die Approximationsqualität, die wir erreichen, wenn wir die Singulärwerte  $\sigma_{k+1}, \dots, \sigma_r$  auf Null setzen, also die Datenmatrix  $A$  durch die in (7.1) definierte Approximation  $\check{A}$  ersetzen.

**Satz 7.17.** (Nullsetzen von Singulärwerten) *Sei  $A \in \mathbb{R}^{n \times d}$  und  $A = U\Sigma V^\top$  eine Singulärwertzerlegung von  $A = (a_{ij})$  mit Singulärwerten  $\sigma_1 \geq \dots \geq \sigma_r > 0$ . Für fixiertes  $1 \leq k \leq r$  sei  $\check{A} = (\check{a}_{ij})$  definiert wie in (7.1). Wir bezeichnen mit  $v_j$  die  $j$ -te Spalte von  $V$  und mit  $a_i$  bzw.  $\check{a}_i$  die  $i$ -te Zeile von  $A$  bzw. von  $\check{A}$ , jeweils gelesen als (Spalten-)vektoren in  $\mathbb{R}^d$ . Dann gilt*

- (i)  $\|A - \check{A}\|_F = \sqrt{\sigma_{k+1}^2 + \dots + \sigma_r^2}$ ,
- (ii)  $\check{a}_i = \pi_{\text{span}\{v_1, \dots, v_k\}}(a_i)$  für alle  $1 \leq i \leq n$ .

*Beweis.* (i) Wegen  $A - \check{A} = U(\Sigma - \check{\Sigma})V^\top = U \text{diag}(0, \dots, 0, \sigma_{k+1}, \dots, \sigma_r)V^\top$  sind nach Proposition 7.9 die  $\sigma_{k+1} \geq \dots \geq \sigma_r$  gerade die Singulärwerte von  $A - \check{A}$ . Damit folgt aus Proposition 7.16 die gewünschte Formel.

(ii) In der Notation von Lemma 7.7 liefert eben dieses eine Formel für die Einträge der  $i$ -ten Zeile von  $A$ , nämlich

$$\begin{aligned}a_i &= (a_{i1}, \dots, a_{id}) = \left( \sum_{\ell=1}^r u_{i\ell} \sigma_\ell v_{1\ell}, \dots, \sum_{\ell=1}^r u_{i\ell} \sigma_\ell v_{d\ell} \right) \\ &= \sum_{\ell=1}^r u_{i\ell} \sigma_\ell (v_{1\ell}, \dots, v_{d\ell}) = \sum_{\ell=1}^r u_{i\ell} \sigma_\ell v_\ell\end{aligned}$$

wobei  $v_\ell = (v_{1\ell}, \dots, v_{d\ell})$ . Für  $\check{A}$  gilt die gleiche Rechnung, außer das die Summation jeweils bei  $k \leq r$  abbricht. Daraus ergibt sich dann

$$\begin{aligned}\pi_{\text{span}\{v_1, \dots, v_k\}}(a_i) &= \sum_{\mu=1}^k \langle a_i, v_\mu \rangle v_\mu = \sum_{\mu=1}^k \left\langle \sum_{\ell=1}^r u_{i\ell} \sigma_\ell v_\ell, v_\mu \right\rangle v_\mu \\ &= \sum_{\mu=1}^k \sum_{\ell=1}^r u_{i\ell} \sigma_\ell \langle v_\ell, v_\mu \rangle v_\mu = \sum_{\ell=1}^k u_{i\ell} \sigma_\ell v_\ell = \check{a}_i\end{aligned}$$

wie behauptet.  $\square$

Schreibt man die Matrix-Matrix-Matrix-Multiplikation  $U\check{\Sigma}V^\top$  in Satz 7.17 aus, vergleiche Lemma 7.7, so erhält man

$$\check{A} = U\check{\Sigma}V^\top = \begin{bmatrix} u_{11} & \dots & u_{1k} \\ \vdots & & \vdots \\ u_{n1} & \dots & u_{nk} \end{bmatrix} \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_k \end{bmatrix} \begin{bmatrix} v_{11} & \dots & v_{d1} \\ \vdots & & \vdots \\ v_{1k} & \dots & v_{dk} \end{bmatrix},$$

was die Berechnung von  $\check{A}$  in Beispielen vereinfacht. Ferner sieht man, dass  $\check{A}$  durch  $(k \cdot n + k + k \cdot d)$ -viele reelle Zahlen rekonstruiert werden kann. Dies verwenden wir in Aufgabe 7.4 zur Bildkompression via Singulärwertzerlegung.

Nach Satz 7.17(i) approximiert das Weglassen von Singulärwerten die Ausgangsmatrix in der Frobeniusnorm, wobei wir den Frobeniusabstand dadurch steuern können wieviele Singulärwerte wir auf Null setzen. Wegen  $|a_{ij} - \check{a}_{ij}| \leq \|A - \check{A}\|_F$  erhalten wir insbesondere eine *punktweise* Approximation der Matrix  $A = (a_{ij})$  durch  $\check{A} = (\check{a}_{ij})$ . Außerdem kontrollieren wir durch die Wahl von  $k$  den Rang der Matrix  $\check{A}$ , erhalten also eine *Approximation von A durch Matrizen niedrigeren Ranges* die im folgenden Sinne optimal ist.

**Korollar 7.18.** (Bestapproximation niedrigeren Ranges) *Unter den Voraussetzungen von Satz 7.17 gilt*

$$\check{A} \in \operatorname{argmin}_{\substack{B \in \mathbb{R}^{n \times d} \\ \operatorname{rk} B = k}} \|A - B\|_F.$$

*Beweis.* Angenommen, es existiert  $B \in \mathbb{R}^{n \times d}$  mit  $\|A - B\|_F < \|A - \check{A}\|_F$ . Wir setzen

$$V_k := \operatorname{span}\{v_1, \dots, v_k\}, \quad W := \operatorname{span}\{b_1, \dots, b_n\} \quad \text{und} \quad \check{V} := \operatorname{span}\{\check{a}_1, \dots, \check{a}_n\}$$

wobei  $b_i \in \mathbb{R}^d$  derjenige Vektor ist, der durch die Einträge der  $i$ -ten Zeile von  $B$  gegeben ist. Nach Satz 7.12 ist  $V_k$  ein  $k$ -bestpassender Unterraum für  $\{a_1, \dots, a_n\}$ . Wir behaupten, dass  $\check{V} = V_k$  gilt. Nach Satz 7.17(ii) haben wir  $\check{a}_i \in V_k$  für  $i = 1, \dots, n$ , also  $\check{V} \subseteq V_k$ . Die Gleichheit folgt dann aus  $\dim \check{V} = \operatorname{rk} \check{A} = k = \dim V_k$ . Wir erhalten damit für jedes  $i = 1, \dots, n$  die Abschätzungen

$$\|a_i - b_i\| \underset{\substack{\uparrow \\ b_i \in W}}{\geq} \|a_i - \pi_W(a_i)\| \underset{\substack{\uparrow \\ V_k \text{ k-best-} \\ \text{passend für} \\ a_1, \dots, a_n}}{\geq} \|a_i - \pi_{V_k}(a_i)\| \underset{\substack{\uparrow \\ V_k = \check{V}}}{=} \|a_i - \pi_{\check{V}}(a_i)\| \underset{\substack{\uparrow \\ \text{Satz} \\ 7.17(ii)}}{=} \|a_i - \check{a}_i\|.$$

Quadrieren und Absummieren liefert dann

$$\|A - \check{A}\|_F^2 = \sum_{i=1}^n \|a_i - \check{a}_i\|^2 \leq \sum_{i=1}^n \|a_i - b_i\|^2 = \|A - B\|_F^2$$

im Widerspruch zur Annahme.  $\square$

Als nächstes beweisen wir eine quantitative Version von Proposition 6.5: Sind uns alle Singulärwerte  $\sigma_1 \geq \dots \geq \sigma_r > 0$  bekannt, so können wir via Korollar 7.19 für eine gewünschte Approximationsqualität das dafür nötige  $k$  finden oder umgekehrt für ein gewünschtes  $k$  den Approximationsfehler berechnen.

**Korollar 7.19.** (Finetuning bestpassender Unterräume) *Unter den Voraussetzungen von Satz 7.17 gilt*

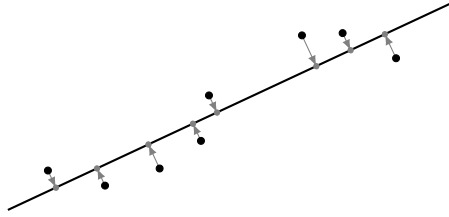
$$\sum_{j=k+1}^r \sigma_j^2 = \sum_{i=1}^n \operatorname{dist}(a_i, \operatorname{span}\{v_1, \dots, v_k\})^2 = \min_{\substack{V \subseteq \mathbb{R}^d \\ \dim V = k}} \sum_{i=1}^n \operatorname{dist}(a_i, V)^2.$$

*Beweis.* Sei  $V_k := \text{span}\{v_1, \dots, v_k\}$ . Nach Satz 7.12 ist  $V_k$  ein  $k$ -bestpassender Unterraum für  $\{a_1, \dots, a_n\}$ . Die erste Gleichung im Korollar folgt dann per

$$\sum_{i=1}^n \text{dist}(x_i, V_k)^2 = \sum_{i=1}^n \|a_i - \pi_{V_k}(a_i)\|^2 \underset{\substack{\uparrow \\ \text{Satz} \\ 7.17(\text{ii})}}{=} \|A - \tilde{A}\|_F^2 \underset{\substack{\uparrow \\ \text{Satz} \\ 7.17(\text{i})}}{=} \sum_{j=k+1}^r \sigma_j^2.$$

Die zweite Gleichung gilt, weil  $V_k$  ein  $k$ -bestpassender Unterraum ist.  $\square$

Ist eine Datenmenge in  $\mathbb{R}^d$  gegeben, und liegen vielleicht sogar alle Datenpunkte nah an einem niedrigdimensionalen Unterraum  $W$  von  $\mathbb{R}^d$ , so kann man die ‘Dimensionalität’ der Datenmenge reduzieren, indem man die Datenpunkte auf  $W$  orthogonalprojiziert. Denkt man z.B. an den in Kapitel 3 diskutierten  $k$ -NN Algorithmus, so ist es wünschenswert, dass durch diese Projektion paarweise Abstände zwischen möglichst vielen Datenpunkten nur leicht verändert werden. Das folgende Bild suggeriert, dass ein bestpassender Unterraum mit nicht zu kleiner Dimension ein guter Kandidat für  $W$  sein könnte.



Um nicht nur die Dimension, sondern auch die Anzahl der Koordinaten zu reduzieren, verwenden wir nicht die Projektion  $\pi_{V_k}$  selbst, sondern die Abbildung

$$T_{V_k}: \mathbb{R}^d \rightarrow \mathbb{R}^k, \quad x \mapsto \begin{bmatrix} v_{11} & \dots & v_{d1} \\ \vdots & & \vdots \\ v_{1k} & \dots & v_{dk} \end{bmatrix} x. \quad (7.2)$$

Wegen  $T_{V_k}x = (\langle v_1, x \rangle, \dots, \langle v_k, x \rangle)$ , d.h.  $T_{V_k}x$  ist der Koordinatenvektor der orthogonalen Projektion von  $x$  auf  $V_k = \text{span}\{v_1, \dots, v_k\}$  bezüglich der Basis  $\mathcal{V}_k = \{v_1, \dots, v_k\}$ .

**Korollar 7.20.** (Dimensionalitätsreduktion via SVD) *Unter den Voraussetzungen von Satz 7.19 sei  $T_{V_k}$  definiert wie in (7.2). Dann gilt*

$$\forall i, j \in \{1, \dots, n\}: \|T_{V_k}a_i - T_{V_k}a_j\| \leq 2 \left( \sum_{\ell=k+1}^r \sigma_\ell \right)^{1/2}.$$

*Beweis.* Wir bezeichnen mit  $\mathcal{V}_k$  die Orthonormalbasis  $\{v_1, \dots, v_k\}$  des Raums  $V_k = \text{span}\{v_1, \dots, v_k\}$ . Nach Satz 7.17(ii) haben wir dann

$$(\tilde{a}_i)_{\mathcal{V}_k} = (\pi_{V_k}(a_i))_{\mathcal{V}_k} = (\langle v_1, a_i \rangle, \dots, \langle v_k, a_i \rangle) = T_{V_k}a_i$$

für  $1 \leq i \leq n$ . Für  $i$  und  $j$  wie im Korollar ergibt sich

$$\|T_{V_k} a_i - T_{V_k} a_j\| = \|\check{a}_i - \check{a}_j\| \leq \|\check{a}_i - a_i\| + \|a_i - a_j\| + \|a_j - \check{a}_j\|$$

sowie

$$\begin{aligned} \|a_i - a_j\| &\leq \|a_i - \check{a}_i\| + \|\check{a}_i - \check{a}_j\| + \|\check{a}_j - a_j\| \\ &= \|a_i - \check{a}_i\| + \|T_{V_k} a_i - T_{V_k} a_j\| + \|\check{a}_j - a_j\| \end{aligned}$$

Abziehen von  $\|a_i - a_j\|$  von der ersten Ungleichung und  $\|T_{V_k} a_i - T_{V_k} a_j\|$  von der zweiten Ungleichung liefert

$$\begin{aligned} \left| \|T_{V_k} a_i - T_{V_k} a_j\| - \|a_i - a_j\| \right| &\leq \|a_i - \check{a}_i\| + \|a_j - \check{a}_j\| \\ &= \left( \sum_{\ell=1}^d |a_{j\ell} - \check{a}_{j\ell}|^2 \right)^{1/2} + \left( \sum_{\ell=1}^d |a_{i\ell} - \check{a}_{i\ell}|^2 \right)^{1/2} \\ &\leq 2 \left( \sum_{j=1}^n \sum_{\ell=1}^d |a_{j\ell} - \check{a}_{j\ell}|^2 \right)^{1/2} = 2 \|A - \check{A}\|_F \end{aligned}$$

was mit Satz 7.17(i) den Beweis beendet.  $\square$

Anders als in Korollar 7.18, ist  $V_k$  unter den  $k$ -dimensionalen Unterräumen im Allgemeinen nicht optimal im Sinne, dass die Summe der durch die Projektion entstehenden Abstandsänderungen minimiert würde. Etwas formaler haben wir

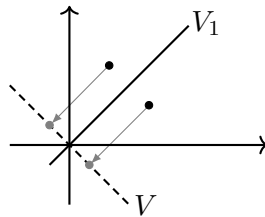
$$V_k \notin \operatorname{argmin}_{\substack{V \subseteq \mathbb{R}^d \\ \dim V = k}} \sum_{i,j=1}^n \left| \|\pi_V(a_i) - \pi_V(a_j)\| - \|a_i - a_j\| \right|$$

im Allgemeinen und zeigen dies mit dem folgenden Beispiel.

**Beispiel 7.21.** Betrachte für  $d = n = 2$  die Matrix  $A = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$  mit Singulärwertzerlegung

$$A = \begin{bmatrix} \sqrt{2}/2 & -\sqrt{2}/2 \\ \sqrt{2}/2 & \sqrt{2}/2 \end{bmatrix} \begin{bmatrix} 3 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \sqrt{2}/2 & \sqrt{2}/2 \\ -\sqrt{2}/2 & \sqrt{2}/2 \end{bmatrix}.$$

Für die Datenmenge  $\{(2, 1), (1, 2)\}$  ergibt sich also  $V_1 = \operatorname{span}\{(\sqrt{2}/2, \sqrt{2}/2)\}$  als ein 1-bestpassender Unterraum. Das folgende Bild zeigt allerdings, dass bei Projektion auf den Unterraum  $V$  die Abstände der zwei Punkte unverändert bleiben, während bei Projektion auf  $V_1$  der Abstand nach Projektion Null ist.



Das folgende Beispiel illustriert einerseits eine klassische Anwendung für die Ap-

proximation einer Datenmatrix durch eine Matrix niedrigeren Ranges per Singulärwertzerlegung und die dadurch sich ergebende Möglichkeit der Kompression von Daten. Andererseits werden wir in den folgenden zwei Unterkapiteln sowohl die *Hauptkomponentenanalyse* als auch das SVD-basierte *kollaborative Filtern* anhand des folgenden Beispiels erklären.

**Beispiel 7.22.** Gegeben sei die folgende Tabelle von Filmbewertungen

	Alien	Casablanca	Star Wars	Titanic	Matrix
Abbey	0	2	0	2	1
Bailey	1	0	1	0	1
Caitlin	5	0	5	0	5
Daisy	0	4	0	4	2
Edith	3	0	3	0	3
Fae	0	5	0	5	0
Gail	4	0	4	0	4

welche wir in der Datenmatrix

$$A := \begin{bmatrix} 0 & 2 & 0 & 2 & 1 \\ 1 & 0 & 1 & 0 & 1 \\ 5 & 0 & 5 & 0 & 5 \\ 0 & 4 & 0 & 4 & 2 \\ 3 & 0 & 3 & 0 & 3 \\ 0 & 5 & 0 & 5 & 0 \\ 4 & 0 & 4 & 0 & 4 \end{bmatrix}$$

organisieren und für die wir per Python eine Singulärwertzerlegung bestimmen:

$$A = \underbrace{\begin{bmatrix} 0.07 & 0.29 & 0.32 & 0.51 & 0.66 & 0.18 & -0.23 \\ 0.13 & -0.02 & -0.01 & -0.79 & 0.59 & -0.02 & -0.06 \\ 0.68 & -0.11 & -0.05 & -0.05 & -0.24 & 0.56 & -0.35 \\ 0.15 & 0.59 & 0.65 & -0.25 & -0.33 & -0.09 & 0.11 \\ 0.41 & -0.07 & -0.03 & 0.10 & -0.02 & -0.78 & -0.43 \\ 0.07 & 0.73 & -0.67 & 0.00 & -0.00 & 0.00 & 0.00 \\ 0.55 & 0.09 & -0.04 & 0.17 & 0.17 & -0.11 & 0.78 \end{bmatrix}}_U \underbrace{\begin{bmatrix} 12.4 & 0.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 9.5 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 1.3 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.0 & 0.0 \end{bmatrix}}_\Sigma \underbrace{\begin{bmatrix} 0.56 & 0.09 & 0.56 & 0.09 & 0.59 \\ -0.12 & 0.69 & -0.12 & 0.69 & 0.02 \\ -0.40 & -0.09 & -0.40 & -0.09 & 0.80 \\ 0.41 & 0.09 & 0.40 & 0.09 & -0.80 \\ 0.51 & 0.48 & -0.51 & -0.48 & -0.00 \\ 0.48 & -0.51 & -0.48 & 0.51 & -0.00 \end{bmatrix}}_{V^T}$$

$$= \begin{bmatrix} 0.07 & 0.29 & 0.32 \\ 0.13 & -0.02 & -0.01 \\ 0.68 & -0.11 & -0.05 \\ 0.15 & 0.59 & 0.65 \\ 0.41 & -0.07 & -0.03 \\ 0.07 & 0.73 & -0.67 \\ 0.55 & 0.09 & -0.04 \end{bmatrix} \begin{bmatrix} 12.4 & & \\ & 9.5 & \\ & & 1.3 \end{bmatrix} \begin{bmatrix} 0.56 & 0.09 & 0.56 & 0.09 & 0.59 \\ -0.12 & 0.69 & -0.12 & 0.69 & 0.02 \\ -0.40 & -0.09 & -0.40 & -0.09 & 0.80 \end{bmatrix}.$$

Wir schreiben oben ‘ $A = \dots$ ’ auch wenn die Einträge von  $A$  nach zwei Nachkommastellen abgeschnitten wurden und daher nicht exakt sind. Nullsetzen von  $\sigma_3$  liefert die Rang-2-Bestapproximation

$$\tilde{A} = \begin{bmatrix} 0.15 & 1.97 & 0.15 & 1.97 & 0.56 \\ 0.92 & 0.01 & 0.92 & 0.01 & 1.94 \\ 4.84 & 0.03 & 4.84 & 0.03 & 4.95 \\ 0.36 & 4.03 & 0.36 & 4.03 & 1.20 \\ 2.92 & 0.00 & 2.92 & 0.00 & 2.98 \\ -0.34 & 4.86 & -0.34 & 4.86 & 0.65 \\ 3.71 & 1.20 & 3.71 & 1.20 & 4.04 \end{bmatrix} = \begin{bmatrix} 0.07 & 0.29 \\ 0.13 & -0.02 \\ 0.68 & -0.11 \\ 0.15 & 0.59 \\ 0.41 & -0.07 \\ 0.07 & 0.73 \\ 0.55 & 0.09 \end{bmatrix} \begin{bmatrix} 12.4 & & \\ & 9.5 & \end{bmatrix} \begin{bmatrix} 0.56 & 0.09 & 0.56 & 0.09 & 0.59 \\ -0.12 & 0.69 & -0.12 & 0.69 & 0.02 \end{bmatrix}$$

mit  $\|A - \tilde{A}\|_F = 1.3$ . Um die Matrix  $\tilde{A}$  anhand der Singulärwertzerlegung zu speichern würden wir  $14 + 2 + 10 = 26$  reelle Zahlen speichern müssen im Gegensatz zu  $5 \cdot 7 = 35$  für die Originalmatrix  $A$ .

Bevor wir jetzt zur Hauptkomponentenanalyse kommen, notieren wir noch, dass in diesem Kapitel der Begriff ‘Dimensionalität’ bisher nicht formal definiert wurde. Im Sinne der Korollare 7.18 und 7.20 kann dies leicht nachgeholt werden, in dem für eine gegebene Datenmenge  $D = \{a_1, \dots, a_n\} \in \mathbb{R}^d$  die Dimension des Spanns von  $D$  oder, äquivalent, den Rang der Datenmatrix  $A = (a_{ij})$  als *Dimensionalität* von  $D$  bzw. von  $A$  formal einführt.

## 7.2 Hauptkomponentenanalyse

Wie angekündigt, bleiben wie bei Beispiel 7.22 und erläutern nun die Methode der *Hauptkomponentenanalyse* oder kurz *PCA*. Ganz konkret betrachten wir zunächst die Bewertung des Films Titanic durch Bewerterin Bailey, also den vorletzten Eintrag in der zweiten Zeile der Matrix:

$$\begin{array}{c} \text{Bailey} \rightarrow \end{array} \begin{array}{c} \text{Titanic} \\ \downarrow \\ \begin{bmatrix} 0 & 2 & 0 & 2 & 1 \\ 1 & 0 & 1 & 0 & 1 \\ 5 & 0 & 5 & 0 & 5 \\ 0 & 4 & 0 & 4 & 2 \\ 3 & 0 & 3 & 0 & 3 \\ 0 & 5 & 0 & 5 & 0 \\ 4 & 0 & 4 & 0 & 4 \end{bmatrix} \end{array}.$$

Diesen können wir herausprojizieren, indem wir die Matrix von links mit einem Zeilenvektor multiplizieren, der im zweiten Eintrag eine Eins und sonst Nullen hat und von rechts mit einem Spaltenvektor, der eine Eins im vorletzten Eintrag hat und sonst Nullen. Die Einträge dieser Zeilen- bzw. Spaltenvektoren können als kartesische Koordinaten von Bailey im *Raum der Bewerterinnen* bzw. von Titanic im *Raum der Filme* verstanden werden. Letztere definieren wir als die freien Vektorräume

$$B := \mathbb{R}\{\text{Abbey, Bailey, } \dots\} \quad \text{und} \quad F := \mathbb{R}\{\text{Alien, Casablanca, } \dots\}$$

über den Mengen  $\mathcal{B} = \{\text{Abbey, Bailey, } \dots\}$  und  $\mathcal{F} = \{\text{Alien, Casablanca, } \dots\}$ . Benutzen wir letztere als Basen, so sind die Einheitsvektoren in  $\mathbb{R}^5$  bzw. in  $\mathbb{R}^7$  gerade die Koordinatenvektoren der Bewerterinnen bzw. der Filme.

Wir nennen

$$R: B \times F \rightarrow \mathbb{R}, (b, f) \mapsto (b)_{\mathcal{B}}^T \cdot A \cdot (f)_{\mathcal{F}}$$

die *Bewertungsabbildung*. In unserem Beispiel von Bailey und Titanic ergibt sich unter Einbeziehung der in Beispiel 7.22 angegebene Singulärwertzerlegung:

$$R(\text{Bailey, Titanic}) = \begin{bmatrix} 0 & 1 & 0 & \dots & 0 \end{bmatrix} \begin{bmatrix} 0 & 2 & 0 & 2 & 1 \\ 1 & 0 & 1 & 0 & 1 \\ 5 & 0 & 5 & 0 & 5 \\ 0 & 4 & 0 & 4 & 2 \\ 3 & 0 & 3 & 0 & 3 \\ 0 & 5 & 0 & 5 & 0 \\ 4 & 0 & 4 & 0 & 4 \end{bmatrix} \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}$$

$\nearrow$  Bailey in kartesischen Koordinaten       $\nwarrow$  Titanic in kartesischen Koordinaten

$$\begin{aligned}
&= \begin{bmatrix} 0 & 1 & 0 & \cdots & 0 \end{bmatrix} \begin{bmatrix} 0.07 & 0.29 & \cdots & -0.23 \\ 0.13 & -0.02 & & -0.06 \\ 0.68 & -0.11 & & -0.35 \\ 0.15 & 0.59 & & 0.11 \\ 0.41 & -0.07 & & -0.43 \\ 0.07 & 0.73 & & 0.00 \\ 0.55 & 0.09 & \cdots & 0.78 \end{bmatrix} \begin{bmatrix} 12.4 & & & \\ & 9.5 & & \\ & & 1.3 & \\ & & & \ddots \end{bmatrix} \begin{bmatrix} 0.56 & 0.09 & 0.56 & 0.09 & 0.59 \\ -0.12 & 0.69 & -0.12 & 0.69 & 0.02 \\ \vdots & & & & \vdots \\ 0.48 & -0.51 & -0.48 & 0.51 & 0.00 \end{bmatrix} \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \end{bmatrix} \\
&= \begin{bmatrix} 0.13 & -0.02 & \cdots & -0.06 \end{bmatrix} \begin{bmatrix} 12.4 & & & \\ & 9.5 & & \\ & & 1.3 & \\ & & & \ddots \end{bmatrix} \begin{bmatrix} 0.09 \\ 0.69 \\ \vdots \\ 0.00 \end{bmatrix} \quad \begin{array}{l} \nearrow \text{Bailey in } \mathcal{U}\text{-} \\ \text{Koordinaten} \end{array} \quad \begin{array}{l} \nwarrow \text{Titanic in } \mathcal{V}\text{-} \\ \text{Koordinaten} \end{array}
\end{aligned}$$

In der letzten Zeile können wir die Einträge des Zeilenvektors links und des Spaltenvektors rechts wieder als Koordinaten von Bailey und Titanic auffassen — jetzt aber bezüglich neuer Basen  $\mathcal{U} = \{u_1, \dots, u_5\}$  von  $B$  und  $\mathcal{V} = \{v_1, \dots, v_7\}$  von  $F$ , die sich per

$$\begin{aligned}
u_1 &= 0.07 \cdot \text{Abbey} + 0.13 \cdot \text{Bailey} + \cdots + 0.55 \cdot \text{Gail}, \\
u_2 &= 0.29 \cdot \text{Abbey} - 0.02 \cdot \text{Bailey} + \cdots + 0.09 \cdot \text{Gail}, \\
&\vdots \\
u_5 &= -0.23 \cdot \text{Abbey} - 0.06 \cdot \text{Bailey} + \cdots + 0.78 \cdot \text{Gail}, \\
v_1 &= 0.56 \cdot \text{Alien} + 0.09 \cdot \text{Casablanca} + \cdots + 0.59 \cdot \text{Matrix}, \\
v_2 &= -0.12 \cdot \text{Alien} + 0.69 \cdot \text{Casablanca} + \cdots + 0.02 \cdot \text{Matrix}, \\
&\vdots \\
v_7 &= 0.48 \cdot \text{Alien} - 0.51 \cdot \text{Casablanca} + \cdots + 0.00 \cdot \text{Matrix}
\end{aligned}$$

aus den alten Basisvektoren ergeben. Bezeichnen wir mit  $\mathbb{R}_{\mathcal{B}}^5$  den Raum der Koordinatenvektoren bezüglich der Basis  $\mathcal{B}$  des Raumes  $B$  usw., so erhalten wir das folgende Basiswechseldiagramm:

$$\begin{array}{ccccc}
(\text{Film})_{\mathcal{F}} & \mathbb{R}_{\mathcal{F}}^5 & \xrightarrow{A} & \mathbb{R}_{\mathcal{B}}^7 & (\text{Bewerterin})_{\mathcal{B}} \\
\downarrow & \downarrow V^{\top} & \sim & \sim \downarrow U^{\top} & \downarrow \\
(\text{Film})_{\mathcal{V}} & \mathbb{R}_{\mathcal{V}}^5 & \xrightarrow{\Sigma} & \mathbb{R}_{\mathcal{U}}^7 & (\text{Bewerterin})_{\mathcal{U}}
\end{array}$$

Wir weisen darauf hin, dass wir in zu Beginn des Kapitels mit  $\mathcal{U} = \{u_1, \dots, u_n\}$  und  $\mathcal{V} = \{v_1, \dots, v_n\}$  die Basen von  $\mathbb{R}^n$  bzw. von  $\mathbb{R}^d$  bezeichnet haben, deren Elemente  $u_j$  und  $v_j$  die Spalten der Matrizen  $U$  bzw.  $V$  waren. Hier erhalten wir natürlich genau dasselbe, wenn wir die abstrakten Räume  $B$  und  $F$  mit  $\mathbb{R}_{\mathcal{V}}^5$  bzw.  $\mathbb{R}_{\mathcal{U}}^7$  identifizieren.

Die neuen Basisvektoren  $u_1, \dots, u_7$  und  $v_1, \dots, v_5$  sind per Definition Elemente des Raumes der Bewerterinnen bzw. der Filme. Da sie in unserer Ausgangstabelle nicht vorkommen, bezeichnen wir sie als *künstliche Bewerterinnen* und *künstliche Filme*. Wegen der Basiseigenschaft, können wir die echten Bewerterinnen und Filme aus den künstlichen rekonstruieren, z.B.

$$\text{Abbey} = \begin{bmatrix} u_1 & u_2 & \cdots & u_7 \end{bmatrix} \cdot (\text{erste Zeile von } U)^{\top}$$

$$= 0.07 \cdot u_1 + 0.29 \cdot u_2 + \cdots - 0.23 \cdot u_7,$$

$$\begin{aligned} \text{Alien} &= [v_1 \ v_2 \ \cdots \ v_5] \cdot (\text{erste Zeile von } V)^\top \\ &= [v_1 \ v_2 \ \cdots \ v_5] \cdot \text{erste Spalte von } V^\top \\ &= 0.56 \cdot v_1 - 0.12 \cdot v_2 + \cdots + 0.48 \cdot v_5. \end{aligned}$$

Der springende Punkt ist, dass wir die Basen  $\mathcal{U}$  und  $\mathcal{V}$  gerade so gewählt haben, dass für alle  $b \in B$  und alle  $f \in F$  gilt

$$R(b, f) = (b)_{\mathcal{U}}^\top \cdot \Sigma \cdot (f)_{\mathcal{V}}.$$

Da  $\Sigma$  nur drei Diagonalelemente ungleich Null hat, ist also jede Bewertung, von künstlichen oder echten Bewerterinnen bzw. Filmen, eine Summe dreier Zahlen in denen die Singulärwerte als fallende Gewichte interpretiert werden können. Wegen Satz 7.17(i) wissen wir überdies, dass die Bewertungen aus unserer Ausgangstabelle nur minimal verfälscht werden wenn wir  $\sigma_3 = 1.3$  auf Null setzen, also nur zwei Summanden behalten. Für Baileys Bewertung von Titanic haben wir z.B.

$$\begin{aligned} R(\text{Bailey}, \text{Titanic}) &= [0.13 \ -0.02 \ \cdots \ -0.06] \begin{bmatrix} 12.4 & & & \\ & 9.5 & & \\ & & 1.3 & \\ & & & \ddots \end{bmatrix} \begin{bmatrix} 0.09 \\ 0.69 \\ \vdots \\ 0.00 \end{bmatrix} \\ &= 0.13 \cdot 12.4 \cdot 0.09 - 0.02 \cdot 9.5 \cdot 0.69 + 0.01 \cdot 1.3 \cdot 0.09 \\ &\approx 0.13 \cdot 12.4 \cdot 0.09 - 0.02 \cdot 9.5 \cdot 0.69. \end{aligned}$$

Es folgt, dass jede Bewertung fast komplett durch die ersten zwei Einträge der  $\mathcal{U}$ - bzw.  $\mathcal{V}$ -Koordinaten einer Bewerterin bzw. eines Films zustande kommt. Diese sogenannten *Hauptkomponenten* können sehr einfach an der Rang-2-Bestapproximation von  $A$  aus Beispiel 7.22 abgelesen werden

$$\begin{array}{l} \text{Abbey} \rightarrow \\ \text{Bailey} \rightarrow \\ \text{Caitlin} \rightarrow \\ \text{Daisy} \rightarrow \\ \check{A} = \\ \text{Edith} \rightarrow \\ \text{Fae} \rightarrow \\ \text{Gail} \rightarrow \end{array} \begin{bmatrix} 0.07 & 0.29 \\ 0.13 & -0.02 \\ 0.68 & -0.11 \\ 0.15 & 0.59 \\ 0.41 & -0.07 \\ 0.07 & 0.73 \\ 0.55 & 0.09 \end{bmatrix} \begin{bmatrix} 12.4 \\ 9.5 \end{bmatrix} \begin{bmatrix} \text{Alien} & \text{Casablanca} & \text{Star Wars} & \text{Titanic} & \text{Matrix} \\ \begin{bmatrix} 0.56 & 0.09 & 0.56 & 0.09 & 0.59 \\ -0.12 & 0.69 & -0.12 & 0.69 & 0.02 \end{bmatrix} \end{bmatrix}$$

und korrespondieren in unserem Fall mit zwei *versteckten Konzepten*. Schaut man oben auf die Hauptkomponenten der Filme, so liegt es nicht fern die versteckten Konzepte als ‘SF’ und ‘Romantik’ zu identifizieren: Alien, Star Wars und Matrix haben große erste Hauptkomponenten und kleine zweite während es bei Casablanca und Titanic genau andersherum ist. Wegen der Singulärwerte 12.4 und 9.5 hat dabei ‘SF’ größeren Einfluss auf die Bewertungen als ‘Romantik’. An den Hauptkomponenten der Bewerterinnen kann man ablesen, wie stark sie auf die zwei Konzepte reagieren. Beispielsweise ist Bailey für das erste Konzept wenig und für das zweite gar nicht empfänglich. Für Caitlins Bewertung spielt hingegen das erste Konzept eine



große Rolle und für Faes Bewertung das zweite.

Um das Obige etwas formaler zu machen, beachten wir zunächst, dass wir zwar alle künstlichen Bewerterinnen und Filme brauchen um die echten Bewerterinnen und Filme jeweils einzeln zu rekonstruieren, dass aber jeweils die ersten beiden ausreichen, wenn wir *nur die Bewertungen* rekonstruieren wollen. Wir definieren die *Konzepträume*

$$U_2 := \text{span}\{u_1, u_2\} \subseteq B; \text{ und } V_2 := \text{span}\{v_1, v_2\} \subseteq F$$

der Bewerterinnen bzw. der Filme, und stattdessen sie mit den Orthonormalbasen  $\check{U} = \{u_1, u_2\}$  und  $\check{V} = \{v_1, v_2\}$  aus. Weiter bezeichnen wir mit

$$\check{U} = [u_1 \ u_2], \ \check{V} = [v_1 \ v_2] \text{ und } \check{\Sigma} = \begin{bmatrix} 12.4 & \\ & 9.5 \end{bmatrix}$$

die abgeschnittenen Matrizen. Wegen  $A \approx \check{A} = \check{U}\check{\Sigma}\check{V}^\top$  erhalten wir das fast kommutative Diagramm

$$\begin{array}{ccccc} (\text{Film})_{\mathcal{E}} & \mathbb{R}_{\mathcal{E}}^5 & \xrightarrow{A} & \mathbb{R}_{\mathcal{F}}^7 & (\text{Bewerterin})_{\mathcal{F}} \\ \downarrow & \check{V}^\top \downarrow & & \downarrow \check{U}^\top & \downarrow \\ (\pi_{V_2}(\text{Film}))_{\check{V}} & \mathbb{R}_{\check{V}}^2 & \xrightarrow{\check{\Sigma}} & \mathbb{R}_{\check{U}}^2 & (\pi_{U_2}(\text{Bewerterin}))_{\check{U}} \end{array}$$

in welchem die Einträge der Vektoren  $(\pi_{U_2}(\text{Bewerterin}))_{\check{U}}$  und  $(\pi_{V_2}(\text{Film}))_{\check{V}}$  gerade die Hauptkomponenten sind. Lassen wir zur Schreiberleichterung  $\pi_{V_2}$  und  $\pi_{U_2}$  weg, so suggerieren die Ausdrücke  $(\text{Bewerterin})_{\check{U}}$  und  $(\text{Film})_{\check{V}}$ , dass es sich um die Koordinaten der Bewerterin und des Films jeweils im Konzeptraum handelt. Dies ist zwar nicht korrekt, denn Nachrechnen (!) zeigt, dass z.B. weder

$$\text{Bailey} \approx 0.13 \cdot u_1 - 0.02 \cdot u_2 \quad \text{noch} \quad \text{Titanic} \approx 0.09 \cdot v_1 + 0.69 \cdot v_2$$

gilt, sondern eben nur

$$R(\text{Bailey}, \text{Titanic}) \approx [0.13 \quad -0.02] \begin{bmatrix} 12.4 & \\ & 9.5 \end{bmatrix} \begin{bmatrix} 0.09 \\ 0.69 \end{bmatrix}.$$

Die Vorstellung ist aber zweckdienlich, denn  $\begin{bmatrix} 0.13 \\ -0.02 \end{bmatrix}$  und  $\begin{bmatrix} 0.09 \\ 0.69 \end{bmatrix}$  enthalten offenbar alle Informationen über Bailey und Titanic, die nötig sind, um Bewertungen von Bailey für echte Filme oder Bewertungen echter Bewerterinnen für Titanic zu rekonstruieren.

Die künstlichen Bewerterinnen  $u_1, u_2$  können als *Prototypen von Bewerterinnen* gesehen werden, die auf genau eins der Konzepte ansprechen. Die künstlichen Filme  $v_1, v_2$  kann man sich analog als *Prototypen von Filmen* vorstellen, die auf genau einem der Konzepte basieren. In unserem Beispiel ist  $v_1$  ein *reiner SF-Film* und  $u_1$  eine *reine SF-Fanin*, während  $v_2$  ein *reiner Romantikfilm* ist und  $u_2$  eine *reine*

*Romantik-Fanin.* In der Tat sprengt z.B.

$$R(u_1, \text{Matrix}) \approx (u_1)_u^T \begin{bmatrix} 12.4 & \\ & 9.5 \end{bmatrix} (\text{Matrix})_{\tilde{v}} = [1 \ 0] \begin{bmatrix} 12.4 & \\ & 9.5 \end{bmatrix} \begin{bmatrix} 0.59 \\ 0.02 \end{bmatrix} = 7.19$$

sogar die Bewertungsskala! — Dies zeigt, dass es in Anwendungen nötig sein kann, die Daten erst einem geeigneten Preprocessing zu unterziehen, siehe die Bemerkungen ganz am Ende des Kapitels.

Manche Autoren nennen die  $u_i$  *typische Bewerterinnen* und die  $v_i$  *typische Filme*. Dabei ist dann ‘typisch’ nicht im Sinne von durchschnittlich zu verstehen, sondern im obigen Sinne eines Typs von Film.

Zum Abschluss bemerken wir noch, dass man in unserem niedrigdimensionalen Beispiel die versteckten Konzepte auch direkt an der Bewertungsmatrix ablesen kann, wenn man geeignet umsortiert.

	Alien	Matrix	Star Wars	Titanic	Casablanca
Bailey	<b>1</b>	<b>1</b>	<b>1</b>	0	0
Edith	<b>3</b>	<b>3</b>	<b>3</b>	0	0
Gail	<b>4</b>	<b>4</b>	<b>4</b>	0	0
Caitlin	<b>5</b>	<b>5</b>	<b>5</b>	0	0
Daisy	0	2	0	<b>4</b>	<b>4</b>
Fae	0	0	0	<b>5</b>	<b>5</b>
Abbey	0	1	0	<b>2</b>	<b>2</b>

Entscheidend ist aber, dass die vorgestellte SVD-basierte Methode die Hauptkomponenten und Konzepträume automatisch findet, und zwar auch dann wenn wir a priori keine Vermutungen haben, was für versteckte Konzepte hinter einer Datenmatrix liegen könnten. In der Tat kann es sogar sein, dass wir auch nach Durchführung der Hauptkomponentenanalyse keine gute Interpretation haben (z.B. weil uns ganz und gar unbekannte Produkte A, B, C, D, ... bewertet wurden).

### 7.3 Kollaboratives Filtern

Eine Instanz kollaborativen Filterns haben wir bereits in Kapitel 3 diskutiert und dort auch in Bemerkung 3.13(ii) erläutert, dass die Bezeichnung daher rührt, dass Vorhersagen für einen Bewerter auf Bewertungen anderer Bewerter basieren. Im dortigen Beispiel 3.12 hat dieser Vergleich auch der Kosinusähnlichkeit der Bewertungsvektoren beruht. Im folgenden nutzen wir stattdessen eine lineare Fortsetzung der abgeschnittenen Diagonalmatrix  $\tilde{\Sigma}$ .

Wir beginnen mit dem Szenario, dass zusätzlich zu der Bewertungsmatrix aus Beispiel 7.22 vorstellen wollen, eine weitere Bewerterin mit Bewertungen nur für die Filme Alien und Matrix gegeben ist.

	Alien	Casablanca	Star Wars	Titanic	Matrix
Hazel	4	1	n/a	n/a	?

Um die Bewertung von Hazel für den Film Matrix vorherzusagen, benutzen wir nun zuerst die gegebenen zwei Bewertungen um Hazel in der Basis  $\mathcal{U}_2$  zu entwickeln. Diese führen in der Tat auf das lineare Gleichungssystem

$$\begin{aligned} R(\text{Hazel}, \text{Alien}) &= [x \ y] \begin{bmatrix} 12.4 & \\ & 9.5 \end{bmatrix} \begin{bmatrix} 0.56 \\ -0.12 \end{bmatrix} = 6.94x - 1.14y \stackrel{!}{=} 4 \\ R(\text{Hazel}, \text{Casablanca}) &= [x \ y] \begin{bmatrix} 12.4 & \\ & 9.5 \end{bmatrix} \begin{bmatrix} 0.09 \\ 0.69 \end{bmatrix} = 1.12x + 6.55y \stackrel{!}{=} 1 \end{aligned}$$

was die Lösung  $x = 0.58$  und  $y = 0.04$  hat, woraus folgt  $(\text{Hazel})_{\mathcal{U}_2} = \begin{bmatrix} 0.58 \\ 0.04 \end{bmatrix}$ . Damit können wir nun

$$R(\text{Hazel}, \text{Matrix}) = [0.58 \ 0.04] \begin{bmatrix} 12.4 & \\ & 9.5 \end{bmatrix} \begin{bmatrix} 0.59 \\ 0.02 \end{bmatrix} = 4.24$$

berechnen und auf dieser Grundlage die Bewertung 4 des Films Matrix durch Hazel vorhersagen.

Das SVD-basierte kollaborative Filtern eignet sich gut für sogenannte *out-of-sample-Vorhersagen* wie wir sie oben durchgeführt haben. Bei ausreichend großer anfänglich bekannter Stichprobe können Bewertungen neuer Bewerterinnen oder auch neuer Filme durch die Bewertungen der Prototypen (=Basisvektoren der Konzepträume) dargestellt werden. Beachte, dass es hierfür (im Gegensatz z.B. zur Vorhersage mit Kosinusähnlichkeit) nicht nötig ist, dass eine neue Bewerterin einen sehr ähnlichen Filmgeschmack hat wie eine existierende Bewerterin. Oben war das der Fall, liegt aber daran, dass wir eine sehr kleine Matrix und nur zwei Konzepte hatten.

Wir haben oben schon gesehen, dass die Methode so wie wir sie oben erklärt haben, nicht automatisch Vorhersagen im Intervall  $[0, 5]$  liefert. Um solche Effekte zu korrigieren muss man nachsteuern, z.B. vorher normalisieren und danach runden. Darüber hinaus können Basisvektoren auftauchen, die schwieriger zu interpretieren sind als das in unserem künstlichen Beispiel der Fall war.

Sind bereits alle Bewerterinnen und Filme zu Beginn gegeben, aber nicht alle Bewertungen bekannt, so spricht man von *in-sample-Vorhersagen*, d.h. man will die fehlenden Einträge der Matrix

$$A = \begin{bmatrix} 1 & 1 & \square & 0 & 5 & \square \\ 2 & \square & 3 & \square & 4 & \square \\ \square & \square & 5 & 0 & 0 & \square \\ \square & 2 & 2 & \square & \square & 5 \end{bmatrix}$$

vorhersagen. Analog zu Beispiel 3.11 kann man hier provisorisch mit Zeilen- oder Spaltenmittelwert auffüllen, die Singulärwertzerlegung berechnen, dann das aufgefüllte wieder löschen und neu mit der SVD-Methode vorhersagen. Alternativ kann man anstelle der Faktorisierung  $A = U\Sigma V^\top$  auch  $A = PQ^\top$  verwenden, wobei  $P \in \mathbb{R}^{n \times k}$  und  $\mathbb{R}^{d \times k}$  orthogonale Spalten haben. Bei der letzteren Faktorisierung sieht man die Singulärwerte nicht mehr, aber diese explizit zu kennen ist für die Be-

rechnung von Vorhersagen auch nicht unbedingt nötig—solange man der Methode vertraut. Es gilt dann

$$(P, Q) \in \underset{\substack{(P, Q) \text{ haben} \\ \text{orthogonale} \\ \text{Spalten}}}{\operatorname{argmin}} \|A - PQ^\top\|_F^2 = \underset{\substack{(P, Q) \text{ haben} \\ \text{orthogonale} \\ \text{Spalten}}}{\operatorname{argmin}} \sum_{i,j} (a_{ij} - \langle p_i, q_j \rangle)^2$$

wobei  $p_i$  die  $i$ -te Zeile von  $P$  und  $q_j$  die  $j$ -te Zeile von  $Q$  bezeichnet. Fehlen jetzt einige der  $a_{ij}$ 's, so lässt man diese Summanden weg und optimiert über den Rest. Mit dieser Methode hat Simon Funk bei der Netflix Challenge 2006 einen vorderen Platz erreicht.

## Referenzen

Die Singulärwertzerlegung ist von erheblicher theoretischer und praktischer Bedeutung in vielen Anwendungen [Kal96]. Unser Zugang zu Beginn dieses Kapitels folgt [SSBD14, Appendix C.4]. Der algorithmische Beweis für die Existenz der SVD in Bemerkung 7.13 folgt [DH08, Satz 5.15]. Der wagemutige Vorlesende einer Linearen Algebra Grundvorlesung kann mit letzterem beginnen und dann, als Spezialfall, den Satz über die orthogonale Diagonalisierbarkeit bringen. Die Anwendungen in den Unterkapiteln 7.3–7.1 folgen [LRU12]. Wir verweisen außerdem auf [Agg16] für eine sehr ausführliche Diskussion des kollaborativen Filterns. Mehr Informationen zur Netflix Challenge finden sich in [Gow14] und unter [Fun06] kann ein originaler Blogpost von Simon Funk nachgelesen werden.

## Aufgaben

**Aufgabe 7.1.** Finden Sie Beispiele für Matrizen  $A \in \mathbb{R}^{n \times d}$  mit

- (i)  $\sigma(A^\top A) \neq \sigma(AA^\top)$ ,
- (ii)  $\dim \ker(A^\top A - 0) \neq \dim \ker(AA^\top - 0)$ .

**Aufgabe 7.2.** Sei  $A \in \mathbb{R}^{n \times d}$  und  $\sigma > 0$ . Zeigen Sie, dass  $v \in \mathbb{R}^d$  Eigenvektor von  $A^\top A$  und  $u \in \mathbb{R}^n$  Eigenvektor von  $AA^\top$  jeweils mit Eigenwert  $\sigma^2$  sein kann, ohne dass  $v$  und  $u$  (ein zusammengehöriges Paar von) Singulärvektoren mit Singulärwert  $\sigma$  sind.

**Aufgabe 7.3.** Finden Sie 1-bestpassende Unterräume jeweils für die im folgenden angegebenen Punkte mithilfe von Methode 1 oder  $1\frac{1}{2}$ :

- (i)  $\begin{bmatrix} 4 \\ 0 \end{bmatrix}, \begin{bmatrix} 3 \\ -5 \end{bmatrix} \in \mathbb{R}^2$       (ii)  $\begin{bmatrix} 2 \\ 2 \end{bmatrix}, \begin{bmatrix} -1 \\ 1 \end{bmatrix} \in \mathbb{R}^2$       (iii)  $\begin{bmatrix} 3 \\ 2 \end{bmatrix}, \begin{bmatrix} 2 \\ 6 \end{bmatrix} \in \mathbb{R}^2$
- (iv)  $\begin{bmatrix} 3 \\ 2 \end{bmatrix}, \begin{bmatrix} 2 \\ 3 \\ -2 \end{bmatrix} \in \mathbb{R}^3$ ,      (v)  $\begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix} \in \mathbb{R}^4$ .

Malen Sie Bilder für (i)–(iv).

**Aufgabe 7.4.** Das folgende Graustufenbild



wird durch die [hier](#) abrufbare Matrix  $A \in \mathbb{R}^{320 \times 240}$  dargestellt. Berechnen Sie für  $A$  eine Singulärwertzerlegung. Betrachten Sie dann für  $k = 1, 5, 10, 50$  jeweils  $\tilde{A}$  entsprechend (7.1).

- (i) Plotten Sie jeweils das durch  $\tilde{A}$  gegebene Graustufenbild.
- (ii) Berechnen Sie jeweils  $\left\| \begin{bmatrix} \sigma_1 \\ \vdots \\ \sigma_k \end{bmatrix} \right\| / \left\| \begin{bmatrix} \sigma_1 \\ \vdots \\ \sigma_r \end{bmatrix} \right\|$  und interpretieren Sie diese Zahl.
- (iii) Berechnen Sie jeweils, wieviele reelle Zahlen gespeichert werden müssen, um das  $\tilde{A}$  entsprechende Bild aus diesen zu rekonstruieren.

**Aufgabe 7.5.** Nehmen Sie an, dass wir zusätzlich zu den Informationen in Beispiel 7.22 noch wissen, dass der Film Notting Hill durch Abbey mit 4 und durch Bailey mit 1 bewertet wurde. Welche Bewertung durch Caitlin und welche durch Fae sagt dann die Methode aus Kapitel 7.3 vorher?

**Aufgabe 7.6.** Betrachte die folgende Tabelle von Buchbewertungen.

	Hunger Games	Jane Eyre	Twilight	Animal Farm	Da Vinci Code
Aaron	1	0	1	0	1
Bobby	0	2	0	2	1
Charles	5	1	5	0	5
Dustin	1	4	1	4	3
Ethan	3	1	3	0	3
Finn	0	5	0	5	0
Gareth	4	1	4	0	4

- (i) Angenommen, Leser Harvey bewertet Hunger Games mit 2 und Jane Eyre mit 5. Welche der anderen Bücher würden Sie Harvey als Leseempfehlung geben?
- (ii) Was könnten die versteckten Konzepte hinter der Bewertungstabelle sein?

*Hinweis:* Nutzen Sie ohne Beweis, dass die Bewertungsmatrix  $A$  die folgende SVD hat:

$$A = \begin{bmatrix} 0.12 & 0.07 & 0.01 & -0.54 & 0.57 & 0.57 & -0.14 \\ 0.12 & -0.27 & 0.32 & 0.22 & 0.00 & 0.00 & -0.86 \\ 0.62 & 0.28 & -0.09 & -0.40 & -0.57 & 0.00 & -0.14 \\ 0.36 & -0.48 & 0.66 & -0.08 & 0.00 & 0.00 & 0.43 \\ 0.38 & 0.14 & -0.11 & 0.68 & 0.00 & 0.57 & 0.14 \\ 0.19 & -0.72 & -0.65 & -0.06 & 0.00 & 0.00 & 0.00 \\ 0.50 & 0.21 & -0.10 & 0.13 & 0.57 & -0.57 & 0.00 \end{bmatrix} \begin{bmatrix} 13.2 & 0.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 8.9 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 1.3 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.3 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.0 & 0.0 \end{bmatrix} \begin{bmatrix} 0.51 & 0.31 & 0.51 & 0.20 & 0.57 \\ 0.25 & -0.61 & 0.25 & -0.68 & 0.11 \\ -0.39 & -0.20 & -0.39 & 0.02 & 0.80 \\ -0.11 & 0.69 & -0.11 & -0.69 & 0.07 \\ -0.70 & 0.00 & 0.70 & 0.00 & 0.00 \end{bmatrix}.$$

## Kapitel 8

# Fluch und Segen der hohen Dimension

Datenmengen bestehen oft nicht nur aus sehr vielen Datenpunkten, sondern oft besteht auch jeder Datenpunkt aus sehr vielen Features und wir haben es dann mit Mengen

$$D = \{(x^{(i)}, y^{(i)}) \mid i = 1, \dots, n\} \subseteq \mathbb{R}^d \times \mathbb{R}$$

zu tun, bei denen  $n$  und  $d$  beide sehr groß sind. Als konkretes Beispiel rufen wir die deutsche Wikipedia ins Gedächtnis. Diese besteht aus ca. 2.8 Millionen Artikeln. Vektorisiert man diese via der Grundformen der deutschen Sprache, so kann jeder Artikel als Element des  $\mathbb{R}^{148000}$  aufgefasst werden. In der Tat kann aber  $d$  sogar sehr viel größer sein als  $n$ . Betrachte z.B. alle Informationen die für die Bewohner eines Landes wie Deutschland existieren. Schließt man hier neben biographischen und medizinischen Daten, Socialmedia Posts und jemals verfassten Textnachrichten insbesondere jedes Foto und jedes Video mit ein, so kann man sich vorstellen, dass man bei der Dimension schnell in die Milliarden gerät, während die Anzahl der Datenpunkte in den Millionen verbleibt.

In den vorhergehenden Kapiteln haben wir teils mehr und teils weniger weit entwickelte mathematische Methoden gesehen, mit denen Datenmengen der Form  $D$ , oder auch solche mit mehrdimensionalen Labels  $y \in \mathbb{R}^m$ , untersucht werden können. Bei diesen Methoden scheint es auf den ersten Blick unerheblich, wie groß  $n$  und  $d$  am Ende sind, da die benutzte Lineare Algebra, Analysis oder Wahrscheinlichkeitstheorie unabhängig davon funktioniert. Auf den zweiten Blick wird klar, dass bei der konkreten Berechnung eines nächsten Nachbarn, einer Kosinusähnlichkeit, einer Singulärwertzerlung oder des Minimierers einer Kostenfunktion, die Größe von  $n$  und  $d$  sehr wohl relevant ist: je größer letztere sind, desto aufwendiger wird die Berechnung. Dies wird oft als der *Fluch der hohen Dimension* bezeichnet. Aus diesem Grund sind z.B. Techniken zur Dimensionalitätsreduktion erstrebenswert und zwar insbesondere solche die nicht selbst wieder auf einer mathematischen Methode basieren die in hohen Dimensionen Schwierigkeiten machen kann. Nach einigen Vorbereitungen werden wir eine solche in Kapitel 11 vorstellen.

Neben dem Problem der praktischen Berechenbarkeit haben Räume hoher Dimension einige zunächst sehr unintuitiv erscheinende geometrische Eigenschaften, die sich insbesondere dann zeigen, wenn man die Verteilung zufällig gewählter Punkte betrachtet. Obwohl gewöhnungsbedürftig, führen diese hochdimensionalen Effekte zu Methoden, die in mäßig großen Dimensionen nicht zur Verfügung stehen. Dies bezeichnet man als *Segen der hohen Dimension*. In diesem ersten von mehreren Kapiteln zu hochdimensionalen Räumen werden wir einige der unintuitiven Effekte kennenlernen. Allgemeine Resultate zu den vorgenannten Regelmäßigkeiten werden wir in späteren Kapiteln vorführen und dann auch zu deren Rolle bei der Behandlung hochdimensionaler Daten zurückkommen.

Im Folgenden bezeichnen wir wie immer mit  $\|\cdot\|$  die euklidische Norm, mit  $\mathcal{B}^d$  die  $\sigma$ -Algebra der Borelmengen und mit  $\lambda^d$  das Lebesguemaß auf  $\mathbb{R}^d$ .

**Definition 8.1.** Sei  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum, sei  $B \in \mathcal{B}^d$  eine Borelmenge und sei  $X: \Omega \rightarrow \mathbb{R}^d$  ein Zufallsvektor.

(i) Wir sagen, dass  $X$  auf  $B$  *gleichmäßig verteilt* ist, falls

$$P[X \in A] = \frac{\lambda^d(A \cap B)}{\lambda^d(B)}$$

für jede Borelmenge  $A \in \mathcal{B}^d$  gilt. Ist dies der Fall, so schreiben wir  $X \sim \mathcal{U}(B)$  und notieren bereits hier, dass wir später insbesondere den Fall des *Hypercubes*  $H_d := [-1, 1]^d$  und den der abgeschlossenen *Einheitskugel*  $B_{\mathbb{R}^d} := \bar{B}_1(0) \subseteq \mathbb{R}^d$  betrachten werden.

(ii) Wir sagen, dass  $X$  (*sphärisch*) *gaußverteilt* ist mit Mittelwert  $\mu \in \mathbb{R}^d$  und Varianz  $\sigma > 0$ , falls

$$P[X \in A] = \frac{1}{(2\pi\sigma^2)^{d/2}} \int_A e^{-\frac{\|x-\mu\|^2}{2\sigma^2}} d\lambda^d(x)$$

für jede Borelmenge  $A \in \mathcal{B}^d$  gilt. Ist dies der Fall, so schreiben wir  $X \sim \mathcal{N}(0, \sigma^2, \mathbb{R}^d)$ , oder, wenn keine Verwechslungsgefahr besteht, nur  $X \sim \mathcal{N}(0, \sigma^2)$ . Eine einfache Rechnung, siehe Anhang A.2, zeigt, dass  $X = (X_1, \dots, X_d) \sim \mathcal{N}(0, \sigma^2, \mathbb{R}^d)$  genau dann gilt, wenn  $X_i \sim \mathcal{N}(0, \sigma^2, \mathbb{R})$  für alle  $i = 1, \dots, d$  und die  $X_i$  unabhängig sind. Ist  $\mu = 0$  und  $\sigma = 1$ , so sprechen wir von *normalverteilten* Zufallsvektoren.

Wir beginnen nun mit dem  $d$ -dimensionalen Hypercube  $H_d = [-1, 1]^d$ , welcher offenbar eine kompakte und absolutkonvexe Teilmenge von  $\mathbb{R}^d$  ist. Betrachten wir seine ‘Ecken’  $v = (\pm 1, \pm 1, \dots, \pm 1) \in H_d$ , so sehen wir, dass

$$\|v - 0\| = \sqrt{(\pm 1)^2 + (\pm 1)^2 + \dots + (\pm 1)^2} = \sqrt{d} \xrightarrow{d \rightarrow \infty} \infty$$

gilt, und damit deren Abstände zum Mittelpunkt 0 des Hypercubes mit wachsender Dimension immer größer werden. Außerdem gibt es mit wachsendem  $d$  immer mehr, nämlich  $2^d$ -viele, Ecken für die das vorgenannte gilt. Bezeichnen wir mit

$e_i = (0, \dots, 0, 1, 0, \dots, 0)$  den  $i$ -ten Einheitsvektor, so gilt für  $i \neq j$  andererseits

$$\|\pm e_i - (\pm e_j)\| = \|(0, \dots, 0, \pm 1, 0, \dots, 0, \pm 1, 0, \dots, 0)\| = \sqrt{2}$$

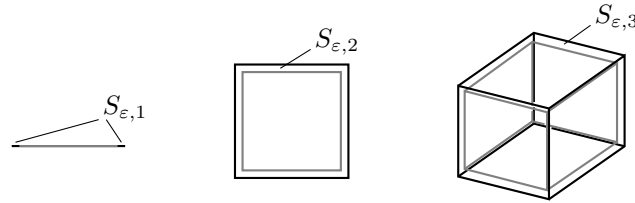
und weiter

$$\|e_i - (-e_i)\| = \|(0, \dots, 0, 2, 0, \dots, 0)\| = 2.$$

Liest man die  $\pm e_i$  als Mittelpunkte zweier verschiedener ‘Seiten’ des Hypercubes, so heißt obiges, dass die Abstände beliebiger verschiedener Seitenmittelpunkte konstant, d.h., unabhängig von der Dimension, sind. Beide Berechnungen legen in Kombination nahe, dass die Ecken des Hypercubes mit wachsender Dimensionen wie immer länger werdende Stacheln nach außen abstecken, während die Seitenmittelpunkte nahe dem Ursprung verbleiben. Wenn sich Abstände bei wachsendem  $d$  nicht homogen verhalten, so steht zu erwarten, dass das Volumen im hochdimensionalen Hypercube inhomogen verteilt ist. Um dies zu sehen, betrachten wir für  $0 < \varepsilon < 1$  die Teilmenge

$$S_{\varepsilon,d} = H_d \setminus (1 - \varepsilon) H_d$$

von  $H_d$ , die wir uns als  $\varepsilon$ -dicke äußere Schale vorstellen können und die in den folgenden drei Bildern für niedrige Dimensionen veranschaulicht ist.



Da  $(1 - \varepsilon) H_d = [-1 + \varepsilon, 1 - \varepsilon]^d$  ein Quader ist, können wir das Verhältnis des Volumens von  $S_{\varepsilon,d}$  zum Volumen des gesamten Hypercubes wie folgt berechnen

$$\frac{\lambda^d(S_{\varepsilon,d})}{\lambda^d(H_d)} = \frac{\lambda^d(H_d) - \lambda^d((1 - \varepsilon) H_d)}{\lambda^d(H_d)} = \frac{2^d - ((2(1 - \varepsilon))^d)}{2^d} = 1 - (1 - \varepsilon)^d \xrightarrow{d \rightarrow \infty} 1$$

wobei  $0 < \varepsilon < 1$  beliebig ist. Letzteres bedeutet, dass der überwiegende Teil des Volumens des Hypercubes nah an dessen Oberfläche verortet ist.

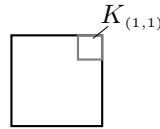
Das Volumen des gesamten Hypercubes  $H_d = [-1, 1]^d$ , so wie wir ihn bisher betrachtet haben, geht gegen unendlich. Betrachten wir stattdessen einen Hypercube mit Seitenlänge Eins, z.B.  $[0, 1]^d$ , so haben wir  $\lambda^d([0, 1]^d) = 1$  für jede Dimension  $d$ . Das Volumen in einer  $\varepsilon$ -dicken Schale  $\lambda^d([0, 1]^d \setminus [\varepsilon, 1 - \varepsilon]^d) = 1 - (1 - 2\varepsilon)^d$  konvergiert gegen Eins solange  $0 < \varepsilon < 1/2$  ist. Wir haben also auch hier den Effekt, dass für großes  $d$  fast das ganze Volumen des Hypercubes nah an dessen Oberfläche liegt.

Andererseits kann auch nicht viel Volumen von  $[-1, 1]^d$  bzw.  $[0, 1]^d$  in den Ecken liegen: Wir betrachten für eine Ecke  $v = (v_1, \dots, v_d) \in H_d$  mit  $v_i \in \{1, -1\}$  die Menge

$$K_v := \{x \in \mathbb{R}^d \mid x_i \in [1 - \varepsilon, 1] \text{ falls } v_i = 1 \text{ und } x_i \in [-1, -1 + \varepsilon] \text{ falls } v_i = -1\}$$



also einen Quader mit Seitenlänge  $0 < \varepsilon < 1/2$  der in der  $v$ -ten Ecke von  $H_d$  sitzt, wie im folgenden Bild für  $d = 2$  veranschaulicht.

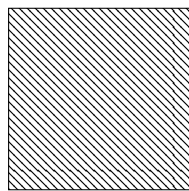


Es folgt  $\lambda^d(K_v) = \varepsilon^d$  und da es  $2^d$  Ecken gibt, geht das Volumen in der Nähe der Ecken gegen Null

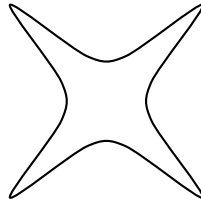
$$\sum_{v \in V} \lambda^d(K_v) = 2^d \cdot \varepsilon^d \xrightarrow{d \rightarrow \infty} 0$$

wobei wir mit  $V$  die Menge der Ecken von  $H_d$  bezeichnen. Wir können uns also vorstellen, dass das Volumen von  $H_d$  einerseits in einer dünnen Schicht nah unter der Oberfläche des Hypercubes verortet ist, und sich gleichzeitig um die Seitenmittelpunkte herum ansammelt.

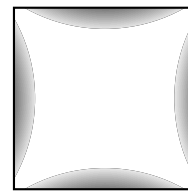
Die folgenden Bilder illustrieren die drei ganz verschiedenen Eigenschaften des Hypercubes, die wir oben diskutiert haben. Wir halten den Leser zum vorsichtigen Umgang mit diesen Bildern an: Neben der Tatsache, dass jeweils nur ein bestimmter Effekt veranschaulicht und die restlichen Eigenschaften ignoriert werden, handelt es sich um 2-dimensionale Bilder eines hochdimensionalen Objektes.



Der Hypercube ist kompakt und konvex.



Die Abstände der Ecken vom Mittelpunkt wachsen mit der Dimension, die der Seitenmittelpunkte nicht.

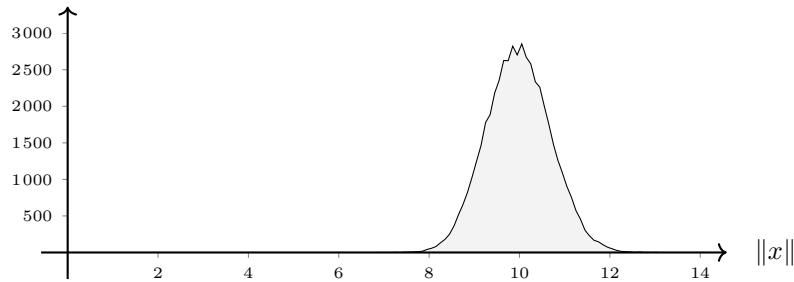


Das Volumen konzentriert sich an der Oberfläche und um den Seitenmittelpunkt.

Interpretieren wir unsere obigen Erkenntnisse über die Volumenverteilung in einer probabilistischen Weise, so ergibt sich, dass aus  $H_d$  gleichmäßig zufällig ausgewählte Punkte mit hoher Wahrscheinlichkeit nah an dessen Oberfläche und nah an der Mitte einer Seite liegen werden. Haben wir also z.B. eine Datenmenge gegeben, deren Features gleichmäßig in  $H_d$  verteilt sind, so legt obiges nahe, dass die paarweisen Abstände von Punkten wie auch deren Winkel alle sehr nah beieinander liegen werden: Das ganz rechte Bild suggeriert etwa  $\angle(x, y) \approx 90^\circ$  für zwei gleichmäßig zufällig gewählte  $x, y \in H_d$  und in Aufgabe 8.1 werden wir experimentell zeigen, dass  $\|x - y\| \approx \sqrt{\frac{3}{2}d}$  zu erwarten ist. Beides zusammen bedeutet, dass man mit Prediktoren, die auf Abständen oder beispielsweise auf der Kosinusähnlichkeit basieren, in Schwierigkeiten gerät.

Als nächstes betrachten wir Punkte, die zufällig aus dem ganzen Raum  $\mathbb{R}^d$  gewählt werden, und zwar entsprechend einer Gaußverteilung von der wir der Einfachheit halber annehmen, dass ihr Mittelwert Null und ihre Varianz Eins sind. Wie in Definition

8.1(ii) erklärt, können wir dies durch einen Zufallsvektor  $X: \Omega \rightarrow \mathbb{R}^d$  formalisieren, dessen Koordinaten unabhängige normalverteilte Zufallsvariablen sind. Das folgende Bild zeigt die Verteilung der Normen von 50 000 Punkten  $x^{(i)} \in \mathbb{R}^{100}$ , die auf diese Weise zufällig gewählt wurden.



Die Simulation zeigt, dass im Durchschnitt  $\|x\| \approx 10$  gilt und sogar fast alle Punkte in der Nähe der Oberfläche einer Kugel mit Radius 10 im  $\mathbb{R}^{100}$  liegen. Führen wir dieses Experiment mehrfach durch und variieren dabei die Dimension  $d$ , so suggerieren die Ergebnisse, dass  $E(\|X\|) \approx \sqrt{d}$  für  $X \sim \mathcal{N}(0, 1, \mathbb{R}^d)$  gilt. Außerdem scheint die Stichprobenvarianz der Normen durch eine von der Dimension unabhängige Konstante beschränkt zu sein, was zur Vermutung führt, dass  $V(X)$  beschränkt ist.

$d$	1	10	100	1 000	10 000	100 000	1 000 000
$\frac{1}{100} \sum_{i=1}^{100} \ x^{(i)}\ $	0.73	3.05	10.10	31.61	100.03	316.21	1000.03
$\sqrt{d}$	1.00	3.16	10.00	31.62	100.00	316.22	1000.00
Varianz	0.33	0.48	0.54	0.45	0.52	0.36	0.44

Wir überlassen es dem Leser als Aufgabe 8.2 die obigen Experimente selbst durchzuführen und schauen uns nun Erwartungswert und Varianz für  $X \sim \mathcal{N}(0, 1, \mathbb{R}^d)$  abstrakt an. Da die Koordinatenfunktionen  $X_i \sim \mathcal{N}(0, 1)$  unabhängig und normalverteilt sind, ergibt sich ohne Schwierigkeiten

$$\begin{aligned}
 E(\|X\|^2) &= E(X_1^2 + \dots + X_d^2) = E(X_1^2) + \dots + E(X_d^2) \\
 &= E((X_1 - E(X_1))^2) + \dots + E((X_d - E(X_d))^2) \\
 &= V(X_1) + \dots + V(X_d) = d
 \end{aligned} \tag{8.1}$$

wobei wir erstmal das Normquadrat betrachten, da a priori nicht klar, ist, wie die Wurfelfunktion mit dem Erwartungswert interagiert. Analog ergibt sich für die Varianz

$$\begin{aligned}
 V(\|X\|^2) &= V(X_1^2) + \dots + V(X_d^2) = d \cdot V(X_1^2) = d \cdot (E(X_1^4) - E(X_1^2)^2) \\
 &= d \cdot \left( \frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} x^4 \exp(-x^2/2) dx - V(X_1) \right) \\
 &= (3 - 1)d = 2d.
 \end{aligned} \tag{8.2}$$

Nach dieser Vorarbeit können wir das folgende genauere Resultat beweisen, wel-

ches das durch die Experimente gezeichnete Bild bestätigen wird: Für einen  $d$ -dimensionalen Zufallsvektor  $X$  ist  $E(X) \approx \sqrt{d}$  für  $d \rightarrow \infty$  und die Varianz ist durch eine von  $d$  unabhängige Konstante beschränkt.

**Satz 8.2.** Sei  $X \sim \mathcal{N}(0, 1, \mathbb{R}^d)$ . Dann gilt

$$(i) \quad \forall d \in \mathbb{N}: |E(\|X\| - \sqrt{d})| \leq 1/\sqrt{d},$$

$$(ii) \quad \forall d \in \mathbb{N}: V(\|X\|) \leq 2.$$

*Beweis.* (i) Wir beginnen mit der folgenden Gleichung

$$\|X\| - \sqrt{d} = \frac{\|X\|^2 - d}{2\sqrt{d}} - \frac{(\|X\|^2 - d)^2}{2\sqrt{d}(\|X\| + \sqrt{d})^2} =: S_d - R_d$$

die aus  $\|X\|^2 - d = (\|X\| - \sqrt{d})(\|X\| + \sqrt{d})$  folgt. Dann benutzen wir (8.2) und  $\|X\| \geq 0$  um wie folgt abzuschätzen

$$0 \leq E(R_d) \leq \frac{E((\|X\|^2 - d)^2)}{2d^{3/2}} = \frac{V(\|X\|^2)}{2d^{3/2}} = \frac{2d}{2d^{3/2}} = \frac{1}{\sqrt{d}}.$$

Es folgt  $E(R_d) \rightarrow 0$  für  $d \rightarrow \infty$ . Da  $E(\|X\|^2) = d$  gilt, erhalten wir  $E(S_d) = 0$  und folglich

$$|E(\|X\| - \sqrt{d})| = |E(S_d - R_d)| = |-E(R_d)| \leq \frac{1}{\sqrt{d}}$$

wie behauptet.

(ii) Für die Varianz berechnen wir

$$\begin{aligned} V(\|X\|) &= V(\|X\| - \sqrt{d}) = E((\|X\| - \sqrt{d})^2) - (E(\|X\| - \sqrt{d}))^2 \\ &\leq E((\|X\| - \sqrt{d})^2) = E(\|X\|^2 - 2\|X\|\sqrt{d} + d) \\ &= E(\|X\|^2) - 2\sqrt{d} E(\|X\|) + d = 2d - 2\sqrt{d} E(\|X\| - \sqrt{d} + \sqrt{d}) \\ &= 2\sqrt{d} E(R_d) \\ &\leq 2 \end{aligned}$$

womit (ii) gezeigt ist. □

Als nächstes wollen wir den Abstand zweier unabhängiger normalverteilter Zufallsvektoren  $X, Y \sim \mathcal{N}(0, 1, \mathbb{R}^d)$  untersuchen. Für den quadrierten Abstand ergibt sich wieder durch direkte Rechnung

$$\begin{aligned} E(\|X - Y\|^2) &= \sum_{i=1}^d E((X_i - Y_i)^2) \\ &= \sum_{i=1}^d E(X_i^2) + 2E(X_i)E(Y_i) + E(Y_i^2) \\ &= \sum_{i=1}^d (1 + 2 \cdot 0 \cdot 0 + 1) = 2d \end{aligned}$$

was  $E(\|X - Y\|) \approx \sqrt{2d}$  nahelegt. Experimente bestätigen dies, siehe Aufgabe 8.2, und in der Tat zeigen ähnliche Argumente wie die im Beweis von Satz 8.2 das folgende.

**Satz 8.3.** *Seien  $X, Y \sim \mathcal{N}(0, 1, \mathbb{R}^d)$ . Dann gilt*

$$(i) \quad \forall d \in \mathbb{N}: E(\|X - Y\| - \sqrt{2d}) \leq 1/\sqrt{2d},$$

$$(ii) \quad \forall d \in \mathbb{N}: V(\|X - Y\|) \leq 3. \quad \square$$

Die Details des Beweises lassen wir als Aufgabe 8.3.

Um den Erwartungswert des Winkels  $\angle(X, Y)$  zu bestimmen, können wir auf die bereits bekannten Gleichungen  $E(\|X\|^2) = E(\|Y\|^2) = d$  und  $E(\|X - Y\|^2) = 2d$  zurückgreifen und erhalten damit das folgende Resultat.

**Satz 8.4.** *Seien  $X, Y \sim \mathcal{N}(0, 1, \mathbb{R}^d)$  und  $\xi \in \mathbb{R}^d$  ein konstanter Vektor. Dann gelten*

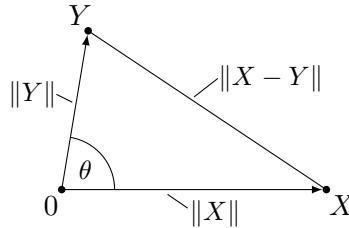
$$(i) \quad E(\langle X, Y \rangle) = 0 \text{ und } V(\langle X, Y \rangle) = d,$$

$$(ii) \quad E(\langle X, \xi \rangle) = 0 \text{ und } V(\langle X, \xi \rangle) = \|\xi\|^2.$$

*Beweis.* (i) Wir verwenden das Kosinusetz

$$\|X - Y\|^2 = \|X\|^2 + \|Y\|^2 - 2\|X\|\|Y\|\cos(\theta)$$

wobei  $\theta$  wie im folgenden Bild der Winkel zwischen  $X$  und  $Y$  ist.



Dann benutzen wir  $\cos \theta = \langle X/\|X\|, Y/\|Y\| \rangle$  und erhalten

$$\langle X, Y \rangle = \frac{1}{2}(\|X\|^2 + \|Y\|^2 - \|X - Y\|^2).$$

Den Erwartungswert dieses Skalarproduktes können wir nun durch Ausnutzung unserer Vorarbeiten per

$$E(\langle X, Y \rangle) = \frac{1}{2}(E(\|X\|^2) + E(\|Y\|^2) - E(\|X - Y\|^2)) = \frac{1}{2}(d + d - 2d) = 0$$

berechnen. Für die Varianz erhalten wir

$$\begin{aligned} V(\langle X, Y \rangle) &= V(X_1 Y_1 + \dots + X_d Y_d) = V(X_1 Y_1) + \dots + V(X_d Y_d) \\ &= d \cdot V(X_1 Y_1) = d \cdot (E(X_1^2) E(Y_1^2) - E(X_1)^2 E(Y_1)^2) \\ &= d \cdot (E(X_1^2) - E(X_1))^2 = d \cdot V(X_1) V(Y_1) = d. \end{aligned}$$

(ii) Mit  $\xi = (\xi_1, \dots, \xi_d)$  sieht man sofort

$$\mathbb{E}(\langle X, \xi \rangle) = \sum_{i=1}^d \mathbb{E}(X_i \xi_i) = \sum_{i=1}^d \xi_i \mathbb{E}(X_i) = 0$$

und

$$\mathbb{V}(\langle X, \xi \rangle) = \sum_{i=1}^d \mathbb{V}(\xi_i X_i) = \sum_{i=1}^d \xi_i^2 \mathbb{V}(X_i) = \|\xi\|^2$$

wie behauptet.  $\square$

Nach Satz 8.4 ist zu erwarten, dass zwei normalverteilt zufällig gewählte Punkte  $x, y \in \mathbb{R}^d$  orthogonal sind. Im Gegensatz zu allen vorherigen Resultaten sind aber die Stichprobenvarianzen unbeschränkt und man kann daher nicht schließen, dass tatsächlich für die meisten der Punkte  $x, y$  in einer Stichprobe  $\langle x, y \rangle \approx 0$  gilt. Normieren wir allerdings die Vektoren vor der Skalarproduktbildung, d.h. betrachten wir  $\langle x/\|x\|, y/\|y\| \rangle$ , so suggerieren Simulationen, vergleiche Aufgabe 8.2(iii)–(iv), dass in der Tat

$$\mathbb{E}(\langle \frac{X}{\|X\|}, \frac{Y}{\|Y\|} \rangle) \approx 0 \quad \text{und} \quad \mathbb{V}(\langle \frac{X}{\|X\|}, \frac{Y}{\|Y\|} \rangle) \xrightarrow{d \rightarrow 0} 0$$

gelten. Beachte, dass eine direkte Berechnung von Erwartungswert und Varianz nicht-trivial ist, insbesondere da  $X$  und  $\|X\|$  nicht unabhängig sind.

Bisher haben wir mit Erwartungswert und Varianz gearbeitet und argumentiert, dass eine kleine und beschränkte Varianz impliziert, dass für die Elemente einer Stichprobe die uns interessierenden Größen (Norm, Abstand, Winkel) nah beim Erwartungswert liegen. Experimente haben dies bestätigt. Im folgenden Kapitel 10 werden wir explizite Abschätzungen für die Wahrscheinlichkeiten

$$\mathbb{P}[|\|X\| - \sqrt{d}| \geq \varepsilon] \quad \text{und} \quad \mathbb{P}[|\langle X, Y \rangle| \geq \varepsilon]$$

beweisen unter der Voraussetzung  $X, Y \sim \mathcal{N}(0, 1, \mathbb{R}^d)$ . Im Fall der Gleichverteilung werden wir statt des Hypercubes die abgeschlossene Einheitskugel betrachten, für  $X, Y \sim \mathcal{U}(\bar{B}_1(0))$  ebenfalls explizite Abschätzungen des obigen Typs beweisen, und somit die Ergebnisse dieses Kapitels quantifizieren. Wir weisen jetzt schonmal vorsorglich darauf hin, dass bei diesen Abschätzungen genau darauf geachtet werden muss, ob und wenn ja in welcher Weise,  $\varepsilon \searrow 0$  und  $d \rightarrow \infty$  gehen dürfen bzw. müssen. Wir werden in den Kapiteln 9–10 zeigen, dass

Verteilung	$\ X\ $	$\ X - Y\ $	$\langle X, Y \rangle$
$X, Y \sim \mathcal{N}(0, 1, \mathbb{R}^d)$	$\approx \sqrt{d}$	$\approx \sqrt{2d}$	$\approx 0$
$X, Y \sim \mathcal{U}(\bar{B}_1(0))$	$\approx 1$	$\approx \sqrt{2}$	$\approx 0$

gilt. Hierbei ist z.B. mit  $\|X\| \approx \sqrt{d}$  gemeint, dass die Wahrscheinlichkeit dafür, dass  $|\|X\| - \sqrt{d}|$  größer als ein Schwellwert ist, durch einen Ausdruck abgeschätzt werden kann, der klein ist, wenn die Dimension groß genug ist. Idealerweise ist der

Schwellwert eine Konstante  $\varepsilon > 0$ , die wir beliebig vorgeben können und die Schranke für die Wahrscheinlichkeit hängt allein von  $d$  ab, in der Tat müssen wir aber damit leben, dass oft beide Ausdrücke von  $\varepsilon$  und  $d$  abhängen. Wir werden dann Grenzwerte vermeiden und Sätze beweisen, die garantieren, dass für ‘fast alle’ Datenpunkte einer Menge  $D$ , z.B. mehr als der  $(1 - 1/N)$ -te Teil, eine Eigenschaft gilt, oder das für jeden Punkt von  $D$  die Eigenschaft ‘mit hoher Wahrscheinlichkeit’, z.B.  $P[\cdot]$  größer als  $1 - 1/N$ , gilt. Hierbei kann  $N$  von  $n = \#D$  und  $d$  abhängen. Diese Herangehensweise wird als *nicht-asymptotische Analysis* bezeichnet.

## Referenzen

Dieses Kapitel orientiert sich hauptsächlich an [BHK20] und [Köp13]. Der Begriff *Segen der hohen Dimension* ist [Don04] entnommen und der Beweis von Satz 8.2 folgt [Pin20].

## Aufgaben

**Aufgabe 8.1.** Schreiben Sie ein Programm, z.B. in Python, das Punkte  $x^{(1)}, \dots, x^{(n)}$  erzeugt, die gleichmäßig zufällig aus dem Hypercube  $H^d(1)$  gewählt werden, indem Sie deren Koordinaten gleichmäßig zufällig aus  $[-1, 1]$  wählen lassen. Verifizieren Sie dann, dass für geeignet große  $d$

$$\begin{aligned} \text{(i)} \quad & \frac{1}{n} \sum_{i=1}^n \|x^{(i)}\|_{\infty} \approx 1, & \text{(ii)} \quad & \frac{1}{n} \sum_{i=1}^n \|x^{(i)}\|_2 \approx \sqrt{\frac{d}{3}}, \\ \text{(iii)} \quad & \frac{1}{n(n-1)} \sum_{i \neq j} \|x^{(i)} - x^{(j)}\|_2 \approx \sqrt{\frac{2}{3}d}, & \text{(iv)} \quad & \frac{1}{n(n-1)} \sum_{i \neq j} \angle(x^{(i)}, x^{(j)}) \approx \frac{\pi}{2} \end{aligned}$$

gelten. Hier bei bezeichnet  $\|\cdot\|_{\infty}$  die Maximumsnorm,  $\|\cdot\|_2$  die euklidische Norm und der Winkel  $\theta = \angle(x, y)$  ist per  $\cos \theta = \langle x/\|x\|, y/\|y\| \rangle$  gegeben. Plotten Sie außerdem die Verteilung der euklidischen Normen. Wächst jeweils die Varianz mit  $d$ ?

*Hinweis:* Testen Sie ihr Programm erst z.B. mit  $n = 10$  und  $d = 5$  und erhöhen Sie beide Parameter später. Nutzen Sie Funktionen wie `average`, `maximum`, `acos` etc. Für den Plot müssen Sie  $[\min_i \|x^{(i)}\|, \max_i \|x^{(i)}\|]$  in gleichlange Intervalle teilen. Deren Anzahl ist ein zusätzlicher Parameter für den Sie einen geeigneten Wert durch Ausprobieren finden müssen.

**Aufgabe 8.2.** Schreiben Sie ein Programm, z.B. in Python, das Punkte  $x^{(1)}, \dots, x^{(n)} \in \mathbb{R}^d$  erzeugt, deren Koordinaten normalverteilt sind.

- (i) Plotten Sie für  $d = 100$  und  $n = 50\,000$  die Verteilung der Normen. Es sollte sich ein Bild wie auf Seite 114 ergeben.
- (ii) Verifizieren Sie, dass für geeignet große  $d$

$$\frac{1}{n} \sum_{j=1}^n \|x^{(j)}\| \approx \sqrt{d} \quad \text{und} \quad \frac{1}{n(n-1)} \sum_{k \neq j} \|x^{(j)} - x^{(k)}\| \approx \sqrt{2d}$$

gelten und dass die Stichprobenvarianzen jeweils beschränkt sind. Vergleichen Sie Ihre Ergebnisse des ersten Teils mit der Tabelle auf Seite 114. Für den zweiten Teil kommt etwas sehr ähnliches heraus.

(iii) Verifizieren Sie, dass für geeignet große  $d$

$$\frac{1}{n(n-1)} \sum_{i \neq j} \langle x^{(i)}, x^{(j)} \rangle \approx 0$$

gilt, dass aber die Stichprobenvarianz gegen Unendlich geht, vergleiche Satz 8.4, bzw. Aufgabe 8.3 für eine quantitative und nicht-experimentelle Version hiervon.

(iv) Verifizieren Sie, dass für geeignet große  $d$

$$\frac{1}{n(n-1)} \sum_{i \neq j} \left\langle \frac{x^{(i)}}{\|x^{(i)}\|}, \frac{x^{(j)}}{\|x^{(j)}\|} \right\rangle \approx 0$$

gilt und dass nun die Stichprobenvarianz nicht nur beschränkt, sondern für  $d \rightarrow \infty$  sogar gegen Null konvergieren.

**Aufgabe 8.3.** Seien  $X, Y \sim \mathcal{N}(0, 1, \mathbb{R}^d)$ . Zeigen Sie:

- (i)  $\forall d \in \mathbb{N}: \mathbb{E}(\|X - Y\| - \sqrt{2d}) \leq 1/\sqrt{2d}$ .
- (ii)  $\forall d \in \mathbb{N}: \mathbb{V}(\|X - Y\|) \leq 3$ .

*Hinweis:* Berechnen Sie zuerst  $\mathbb{V}((X_i - Y_i)^2) = 3$  durch Nachweis von  $X_i - Y_i \sim \mathcal{N}(0, 2, \mathbb{R})$  und Anwendung einer geeigneten Formel zur Berechnung des vierten Moments. Folgern Sie dann  $\mathbb{V}(\|X - Y\|^2) \leq 3d$ . Adaptieren Sie schließlich die Argumente, mit denen wir  $\mathbb{E}(\|X\| - \sqrt{d})$  und  $\mathbb{V}(\|X\|)$  in Beweis von Satz 8.2 berechnet bzw. abgeschätzt haben.

**Aufgabe 8.4.** Falls Sie Horrorfilme mögen, schauen Sie sich Andrzej Sekuła's Film *Hypercube* von 2002 an.

## Kapitel 9

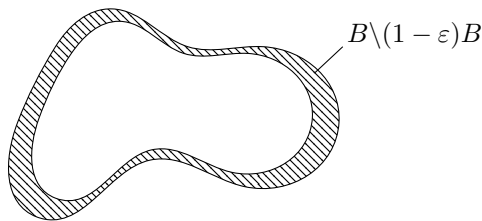
# Maßkonzentration

Wie auch schon im vorherigen Kapitel 8, bezeichnet im folgenden  $\|\cdot\|$  die euklidische Norm,  $\mathcal{B}^d$  die  $\sigma$ -Algebra der Borelmengen und  $\lambda^d$  das Lebesguemaß auf  $\mathbb{R}^d$ . Weiter bezeichnen wir mit

$$B_r(x_0) = \{x \in \mathbb{R}^d \mid \|x - x_0\| < r\} \quad \text{und} \quad \bar{B}_r(x_0) = \{x \in \mathbb{R}^d \mid \|x - x_0\| \leq r\}$$

die offene, bzw. abgeschlossene, Kugel mit Mittelpunkt  $x_0$  und Radius  $r > 0$  in  $\mathbb{R}^d$ . Die abgeschlossene Einheitskugel kürzen wir mit  $B_{\mathbb{R}^d} := \bar{B}_1(0)$  ab und wir erinnern daran, dass  $r B_{\mathbb{R}^d} = \bar{B}_r(x_0)$  gilt.

Wir beginnen nun mit der folgenden Beobachtung. Sei  $B \subseteq \mathbb{R}^d$  eine messbare Menge sodass  $\mu \cdot B \subseteq B$  für  $0 < \mu < 1$  gilt; also z.B. ein Sterngebiet mit Mittelpunkt Null, oder eine kreisförmige Menge. Dann können wir  $B \setminus (1 - \varepsilon)B \subseteq B$  als den Teil von  $B$  interpretieren, der nah an der ‘Oberfläche’ von  $B$  liegt.



Analog zum Hypercube in Kapitel 8, können wir das Verhältnis

$$\frac{\lambda^d(B \setminus (1 - \varepsilon)B)}{\lambda^d(B)} = \frac{\lambda^d(B) - \lambda^d((1 - \varepsilon)B)}{\lambda^d(B)} = \frac{\lambda^d(B) - (1 - \varepsilon)^d \lambda^d(B)}{\lambda^d(B)} = 1 - (1 - \varepsilon)^d$$

berechnen und sehen, dass dieses nah bei Eins liegt, wenn  $\varepsilon$  klein und  $d$  groß ist. Dies verallgemeinert den Fall des Hypercubes und zeigt, dass sich fast das ganze Maß von  $B$  knapp unter dessen ‘Oberfläche’ befindet, während das ‘Innere’ wenig zum Maß beiträgt. Mit der wohlbekannten Abschätzung  $e^t \leq \frac{1}{1-t}$  liefert das obige für den Spezialfall der Einheitskugel  $B = B_{\mathbb{R}^d}$  den folgenden sogenannten Satz von der Oberflächenkonzentration.



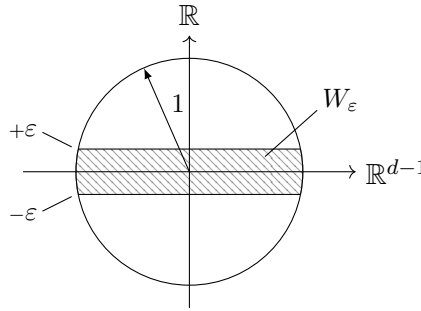
**Satz 9.1.** (Oberflächenkonzentration) *Sei  $d \in \mathbb{N}$  und  $0 < \varepsilon \leq 1$ . Dann gilt*

$$\lambda^d(\bar{B}_1(0) \setminus B_{1-\varepsilon}(0)) \geq (1 - e^{-\varepsilon^d}) \cdot \lambda^d(\bar{B}_1(0))$$

*d.h. mindestens der  $(1 - e^{-\varepsilon^d})$ -te Teil des Volumens der  $d$ -dimensionalen Einheitskugel befindet sich  $\varepsilon$ -nah an deren Oberfläche.*  $\square$

Gerade im Fall der Einheitskugel wird dieser Effekt der Maßkonzentration anschaulich klarer: Teilt man die Kugel in Schalen auf, die alle die gleiche Dicke, aber entsprechend unterschiedliche Radien haben, so werden Schalen mit kleinem Radius weniger Volumen beitragen als Schalen mit großem Radius. Letzteres ist ab Dimensionen 2 der Fall, aber je mehr Dimensionen zur Verfügung stehen, desto stärker wirkt sich der Radius auf das Volumen aus.

Um den nächsten Satz über die Tailenkonzentration zu formulieren, fassen wir die erste Koordinate des  $\mathbb{R}^d$  als ‘Norden’ auf und betrachten dann eine um den ‘Äquator’ herum aus der Einheitskugel herausgeschnittene  $2\varepsilon$ -breite Scheibe  $W_\varepsilon$ .



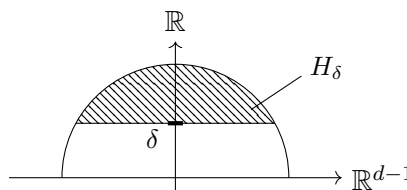
Satz 9.2 wird zeigen, dass sich der Großteil des Volumens der Einheitskugel in dieser Scheibe befindet, oder mit anderen Worten, an dessen Taille. Ebenso wie im Fall der Oberflächenkonzentration ist das nur auf den ersten Blick überraschend: Gehen wir vom Ursprung aus in ‘Nordrichtung’, so wird die Kugel dünner und zwar in allen der verbleibenden  $(d-1)$ -vielen Richtungen.

**Satz 9.2.** (Tailenkonzentration) *Sei  $d \in \mathbb{N}_{\geq 3}$  und  $\varepsilon > 0$ . Dann gilt*

$$\lambda^d(W_\varepsilon) \geq \left(1 - \frac{2}{\varepsilon \sqrt{d-1}} e^{-\frac{\varepsilon^2(d-1)}{2}}\right) \cdot \lambda^d(B_{\mathbb{R}^d})$$

wobei  $W_\varepsilon = \{(x_1, \dots, x_d) \in B_{\mathbb{R}^d} \mid |x_1| \leq \varepsilon\}$ .

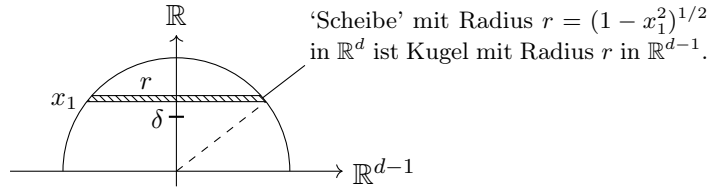
*Beweis.* Wir stellen zuerst fest, dass die Aussage für  $\varepsilon \geq 1$  trivial ist und nehmen daher  $\varepsilon < 1$  an. Für  $\delta \geq 0$  definieren wir  $H_\delta = \{(x_1, \dots, x_d) \in \bar{B}_1(0) \mid x_1 \geq \delta\}$ , siehe das folgende Bild, und verwenden, dass offenbar  $\frac{1}{2} \cdot \lambda^d(W_\varepsilon) = \lambda^d(H_0 \setminus H_\varepsilon)$  gilt.



Wir berechnen also als erstes das Volumen von  $H_\delta$ . Für  $x_1 \in [\delta, 1]$  schneiden wir eine Scheibe in Höhe  $x_1$  aus  $H_\delta$  heraus und erhalten

$$\begin{aligned} H_{\delta, x_1} &= \{(x_2, \dots, x_d) \in \mathbb{R}^{d-1} \mid (x_1, x_2, \dots, x_d) \in H^\delta\} \\ &= \{(x_2, \dots, x_d) \in \mathbb{R}^{d-1} \mid \|(x_1, x_2, \dots, x_d)\| \leq 1\} \\ &= \{y \in \mathbb{R}^{d-1} \mid x_1^2 + \|y\|^2 \leq 1\} \\ &= \bar{B}_{(1-x_1^2)^{1/2}}(0) \end{aligned}$$

wobei die Kugel in der letzten Zeile in  $\mathbb{R}^{d-1}$  genommen wird.



Jetzt wenden wir das Cavalieri'sche Prinzip an und erhalten

$$\begin{aligned} \lambda^d(H_\delta) &= \int_\delta^1 \lambda^{d-1}(H_{\delta, x_1}) d\lambda(x_1) \\ &= \int_\delta^1 \lambda^{d-1}(\bar{B}_{(1-x_1^2)^{1/2}}(0)) d\lambda(x_1) \\ &= \int_\delta^1 (1-x_1^2)^{\frac{d-1}{2}} \lambda^{d-1}(\bar{B}_1(0)) d\lambda(x_1) \\ &= \lambda^{d-1}(B_{\mathbb{R}^{d-1}}) \cdot \int_\delta^1 (1-x_1^2)^{\frac{d-1}{2}} d\lambda(x_1). \end{aligned} \tag{9.1}$$

Da wir  $\lambda^d(W_\varepsilon) = 2 \cdot (\lambda^d(H_0 \setminus H_\varepsilon)) = 2 \cdot (\lambda^d(H_0) - \lambda^d(H_\varepsilon))$  nach unten abschätzen wollen, benötigen wir einerseits eine Abschätzung des Integrals in der letzten Zeile von (9.1) nach unten für  $\delta = 0$  und andererseits eine Abschätzung nach oben für  $\delta = \varepsilon$ .

① Wir beginnen mit  $\delta = 0$ . Durch Verkleinerung der oberen Grenze im Integral der letzten Zeile von (9.1) schätzen wir sicher nach unten ab. Für  $0 \leq x_1 \leq \frac{1}{\sqrt{d-1}}$  ergeben sich die folgenden Umformungen

$$x_1^2 \leq \frac{1}{d-1} \implies 1 - \frac{1}{d-1} \leq 1 - x_1^2 \implies \left(1 - \frac{1}{d-1}\right)^{\frac{d-1}{2}} \leq (1 - x_1^2)^{\frac{d-1}{2}}.$$

Zur Schreiberleichterung benennen wir ab jetzt die Integrationsvariable in  $x$  um. Mit der Änderung der Integrationsgrenze und der letzten Ungleichung oben erhalten wir dann

$$\int_0^1 (1-x^2)^{\frac{d-1}{2}} d\lambda(x) \geq \int_0^{\frac{1}{\sqrt{d-1}}} (1-x^2)^{\frac{d-1}{2}} d\lambda(x) \geq \int_0^{\frac{1}{\sqrt{d-1}}} \left(1 - \frac{1}{d-1}\right)^{\frac{d-1}{2}} d\lambda(x)$$

$$\begin{aligned}
&= \left(1 - \frac{1}{d-1}\right)^{\frac{d-1}{2}} \cdot \frac{1}{\sqrt{d-1}} \underset{\substack{\uparrow \\ \text{Bernoulli-} \\ \text{Ungl.}}}{\geq} \left(1 - \frac{d-1}{2} \cdot \frac{1}{d-1}\right) \cdot \frac{1}{\sqrt{d-1}} \\
&= \frac{1}{2\sqrt{d-1}}.
\end{aligned}$$

Hierbei konnten wir die Bernoulli-Ungleichung anwenden, da  $\frac{d-1}{2} \geq 1$  wegen unserer im Satz gemachten Voraussetzung  $d \geq 3$  gilt.

② Als nächstes schätzen wir das Integral für  $\delta = \varepsilon > 0$  nach oben ab. Aus der elementaren Ungleichung  $e^x \leq \frac{1}{1-x}$  folgt  $1 - x^2 \leq e^{-x^2}$  und damit erhalten wir

$$\begin{aligned}
\int_{\varepsilon}^1 (1 - x^2)^{\frac{d-1}{2}} d\lambda(x) &\leq \int_{\varepsilon}^1 (e^{-x^2})^{\frac{d-1}{2}} \frac{x}{\varepsilon} dx \underset{\substack{\uparrow \\ x/\varepsilon \geq 1}}{\leq} \int_{\varepsilon}^1 (e^{-x^2})^{\frac{d-1}{2}} \frac{x}{\varepsilon} d\lambda(x) \\
&= \frac{1}{\varepsilon} \cdot \int_{\varepsilon}^1 x e^{\frac{d-1}{2} x^2} x^2 dx = \frac{1}{\varepsilon} \cdot \frac{e^{-\frac{d-1}{2}} - e^{-\frac{\varepsilon^2(d-1)}{2}}}{-(d-1)} \\
&\leq \frac{1}{\varepsilon(d-1)} e^{-\frac{\varepsilon^2(d-1)}{2}}.
\end{aligned}$$

Einsetzen der Ergebnisse von ① und ② in (9.1) liefert die zwei Abschätzungen

$$\lambda^d(H_0) \geq \lambda^d(B_{\mathbb{R}^{d-1}}) \frac{1}{2\sqrt{d-1}} \quad \text{und} \quad \lambda^d(H_{\varepsilon}) \leq \lambda^d(B_{\mathbb{R}^{d-1}}) \frac{1}{\varepsilon(d-1)} e^{-\frac{\varepsilon^2(d-1)}{2}}$$

mit denen schließlich

$$\begin{aligned}
\frac{\lambda^d(W_{\varepsilon})}{\lambda^d(B_{\mathbb{R}})} &= \frac{2 \cdot (\lambda^d(H_0) - \lambda^d(H_{\varepsilon}))}{2 \cdot \lambda^d(H_0)} = 1 - \frac{\lambda^d(H_{\varepsilon})}{\lambda^d(H_0)} \\
&\geq 1 - \frac{2\sqrt{d-1}}{\varepsilon(d-1)} e^{-\frac{\varepsilon^2(d-1)}{2}} = 1 - \frac{2}{\varepsilon\sqrt{d-1}} e^{-\frac{\varepsilon^2(d-1)}{2}}
\end{aligned}$$

folgt wie gewünscht.  $\square$

Am Ende von Kapitel 8 hatten wir bereits darauf aufmerksam gemacht, dass im aktuellen Kontext mitunter Vorsicht angesagt ist bei der asymptotischen Interpretation von Abschätzungen. Wir illustrieren dies nun am Satz über die Tailkonzentration. Dessen Name suggeriert bereits, dass in einem *hochdimensionalen* Raum das Volumen der Einheitskugel *nah* an deren Taille konzentriert ist. Die letzten zwei Größenangaben werden im Satz durch (kleines)  $\varepsilon$  und (großes)  $d$  abgebildet. Hierbei ist in folgendem Sinne Vorsicht geboten.

**Bemerkung 9.3.** (i) Bezeichnen wir die Konstante in Satz 9.2 mit

$$K_{\varepsilon,d} = 1 - \frac{2}{\varepsilon\sqrt{d-1}} e^{-\frac{\varepsilon^2(d-1)}{2}} \in (-\infty, 1)$$

so gilt  $K_{\varepsilon,d} \xrightarrow{d \rightarrow \infty} 1$  für jedes feste  $\varepsilon > 0$ . In diesem Sinne ist die oben erläuterte Interpretation des Satzes zu verstehen.

(ii) Fixiert man hingegen die Dimension  $d$ , so gilt  $K_{\varepsilon,d} \xrightarrow{\varepsilon \rightarrow 0} -\infty$ . Auf diese Weise

darf man den Satz folglich *nicht* interpretieren.

(iii) Betrachtet man  $\varepsilon > 0$  und  $d \in \mathbb{N}$  beide als fixiert, so ist die Abschätzung  $\lambda^d(W_\varepsilon) \geq K_{\varepsilon,d} \lambda^d(B_{\mathbb{R}^d})$  aus Satz 9.2 nur dann nicht-trivial, wenn  $K_{\varepsilon,d} > 0$  ausfällt. Per Substitution  $a := \varepsilon\sqrt{d-1}$  sieht man, dass es ein  $a_0 \in (1, 2)$  gibt, sodass

$$\forall \varepsilon > 0, d \in \mathbb{N}_{\geq 2}: \varepsilon > \frac{a_0}{\sqrt{d-1}} \implies K_{\varepsilon,d} > 0$$

gilt, vergleiche Aufgabe 9.3. Man kann also durchaus in Satz 9.2 beliebig kleine  $\varepsilon$  betrachten — wenn man dies durch entsprechend große  $d$  kompensiert.

(iv) Auch im Satz 9.1 zur Oberflächenkonzentration hängt die Konstante  $1 - e^{-\varepsilon d}$  von  $\varepsilon$  und  $d$  ab, ist aber stets positiv. Dennoch ist der Satz natürlich nur dann interessant, wenn  $1 - e^{-\varepsilon d}$  nah bei 1 liegt, und man sieht, dass auch hier ein Tradeoff zwischen kleinem  $\varepsilon$  und großem  $d$  möglich ist.

Die obigen Bemerkungen erklären jetzt sehr genau, warum man vom *Segen der hohen Dimension* spricht: Ist die Dimension  $d$  groß genug, so kann oben der Parameter  $\varepsilon > 0$  derart klein gewählt werden, dass sich eine nützliche Aussage ergibt. Ist  $d$  hingegen zu klein, so sind die Aussagen entweder trivial oder nicht von Nutzen, vergleiche auch Aufgabe 9.4.

Bevor wir zur probabilistischen Interpretation der Maßkonzentrationssätze kommen, notieren die folgende bloße Umformulierung von Satz 9.2.

**Korollar 9.4.** *Sei  $d \in \mathbb{N}_{\geq 3}$ . Dann gilt*

$$\frac{\lambda^d(\{(x_1, \dots, x_d) \in B_{\mathbb{R}^d} \mid |x_1| > \varepsilon\})}{\lambda^d(B_{\mathbb{R}^d})} \leq \frac{2}{\varepsilon\sqrt{d-1}} e^{-\frac{\varepsilon^2(d-1)}{2}}$$

für jedes  $\varepsilon > 0$ . □

Im Rest des Kapitels werden wir die zwei vorhergehenden Sätze und das obige Korollar nutzen, um Aussagen über gleichmäßig auf der Einheitskugel verteilte Zufallsvektoren, d.h.  $X: \Omega \rightarrow \mathbb{R}^d$  mit  $X \sim \mathcal{U}(B_{\mathbb{R}^d})$ , zu gewinnen, vergleiche Definition 8.1. Für eine solche Zufallsvariable erhält man sofort

$$\mathbb{P}[\|X\| \geq 1 - \varepsilon] = \frac{\lambda^d(\bar{B}_1(0) \setminus B_{1-\varepsilon}(0))}{\lambda^d(\bar{B}_1(0))} \geq 1 - e^{-\varepsilon d} \quad (9.2)$$

und kann interpretieren, dass ein aus  $B_{\mathbb{R}^d}$  zufällig gewählter Punkt mit hoher Wahrscheinlichkeit nah an dessen Oberfläche liegt. Für den folgenden Satz wählen wir nun nicht nur einen Punkt zufällig aus, sondern mehrere auf einmal.

**Satz 9.5.** *Sei  $(\Omega, \Sigma, \mathbb{P})$  ein Wahrscheinlichkeitsraum, sei  $d \in \mathbb{N}_{\geq 3}$  und für  $n \in \mathbb{N}_{\geq 2}$  seien  $X^{(1)}, \dots, X^{(n)}: \Omega \rightarrow \mathbb{R}^d$  unabhängig mit  $X^{(i)} \sim \mathcal{U}(B_{\mathbb{R}^d})$ . Dann gilt*

$$\mathbb{P}\left[\|X^{(i)}\| \geq 1 - \frac{2 \log n}{d} \text{ für alle } i = 1, \dots, n\right] \geq 1 - \frac{1}{n}.$$

*Beweis.* Wir setzen  $\varepsilon = \frac{2 \log n}{d}$  und fixieren  $1 \leq i \leq n$ . Aus dem Satz 9.1 über die

Oberflächenkonzentration folgt dann

$$\mathbb{P}[\|X^{(i)}\| < 1 - \varepsilon] = \frac{\lambda^d(B_{1-\varepsilon}(0))}{\lambda^d(B_1(0))} \leq e^{-\varepsilon d} = e^{-2 \log n} = \frac{1}{n^2}.$$

Daraus ergibt sich

$$\begin{aligned} \mathbb{P}[\forall i: \|X^{(i)}\| \geq 1 - \varepsilon] &= 1 - \mathbb{P}[\exists i: \|X^{(i)}\| < 1 - \varepsilon] \\ &= 1 - \sum_{i=1}^n 1/n^2 = 1 - \frac{1}{n} \end{aligned}$$

wie behauptet.  $\square$

**Bemerkung 9.6.** Satz 9.5 kann so interpretiert werden, als dass  $n$ -viele aus der Einheitskugel des  $\mathbb{R}^d$  gleichmäßig zufällig ausgewählte Punkte mit hoher Wahrscheinlichkeit nah an deren Oberfläche liegen werden. Dabei heben wir das folgende nochmal besonders hervor:

- (i) Die Normabschätzung gilt für alle Punkte — und nicht nur für alle bis auf eine in irgendeinem Sinne kleine Ausnahmemenge.
- (ii) Die untere Schranke in der Normabschätzung wächst in der Dimension  $d$  und fällt in der Anzahl der Punkte  $n$  und ist unabhängig von  $d$ .
- (iii) Die untere Schranke in der Wahrscheinlichkeitsabschätzung wächst in der Anzahl der Punkte  $n$ .
- (iv) Dadurch, dass  $n$  in die Normabschätzung mit  $\log n$  und in die Wahrscheinlichkeitsabschätzung mit  $1/n$  eingeht, führt eine Erhöhung des Stichprobenumfangs zu einer viel stärkeren Verbesserung der Wahrscheinlichkeitsabschätzung, als dadurch die Normabschätzung verschlechtert wird.

Für großes  $d$  und moderat großes  $n$  erhält man also  $1 - \frac{2 \log n}{d} \approx 1$  und auch  $1 - \frac{1}{n} \approx 1$ . In diesem Sinne ist die am Beginn dieser Bemerkung gegebene Interpretation zu verstehen.

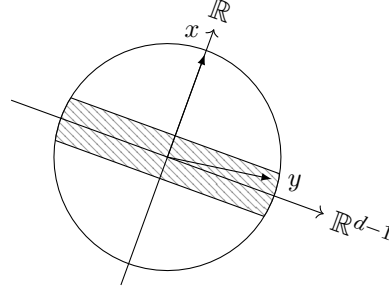
Als nächstes kommen wir zum Winkel zwischen zufällig aus  $B_{\mathbb{R}^d}$  gewählten Punkten. Da nach Satz 9.5 deren Normen mit hoher Wahrscheinlichkeit nah bei Eins liegen, verzichten wir auf die Normierung und betrachten unten nur die paarweisen Skalarprodukte. Vergleiche aber Aufgabe 9.3.

**Satz 9.7.** Sei  $(\Omega, \Sigma, \mathbb{P})$  ein Wahrscheinlichkeitsraum, sei  $d \in \mathbb{N}_{\geq 3}$  und für  $n \in \mathbb{N}_{\geq 2}$  seien  $X^{(1)}, \dots, X^{(n)}: \Omega \rightarrow \mathbb{R}^d$  unabhängig mit  $X^{(i)} \sim \mathcal{U}(B_{\mathbb{R}^d})$ . Dann gilt

$$\mathbb{P}[\langle X^{(j)}, X^{(k)} \rangle \leq \frac{\sqrt{6 \log n}}{\sqrt{d-1}} \text{ für alle } j \neq k] \geq 1 - \frac{1}{n}.$$

*Beweis.* Wir erläutern zuerst die Beweisidee anschaulich und beschränken uns auf den Fall, dass zwei Punkte  $x, y \in B_{\mathbb{R}^d}$  gleichmäßig zufällig und unabhängig gewählt werden. Nach dem Satz von der Oberflächenkonzentration gilt  $\|x\| \approx 1$  und  $\|y\| \approx 1$ . Die Punkte  $x$  und  $y$  liegen also knapp unter der Oberfläche der Einheitskugel. Ändern

wir das Koordinatensystem derart, dass  $x$  ‘in Nordrichtung’ zeigt, so muss dann der unabhängig von  $x$  gewählte Punkt  $y$  nach dem Satz von der Tailenkonzentration nah am Äquator liegen:



Folglich wird der Winkel der zwei Punkte ungefähr  $90^\circ$  sein, und daher das Skalarprodukt  $\langle x, y \rangle \approx 0$  ausfallen. Wir werden jetzt zuerst diese Heuristik für zwei Punkte formalisieren und dann mithilfe der üblichen Rechenregeln für Wahrscheinlichkeiten die im Satz behauptete Abschätzung beweisen.

Wir setzen  $\varepsilon = \frac{\sqrt{6 \log n}}{\sqrt{d-1}}$  und betrachten zunächst zwei Zufallsvariablen

$$X, Y: \Omega \rightarrow \mathbb{R}^d \text{ mit } X, Y \sim \mathcal{U}(\mathbb{B}_{\mathbb{R}^d}).$$

Sei  $\omega \in \Omega$  derart dass  $x := X(\omega) \neq 0$ . Wir setzen  $b_1 := \frac{x}{\|x\|}$  und ergänzen zu einer Orthonormalbasis  $\{b_1, \dots, b_d\}$  von  $\mathbb{R}^d$ . Weiter bezeichnen wir mit  $\{e_1, \dots, e_d\}$  die Standardbasis von  $\mathbb{R}^d$ . Dann gibt es genau eine lineare Abbildung

$$T: \mathbb{R}^d \rightarrow \mathbb{R}^d \text{ mit } Tb_i = e_i \text{ für } i = 1, \dots, d$$

und mit der Abkürzung  $y := Y(\omega)$  folgt

$$\begin{aligned} |\langle X(\omega), Y(\omega) \rangle| &= |\langle x, y \rangle| = |\langle Tx, Ty \rangle| = |\langle \|x\| b_1, Ty \rangle| \\ &\quad \begin{array}{c} \uparrow \\ T \text{ orthog.} \end{array} \quad \begin{array}{c} \uparrow \\ x = \|x\| b_1 \end{array} \\ &= \|x\| \cdot |\langle e_1, Ty \rangle| = \|X(\omega)\| \cdot |\langle e_1, TX(\omega) \rangle|. \end{aligned}$$

Da die Einheitskugel invariant unter orthogonalen Abbildungen ist, folgt

$$\begin{aligned} \mathbb{P}[|\langle e_1, TY \rangle| > \varepsilon] &= \frac{\lambda^d(\{y \in \mathbb{B}_{\mathbb{R}^d} \mid |\langle e_1, Ty \rangle| > \varepsilon\})}{\lambda^d(\mathbb{B}_{\mathbb{R}^d})} = \frac{\lambda^d(\{T^{-1}z \in \mathbb{B}_{\mathbb{R}^d} \mid |\langle e_1, z \rangle| > \varepsilon\})}{\lambda^d(\mathbb{B}_{\mathbb{R}^d})} \\ &= \frac{\lambda^d(\{(z_1, \dots, z_d) \in \mathbb{B}_{\mathbb{R}^d} \mid |z_1| > \varepsilon\})}{\lambda^d(\mathbb{B}_{\mathbb{R}^d})} \leq \frac{2}{\varepsilon \sqrt{d-1}} e^{-\frac{\varepsilon^2(d-1)}{2}} \\ &\quad \begin{array}{c} \uparrow \\ \text{Korollar} \\ 9.4 \end{array} \\ &= \frac{2}{\sqrt{6 \log n}} e^{-\frac{6 \log n}{2}} \leq \frac{1}{n^3} \\ &\quad \begin{array}{c} \uparrow \\ \frac{2}{\sqrt{6 \log n}} \leq 1 \end{array} \end{aligned}$$

wobei wir insbesondere beachten, dass die Abbildung  $T$  von unserem vorgewählten  $\omega \in \Omega$  mit  $X(\omega) \neq 0$  abhängt, aber die oben gefunden Schranke von diesem unab-

hängig ist. Mit letzterem im Hinterkopf berechnen wir

$$\begin{aligned}
 \mathbb{P}[|\langle X, Y \rangle| > \varepsilon] &= \mathbb{P}(\{\omega \in \Omega \mid |\langle X(\omega), Y(\omega) \rangle| > \varepsilon\}) \\
 &\stackrel{\substack{\uparrow \\ \text{d.h. } X(\omega) \neq 0 \\ \text{automatisch}}}{=} \mathbb{P}(\{\omega \in \Omega \mid \|X(\omega)\| \cdot |\langle e_1, TY(\omega) \rangle| > \varepsilon\}) \\
 &\leq \mathbb{P}(\{\omega \in \Omega \mid |\langle e_1, TY(\omega) \rangle| > \varepsilon\}) \stackrel{\substack{\uparrow \\ X(\omega) \in B_{\mathbb{R}^d}}}{\leq} \frac{1}{n^3} \stackrel{\substack{\uparrow \\ \text{s.o.}}}{\leq} \frac{1}{n^3}.
 \end{aligned}$$

Seien nun schließlich  $X^{(1)}, \dots, X^{(n)} \sim \mathcal{U}(B_{\mathbb{R}^d})$  gegeben. Dann gilt

$$\begin{aligned}
 \mathbb{P}[\forall j \neq k: |\langle X^{(j)}, X^{(k)} \rangle| \leq \varepsilon] &= 1 - \mathbb{P}[\exists j \neq k: |\langle X^{(j)}, X^{(k)} \rangle| > \varepsilon] \\
 &= 1 - \binom{n}{2} \cdot \mathbb{P}[|\langle X, Y \rangle| > \varepsilon] \\
 &\geq 1 - \frac{n^2 - n}{n^3} \\
 &\geq 1 - \frac{1}{n}
 \end{aligned}$$

wie behauptet.  $\square$

Für  $n = 2$  liefert der Satz nur eine Abschätzung der Wahrscheinlichkeit nach unten durch  $1/2$  dafür dass der Betrag des Skalarproduktes zweier Zufallsvektoren kleiner gleich  $\frac{\sqrt{6 \log 2}}{\sqrt{d-1}}$  ausfällt. Der Beweis zeigt in der Tat aber

$$\mathbb{P}[|\langle X, Y \rangle| \leq \varepsilon] \geq 1 - \frac{2}{\varepsilon \sqrt{d-1}} e^{-\frac{\varepsilon^2(d-1)}{2}} \quad (9.3)$$

für  $X, Y \sim \mathcal{U}(B_{\mathbb{R}^d})$ .

**Bemerkung 9.8.** Zusammengenommen führen die Sätze 9.5 und 9.7 zu der Interpretation, dass  $n$ -viele gleichmäßig zufällig aus der  $d$ -dimensionalen Einheitskugel gewählte Punkte mit hoher Wahrscheinlichkeit paarweise orthogonal sind: Satz 9.5 liefert  $\|x^{(j)}\| \approx 1$  und  $\|x^{(k)}\| \approx 1$  woraus mit Satz 9.7 dann  $\langle \frac{x^{(j)}}{\|x^{(j)}\|}, \frac{x^{(k)}}{\|x^{(k)}\|} \rangle \approx 0$  folgt, vergleiche Aufgabe 9.3 für eine formale Behandlung. Zu der Frage, unter welcher Konfiguration von  $n$  und  $d$ , bzw. im Fall von nur zwei Punkten,  $d$  und  $\varepsilon$ , diese Interpretation gültig ist, ermuntern wir den Leser dazu, Überlegungen analog zu Bemerkung 9.6 anzustellen und selbständig zu verschriftlichen. Insbesondere empfiehlt es sich hier auch konkrete Zahlbeispiele auszuprobieren, siehe Aufgabe 9.4.

Jetzt betrachten wir den Abstand von gleichmäßig zufällig aus  $B_{\mathbb{R}^d}$  gewählten Punkten. Wir beschränken uns hier auf zwei Punkte und überlassen es dem Leser, Resultate für  $n$ -viele Punkte ähnlich den vorhergehenden Sätzen abzuleiten.

**Satz 9.9.** Sei  $(\Omega, \Sigma, \mathbb{P})$  ein Wahrscheinlichkeitsraum,  $X, Y: \Omega \rightarrow \mathbb{R}^d$  unabhängige Zufallsvektoren mit  $X, Y \sim \mathcal{U}(B_{\mathbb{R}^d})$ . Dann gilt

$$\mathbb{P}[|\|X - Y\| - \sqrt{2}| \leq \varepsilon] \geq 1 - 2e^{-\varepsilon d/5} - \frac{9}{\varepsilon \sqrt{d-1}} e^{-\frac{\varepsilon^2(d-1)}{36}}$$

für  $0 < \varepsilon < 1$  und  $d \in \mathbb{N}_{\geq 3}$ .

*Beweis.* Wir drehen die Abstandsabschätzung zunächst um und multiplizieren dann mit der Ungleichung  $\|X - Y\| + \sqrt{2} \geq \sqrt{2}$ . Auf diese Weise erhalten wir

$$\begin{aligned} \mathbb{P}[|\|X - Y\| - \sqrt{2}| \leq \varepsilon] &= 1 - \mathbb{P}[|\|X - Y\| - \sqrt{2}| \geq \varepsilon] \\ &\geq 1 - \mathbb{P}[|\|X - Y\| - \sqrt{2}|(\|X - Y\| + \sqrt{2}) \geq \varepsilon\sqrt{2}] \\ &= 1 - \mathbb{P}[|\|X - Y\|^2 - 2| \geq \varepsilon\sqrt{2}] \\ &= \mathbb{P}[-\varepsilon\sqrt{2} \leq \|X\|^2 + \|Y\|^2 - 2\langle X, Y \rangle - 2 \leq \varepsilon\sqrt{2}] \\ &\geq \mathbb{P}[|\|X\|^2 - 1| \leq \frac{\varepsilon\sqrt{2}}{3}]^2 \cdot \mathbb{P}[|\langle X, Y \rangle| \leq \frac{\varepsilon\sqrt{2}}{6}]. \end{aligned}$$

Dann schätzen wir den ersten Term mithilfe der Bernoulli-Ungleichung und (9.2) ab, was auf das folgende führt

$$\begin{aligned} \mathbb{P}[|\|X\|^2 - 1| \leq \frac{\varepsilon\sqrt{2}}{3}] &= \mathbb{P}[1 - \|X\|^2 \leq \frac{\varepsilon\sqrt{2}}{3}] = \mathbb{P}[\|X\| \geq (1 - \frac{\varepsilon\sqrt{2}}{3})^{1/2}] \\ &\geq \mathbb{P}[\|X\| \geq 1 - \frac{1}{2} \cdot \frac{\varepsilon\sqrt{2}}{3}] \geq 1 - e^{-\frac{\varepsilon\sqrt{2}}{6}d}. \end{aligned}$$

Den zweiten Term erledigen wir mit (9.3) wie folgt

$$\mathbb{P}[|\langle X, Y \rangle| \leq \frac{\varepsilon\sqrt{2}}{6}] \geq 1 - \frac{6\sqrt{2}}{\varepsilon\sqrt{d-1}} e^{-\frac{\varepsilon^2(d-1)}{36}}.$$

Ausmultiplizieren liefert

$$\begin{aligned} \mathbb{P}[|\|X - Y\| - \sqrt{2}| \leq \varepsilon] &\geq (1 - e^{-\frac{\varepsilon\sqrt{2}}{6}d})^2 (1 - \frac{6\sqrt{2}}{\varepsilon\sqrt{d-1}} e^{-\frac{\varepsilon^2(d-1)}{36}}) \\ &= (1 - 2e^{-\frac{\varepsilon\sqrt{2}}{6}d} + e^{-\frac{\varepsilon\sqrt{2}}{3}d}) (1 - \frac{6\sqrt{2}}{\varepsilon\sqrt{d-1}} e^{-\frac{\varepsilon^2(d-1)}{36}}) \\ &\geq 1 - 2e^{-\frac{\varepsilon\sqrt{2}}{6}d} - \frac{6\sqrt{2}}{\varepsilon\sqrt{d-1}} e^{-\frac{\varepsilon^2(d-1)}{36}} \end{aligned}$$

und damit das Gewünschte, wenn man  $\frac{\sqrt{2}}{6} \geq \frac{1}{5}$  und  $6\sqrt{2} \leq 9$  benutzt.  $\square$

**Bemerkung 9.10.** Die Resultate oben suggerieren, dass für  $X, Y \sim \mathcal{U}(B_{\mathbb{R}^d})$

$$\mathbb{E}(\|X\|) \approx 1, \quad \mathbb{E}(\langle X, Y \rangle) \approx 0 \quad \text{und} \quad \mathbb{E}(\|X - Y\|) \approx \sqrt{2}$$

gelten, in Analogie zu den Sätzen 8.2 und 8.4 für normalverteilte Zufallsvektoren. Beachte allerdings, dass wir einerseits obiges nicht bewiesen haben, und andererseits, dass unsere Aussagen ‘mit hoher Wahrscheinlichkeit’ quantitativ, und daher in gewissem Sinne sogar besser sind, vergleiche auch die Diskussion am Ende von Kapitel 8.



## Referenzen

Dieses Kapitel basiert hauptsächlich auf [BHK20].

## Aufgaben

**Aufgabe 9.1.** Sei  $\lambda^d$  das Lebesguemaß und  $B_{\mathbb{R}^d}$  die  $d$ -dimensionale Einheitskugel.

- (i) Zeigen Sie, dass  $\lambda^d(B_{\mathbb{R}^d}) = \frac{2\pi}{d} \cdot \lambda^{d-2}(B_{\mathbb{R}^{d-2}})$  für  $d \geq 3$  gilt.
- (ii) Benutzen Sie die Rekursionformel aus (i) und die bekannten Fälle  $d = 1, 2$  um  $\lambda^d(B_{\mathbb{R}^d})$  für  $d = 3, \dots, 10$  zu berechnen und die Funktion  $d \mapsto \lambda^d(B_{\mathbb{R}^d})$  zu skizzieren.
- (iii) Berechnen Sie  $\lim_{d \rightarrow \infty} \lambda^d(B_{\mathbb{R}^d})$ .
- (iv) Zeigen Sie, dass  $\lambda^d(B_{\mathbb{R}^d}) = \frac{\pi^{d/2}}{\Gamma(d/2+1)}$  gilt, wobei  $\Gamma: (0, \infty) \rightarrow \mathbb{R}$  die Gammafunktion bezeichnet.

*Hinweis:* Überlegen Sie sich für (i) zunächst, dass

$$\lambda^d(B_{\mathbb{R}^d}) = \int_{x_1^2 + x_2^2 \leq 1} \left( \int_{x_3^2 + \dots + x_d^2 \leq 1 - x_1^2 - x_2^2} 1 \, d\lambda(x_3, \dots, x_d) \right) d\lambda(x_1, x_2)$$

nach dem Satz von Fubini gilt. Schreiben Sie als nächstes das innere Integral als Vielfaches von  $\lambda^{d-2}(B_{\mathbb{R}^{d-2}})$  um. Es verbleibt dann ein 2-dimensionales Integral über den Vorfaktor, welches durch Transformation in Polarkoordinaten gelöst werden kann.

**Aufgabe 9.2.** Zeigen Sie, dass es ein  $a_0 \in (1, 2)$  gibt, sodass gilt

$$\forall \varepsilon > 0, d \in \mathbb{N}_{\geq 2}: \varepsilon > \frac{a_0}{\sqrt{d-1}} \implies 1 - \frac{2}{\varepsilon\sqrt{d-1}} e^{-\frac{\varepsilon^2(d-1)}{2}} > 0,$$

vergleiche Bemerkung 9.3(iii). In der Tat gilt sogar  $a_0 < \sqrt{5} - 1$ .

*Hinweis:* Betrachten Sie  $f: (0, \infty) \rightarrow \mathbb{R}$ ,  $f(a) = \frac{2}{a} e^{a^2/2}$ , und zeigen Sie erstmal, dass es ein  $a_0 > 0$  gibt, sodass  $f(a) < 1$  gilt für  $a > a_0$ . Folgern Sie dann, dass  $a_0 \in (1, 2)$  gelten muss. Substituieren Sie schließlich  $a = \varepsilon\sqrt{d-1}$ .

**Aufgabe 9.3.** Sei  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum, sei  $d \in \mathbb{N}_{\geq 3}$  und für  $n \in \mathbb{N}$ , derart dass  $2 \log(n) \leq d$ , seien  $X^{(1)}, \dots, X^{(n)}: \Omega \rightarrow \mathbb{R}^d$  unabhängig mit  $X^{(j)} \sim \mathcal{U}(B_{\mathbb{R}^d})$ . Zeigen Sie, dass gilt:

$$P\left[\left|\left\langle \frac{X^{(j)}}{\|X^{(j)}\|}, \frac{X^{(k)}}{\|X^{(k)}\|} \right\rangle\right| \leq \frac{\sqrt{6 \log n}}{\sqrt{d-1}} \text{ for all } j \neq k\right] \geq 1 - \frac{1}{n}.$$

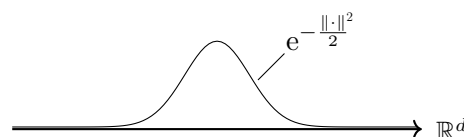
*Hinweis:* Benutzen Sie  $1 \leq \frac{d^2}{(d-2 \log n)^2}$  und den Satz von der totalen Wahrscheinlichkeit.

**Aufgabe 9.4.** Gehen Sie ein paar Zahlbeispiele durch, um ein Gefühl dafür zu bekommen, für welche  $n$ ,  $d$  und  $\varepsilon$  die Aussagen in diesem Kapitel tatsächlich die angegebene Interpretationen zulassen, also z.B. Satz 9.5 wirklich so gelesen werden kann, als dass  $\|X\| \approx 1$  mit einer Wahrscheinlichkeit nah bei 1 gilt, usw.

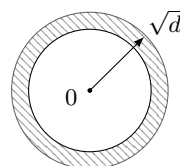
## Kapitel 10

# Gaußsche Zufallsvektoren in hohen Dimensionen

Wir wenden uns nun gaußverteilten Zufallsvariablen zu. Der Einfachheit halber betrachten wir die meiste Zeit solche mit Mittelwert Null und Varianz Eins, also  $X \sim \mathcal{N}(0, 1, \mathbb{R}^d)$ . In Satz 8.2 hatten wir bereits gezeigt, dass dann  $|\mathbb{E}(X) - \sqrt{d}| \leq \frac{1}{\sqrt{d}}$  und  $V(\|X\|) \leq 2$  gelten. Experimente, vgl. Seite 114, hatten ebenfalls nahegelegt, dass sich fast alle Stichproben eines normalverteilten Zufallsvektors auf einem relativ schmalen ringförmigen Gebiet mit Radius  $\sqrt{d}$  wiederfinden. Auf den ersten Blick mag dies überraschen, ist man doch aus niedrigen Dimensionen ein anderes Bild gewöhnt und sieht, dass auch in hoher Dimension die Dichtefunktion um Null herum die größten Werte annimmt und nach außen hin schnell abfällt.



Stellt man sich aber den  $\mathbb{R}^d$  aufgeteilt in Ringe  $R$  gleicher Dicke mit wachsenden Radien  $r$  vor, und erinnert sich an Kapitel 9, so sieht man, dass das Maß der Ringe von innen nach außen wächst. Integriert man nun die Dichte über einen solchen Ring, so ist bei kleinem Radius zwar der Integrand groß, aber das Maß des Rings ist klein und daher das Integral klein. Bei großen Radien ist es gerade umgekehrt. Bei mittleren Radien halten sich beide Effekte die Waage, sodass sich hier — wie wir zeigen werden in einem Ring um den Radius  $r \approx \sqrt{d}$  herum — die Wahrscheinlichkeit konzentriert.



Unser erstes Ziel in diesem Kapitel ist es, das Obige zu quantifizieren und zu beweisen. Dies führt auf den *Gaußschen Ringsatz* 10.4. Wir benötigen hierfür einige

Vorbereitungen und insbesondere zwei wichtige Eigenschaften gaußscher Zufallsvektoren. Diese notieren wir unten als Fakt; ausführliche Beweise finden sich in A.2.

**Fakt 10.1.** *Sei  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum.*

- (i) *Für einen Zufallsvektor  $X: \Omega \rightarrow \mathbb{R}^d$  gilt  $X \sim \mathcal{N}(0, 1, \mathbb{R}^d)$  genau dann wenn seine Koordinaten  $X_1, \dots, X_d$  unabhängig sind und  $X_i \sim \mathcal{N}(0, 1)$  für alle  $i$  gilt.*
- (ii) *Sind  $X_1, \dots, X_n: \Omega \rightarrow \mathbb{R}$  unabhängige Zufallsvariablen mit  $X_i \sim \mathcal{N}(0, 1)$  für alle  $i$  und  $\lambda_1, \dots, \lambda_n \in \mathbb{R}$ , so ist  $\lambda_1 X_1 + \dots + \lambda_n X_n \sim \mathcal{N}(0, \lambda_1^2 + \dots + \lambda_n^2)$ .*

Wir starten nun mit dem ersten von zwei vorbereitenden Resultaten, welches auf der sogenannten *Chernoff-Methode* basiert.

**Lemma 10.2.** (Bernstein Tailbound) *Sei  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum und seien  $Y_1, \dots, Y_d: \Omega \rightarrow \mathbb{R}$  unabhängige Zufallsvariablen mit  $E(Y_i) = 0$  und  $|E(Y_i^k)| \leq k!/2$  für  $i = 1, \dots, d$  und  $k \geq 2$ . Dann gilt für  $a > 0$*

$$P[|Y_1 + \dots + Y_d| \geq a] \leq 2e^{-\frac{1}{4} \min(\frac{a^2}{d}, a)}.$$

*Beweis.* Wir setzen  $Y = Y_1 + \dots + Y_d$  und aus Gründen die wir gleich sehen werden, betrachten wir für  $0 < t \leq 1/2$ :

$$\begin{aligned} P[Y \geq a] & \underset{t > 0}{=} P[e^{tY} \geq e^{ta}] \underset{\text{Markov}}{\leq} \frac{E(e^{tY})}{e^{ta}} \\ & = e^{-ta} E(e^{t(Y_1 + \dots + Y_d)}) \underset{Y_i \text{ unabh.}}{=} e^{-ta} \prod_{i=1}^d E(e^{tY_i}). \end{aligned} \tag{10.1}$$

Als nächstes schätzen wir die Faktoren von oben einzeln wie folgt ab:

$$\begin{aligned} E(e^{tY_i}) & = |E(1 + tY + \sum_{k=2}^{\infty} \frac{(tY_i)^k}{k!})| \leq 1 + \sum_{k=2}^{\infty} \frac{t^k |E(Y_i^k)|}{k!} \underset{|E(Y_i^k)| \leq k!/2}{\leq} 1 + \sum_{k=2}^{\infty} \frac{t^k}{2} \\ & = 1 + \frac{1}{2} \left( \sum_{k=0}^{\infty} t^k - t - 1 \right) \underset{\substack{t < 1 \text{ und} \\ \text{geom. R.}}}{=} 1 + \frac{1}{2} \left( \frac{1}{1-t} - \frac{(t+1)(t-1)}{t-1} \right) \\ & = 1 + \frac{1}{2} \frac{t^2}{1-t} \underset{\substack{\uparrow \\ 1-t \leq 2}}{\leq} 1 + t^2 \leq e^{t^2}. \end{aligned}$$

Da obiges für alle  $0 < t \leq 1/2$  gilt und in (10.1) die linke Seite nicht von  $t$  abhängt, minimieren wir dort jetzt die rechte Seite über diese  $t$  und erhalten

$$P[Y \geq a] \leq \inf_{0 < t \leq 1/2} e^{-ta} \prod_{i=1}^d e^{t^2} \leq \inf_{0 < t \leq 1/2} e^{-ta + dt^2} \tag{10.2}$$

wobei sich gleich zeigen wird, dass das Infimum in Wirklichkeit ein Minimum ist. Wir betrachten die Funktion  $f: \mathbb{R} \rightarrow \mathbb{R}$ ,  $f(t) = e^{-ta + dt^2}$ . Ableiten und Nullsetzen

liefert

$$\frac{d}{dt}f(t) = (-a + 2dt)e^{-ta+dt^2} \stackrel{!}{=} 0 \implies t = \frac{a}{2d}$$

und da  $f(t) \rightarrow \infty$  für  $t \rightarrow \pm\infty$  gilt, ist  $t_0 := \frac{a}{2d} > 0$  der einzige Minimierer von  $f$ . Ist nun  $t_0 \leq 1/2$ , so wird das Infimum von  $f|_{(0,1/2]}$  an dieser Stelle angenommen. Ist  $t_0 > 1/2$ , so ist  $f|_{(0,1/2]}$  monoton fallend und das Infimum wird in  $t_1 := 1/2$  angenommen. Insbesondere haben wir in letzterem Fall  $a > d$  und es ergibt sich für die jeweilige Minimalstelle

$$-t_0a + dt_0^2 = -\frac{a^2}{2d} + \frac{da^2}{4d^2} = -\frac{a^2}{4d} \quad \text{bzw.} \quad -t_1a + dt_1^2 = -\frac{a}{2} + \frac{d}{4} \underset{d < a}{<} -\frac{a}{4}.$$

Damit können wir die Abschätzung in (10.2) fortsetzen zu

$$\inf_{0 < t \leq 1/2} e^{-ta+dt^2} = e^{-ta+dt^2} \Big|_{t=\min(\frac{a}{2d}, \frac{1}{2})} \leq e^{-\min(\frac{a^2}{4d}, \frac{a}{4})} \quad (10.3)$$

und erhalten

$$P[Y_1 + \dots + Y_d \geq a] \leq e^{-\frac{1}{4} \min(\frac{a^2}{d}, a)}.$$

Wiederholen wir das Bisherige, ersetzen aber überall  $Y_i$  mit  $-Y_i$ , so führen die gleichen Rechnungen auf

$$P[Y_1 + \dots + Y_d \leq -a] = P[-Y_1 - \dots - Y_d \geq a] \leq e^{-\frac{1}{4} \min(\frac{a^2}{d}, a)}$$

und beides zusammen beendet den Beweis.  $\square$

Jetzt kommen wir zum zweiten vorbereitenden Resultat, welches wir gleich benötigen um das  $sk$ -te Moment einer gaußverteilten Zufallsvariable zu berechnen.

**Lemma 10.3.** Für  $k \in \mathbb{N}_{\geq 1}$  gilt  $\frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} t^{2k} e^{-t^2/2} dt = \frac{(2k)!}{2^k k!}.$

*Beweis.* ① Mithilfe des Lebesgue'schen Konvergenzsatzes können wir in der folgenden Rechnung Ableitung und Integral vertauschen

$$\frac{d^k}{da^k} \left( \int_{\mathbb{R}} e^{-at^2} dt \right) \Big|_{a=\frac{1}{2}} = \left( \int_{\mathbb{R}} \frac{d^k}{da^k} e^{-at^2} dt \right) \Big|_{a=\frac{1}{2}} = (-1)^k \int_{\mathbb{R}} t^{2k} e^{-t^2/2} dt.$$

② Andererseits erhalten wir durch Substitution  $u := \sqrt{at}$  für  $a > 0$

$$\int_{\mathbb{R}} e^{-at^2} dt = \frac{1}{\sqrt{a}} \int_{\mathbb{R}} e^{-u^2} du = \sqrt{\frac{\pi}{a}}.$$

③ Lesen wir nun zuerst ① rückwärts und setzen dann ② ein, so erhalten wir

$$\int_{\mathbb{R}} t^{2k} e^{-t^2/2} dt = (-1)^k \frac{d^k}{da^k} \left( \int_{\mathbb{R}} e^{-at^2} dt \right) \Big|_{a=\frac{1}{2}} = (-1)^k \frac{d^k}{da^k} \left( \sqrt{\frac{\pi}{a}} \right) \Big|_{a=\frac{1}{2}}$$

und müssen als nächstes die Ableitung auf der rechten Seite ausrechnen

$$\begin{aligned}
 \frac{d^k}{da^k}(a^{-1/2}) &= (-1/2) \cdot (-1/2 - 1) \cdots (-1/2 - (k-1)) \cdot a^{-1/2-k} \\
 &= \frac{(-1)^k}{2^k} \cdot (1 \cdot 3 \cdots (2k-3) \cdot (2k-1)) \cdot a^{-1/2-k} \\
 &= \frac{(-1)^k}{2^k} \cdot \frac{1 \cdot 2 \cdot 3 \cdots (2k-3) \cdot (2k-2) \cdot (2k-1) \cdot (2k)}{2 \cdot 4 \cdots (2k-2) \cdot (2k)} \cdot a^{-1/2-k} \\
 &= \frac{(-1)^k}{2^k} \cdot \frac{(2k)!}{2^k k!} \cdot a^{-1/2-k}.
 \end{aligned}$$

④ Einsetzen der Ableitung in ③ und Division durch  $\sqrt{2\pi}$  führt auf

$$\frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} t^{2k} e^{-t^2/2} dt = \frac{(-1)^k}{\sqrt{2}} \frac{(-1)^k}{2^k} \frac{(2k)!}{2^k k!} a^{-1/2-k} \Big|_{a=\frac{1}{2}} = \frac{1}{\sqrt{2}} \frac{1}{2^k} \frac{(2k)!}{2^k k!} 2^{1/2+k}$$

was sich nach Kürzen genau zum im Lemma angegebenen Ausdruck vereinfacht.  $\square$

Jetzt sind wir bereit dafür den Gaußschen Ringsatz zu beweisen. Bevor wir beginnen, machen wir darauf aufmerksam, dass in der folgenden Ungleichung die linke Seite gleich Eins ist für  $\varepsilon = 4 \log 2 \approx 2.8$ . Interessant wird die Ungleichung also erst für entsprechend größere  $\varepsilon$ , vergleiche Beispiel 10.5. Es ist hier auch nicht möglich durch Erhöhung der Dimension kleinere  $\varepsilon$ 's zuzulassen. Anschaulich bedeutet dies, dass die ‘absolute Dicke’ des Ringgebiets im Bild auf Seite 130 konstant ist — wenn wir von einem *schmalen Ring* sprechen, so ist damit die ‘Dicke relativ zu dessen Radius’ gemeint, denn dieser wächst in der Tat mit  $\sqrt{d}$ .

**Satz 10.4.** (*Gaußscher Ringsatz*) Sei  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum, sei  $X: \Omega \rightarrow \mathbb{R}^d$  ein Zufallsvektor mit  $X \sim \mathcal{N}(0, 1, \mathbb{R}^d)$  und  $0 \leq \varepsilon \leq \sqrt{d}$ . Dann gilt

$$P[||X|| - \sqrt{d} \leq \varepsilon] \geq 1 - 2e^{-\varepsilon^2/16}.$$

*Beweis.* Wir multiplizieren mit der Ungleichung  $||X|| + \sqrt{d} \geq \sqrt{d}$  und erhalten

$$\begin{aligned}
 P[||X|| - \sqrt{d} \geq \varepsilon] &\leq P[||X|| - \sqrt{d} \cdot (||X|| + \sqrt{d}) > \varepsilon \cdot \sqrt{d}] \\
 &= P[|X_1^2 + \cdots + X_d^2 - d| > \varepsilon \sqrt{d}] \\
 &= P[|(X_1^2 - 1) + \cdots + (X_d^2 - 1)| > \varepsilon \sqrt{d}] \\
 &= P\left[\left|\frac{X_1^2 - 1}{2} + \cdots + \frac{X_d^2 - 1}{2}\right| > \frac{\varepsilon \sqrt{d}}{2}\right] \\
 &= P[|Y_1 + \cdots + Y_d| > a]
 \end{aligned}$$

wobei wir im letzten Schritt  $Y_i := \frac{X_i^2 - 1}{2}$  und  $a := \frac{\varepsilon \sqrt{d}}{2}$  abgekürzt haben und natürlich nun die Bernstein Tailbound aus Lemma 10.2 anwenden wollen. Hierfür prüfen wir erst die drei dort geforderten Voraussetzungen.

- ① Die  $Y_i$  sind paarweise unabhängig nach Fakt 10.1(i).
- ② Es gilt  $E(Y_i) = \frac{1}{2}(E(X_i^2) - 1) = \frac{1}{2}(E((X_i - E(X_i))^2) - 1) = 0$ , wobei wir erst

$E(X_i) = 0$  eingefügt und dann  $E((X_i - E(X_i))^2) = V(X_i) = 1$  eingesetzt haben.

③ Für  $k \geq 2$  gilt

$$|E(Y_i^k)| = |E((\frac{X_i^2 - 1}{2})^k)| = \frac{1}{2^k} |E((X_i^2 - 1)^k)| \leq \frac{1}{2^k} E(X_i^{2k} + 1)$$

wobei die Ungleichung aus der Monotonie des Erwartungswertes folgt und aus der (punktweisen) Abschätzung  $|X_i^2 - 1|^k \leq X_i^{2k} + 1$ : Ist nämlich  $|X_i| \leq 1$ , so haben wir  $0 \leq X_i^2 \leq 1$ , also  $|X_i^2 - 1|^k \leq 1$ . Ist andererseits  $|X_i| > 1$ , so haben wir  $X_i^2 - 1 > 0$ , also  $|X_i^2 - 1|^k = (X_i^2 - 1)^k \leq (X_i^2)^k = X_i^{2k}$ . Jetzt verwenden wir das Gesetz des unbewussten Statistikers und Lemma 10.3, um das  $2k$ -te Moment von  $X_i \sim \mathcal{N}(0, 1)$  zu berechnen.

$$\begin{aligned} \frac{1}{2^k} E(X_i^{2k} + 1) &= \frac{1}{2^k} \left( \frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} t^{2k} e^{-t^2/2} dt + 1 \right) = \frac{1}{2^k} \left( \frac{(2k)!}{2^k k!} + 1 \right) \\ &= \frac{(2k) \cdot (2k-1) \cdot (2k-2) \cdot (2k-3) \cdots 5 \cdot 4 \cdot 3 \cdot 2 \cdot 1}{(2k)^2 \cdot (2k-2)^2 \cdots (2 \cdot 3)^2 \cdot (2 \cdot 2)^2 \cdot (2 \cdot 1)^2} \cdot k! + \frac{1}{2^k} \\ &= \frac{(2k-1) \cdot (2k-3) \cdots 5 \cdot 3 \cdot 1}{(2k) \cdot (2k-2) \cdots 6 \cdot 4 \cdot 2} \cdot k! + \frac{1}{2^k} \\ &\leq 1 \cdot 1 \cdots 1 \cdot \frac{3}{4} \cdot \frac{1}{2} \cdot k! + \frac{k!}{4 \cdot 2} \\ &= \left( \frac{3}{4} + \frac{1}{4} \right) \frac{k!}{2} \end{aligned}$$

woraus sich  $|E(Y_i^k)| \leq k!/2$  ergibt und damit auch die letzte Voraussetzung von Lemma 10.2 erfüllt ist. Anwendung von Lemma 10.2, zusammen mit unserer Vorarbeit zu Beginn dieses Beweises, liefert

$$\begin{aligned} P[||x|| - \sqrt{d} \geq \varepsilon] &\leq P[|Y_1 + \cdots + Y_d| \geq \frac{\varepsilon\sqrt{d}}{2}] \\ &\leq 2e^{-\frac{1}{4} \min(\frac{(\varepsilon\sqrt{d}/2)^2}{d}, \frac{\varepsilon\sqrt{d}}{2})} \\ &\leq 2e^{-\frac{1}{4} \min(\frac{\varepsilon^2}{4}, \frac{\varepsilon^2}{2})} \\ &= 2e^{-\varepsilon^2/16} \end{aligned}$$

wobei wir die Voraussetzung  $\varepsilon \leq \sqrt{d}$  benutzt haben um  $\frac{\varepsilon\sqrt{d}}{2} \geq \frac{\varepsilon^2}{2}$  abzuschätzen.  $\square$

Wir setzen die bereits vor dem Satz begonnenen Bemerkungen fort und illustrieren das obige Ergebnis durch ein Zahlenbeispiel.

**Beispiel 10.5.** Sei  $\varepsilon = 10$  und  $d \geq 100$  und bezeichne mit

$$R_d := \bar{B}_{\sqrt{d}+10}(0) \setminus B_{\sqrt{d}-10}(0)$$

das Ringgebiet mit Mittelpunkt Null und Radien  $\sqrt{d} \pm 10$ . Dann gilt für  $X \sim \mathcal{N}(0, 1, \mathbb{R}^d)$  nach Satz 10.4

$$P[X \in R_d] = P[\sqrt{d} - 10 \leq \|X\| \leq \sqrt{d} + 10] \geq 0.99.$$

Dabei ist  $R_d = \bar{B}_{20}(0)$  für  $d = 100$  eine Vollkugel, aber für  $d > 100$  wird es ein Ring und der Ring wird für sehr große  $d$  im Verhältnis zum Radius schmal, da seine ‘Dicke’ konstant gleich 20 ist.

Mit denselben Argumenten die im Beweis von Satz 10.4 zur Anwendung kamen, erhalten wir eine Abschätzung für  $P[||X||^2 - d| \geq \varepsilon]$ , vergleiche unsere Berechnungen von Erwartungswert und Varianz (8.1) und (8.2) auf Seite 114. Wir überlassen den Beweis, sowie die Frage für welche  $\varepsilon$  und  $d$  die Abschätzungen in Korollar 10.6 interessant sind, als Aufgabe 10.2 dem Leser.

**Korollar 10.6.** *Sei  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum,  $X: \Omega \rightarrow \mathbb{R}^d$  ein Zufallsvektor mit  $X \sim \mathcal{N}(0, 1, \mathbb{R}^d)$ . Dann gilt*

$$P[||X||^2 - d| \leq \varepsilon] \geq 1 - 2e^{-\frac{1}{8} \min(\frac{\varepsilon^2}{2d}, \varepsilon)}$$

für  $0 < \varepsilon \leq \sqrt{d}$ . □

Als nächstes betrachten wir Winkel zwischen normalverteilten unabhängigen Zufallsvariablen, vergleiche auch hier unsere Ergebnisse zu Erwartungswert und Varianz in Satz 8.4, sowie die Simulationen in Aufgaben 8.2. Wir weisen wieder darauf hin, dass die Schranke im folgenden Gaußschen Orthogonalitätssatz 10.7 für  $\varepsilon > \frac{2}{\sqrt{d}-7}$  nicht-trivial ist. Im Gegensatz zum Ringsatz ist hier also wieder ein Tradeoff möglich im Sinne, dass umso kleinere  $\varepsilon$  betrachtet werden können, je größer die Dimension ist.

**Satz 10.7.** (Gaußscher Orthogonalitätssatz) *Sei  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum,  $X, Y: \Omega \rightarrow \mathbb{R}^d$  unabhängige Zufallsvektoren mit  $X, Y \sim \mathcal{N}(0, 1, \mathbb{R}^d)$ . Für  $d \in \mathbb{N}$  und  $\varepsilon > 0$  gilt dann*

$$P[|\langle \frac{X}{||X||}, \frac{Y}{||Y||} \rangle| \leq \varepsilon] \geq 1 - \frac{2/\varepsilon + 7}{\sqrt{d}}.$$

*Hierbei lesen wir den linken Ausdruck stillschweigend als bedingte Wahrscheinlichkeit unter der Bedingung, dass  $||X||, ||Y|| \neq 0$  gilt — was natürlich mit Wahrscheinlichkeit Eins der Fall ist.*

*Beweis.* Da  $X$  und  $Y$  unabhängig sind, genügt es, den Beweis für  $Y \equiv y \in \mathbb{R}^d \setminus \{0\}$  zu machen. Wir setzen  $\lambda_i := \frac{y_i}{||y||}$ , und definieren die Zufallsvariable  $U: \Omega \rightarrow \mathbb{R}$  per

$$U := \langle X, \frac{y}{||y||} \rangle = \sum_{i=1}^d \frac{y_i}{||y||} X_i = \sum_{i=1}^d \lambda_i X_i.$$

Mit Fakt 10.1(ii) folgt dann  $U \sim \mathcal{N}(0, \lambda_1^2 + \dots + \lambda_d^2) = \mathcal{N}(0, 1)$ . Jetzt benutzen wir erst den Satz von der totalen Wahrscheinlichkeit, notieren dabei aber den sich ergebenden zweiten Summanden erst gar nicht, sondern schätzen ihn nach unten mit Null ab. Dies mag sehr grob erscheinen, aber wegen des Ringsatzes verlieren wir in der Tat fast nichts, da nämlich  $||X|| \leq \frac{\sqrt{d}}{2}$  nur mit kleiner Wahrscheinlichkeit passiert

wenn  $d$  geeignet groß ist. Wir erhalten

$$\begin{aligned}
 \mathbb{P}[|\langle \frac{X}{\|X\|}, \frac{y}{\|y\|} \rangle| \leq \varepsilon] &= \mathbb{P}[|\langle \frac{X}{\|X\|}, \frac{y}{\|y\|} \rangle| \leq \varepsilon \mid \frac{\sqrt{d}}{2} \leq \|X\|] \cdot \mathbb{P}[\frac{\sqrt{d}}{2} \leq \|X\|] + \dots \\
 &\geq \mathbb{P}[|U| \leq \varepsilon \|X\| \mid \frac{\sqrt{d}}{2} \leq \|X\|] \cdot \mathbb{P}[\sqrt{d} - \frac{\sqrt{d}}{2} \leq \|X\| \leq \sqrt{d} + \frac{\sqrt{d}}{2}] \\
 &\geq \mathbb{P}[|U| \leq \frac{\varepsilon\sqrt{d}}{2}] \cdot \mathbb{P}[|\|X\| - \sqrt{d}| \leq \frac{\sqrt{d}}{2}] \\
 &\geq \mathbb{P}[|U| \leq \frac{\varepsilon\sqrt{d}}{2}] \cdot (1 - 2e^{-(\frac{\sqrt{d}}{2})^2/16})
 \end{aligned}$$

wobei die letzte Abschätzung des zweiten Faktors gerade der Ringsatz 10.4 mit  $\varepsilon = \frac{\sqrt{d}}{2}$  ist. Da  $U \sim \mathcal{N}(0, 1)$  gilt, erhalten wir für den ersten Faktor

$$\begin{aligned}
 \mathbb{P}[|U| \leq \frac{\varepsilon\sqrt{d}}{2}] &= \frac{1}{\sqrt{2\pi}} \int_{-\frac{\varepsilon\sqrt{d}}{2}}^{\frac{\varepsilon\sqrt{d}}{2}} e^{-t^2/2} dt \geq 1 - \frac{2}{\sqrt{2\pi}} \int_{\frac{\varepsilon\sqrt{d}}{2}}^{\infty} \frac{1}{t^2} dt \\
 &= 1 - \frac{2}{\sqrt{2\pi}} \left( -\frac{1}{t} \Big|_{\frac{\varepsilon\sqrt{d}}{2}}^{\infty} \right) \geq 1 - 1 \cdot \frac{2}{\varepsilon\sqrt{d}}
 \end{aligned}$$

und wenn wir dies oben einsetzen, ergibt sich

$$\mathbb{P}[|\langle \frac{X}{\|X\|}, \frac{y}{\|y\|} \rangle| \leq \varepsilon] \geq \left(1 - \frac{2}{\varepsilon\sqrt{d}}\right) (1 - 2e^{-\frac{d}{64}}) = 1 - \frac{2}{\varepsilon\sqrt{d}} - 2e^{-\frac{d}{64}} + \frac{4e^{-\frac{d}{64}}}{\varepsilon\sqrt{d}}.$$

Jetzt schätzen wir den positiven Summand am Ende nach unten mit Null ab und ferner  $2e^{-\frac{d}{64}} \leq \frac{7}{\sqrt{d}}$  nach oben. Dies liefert die gewünschte untere Schranke.  $\square$

Wir notieren die folgende Variante, bei der  $Y$  konstant ist und bei der auf die Normierung von  $X$  verzichtet wird. Diese werden wir in Kapitel 12 benutzen.

**Korollar 10.8.** Sei  $(\Omega, \Sigma, \mathbb{P})$  ein Wahrscheinlichkeitsraum,  $Z: \Omega \rightarrow \mathbb{R}^d$  ein Zufallsvektor mit  $Z \sim \mathcal{N}(0, 1, \mathbb{R}^d)$  und  $\xi \in \mathbb{R}^d$  fest. Für  $d \in \mathbb{N}$  und  $\varepsilon > 0$  gilt dann

$$\mathbb{P}[|\langle Z, \xi \rangle| \leq \varepsilon] \geq 1 - \frac{\|\xi\|}{\varepsilon}.$$

*Beweis.* Für  $\xi = 0$  ist die Aussage trivial. Andernfalls berechnen wir mit  $U := \langle Z, \frac{\xi}{\|\xi\|} \rangle \sim \mathcal{N}(0, 1)$  wie im vorhergehenden Beweis

$$\begin{aligned}
 \mathbb{P}[|\langle Z, \xi \rangle| \leq \varepsilon] &= \mathbb{P}[|\langle Z, \frac{\xi}{\|\xi\|} \rangle| \leq \frac{\varepsilon}{\|\xi\|}] = \mathbb{P}[|U| \leq \frac{\varepsilon}{\|\xi\|}] \\
 &\geq 1 - \frac{2}{\sqrt{2\pi}} \int_{\frac{\varepsilon}{\|\xi\|}}^{\infty} \frac{1}{t^2} dt \geq 1 - 1 \cdot \frac{\|\xi\|}{\varepsilon}. \quad \square
 \end{aligned}$$

Wir betrachten nun ein Zahlenbeispiel um ein Gefühl dafür zu bekommen, für welche Dimensionen der Gaußsche Orthogonalitätssatz interessante Abschätzungen liefert.



**Beispiel 10.9.** Sei  $\varepsilon = 0.1$ . Dann gilt für  $X, Y \sim \mathcal{N}(0, 1, \mathbb{R}^d)$

$$\mathbb{P}\left[\left|\left\langle \frac{X}{\|X\|}, \frac{Y}{\|Y\|} \right\rangle\right| \leq 0.1\right] \geq 0.9$$

wenn  $d \geq 100\,000$  ist. Letzteres heißt, dass der Winkel  $\angle(X, Y)$  im Gradmaß  $90^\circ \pm 6^\circ$  beträgt mit einer Wahrscheinlichkeit von mehr als 0.9. Beachte, dass wir zur Anwendung des Orthogonalitätssatzes deutlich höhere Dimensionen benötigen als für den Gaußschen Ringsatzes. Die folgende Tabelle enthält die mittleren normalisierten paarweisen Skalarprodukte von 100 bezüglich Normalverteilung zufällig gewählten Punkten in  $\mathbb{R}^d$ , vergleiche Aufgabe 8.2.

$d$	1	10	100	1 000	10 000	100 000
$\frac{1}{100 \cdot 99} \sum_{i \neq j} \left\langle \frac{x^{(i)}}{\ x^{(i)}\ }, \frac{x^{(j)}}{\ x^{(j)}\ } \right\rangle$	-0.0097	0.0007	-0.0023	-0.0006	0.0001	-0.00004
Varianz	0.9999	0.1018	0.0098	0.0010	0.0001	0.00001

Die Simulation deutet darauf hin, dass eventuell auch in niedrigeren Dimensionen bereits mit fast orthogonalen Vektoren zu rechnen ist.

Als letztes kommen wir nun zum Abstand unabhängiger Zufallsvektoren.

**Satz 10.10.** (Erster Gaußscher Abstandssatz) *Sei  $(\Omega, \Sigma, \mathbb{P})$  ein Wahrscheinlichkeitsraum,  $X, Y: \Omega \rightarrow \mathbb{R}^d$  unabhängige Zufallsvektoren mit  $X, Y \sim \mathcal{N}(0, 1, \mathbb{R}^d)$  und sei  $0 < \varepsilon \leq \sqrt{d}$ . Dann gilt*

$$\mathbb{P}\left[\left|\|X - Y\| - \sqrt{2d}\right| \leq \varepsilon\right] \geq 1 - 4e^{-\varepsilon^2/72} - \frac{6\sqrt{2}}{\varepsilon d} - \frac{7}{\sqrt{d}}.$$

*Beweis.* Wir gehen analog zum Beweis des Ringsatzes vor, bzw. zum Beweis von Satz 9.9. D.h. wir drehen die Ungleichung um und multiplizieren mit  $\|X - Y\| + \sqrt{2d} \geq \sqrt{2d}$ . Dies führt auf

$$\begin{aligned} \mathbb{P}\left[\left|\|X - Y\| - \sqrt{2d}\right| \leq \varepsilon\right] &= 1 - \mathbb{P}\left[\left|\|X - Y\| - \sqrt{2d}\right| \geq \varepsilon\right] \\ &\geq 1 - \mathbb{P}\left[\left|\|X - Y\| - \sqrt{2d}\right|(\|X - Y\| + \sqrt{2d}) \geq \varepsilon\sqrt{2d}\right] \\ &= 1 - \mathbb{P}\left[\left|\|X - Y\|^2 - 2d\right| \geq \varepsilon\sqrt{2d}\right] \\ &= \mathbb{P}\left[-\varepsilon\sqrt{2d} \leq \|X\|^2 + \|Y\|^2 - 2\langle X, Y \rangle - 2d \leq \varepsilon\sqrt{2d}\right] \\ &\geq \mathbb{P}\left[\left|\|X\|^2 - d\right| \leq \frac{\varepsilon\sqrt{2d}}{3}\right]^2 \cdot \mathbb{P}\left[|\langle X, Y \rangle| \leq \frac{\varepsilon\sqrt{2d}}{6}\right] \end{aligned}$$

und wieder schätzen wir jetzt beide der Faktoren einzeln ab. Dann ergeben sich

$$\mathbb{P}\left[\left|\|X\|^2 - d\right| \leq \frac{\varepsilon\sqrt{2d}}{3}\right] \underset{\substack{\uparrow \\ \text{Kor. 10.6}}}{\geq} 1 - 2e^{-\frac{1}{8} \min(\frac{\varepsilon^2}{9}, \frac{\varepsilon\sqrt{2d}}{3})} \underset{\substack{\uparrow \\ \varepsilon \leq \sqrt{d}}}{\geq} 1 - 2e^{-\varepsilon^2/72}$$

und

$$\mathbb{P}\left[|\langle X, Y \rangle| \leq \frac{\varepsilon\sqrt{2d}}{6}\right] \underset{\substack{\uparrow \\ \text{Satz 10.7}}}{\geq} 1 - \frac{\frac{12}{\varepsilon\sqrt{2d}} + 7}{\sqrt{d}} = 1 - \left(\frac{6\sqrt{2}}{\varepsilon d} + \frac{7}{\sqrt{d}}\right).$$

Ausmultiplizieren und weiteres Abschätzen liefert

$$\begin{aligned} (1 - 2e^{-\varepsilon^2/72})^2 \left(1 - \left(\frac{6\sqrt{2}}{\varepsilon d} + \frac{7}{\sqrt{d}}\right)\right) &= (1 - 4e^{-\varepsilon^2/72} + 4e^{-\varepsilon^2/36}) \left(1 - \left(\frac{6\sqrt{2}}{\varepsilon d} + \frac{7}{\sqrt{d}}\right)\right) \\ &\geq 1 - 4e^{-\varepsilon^2/72} - \frac{6\sqrt{2}}{\varepsilon d} - \frac{7}{\sqrt{d}} \end{aligned}$$

wie behauptet.  $\square$

Wir überlassen es dem Leser, zu prüfen für welche  $\varepsilon$  und  $d$  der Satz tatsächlich die Interpretation erlaubt, dass mit hoher Wahrscheinlichkeit  $\|X - Y\| \approx \sqrt{2d}$  gilt, sowie den Vergleich mit den experimentellen Daten in der folgenden Tabelle

$d$	1	10	100	1000	10 000	100 000
$\frac{1}{100.99} \sum_{i \neq j} \ x^{(i)} - x^{(j)}\ $	1.06	4.41	14.05	44.65	141.50	447.35
$\sqrt{2d}$	1.41	4.47	14.14	44.72	141.42	447.21
Varianz	0.74	1.10	0.92	0.89	0.96	1.08

in welcher wieder für 100 bezüglich Normalverteilung zufällig gewählte Punkte im  $\mathbb{R}^d$  nun allerdings deren paarweise mittlere Abstände angegeben sind.

Zum Abschluss notieren wir noch das folgende Korollar über die Quadrate der Abstände von zufällig gaußverteilten Punkten.

**Korollar 10.11.** *Es gilt*

$$\mathbb{P}\left[|\|X - Y\|^2 - 2d| \leq \varepsilon\right] \geq 1 - 4e^{-\varepsilon^2/72} - \frac{6\sqrt{2}}{\varepsilon d} - \frac{7}{\sqrt{d}}$$

unter den Voraussetzungen von Satz 10.10.  $\square$

## Referenzen

Dieses Kapitel basiert hauptsächlich auf [BHK20] und [Ver18], enthält aber auch einige Umformulierungen und kleine Erweiterungen der dort diskutierten Ergebnisse.

## Aufgaben

**Aufgabe 10.1.** (Klassisches Beispiel der Chernoff-Methode) Seien  $Y_1, \dots, Y_n$  unabhängige Bernoulli Zufallsvariablen mit  $\mathbb{P}[Y_i = 1] = p \in [0, 1]$  und sei  $Y := Y_1 + \dots + Y_n$ . Sei  $\delta > 0$ .

- (i) Zeigen Sie, dass  $\mathbb{E}(e^{tY_i}) \leq e^{p(e^t - 1)}$  für  $t > 0$  gilt.

- (ii) Gehen Sie wie im Beweis von Lemma 10.2 vor um die folgende Abschätzung zu zeigen

$$\mathbb{P}[X \geq (1 + \delta)np] \leq \left( \frac{e^\delta}{(1 + \delta)^{1 + \delta}} \right)^{np}.$$

*Hinweis:* Oft ist es bei der Chernoff-Methode gar nicht nötig, das Infimum, wie in Lemma 10.2 geschehen, explizit auszurechnen. In aktuellen Beispiel genügt es etwa  $t = \log(1 + \delta)$  zu wählen.

- (iii) Wir betrachten jetzt ein Experiment bei dem ein fairer Würfel  $n$ -mal gewürfelt wird. Wenden Sie (ii) an, um die Wahrscheinlichkeit dafür abzuschätzen, dass in mindestens 70% der  $n$ -vielen Würfe eine Sechs fällt.
- (iv) Vergleiche die Schranke aus (iii) mit den Schranken die sich durch Anwendung der Markov- bzw. der Tschebyscheff-Ungleichung ergeben. Simulieren Sie das Experiment und testen Sie, wie scharf die drei theoretischen Schranken sind.

**Aufgabe 10.2.** Beweisen Sie Korollar 10.6.

**Aufgabe 10.3.** Nach dem Gaußschen Orthogonalitätssatz sind zwei Punkte, die im hochdimensionalen Raum gemäß Normalverteilung zufällig gewählt werden, mit hoher Wahrscheinlichkeit fast orthogonal. Andererseits ist klar, dass die Wahrscheinlichkeit dafür, dass die Punkte exakt orthogonal sind, Null ist. Man kann also erwarten, dass auch die Wahrscheinlichkeit dafür, dass das Skalarprodukt sehr klein ist, selbst klein sein wird. Quantifizieren Sie dies, indem Sie zeigen, dass

$$\mathbb{P}\left[\left|\left\langle \frac{X}{\|X\|}, \frac{Y}{\|Y\|} \right\rangle\right| \leq \frac{\varepsilon}{2\sqrt{d}}\right] \leq 2e^{-\frac{d}{16}} + \varepsilon$$

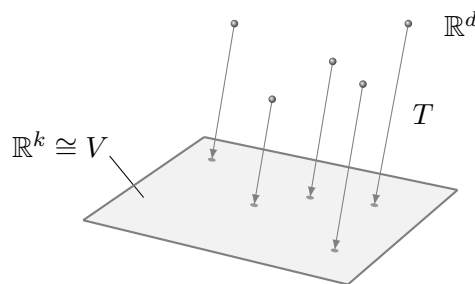
für  $\varepsilon > 0$  gilt wenn  $X, Y \sim \mathcal{N}(0, 1)$ . Finden Sie überdies ein Zahlenbeispiel, welches illustriert, was im Vortext mit ‘sehr kleinem Skalarprodukt’ gemeint ist.

*Hinweis:* Bringen Sie die Zufallsvariable  $U$  wie im Beweis von Satz 10.7 ins Spiel und überlegen Sie sich dann, dass für Zufallsvariablen  $A, B$  und Konstanten  $a, b$  gilt  $\mathbb{P}[A \cdot B \leq a \cdot b] \leq \mathbb{P}[A \leq a \text{ oder } B \leq b]$ .

## Kapitel 11

# Dimensionalitätsreduktion à la Johnson-Lindenstrauss

Zu Beginn von Kapitel 8 hatten wir Situationen diskutiert, in denen die Dimension  $d$  der Features einer natürlich gegebenen Datenmenge  $D$  sehr groß ist, und sogar sehr viel größer sein kann, als die Anzahl  $n$  der Datenpunkte. Besonders im letzteren Fall stellt sich die Frage, ob tatsächlich alle Dimensionen nötig sind, um die in der Datenmenge enthaltenen Informationen darzustellen. Strategien um die Dimensionalität (formal:  $\dim \text{span } D$ ) ohne, oder zumindest mit kontrollierbarem, Informationsverlust zu verringern, laufen unter der Bezeichnung *Dimensionalitätsreduktion*. Eine naheliegende Idee besteht darin, die Datenmenge auf einen  $k$ -dimensionalen Teilraum  $V \subset \mathbb{R}^d$  zu *projizieren*, wobei  $k$  deutlich kleiner als  $d$  ist.



Hierbei kann man zuerst einmal nach einer Projektion  $T: \mathbb{R}^d \rightarrow \mathbb{R}^d$  im Sinne der Linearen Algebra Ausschau halten, d.h.  $T \in L(\mathbb{R}^d)$  mit  $T^2 = T$  fordern und  $V := \text{ran } T$  setzen. Außerdem wäre es gut, wenn  $T$  unabhängig von der gegebenen Datenmenge ist, und idealerweise die paarweisen Abstände zwischen Datenpunkten durch  $T$  nicht verändert würden.

Man merkt nun schnell, dass diese Wunschliste an Eigenschaften für  $T$  etwas zu viel des Guten ist: will man die Erhaltung aller Abstände unabhängig von der Datenmenge, so müsste  $T$  in der Tat eine Isometrie sein und wäre dann insbesondere injektiv, was sich mit der Idee einer *Dimensionalitätsreduktion* nicht verträgt.

Aus der Linearen Algebra wissen wir andererseits, dass eine Isometrie durch Mul-

tiplikation mit einer orthogonalen Matrix gegeben ist. Es liegt also nahe, es für  $k < d$  mit einer orthogonalen Projektion zu versuchen, bzw. mit

$$T_A: \mathbb{R}^d \rightarrow \mathbb{R}^k, \quad x \mapsto Ax,$$

wobei  $A \in \mathbb{R}^{k \times d}$  orthonormale Zeilen hat. Letzteres garantiert  $AA^\top = \text{id}_{\mathbb{R}^k}$  und daher, dass  $T_A$  orthogonal auf den  $k$ -dimensionalen Unterraum  $V := \text{ran } T \subset \mathbb{R}^d$  abbildet. Wir benötigen jetzt also eine *beliebige aber konkrete* Matrix  $A$ . Sei dazu  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum. Wir betrachten eine messbare Abbildung

$$U: \Omega \rightarrow \mathbb{R}^{k \times d}, \quad U = \begin{bmatrix} u_{11} & \cdots & u_{1d} \\ \vdots & & \vdots \\ u_{k1} & \cdots & u_{kd} \end{bmatrix}$$

mit  $U \sim \mathcal{N}(0, 1, \mathbb{R}^{k \times d})$ . Letzteres bedeutet, dass die Einträge  $u_{ij} \sim \mathcal{N}(0, 1)$  unabhängige Zufallsvariablen sind und wir uns daher Realisierungen von  $U$  leicht verschaffen können. Wir halten die folgenden Eigenschaften einer solchen *Zufallsmatrix* fest.

**Bemerkung 11.1.** Sei  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum und für  $k < d$  sei  $U: \Omega \rightarrow \mathbb{R}^{k \times d}$  eine Zufallsmatrix mit  $U \sim \mathcal{N}(0, 1, \mathbb{R}^{k \times d})$ .

- (i) Bezeichnen wir mit  $u_i = [u_{i1} \cdots u_{id}]$  die Zeilen von  $U$ , so sind diese, bzw. deren Transponate, für  $i = 1, \dots, k$  unabhängige  $\mathbb{R}^d$ -wertige Zufallsvektoren mit  $u_i \sim \mathcal{N}(0, 1, \mathbb{R}^d)$ .
- (ii) Der Gaußsche Orthogonalitätssatz 10.7 liefert dann für  $i \neq j$  und  $\varepsilon > 0$

$$P[|\langle \frac{u_i}{\|u_i\|}, \frac{u_j}{\|u_j\|} \rangle| \leq \varepsilon] \geq 1 - \frac{2/\varepsilon + 7}{\sqrt{d}}.$$

Wir erhalten also, dass die Zeilen von  $U$  mit hoher Wahrscheinlichkeit fast orthogonal sind. Würden wir selbige jetzt noch normieren, so hätten wir  $UU^\top \approx \text{id}_{\mathbb{R}^k}$  und kämen damit einer Matrix  $A$  wie wir sie uns anfangs gewünscht hatten, sehr nahe.

**Satz 11.2.** (Zufallsprojektion) Sei  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum und für  $k < d$  sei  $U: \Omega \rightarrow \mathbb{R}^{k \times d}$  eine Zufallsmatrix mit  $U \sim \mathcal{N}(0, 1, \mathbb{R}^{k \times d})$ . Für  $x \in \mathbb{R}^d \setminus \{0\}$  und  $0 < \varepsilon \leq 1$  gilt dann

$$P[|\|Ux\| - \sqrt{k}\|x\|| \leq \varepsilon\sqrt{k}\|x\|] \geq 1 - 2e^{-k\varepsilon^2/16}.$$

*Beweis.* Wir notieren zuert, dass für beliebiges  $x \in \mathbb{R}^d$  die Anwendung von  $U$  auf  $x$

$$Ux = \begin{bmatrix} u_{11} & \cdots & u_{1d} \\ \vdots & & \vdots \\ u_{k1} & \cdots & u_{kd} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix} = \begin{bmatrix} \langle u_1, x \rangle \\ \vdots \\ \langle u_k, x \rangle \end{bmatrix}$$

gerade den obigen Vektor von Skalarprodukten ergibt. Jetzt fixieren wir  $x \in \mathbb{R}^d \setminus \{0\}$  und betrachten den  $k$ -dimensionalen Zufallsvektor

$$U(\cdot) \frac{x}{\|x\|}: \Omega \rightarrow \mathbb{R}^k.$$

Dessen  $i$ -te Koordinatenfunktion ist durch

$$(U(\cdot) \frac{x}{\|x\|})_i = \langle u_i, \frac{x}{\|x\|} \rangle = \sum_{j=1}^d u_{ij} \frac{x_j}{\|x\|}$$

gegeben und wir erhalten mit Fakt 10.1(ii), dass

$$(U(\cdot) \frac{x}{\|x\|})_i \sim \mathcal{N}(0, \frac{x_1^2}{\|x\|^2} + \dots + \frac{x_d^2}{\|x\|^2}) = \mathcal{N}(0, 1)$$

für alle  $i = 1, \dots, k$  gilt und daher  $U(\cdot) \frac{x}{\|x\|} \sim \mathcal{N}(0, 1, \mathbb{R}^k)$  mit Fakt 10.1(i). Es folgt

$$\begin{aligned} \mathbb{P}[\|Ux\| - \sqrt{k}\|x\| \leq \varepsilon\sqrt{k}\|x\|] &= \mathbb{P}[|\|U \frac{x}{\|x\|}\| - \sqrt{k}| \leq \varepsilon\sqrt{k}] \\ &\leq 1 - 2e^{-(\varepsilon\sqrt{k})^2/16} \\ &= 2e^{-\varepsilon^2 k/16} \end{aligned}$$

wobei wir für die Abschätzung den Ringsatz 10.4 angewandt haben und zwar mit  $k$  statt  $d$  und  $\varepsilon\sqrt{k}$  statt  $\varepsilon$ . Die dafür nötige Voraussetzung  $0 < \varepsilon\sqrt{k} \leq \sqrt{k}$  gilt wegen unserer Annahme  $\varepsilon \leq 1$ .  $\square$

Da wir  $k$  vorgeben, macht es Sinn durch  $\sqrt{k}$  zu dividieren und dann  $\frac{1}{\sqrt{k}}$  in die Matrix  $U$  zu ‘absorbieren’.

**Definition 11.3.** Sei  $(\Omega, \Sigma, \mathbb{P})$  ein Wahrscheinlichkeitsraum und für  $k < d$  sei  $U: \Omega \rightarrow \mathbb{R}^{k \times d}$  eine Zufallsmatrix mit  $U \sim \mathcal{N}(0, 1, \mathbb{R}^{k \times d})$  und  $\omega \in \Omega$ . Die folgende Abbildung

$$T_{U(\omega)}: \mathbb{R}^d \rightarrow \mathbb{R}^k, \quad T_{U(\omega)}x := \frac{1}{\sqrt{k}}U(\omega)x.$$

heißt *Johnson-Lindenstrauss-Projektion*.

Hierbei ist zu beachten, dass der Begriff ‘Projektion’ formal nicht konsistent mit dessen Bedeutung im Sinne der Linearen Algebra ist — er ist es aber fast, nämlich wenn wir  $\mathbb{R}^k$  per  $U(\omega)^\top$  mit einem Unterraum von  $\mathbb{R}^d$  identifizieren und dabei ignorieren, dass  $U(\omega)^\top$  eventuell nicht ganz vollen Rang hat.

**Bemerkung 11.4.** (i) Mit der obigen Notation wird Satz 11.2 zu

$$\forall x \in \mathbb{R}^d \setminus \{0\}, \quad 0 < \varepsilon \leq 1: \quad \mathbb{P}[\|T_U x\| - \|x\| \leq \varepsilon\|x\|] \geq 1 - 2e^{-k\varepsilon^2/16}$$

wobei  $U \sim \mathcal{N}(0, 1, \mathbb{R}^{k \times d})$  ist. Die Heuristik ist jetzt natürlich, dass dann

$$\left| \frac{\|T_U x\|}{\|x\|} - 1 \right| \leq \varepsilon \quad \text{bzw.} \quad \frac{\|T_U x\|}{\|x\|} \approx 1$$

mit hoher Wahrscheinlichkeit gilt.

(ii) Der Nachteil bei obigem ist, dass nur die relative Änderung der Längen nah bei Eins ist, aber deren absolute Änderung durchaus groß sein kann. Wählen wir unsere Punkte oben allerdings nur aus einer beschränkte Menge  $B \subseteq \bar{B}_r(0)$  mit

festem  $r > 0$ , so können wir durch passende Wahl von  $\varepsilon$  und  $k$  bei gegebenem  $r$  erreichen, dass

$$|\|T_U x\| - \|x\|| \approx 0$$

mit hoher Wahrscheinlichkeit gilt. Als Zahlenbeispiel betrachte etwa  $r = 1$ . Wählen wir  $\varepsilon = 0.2$  und  $k = 1000$ , so folgt  $P[|\|T_U x\| - \|x\|| < 0.2] \geq 0.8$  für jedes  $x \in \mathbb{R}^d$  mit  $\|x\| \leq 1$  — und dies sogar unabhängig von  $d \geq 1000$ .

Als Nächstes wollen wir untersuchen, was mit paarweisen Abständen passiert wenn wir diese per  $T_U$  projizieren.

**Satz 11.5.** (Johnson-Lindenstrauss Lemma) *Sei  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum und weiter  $0 < \varepsilon < 1$ ,  $n \geq 1$  und  $k \geq \frac{48}{\varepsilon^2} \log n$ . Sei  $U: \Omega \rightarrow \mathbb{R}^{k \times d}$  eine Zufallsmatrix mit  $U \sim \mathcal{N}(0, 1, \mathbb{R}^{k \times d})$ . Dann gilt für je  $n$ -viele Punkte  $x^{(1)}, \dots, x^{(n)} \in \mathbb{R}^d$  die Abschätzung*

$$P[(1 - \varepsilon)\|x_i - x_j\| \leq \|T_U x_i - T_U x_j\| \leq (1 + \varepsilon)\|x_i - x_j\| \text{ für alle } i, j] \geq 1 - \frac{1}{n}.$$

*Beweis.* Für  $i \neq j$  definieren wir  $x := x_i - x_j$  und berechnen mithilfe des vorhergehenden Satzes 11.2

$$\begin{aligned} P_{ij} &:= P[(1 - \varepsilon)\|x_i - x_j\| \leq \|T_U x_i - T_U x_j\| \leq (1 + \varepsilon)\|x_i - x_j\|] \\ &= P[(1 - \varepsilon)\|x\| \leq \|T_U x\| \leq (1 + \varepsilon)\|x\|] \\ &= P[|\|T_U x\| - \|x\|| \leq \varepsilon\|x\|] \\ &= P[|\|Ux\| - \sqrt{k}\|x\|| \leq \varepsilon\sqrt{k}\|x\|] \\ &\geq 1 - 2e^{-k\varepsilon^2/16}. \end{aligned}$$

Da es  $\binom{n}{2} \leq \frac{n^2}{2}$  viele Möglichkeiten für  $1 \leq i, j \leq n$  mit  $i \neq j$  gibt, erhalten wir für die im Satz angegebene Wahrscheinlichkeit

$$\begin{aligned} P[\dots] &= 1 - P[\exists i \neq j: \|T_U x_i - T_U x_j\| \notin ((1 - \varepsilon)\|x_i - x_j\|, (1 + \varepsilon)\|x_i - x_j\|)] \\ &= 1 - \sum_{i \neq j} (1 - P_{ij}) \geq 1 - \sum_{i \neq j} 2e^{-k\varepsilon^2/16} \geq 1 - \frac{n^2}{2} \cdot 2e^{-k\varepsilon^2/16} \\ &\geq 1 - n^2 e^{-(\frac{48}{\varepsilon^2} \log n) \varepsilon^2/16} \geq 1 - n^2 e^{\log(n^{-3})} \\ &= 1 - \frac{1}{n} \end{aligned}$$

wie behauptet. □

Wir notieren zunächst die folgenden bemerkenswerten Punkte zu obigem Satz.

- (i) Die Dimension  $d$  des Ausgangsraumes kommt weder in der Normabschätzung noch in der Wahrscheinlichkeitsabschätzung vor.
- (ii) Die Anzahl  $n$  der Punkte geht in die Wahrscheinlichkeitsabschätzung per  $1/n$  ein, in die Dimension  $k$  des niedrigdimensionalen Raumes aber nur mit  $\log n$ .

- (iii) Die Normabschätzung gilt für alle Punktpaare mit hoher Wahrscheinlichkeit und nicht nur mit hoher Wahrscheinlichkeit für alle Punktpaare.

Jetzt knüpfen wir nochmal an Bemerkung 11.4(ii) an.

**Bemerkung 11.6.** Wir betrachten eine Datenmenge

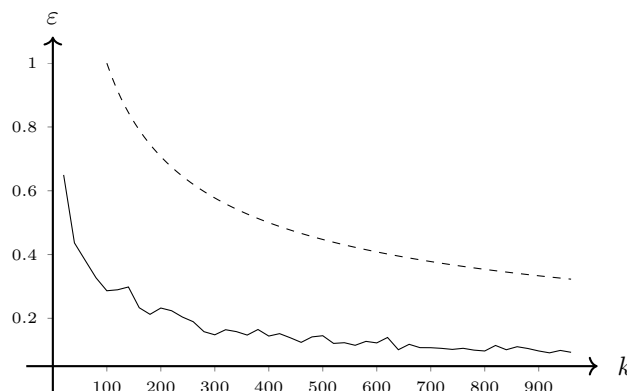
$$D = \{(x^{(i)}, y^{(i)}) \mid i = 1, \dots, n_1\} \subseteq \mathbb{R}^d \times \mathbb{R}$$

bei der die  $x^{(i)}$  in einer beschränkten Menge  $B \subseteq \mathbb{R}^d$  liegen und wir nehmen an, dass wir daran interessiert sind, einen auf  $B$  definierten Prediktor zu ermitteln. Dazu wählen wir  $n_2$ -viele ungelabelte Punkte in  $B$ , projizieren die gelabelten Punkte (Trainingsdaten) sowie die ungelabelten Punkte (Testdaten) auf  $\mathbb{R}^k$  und führen dort einen geeigneten Algorithmus aus, welcher den ungelabelten Punkten Label zuweist. Durch Wahl einer geeignet großen Test- oder Trainingsdatenmenge erreichen wir, dass  $n = n_1 + n_2$  groß und damit schonmal die Wahrscheinlichkeit in Satz 11.5 nah bei Eins liegen wird. Da  $n$  in  $k$  nur logarithmisch eingeht, hängt  $k$  praktisch nur von  $\varepsilon$  ab. Da alle Datenpunkte aus einer beschränkten Menge kommen, können wir wie in Bemerkung 11.4(ii) erreichen, dass

$$\|T_U x^{(i)} - T_U x^{(j)}\| \approx 0 \quad \text{und} \quad k \ll d$$

gilt und zwar für alle Punkte  $x^{(i)}, x^{(j)}$  aus Trainingsmenge oder Testmenge.

In Aufgabe 11.3 werden wir zufällig erzeugte Daten mit der Johnson-Lindenstrauss Projektion behandeln, dabei den tatsächlichen Fehler der Abstände ausrechnen und mit der Schranke aus dem Johnson-Lindenstrauss Lemma vergleichen. Bei  $n = 300$  unabhängigen Stichproben von  $X \sim \mathcal{N}(0, 1, \mathbb{R}^{1000})$  ergibt sich das folgende Bild für den tatsächlichen multiplikativen Fehler (durchgezogene Linie) im Vergleich zur Schranke  $\varepsilon = \sqrt{48 \log(n)/k}$  (gestrichelte Linie) des Johnson-Lindenstrauss-Lemmas.



## Referenzen

Dieses Kapitel basiert ebenfalls auf [BHK20]. Wir weisen auch auf die Originalarbeit [JL84] hin. Einige der Aufgaben in diesem und den vorherigen Kapiteln sind [For] entnommen.



Außerdem notieren wir, dass im Internet Versionen des Johnson-Lindenstrauss-Lemmas zu finden sind, bei denen die Konstante 8 (statt 48) beträgt; dem Autor ist unbekannt, ob diese Version korrekt ist.

## Aufgaben

**Aufgabe 11.1.** Sei  $X = X_1 + \dots + X_d$  mit  $X_i \sim \mathcal{N}(0, 1)$  gegeben.

- (i) Zeigen Sie  $E(e^{tX_i}) = (1 - 2t)^{-d/2}$  für  $t \in (0, \frac{1}{2})$ .
- (ii) Zeigen Sie  $P[X \geq a] \leq \inf_{t \in (0, \frac{1}{2})} \frac{e^{-ta}}{(1-2t)^{d/2}}$  für  $a > 0$ .

*Hinweis:* Für (i) verwende das Gesetz des unbewussten Statistikers  $E(f(X_i)) = \frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} f(t) e^{-\frac{t^2}{2}} dt$ .

**Aufgabe 11.2.** Benutzen Sie Aufgabe 11.1 um einen alternativen Beweis des Johnson-Lindenstrauss-Lemmas zu geben, welcher *ohne* den Gaußschen Rungssatz auskommt.

**Aufgabe 11.3.** Sei  $d > k$  gegeben.

- (i) Implementieren Sie die Johnson-Lindenstrauss-Projektion  $T_U: \mathbb{R}^d \rightarrow \mathbb{R}^k$  und testen Sie diese für kleine  $d$  und  $k$ .
- (ii) Setzen Sie dann  $d = 1000$ , erzeugen Sie 300 Punkte in  $\mathbb{R}^d$  zufällig nach einer Verteilung Ihrer Wahl, projizieren Sie diese via  $T_U$  auf  $\mathbb{R}^k$  und berechnen Sie für unterschiedliche  $k$  die größte auftretende Abstandsänderung

$$\varepsilon := \max \left( 1 - \min_{x \neq y} \frac{\|T_U x - T_U y\|}{\|x - y\|}, \max_{x \neq y} \frac{\|T_U x - T_U y\|}{\|x - y\|} - 1 \right).$$

- (iii) Vergleichen Sie in einem Plot die experimentellen Werte für  $\varepsilon$  mit der Schranke für  $\varepsilon$  die das Johnson-Lindenstrauss Lemma liefert. Dies sollte ein Bild wie auf Seite 144 ergeben.

## Kapitel 12

# Trennung hochdimensionaler Gaußiane und Parameteranpassung

Nachdem wir in Kapitel 10 Eigenschaften der Gaußverteilung im hochdimensionalen Raum untersucht haben, werden wir in diesem Kapitel diese Erkenntnisse auf Datenmengen anwenden die von einer oder mehrerer Gaußverteilungen stammen. Für Datenmengen dieser Gestalt führen wir den folgenden Begriff ein.

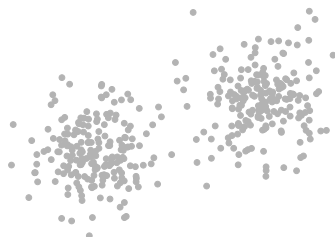
**Definition 12.1.** Sei  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum.

- (i) Sei  $X: \Omega \rightarrow \mathbb{R}^d$  ein Zufallsvektor mit  $X \sim \mathcal{N}(\mu, \sigma^2, \mathbb{R}^d)$  und seien  $\omega_1, \dots, \omega_n \in \Omega$ . Dann nennen wir die Menge  $G = \{X(\omega_i) \mid i = 1, \dots, n\} \subseteq \mathbb{R}^d$  einen (*d-dimensionalen*) *Gaußian*.
- (ii) Seien  $X_1, \dots, X_m: \Omega \rightarrow \mathbb{R}^d$  unabhängige Zufallsvektoren, derart dass  $X_i \sim \mathcal{N}(\mu_i, \sigma_i^2, \mathbb{R}^d)$  und seien weiter Punkte  $\omega_1, \dots, \omega_n \in \Omega$  gegeben. Dann nennen wir die Gaußiane  $G_j = \{X_j(\omega_i) \mid i = 1, \dots, n\}$  für  $j = 1, \dots, m$  *unabhängig*.

Seien jetzt unabhängige Gaußiane  $G_1, \dots, G_m$  gegeben. Wir setzen

$$D := G_1 \cup \dots \cup G_m \subseteq \mathbb{R}^d.$$

Das folgende Beispiel zeigt die Vereinigung  $D = G_1 \cup G_2$  von zwei Gaußianen mit verschiedenen Mittelwerten und jeweils 200 Punkten.



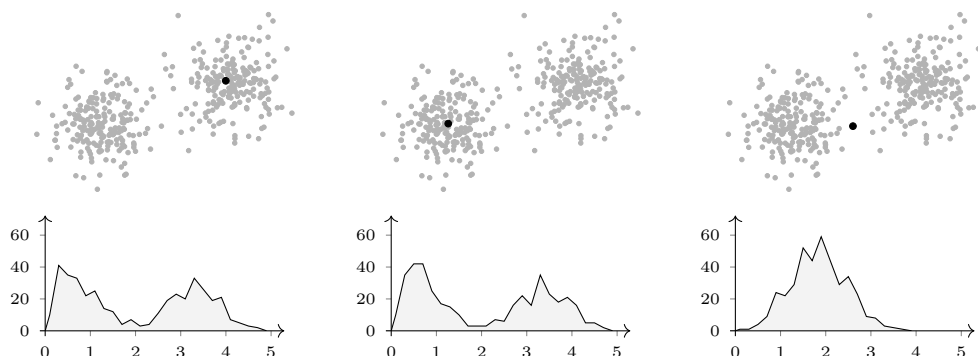
Wie das einfarbige Bild bereits suggeriert, vergessen wir nun, welcher Datenpunkt von welchem Gaußian kommt. Die erste sich natürlich ergebende Aufgabe besteht

darin, diese Information anhand von  $D$  wiederzuerlangen, d.h. die Gaußiane zu *trennen*. Ist dies geschehen, so besteht die zweite natürliche Aufgabe darin, für jeden Gaußian seine Parameter  $\mu_i$  und  $\sigma_i$  zu bestimmen — wobei wir ebenfalls annehmen, dass diese uns nicht a priori bekannt sind. Dies bezeichnet man als *Parameteranpassung*, *Parameterschätzung* oder auch als *Fitting*. Bevor wir mit der Trennungsaufgabe beginnen, notieren wir noch, dass im Fall einer echten Datenmenge  $D$  beide Aufgaben natürlich nicht in einem ‘definitiven’ Sinn gelöst werden können: Bereits im obigen Bild sieht man, dass es durchaus passieren kann, dass ein Punkt, der zwar vom Gaußian links unten stammt trotzdem rechts oben liegt. Jeder sinnvolle Trennungsalgorithmus wird diesen Punkt dann als zum rechten Gaußian zugehörig klassifizieren. Bei gegebenen Daten haben wir keine Möglichkeit diese Fehlklassifikation zu bemerken und streben demnach nicht wirklich eine ‘korrekte’ Klassifizierung sondern eher eine Klassifizierung mit ‘hoher Plausibilität’ an. Wir werden das letztere später, in den Erläuterungen vor Definition 12.9, mithilfe des Begriffs der bedingten Wahrscheinlichkeit noch genauer fassen.

## 12.1 Trennung von Gaußianen

Wir beginnen mit der oben beschriebenen Trennungsaufgabe und beschränken uns im folgenden auf zwei Gaußiane deren Varianzen beide Eins sind, vergleiche aber Bemerkung 12.13 Das Bild auf Seite 146 legt nun die folgende Idee zur Konstruktion eines Klassifizierers sofort nahe: Zuerst wählt man einen Punkt  $z \in D$  zufällig aus, und bemerkt, dass dieser mit hoher Wahrscheinlichkeit in der Nähe eines der Mittelwerte der zwei Gaußiane liegen wird. Dann weist man  $z$  das Label 1 zu und ermittelt für alle weiteren Punkte in  $D$  die Abstände zu  $z$ . Fällt ein solcher Abstand klein aus, so erhält der Punkt ebenfalls Label 1. Andernfalls erhält er Label 2. Ist hierbei der Abstand der Mittelwerte  $\Delta := \|\mu_1 - \mu_2\|$  der Gaußiane groß genug, so ist zu erwarten, dass es nach Berechnung aller Abstände  $\|z - x\|$ ,  $x \in D$ , sinnvoll ist von ‘kleinen’ und ‘großen’ Abständen zu sprechen. Wir betrachten hierzu wieder die auf Seite 146 abgebildete Datenmenge  $D$ .

**Beispiel 12.2.** In den folgenden Bildern ist oben jeweils der Punkt  $z$  markiert. Darunter ist dann die Verteilung der Abstände aller anderen Punkte zu  $z$  abgebildet.



Wie wir oben bereits erwähnt haben, tritt hierbei eine Situation wie ganz rechts, bei der  $z$  ‘auf halber Strecke’ zwischen den Gaußianen liegt, nur mit geringer Wahrscheinlichkeit ein, vergleiche Aufgabe 12.3.

Eine naheliegende Heuristik, nach welcher man nun von den Abständen zum Klassifizierer gelangt, besteht darin, dass man bei  $n$ -vielen Datenpunkten die  $n/2$ -nächsten Nachbarn von  $z$  mit Label 1 versieht und den Rest mit Label 2. Wir formulieren dies nun in Pseudocode.

**Algorithmus 12.3.** *Gegeben sei eine Datenmenge  $D = \{x^{(1)}, \dots, x^{(n)}\} \subseteq \mathbb{R}^d$  die Vereinigung zweier Gaußiane ist, welche jeweils aus  $n/2$ -vielen Punkten bestehen.*

```

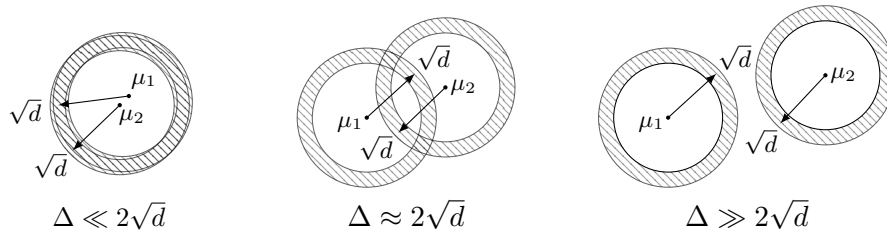
1: function TRENNUNG( $D$ )
2:    $z \leftarrow$  zufälliger Punkt in  $D$ 
3:   for  $i \leftarrow 1$  to  $n$  do
4:      $d_i \leftarrow \|z - x^{(i)}\|$ 
5:    $I \leftarrow (i_1, \dots, i_n)$  mit  $d_{i_j} \leq d_{i_{j+1}}$  für  $j = 1, \dots, n$ 
6:   for  $j \leftarrow 1$  to  $n/2$  do
7:      $\ell_{i_j} \leftarrow 1$ 
8:   for  $j \leftarrow n/2 + 1$  to  $n$  do
9:      $\ell_{i_j} \leftarrow 2$ 
10:  return  $(\ell_1, \dots, \ell_n)$ 

```

Hierbei kann in Zeile 5 für die Sortierung der Datenpunkte nach ihren Abstände von  $z$  ein geeigneter Sortieralgorithmus benutzt werden.

In den Aufgaben werden wir Algorithmus 12.3 sowohl für niedrigdimensionale als auch für hochdimensionale Datenmengen implementieren. Wir notieren, dass wir den Vektor  $\ell$  natürlich auch als Abbildung  $\ell: D \rightarrow \{1, 2\}$  per  $\ell(x^{(i)}) := \ell_i$ , lesen können. Weiter bemerken wir noch, dass es sich bei Obigem um ‘unüberwachtes Lernen’ handelt, da keine gelabelte Trainingsmenge gegeben ist.

Wir wenden uns jetzt der Situation hochdimensionaler Daten zu. Hier wissen wir, dass die Verteilung nicht wie im Bild auf Seite 146 aussieht, sondern dass sich ein Großteil der Datenpunkte eines Gaußians auf einem schmalen Ring mit Radius  $\sqrt{d}$  um den Mittelpunkt  $\mu$  des Gaußians ansammeln wird. Bleiben wir bei zwei Gaußianen, so müssen wir analog zu oben voraussetzen, dass deren Mittelpunkte  $\mu_1$  und  $\mu_2$  weit genug auseinander liegen, damit eine Trennung gelingen kann. Sicher hinreichend wäre es zu verlangen, dass  $\Delta := \|\mu_1 - \mu_2\|$  geeignet größer als  $2\sqrt{d}$  ist — denn dann dürften wir eine ähnliche Verteilung der Abstände zu einem zufällig gewählten Punkt wie in Beispiel 12.2 erwarten. In der Tat kommen wir aber auch mit einem kleineren Abstand aus, solange dieser nicht zu klein ist im Vergleich zu  $\sqrt{d}$ . Die folgenden drei Bilder

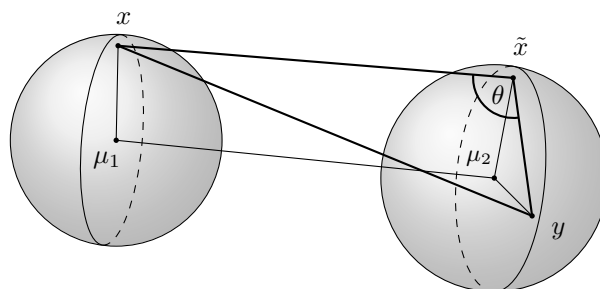


von jeweils zwei Gaußianen, suggerieren schließlich, dass auch falls  $\Delta$  ‘etwas kleiner’ als  $2\sqrt{d}$  ist, die beiden Ringe, in welchen sich der Hauptteil der Punkte der Gaußiane befindet, nur eine relativ kleine Schnittmenge haben. Wie hierbei die Symbole ‘ $\ll$ ’, ‘ $\approx$ ’ und ‘ $\gg$ ’ zu verstehen sind, heben wir in der folgenden Bemerkung hervor.

**Bemerkung 12.4.** Bei der Beschreibung ‘wie groß’ der Abstand  $\Delta = \|\mu_1 - \mu_2\|$  der Mittelpunkte der zwei zu trennenden Gaußiane sein muss, sind wir, analog zu den vorhergehenden Kapiteln 8–10, an Aussagen interessiert, die *für alle geeignet großen Dimensionen* gelten. Wir stellen uns daher im Folgenden  $\Delta = \Delta(d)$  als eine Funktion der Dimension vor.

- (i) Wächst  $\Delta$  im Vergleich zu  $\sqrt{d}$  zu langsam, so geraten wir für große  $d$  notwendigerweise in eine Situation in welcher die Ringgebiete stark überlappen, siehe das ganz linke Bild auf Seite 148. Dies ist mit  $\Delta \ll 2\sqrt{d}$  gemeint.
- (ii) Wächst  $\Delta$  derart, dass für große  $d$  stets  $\Delta(d) > 2\sqrt{d} + \varepsilon$  gilt mit einem festen  $\varepsilon$ , das größer als die nach Satz 10.4 konstante ‘Ringbreite’ ist, vergleiche auch die Diskussion vor Satz 10.4 und Beispiel 10.5, so sollte eine Trennung sicher möglich sein. Auf diese Weise ist  $\Delta \gg 2\sqrt{d}$  zu verstehen.
- (iii) Wächst  $\Delta$  genauso schnell wie  $2\sqrt{d}$ , so suggeriert das mittlere Bild auf Seite 148, dass eine Trennung möglich ist. Wir werden im Folgenden sehen, dass letztere sogar dann gelingt, wenn  $\Delta(d)$  etwas langsamer wächst als  $2\sqrt{d}$ . Dies ist mit  $\Delta \approx 2\sqrt{d}$  gemeint.

Unser Ziel besteht jetzt also darin herauszufinden, wie viel langsamer  $\Delta(d)$  im Vergleich zu  $2\sqrt{d}$  wachsen darf, sodass eine Trennung via des oben bereits diskutierten Algorithmus 12.3 immer noch gelingt. Wir betrachten  $x, y \in D$  aus einer Datenmenge wie dort angegeben. Nach Kapitel 10 erwarten wir, dass  $\|x - y\| \approx \sqrt{2d}$  für Punkte gilt, die zum gleichen Gaußian gehören, und wir müssen uns folglich überlegen, wie  $\|x - y\|$  ausfällt wenn  $x$  und  $y$  zu verschiedenen Gaußianen gehören. In der Tat ergibt sich für solche Punkte  $x$  und  $y$  das folgende Bild.



Denn:

1. Der Punkt  $\tilde{x}$  wird mit hoher Wahrscheinlichkeit nah an der Oberfläche einer Kugel um  $\mu_1$  mit Radius  $\sqrt{d}$  liegen (Ringsatz 10.4) und es wird mit hoher Wahrscheinlichkeit  $x - \mu_1 \perp \mu_2 - \mu_1$  fast gelten (Orthogonalitätssatz 10.7 mit  $Y \equiv \mu_2 - \mu_1$ ).
2. Der Punkt  $\tilde{x} = \mu_2 + x - \mu_1 = x + (\mu_2 - \mu_1)$  entsteht durch Parallelverschiebung, also wird nach obigem mit hoher Wahrscheinlichkeit  $\tilde{x} - \mu_2 \perp \mu_1 - \mu_2$  fast gelten.
3. Der Punkt  $y$  wird mit hoher Wahrscheinlichkeit nah an der Oberfläche einer Kugel um  $\mu_2$  mit Radius  $\sqrt{d}$  liegen (nochmal Ringsatz 10.4). Gleichzeitig wird der Abstand von  $\tilde{x}$  und  $y$  ungefähr  $\sqrt{2d}$  sein (Abstandssatz 10.10) und schließlich wird  $y - \mu_2 \perp \mu_1 - \mu_2$  fast gelten (nochmal Orthogonalitätssatz 10.7).

In dem entstandenen Dreieck ist nach obigem also  $\theta \approx \pi/2$  und es folgt

$$\|x - y\|^2 \approx \|x - \tilde{x}\|^2 + \|\tilde{x} - y\|^2 \approx \Delta^2 + 2d$$

mit dem Satz des Pythagoras. Der folgende Satz formalisiert die angegebene Heuristik mithilfe von Erwartungswert und Varianz, verwendet dabei allerdings das Kosinusetz anstelle des Satzes von Pythagoras.

**Satz 12.5.** Sei  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum, seien  $\mu_1, \mu_2 \in \mathbb{R}^d$  und seien  $X, Y: \Omega \rightarrow \mathbb{R}^d$  unabhängige Zufallsvektoren mit  $X \sim \mathcal{N}(\mu_1, 1, \mathbb{R}^d)$  und  $Y \sim \mathcal{N}(\mu_2, 1, \mathbb{R}^d)$ . Sei  $\Delta = \|\mu_1 - \mu_2\|$ . Dann gelten

- (i)  $E(\|X - Y\|^2) = \Delta^2 + 2d$ ,
- (ii)  $V(\|X - Y\|^2) \leq 3d + 8\Delta^2 + 5\sqrt{d}\Delta$ .

*Beweis.* (i) Seien  $X, Y: \Omega \rightarrow \mathbb{R}^d$  wie oben und sei  $\tilde{X} := \mu_2 + X - \mu_1$ . D.h. es gilt  $\tilde{X} \sim \mathcal{N}(\mu_2, 1, \mathbb{R}^d)$  und  $\tilde{X}$  und  $Y$  sind unabhängig. Wir wenden punktweise das Kosinusetz mit  $\theta = \angle(X - \tilde{X}, Y - \tilde{X})$  an und erhalten

$$\begin{aligned} \|X - Y\|^2 &= \|X - \tilde{X}\|^2 + \|Y - \tilde{X}\|^2 - 2\|X - \tilde{X}\|\|Y - \tilde{X}\|\cos(\theta) \\ &= \|\mu_1 - \mu_2\|^2 + \|Y - \tilde{X}\|^2 - 2\langle X - \tilde{X}, Y - \tilde{X} \rangle \\ &= \Delta^2 + \|Y - \tilde{X}\|^2 - 2(\langle \mu_1 - \mu_2, Y - \mu_2 \rangle + \langle \mu_2 - \mu_1, X - \mu_1 \rangle). \end{aligned}$$

Wenn wir nun um  $\mu_2$  verschieben, liefert (8.1) auf Seite 114, dass  $E(\|Y - \tilde{X}\|^2) = 2d$  ist. Entsprechende Verschiebung und Anwendung von Satz 8.4 zeigt weiter, dass die Erwartungswerte der beiden obigen Skalarprodukte jeweils Null sind. Zusammengekommen erhalten wir also

$$E(\|X_1 - X_2\|^2) = 2d + \Delta^2 + 2(0 - 0) = \Delta^2 + 2d.$$

(ii) Bei der Abschätzung der Varianz von  $\|X - Y\|^2$  gilt es zu beachten, dass die drei sich ergebenden Summanden nicht unabhängig sind, siehe Proposition A.6. Daher

gilt

$$\begin{aligned}
V(\|X - Y\|^2) &= V(\Delta^2 + \|Y - \tilde{X}\|^2 - 2\langle X - \tilde{X}, Y - \tilde{X} \rangle) \\
&= V(\|Y - \tilde{X}\|^2) + 4V(\langle X - \tilde{X}, Y - \tilde{X} \rangle) \\
&\quad + \text{Cov}(\|Y - \tilde{X}\|^2 - 2\langle X - \tilde{X}, Y - \tilde{X} \rangle) \\
&\leq V(\|Y - \tilde{X}\|^2) + 4V(\langle X - \tilde{X}, Y - \tilde{X} \rangle) \\
&\quad + \sqrt{V(\|Y - \tilde{X}\|^2) \cdot 4V(\langle X - \tilde{X}, Y - \tilde{X} \rangle)}
\end{aligned}$$

wobei wir  $V(\|Y - \tilde{X}\|^2) \leq 3d$  gemäß dem Hinweis in Aufgabe 8.3 abschätzen können. Für die Varianz des Skalarproduktes ergibt sich

$$\begin{aligned}
V(\langle X - \tilde{X}, Y - \tilde{X} \rangle) &= V(\langle X - \mu_1, \mu_2 - \mu_1 \rangle + \langle Y - \mu_2, \mu_1 - \mu_2 \rangle) \\
&= V(\langle X - \mu_1, \mu_2 - \mu_1 \rangle) + V(\langle Y - \mu_2, \mu_1 - \mu_2 \rangle) \\
&\leq 2\|\mu_1 - \mu_2\|^2
\end{aligned}$$

wobei wir Satz 8.4 mit  $Z := X_i - \mu_i \sim \mathcal{N}(0, 1, \mathbb{R}^d)$  und  $\xi := \pm(\mu_1 - \mu_2)$  angewandt haben um  $V(\langle Z, \xi \rangle) = \|\xi\|^2$  zu schließen. Schließlich erhalten wir

$$\begin{aligned}
V(\|X - Y\|^2) &\leq V(\|Y - \tilde{X}\|^2) + 4V(\langle X - \tilde{X}, Y - \tilde{X} \rangle) + \sqrt{\cdots} \\
&\stackrel{\text{s.o.}}{\leq} 3d + 8\|\mu_1 - \mu_2\|^2 + \sqrt{3d \cdot 8\|\mu_1 - \mu_2\|^2} \\
&\leq 3d + 8\Delta^2 + 5\sqrt{d}\Delta
\end{aligned}$$

wie behauptet. □

Mit ein bisschen mehr Arbeit erhalten wir Abschätzungen für Erwartungswert und Varianz des Abstandes selbst, d.h. ohne das Quadrat. Im Fall der Varianz sind diese eher technisch und wir formulieren daher im Satz für die Varianz nur die sich ergebende qualitative asymptotische Aussage.

**Satz 12.6.** Sei  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum, seien  $\mu_1, \mu_2 \in \mathbb{R}^d$  und seien  $X, Y: \Omega \rightarrow \mathbb{R}^d$  unabhängige Zufallsvektoren mit  $X \sim \mathcal{N}(\mu_1, 1, \mathbb{R}^d)$  und  $Y \sim \mathcal{N}(\mu_2, 1, \mathbb{R}^d)$ . Sei  $\Delta = \|\mu_1 - \mu_2\|$ . Dann gelten

- (i)  $\forall d \in \mathbb{N}: |\mathbb{E}(\|X - Y\|) - \sqrt{\Delta^2 + 2d}| \leq 5/\sqrt{d}.$
- (ii)  $\forall d \in \mathbb{N}: V(\|X - Y\|) \leq 5/\sqrt{d} + 14 \leq 19.$

*Beweis.* (i) Wir gehen wie im Beweis von Satz 8.2 vor und zerlegen den Term, von welchem der Erwartungswert gesucht ist, wie folgt

$$\|X - Y\| - \sqrt{\Delta^2 + 2d} = \underbrace{\frac{\|X - Y\|^2 - (\Delta^2 + 2d)}{2\sqrt{\Delta^2 + 2d}}}_{=: S_d} - \underbrace{\frac{(\|X - Y\|^2 - (\Delta^2 + 2d))^2}{2\sqrt{\Delta^2 + 2d}(\|X - Y\| + \sqrt{\Delta^2 + 2d})^2}}_{=: R_d}.$$

Nach Satz 12.5(i) ist  $E(S_d) = 0$ . In  $R_d$  schätzen wir im Nenner  $\|X - Y\| \geq 0$  ab und im Zähler wenden wir erst Satz 12.5(i) an, um den Erwartungswert in eine Varianz umzuschreiben und dann Satz 12.5(ii) um diese abzuschätzen. Dies führt auf

$$0 \leq E(R_d) \leq \frac{V(\|X - Y\|^2)}{2(\Delta^2 + 2d)^{3/2}} \leq \frac{3d + 8\Delta^2 + 5\sqrt{d}\Delta}{2(\Delta^2 + 2d)^{3/2}} =: r_d$$

was wir jetzt per Fallunterscheidung behandeln. Sei zunächst  $\Delta \geq \sqrt{2d}$ . Dann folgt

$$r_d = \frac{3d + 8\Delta^2 + 5\sqrt{d}\Delta}{2(\Delta^2 + 2d)^{3/2}} \leq \frac{\frac{3}{2}\Delta^2 + 8\Delta^2 + \frac{5}{\sqrt{2}}\Delta^2}{2(\Delta^2)^{3/2}} \leq \frac{5}{\sqrt{d}}.$$

$\uparrow$   
 $0 \leq d \leq \frac{1}{2}\Delta^2$

Ist andererseits  $\Delta < \sqrt{2d}$  so folgt

$$r_d = \frac{3d + 8\Delta^2 + 5\sqrt{d}\Delta}{2(\Delta^2 + 2d)^{3/2}} \leq \frac{3d + 16d + 5\sqrt{2}d}{2(2d)^{3/2}} \leq \frac{5}{\sqrt{d}}$$

$\uparrow$   $\uparrow$   
 $0 \leq \Delta^2 \leq 2d$   $1 \leq 2^{3/2}$

und somit  $E(R_d) \leq \frac{5}{\sqrt{d}}$  für alle  $d \in \mathbb{N}$ .

(ii) Um die Varianz zu behandeln, berechnen wir zuerst

$$\begin{aligned} V(\|X - Y\|) &= |E(\|X - Y\|^2) - [E(\|X - Y\|)]^2| \\ &\stackrel{\text{Satz 12.5(i)}}{=} |(\Delta^2 + 2d) - [E(\|X - Y\| - \sqrt{\Delta^2 + 2d}) + \sqrt{\Delta^2 + 2d}]^2| \\ &= |(\Delta^2 + 2d) - [(E(\|X - Y\| - \sqrt{\Delta^2 + 2d}))^2 \\ &\quad + 2\sqrt{\Delta^2 + 2d} \cdot E(\|X - Y\| - \sqrt{\Delta^2 + 2d}) + (\Delta^2 + 2d)]| \\ &\leq r_d + 2\sqrt{\Delta^2 + 2d} \cdot r_d. \\ &\quad \uparrow \\ &\quad \text{s.o.} \end{aligned}$$

Der erste Term kann durch  $5/\sqrt{d}$  abgeschätzt werden, den zweiten Term behandeln wir wieder via Fallunterscheidung. Ist  $\Delta \geq \sqrt{2d}$ , so gilt

$$2\sqrt{\Delta^2 + 2d} \cdot r_d = \frac{3d + 8\Delta^2 + 5\sqrt{d}\Delta}{\Delta^2 + 2d} \leq \frac{\frac{3}{2}\Delta^2 + 8\Delta^2 + \frac{5}{\sqrt{2}}\Delta^2}{\Delta^2} \leq 14,$$

$\uparrow$   
 $0 \leq d \leq \frac{1}{2}\Delta^2$

und ist  $\Delta < \sqrt{2d}$  so haben wir

$$2\sqrt{\Delta^2 + 2d} \cdot r_d = \frac{3d + 8\Delta^2 + 5\sqrt{d}\Delta}{\Delta^2 + 2d} \leq \frac{3d + 16d + 5\sqrt{2}d}{2d} \leq 14$$

$\uparrow$   
 $0 \leq \Delta^2 \leq 2d$

was den Beweis beendet. □



Satz 12.6 sagt uns, dass wir zwischen zufällig gewählten Punkten von verschiedenen Gaußianen einen Abstand von  $\|x - y\| \approx \sqrt{\Delta^2 + 2d}$  zu erwarten haben und dass für große  $d$  sich die Verteilung dieser Abstände, unabhängig von  $d$ , nah um diesen Wert konzentriert. Wir formalisieren dies weiter durch eine explizite Wahrscheinlichkeitsabschätzung.

**Satz 12.7.** (Zweiter Gaußscher Abstandssatz) *Sei  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum, seien  $\mu_1, \mu_2 \in \mathbb{R}^d$  und seien  $X, Y: \Omega \rightarrow \mathbb{R}^d$  unabhängige Zufallsvektoren mit  $X \sim \mathcal{N}(\mu_1, 1, \mathbb{R}^d)$  und  $Y \sim \mathcal{N}(\mu_2, 1, \mathbb{R}^d)$ . Sei  $\Delta = \|\mu_1 - \mu_2\|$ . Für  $\varepsilon \leq \sqrt{d}$  gilt dann*

$$P[|\|X - Y\| - \sqrt{\Delta^2 + 2d}| \leq \varepsilon] \geq 1 - 4e^{-\varepsilon^2/648} - \frac{26}{\varepsilon d} - \frac{7}{\sqrt{d}} - \frac{12}{\varepsilon}.$$

*Beweis.* Wir verwenden dieselbe Technik wie im Beweis des Ringsatzes 10.4, d.h. wir schätzen

$$P := P[|\|X - Y\| - \sqrt{\Delta^2 + 2d}| \geq \varepsilon]$$

nach oben ab, indem wir die Abschätzung in  $P[\dots]$  mit der trivialen Ungleichung  $\|X - Y\| + \sqrt{\Delta^2 + 2d} \geq \sqrt{\Delta^2 + 2d}$  multiplizieren. Dann führen wir, wie im Beweis von Satz 12.5, die Zufallsvariable  $\tilde{X} := \mu_2 + X - \mu_1 \sim \mathcal{N}(\mu_2, 1, \mathbb{R}^d)$  ein und wenden wie dort geschehen das Kosinusetz an. Dies liefert

$$\begin{aligned} P &\leq P[|\|X - Y\|^2 - (\Delta^2 + 2d)| \geq \varepsilon \sqrt{\Delta^2 + 2d}] \\ &= 1 - P[-\varepsilon \sqrt{\Delta^2 + 2d} \leq \Delta^2 + \|Y - \tilde{X}\|^2 - 2(\langle \mu_1 - \mu_2, Y - \mu_2 \rangle \\ &\quad + \langle \mu_2 - \mu_1, X - \mu_1 \rangle) - (\Delta^2 + 2d) \leq \varepsilon \sqrt{\Delta^2 + 2d}] \\ &\leq 1 - P\left[-\frac{\varepsilon \sqrt{\Delta^2 + 2d}}{3} \leq \|Y - \tilde{X}\|^2 - 2d \leq \frac{\varepsilon \sqrt{\Delta^2 + 2d}}{3}\right] \\ &\quad \cdot P\left[-\frac{\varepsilon \sqrt{\Delta^2 + 2d}}{3} \leq -2\langle \mu_1 - \mu_2, Y - \mu_2 \rangle \leq \frac{\varepsilon \sqrt{\Delta^2 + 2d}}{3}\right] \\ &\quad \cdot P\left[-\frac{\varepsilon \sqrt{\Delta^2 + 2d}}{3} \leq -2\langle \mu_2 - \mu_1, X - \mu_1 \rangle \leq \frac{\varepsilon \sqrt{\Delta^2 + 2d}}{3}\right] \\ &\leq 1 - P[|\|Y - \tilde{X}\|^2 - 2d| \leq \frac{\varepsilon \sqrt{2d}}{3}] \cdot P[|\langle Z, \xi \rangle| \leq \frac{\varepsilon \sqrt{\Delta^2}}{6}]^2 \end{aligned}$$

wobei  $Z := X_i - \mu_i \sim \mathcal{N}(0, 1, \mathbb{R}^d)$  und  $\xi := \mu_1 - \mu_2$ . Die zwei Wahrscheinlichkeiten am Ende der obigen Rechnung können wir mithilfe von Korollar 10.11 per

$$P[|\|Y - \tilde{X}\|^2 - 2d| \leq \frac{\varepsilon \sqrt{2d}}{3}] \geq 1 - 4e^{-\varepsilon^2/648} - \frac{18\sqrt{2}}{\varepsilon d} - \frac{7}{\sqrt{d}}$$

und mit Satz 10.8 durch

$$P[|\langle Z, \xi \rangle| \leq \frac{\varepsilon \sqrt{\Delta^2}}{6}] \geq 1 - \frac{6\|\xi\|}{\varepsilon \Delta} = 1 - \frac{6}{\varepsilon}$$

abschätzen. Da uns eigentlich  $1 - P$  interessiert, betrachten wir nun

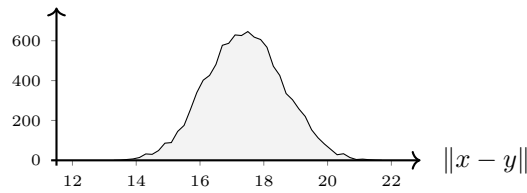
$$1 - P \underset{\text{s.o.}}{\geq} P[|\|Y - \tilde{X}\|^2 - 2d| \leq \frac{\varepsilon \sqrt{2d}}{3}] \cdot P[|\langle Z, \xi \rangle| \leq \frac{\varepsilon \sqrt{\Delta^2}}{6}]^2$$

$$\begin{aligned}
&\geq \left(1 - 4e^{-\varepsilon^2/648} - \frac{18\sqrt{2}}{\varepsilon d} - \frac{7}{\sqrt{d}}\right) \left(1 - \frac{6}{\varepsilon}\right)^2 \\
&\geq 1 - 4e^{-\varepsilon^2/648} - \frac{26}{\varepsilon d} - \frac{7}{\sqrt{d}} - \frac{12}{\varepsilon}
\end{aligned}$$

wobei wir für die letzte Ungleichung erst ausmultipliziert und dann alle positiven Summanden nach unten mit Null abgeschätzt haben.  $\square$

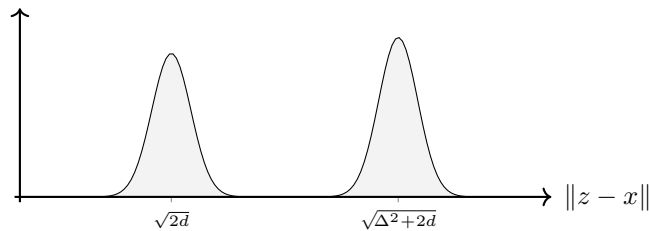
Wir weisen darauf hin, dass  $\varepsilon > 100$  notwendig dafür ist, dass die untere Schranke in der Wahrscheinlichkeitsabschätzung in Satz 12.7 positiv wird aufgrund der zwei von der Dimension unabhängigen Terme auf der rechten Seite. Ähnlich wie in Beispiel 10.9 zum Orthogonalitätssatz, betrachten wir nun eine Simulation und vergleichen diese mit Satz 12.7.

**Beispiel 12.8.** Das folgende Bild zeigt die Verteilung der Abstände  $\|x - y\|$  für  $(x, y) \in G_1 \times G_2$ . Dabei ist  $G_1$  ein Gaußian in  $\mathbb{R}^{100}$  mit Varianz Eins und Mittelpunkt  $\mu_1 = 0$  ist und  $G_2$  ein Gaußian in  $\mathbb{R}^{100}$  mit Varianz Eins und Mittelpunkt  $\mu_2 = (\Delta, 0, \dots, 0)$  und  $\Delta = 10$ .



Es ergibt sich  $\sqrt{\Delta^2 + 2d} \approx 17.32$  und in der Tat gilt in unserer Simulation für *alle* Paare  $(x, y) \in G_1 \times G_2$  die Abschätzung  $|\|x - y\| - \sqrt{\Delta^2 + 2d}| < 5$ .

Haben wir wie in Algorithmus 12.3 einen Punkt  $z \in D = G_1 \cup G_2$  zufällig gewählt, so wissen wir, dass sich die Abstände  $\|z - x\|$  für  $x \in D$  in der Nähe von  $\sqrt{2d}$  und  $\sqrt{\Delta^2 + 2d}$  konzentrieren werden. Ist  $\Delta$  dabei groß genug, so erwarten wir ein Bild der folgenden Form.



Anstatt nun, wie in Algorithmus 12.3 geschehen, erst alle Punkte Ihren Abständen nach zu sortieren, kann man in hohen Dimensionen auch alle  $x$  mit  $\sqrt{2d} - \varepsilon \leq \|z - x\| \leq \sqrt{2d} + \varepsilon$  mit Label 1 versehen und dabei  $\varepsilon$  durch Ausprobieren ermitteln oder sogar ganz *unabhängig von der Datenmenge* schätzen, indem man  $\varepsilon$  z.B. so wählt, dass

$$\begin{aligned}
P[|\|X - Y\| - \sqrt{2d}| \leq \varepsilon] &\geq 1 - 4e^{-\varepsilon^2/72} - \frac{6\sqrt{2}}{\varepsilon d} - \frac{7}{\sqrt{d}} \geq 0.9 \\
&\quad \uparrow \\
&\text{Satz 10.10}
\end{aligned}$$

gilt. Hierfür genügt die Kenntnis von  $d$  und keine weitere Information über  $D$  ist vonnöten. Wir belassen es als Aufgabe 12.4, diese Varianten des Trennungsalgorithmus zu formulieren und zu testen.

Klar ist, dass wir mit jeder Version des Trennungsalgorithmus nur dann eine erfolgreich sein können, wenn  $\Delta := \|\mu_1 - \mu_2\|$  groß genug ist. Die Frage nach dem ‘wie groß’ beantworten wir nun zuerst durch eine asymptotische Wachstumsbedingung an  $\Delta = \Delta(d)$  aufgefasst als *Funktion der Dimension*. Dafür müssen wir die *Chance* quantifizieren, mit der zwei Punkte vom gleichen Gaußian stammen, *unter der Annahme*, dass ihr Abstand nah bei  $\sqrt{2d}$  liegt. In Simulationen kann die entsprechende *Häufigkeit* bestimmt werden indem man bei der Vereinigung zweier Gaußiane  $D = G_1 \cup G_2$  den folgenden Quotient

$$L := \frac{\#\{(x, y) \in D \times D \mid \left| \|x - y\| - \sqrt{2d} \right| < \varepsilon \wedge (x, y \in G_1 \vee x, y \in G_2)\}}{\#\{(x, y) \in D \times D \mid \left| \|x - y\| - \sqrt{2d} \right| < \varepsilon\}}$$

bildet. Betrachtet man für eine feste Datenmenge  $D = G_1 \cup G_2$  wie oben zwei diskrete (!) und unabhängige Zufallsvektoren  $A, B: \Omega \rightarrow \mathbb{R}^d$  die auf  $D$  gleichverteilt sind, also  $A, B \sim \mathcal{U}(G_1 \cup G_2)$ , und notiert man für  $x, y \in \mathbb{R}^d$  als Abkürzung  $x \sim y$ , falls  $x$  und  $y$  beide zu  $G_1$  oder beide zu  $G_2$  gehören, so ergibt sich, dass

$$L = \frac{\text{mit } \#D^2 \text{ erweitern} \uparrow \text{P}[\left| \|A - B\| - \sqrt{2d} \right| < \varepsilon \wedge A \sim B]}{\text{P}[\left| \|A - B\| - \sqrt{2d} \right| < \varepsilon]} = \text{P}[A \sim B \mid \left| \|A - B\| - \sqrt{2d} \right| < \varepsilon]$$

die bedingte Wahrscheinlichkeit dafür ist, dass zwei aus  $D$  gleichmäßig zufällig gewählte Punkte zum selben Gaußian gehören, gegeben dass deren Abstand  $\varepsilon$ -nah an  $\sqrt{2d}$  liegt. Mit  $\#G_1 = \#G_2$  folgt weiterhin

$$\begin{aligned} L &= \frac{\text{Bayes} \uparrow \text{P}[\left| \|A - B\| - \sqrt{2d} \right| < \varepsilon \mid A \sim B] \cdot \text{P}[A \sim B]}{\text{P}[\left| \|A - B\| - \sqrt{2d} \right| < \varepsilon]} \\ &= \frac{\text{totale Wahrsch.} \uparrow \text{P}[\left| \|A - B\| - \sqrt{2d} \right| < \varepsilon \mid A \sim B] \cdot \frac{1}{2}}{\text{P}[\left| \|X - Y\| - \sqrt{2d} \right| < \varepsilon \mid A \sim B] \cdot \frac{1}{2} + \text{P}[\left| \|A - B\| - \sqrt{2d} \right| < \varepsilon \mid A \not\sim B] \cdot \frac{1}{2}}. \end{aligned} \quad (12.1)$$

Im letzteren Ausdruck können wir jede der drei bedingten Wahrscheinlichkeiten für Werte der diskreten Zufallsvektoren  $A$  und  $B$ , inklusive des Faktors  $1/2$ , als (unbedingte) Wahrscheinlichkeit über gaußsche Zufallsvektoren ausdrücken, wenn wir noch berücksichtigen, dass  $x \not\sim y$  für  $x, y \in D$  gerade bedeutet, dass  $x \in G_1$  und  $y \in G_2$  oder umgekehrt gilt.

**Definition 12.9.** Wir definieren die Funktion  $L: \mathbb{N} \times [0, \infty) \times (0, \infty) \rightarrow [0, 1]$  per

$$L(d, \Delta, \varepsilon) := \frac{\text{P}[\left| \|X_1 - Y_1\| - \sqrt{2d} \right| < \varepsilon]}{\text{P}[\left| \|X_1 - Y_1\| - \sqrt{2d} \right| < \varepsilon] + \text{P}[\left| \|X_2 - Y_2\| - \sqrt{2d} \right| < \varepsilon]}$$

wobei  $X_1, Y_1 \sim \mathcal{N}(0, 1, \mathbb{R}^d)$  und  $X_2 \sim \mathcal{N}(\mu_1, 1, \mathbb{R}^d)$ ,  $Y_2 \sim \mathcal{N}(\mu_2, 1, \mathbb{R}^d)$  unabhängige Zufallsvektoren mit  $\Delta = \|\mu_1 - \mu_2\|$  sind.

Wir weisen darauf hin, dass wir oben für jedes  $d$  neue Zufallsvektoren mit den angegebenen Eigenschaften betrachten. Wir dürfen aber annehmen, dass diese alle auf einem einzigen Wahrscheinlichkeitsraum  $(\Omega, \Sigma, P)$  definiert sind, vergleiche Satz A.13.

**Satz 12.10.** (Asymptotischer Trennungssatz) *Sei  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum. Für jedes  $d \geq 1$  seien  $\mu_1, \mu_2 \in \mathbb{R}^d$ , und unabhängige Zufallsvektoren  $X_1, Y_1, X_2, Y_2: \Omega \rightarrow \mathbb{R}^d$  mit  $X_1, Y_1 \sim \mathcal{N}(0, 1, \mathbb{R}^d)$  und  $X_2 \sim \mathcal{N}(\mu_1, 1, \mathbb{R}^d)$ ,  $Y_2 \sim \mathcal{N}(\mu_2, 1, \mathbb{R}^d)$  gegeben. Weiter sei  $\Delta(d) := \|\mu_1 - \mu_2\|$  (gelesen als Funktion der Dimension!) derart dass  $\sqrt{(\Delta(d))^2 - 2d} - \sqrt{2d} \rightarrow \infty$  für  $d \rightarrow \infty$  gilt. Dann haben wir für festes  $\varepsilon > 0$*

$$\liminf_{d \rightarrow \infty} L(d, \Delta(d), \varepsilon) \geq 1 - 4e^{-\varepsilon^2/72}$$

wobei  $L$  die Funktion aus Definition 12.9 ist.

*Beweis.* Sei  $\varepsilon > 0$  fest und sei  $d_0 \in \mathbb{N}$  derart, dass  $\varepsilon < \sqrt{d_0}$  gilt. Mithilfe des ersten Abstandssatzes 10.10 können wir den Zähler von  $L$ , unabhängig von  $\Delta$ , dann durch

$$P[|\|X_1 - Y_1\| - \sqrt{2d}| < \varepsilon] \geq 1 - 4e^{-\varepsilon^2/72} - \frac{6\sqrt{2}}{\varepsilon d} - \frac{7}{\sqrt{d}}$$

nach unten abschätzen. Im Nenner schätzen wir den ersten Summanden nach oben mit Eins ab. Dabei verschenken wir für große  $d$ , wieder aufgrund von Satz 10.10 und unabhängig von  $\Delta$ , nicht viel. Wir setzen

$$\delta(d) := \min(\sqrt{d}, \sqrt{\Delta^2 + 2d} - \sqrt{2d} - \varepsilon) \xrightarrow{d \rightarrow \infty} \infty$$

was per Voraussetzung gegen unendlich geht. Durch eventuelle Vergrößerung von  $d_0$  erreichen wir  $\delta(d) > 0$  für  $d \geq d_0$ . Für jedes solche  $d$  gilt dann  $\sqrt{2d} + \varepsilon \leq \sqrt{\Delta^2 + 2d} - \delta(d)$  und wir erhalten für den zweiten Summanden im Nenner von  $L$  die Abschätzung

$$\begin{aligned} P[|\|X_2 - Y_2\| - \sqrt{2d}| < \varepsilon] &= 1 - P[|\|X_2 - Y_2\| - \sqrt{2d}| \geq \varepsilon] \\ &\leq 1 - P[|\|X_2 - Y_2\| - \sqrt{\Delta^2 + 2d}| \leq \delta(d)] \\ &\leq 4e^{-\delta(d)^2/648} - \frac{26}{\delta(d)d} - \frac{7}{\sqrt{d}} - \frac{12}{\delta(d)} \end{aligned}$$

mithilfe des zweiten Abstandssatzes 12.7. Zusammen folgt also für  $d \geq d_0$

$$L(d, \Delta, \varepsilon) \geq \frac{1 - 4e^{-\varepsilon^2/72} - \frac{6\sqrt{2}}{\varepsilon d} - \frac{7}{\sqrt{d}}}{1 + 4e^{-\delta(d)^2/648} - \frac{26}{\delta(d)d} - \frac{7}{\sqrt{d}} - \frac{12}{\delta(d)}} \xrightarrow{d \rightarrow \infty} 1 - 4e^{-\varepsilon^2/72}$$

und damit die Aussage über den Limes Inferior.  $\square$

Die folgende Proposition gibt explizit an, welche Wachstumsraten für  $\Delta = \Delta(d)$  in Satz 12.10 zugelassen sind.

**Proposition 12.11.** *Für  $\Delta: \mathbb{N} \rightarrow [0, \infty)$  sind die folgenden Aussagen äquivalent.*

- (i)  $\lim_{d \rightarrow \infty} \frac{\Delta(d)}{\sqrt[4]{d}} = \infty.$
- (ii)  $\lim_{d \rightarrow \infty} \sqrt{\Delta(d)^2 - 2d} - \sqrt{2d} = \infty.$

*Beweis.* (i)  $\implies$  (ii) Quadrieren in Bedingung (ii) führt zu

$$\forall R > 0 \exists d_0 \in \mathbb{N} \forall d \geq d_0: \Delta(d)^2 \geq R \cdot \sqrt{d}.$$

Ist letzteres gegeben, so folgt (ii) aus

$$\sqrt{\Delta(d)^2 - 2d} - \sqrt{2d} \geq \sqrt{R\sqrt{d} - 2d} \xrightarrow{d \rightarrow \infty} \frac{R}{2\sqrt{2}}.$$

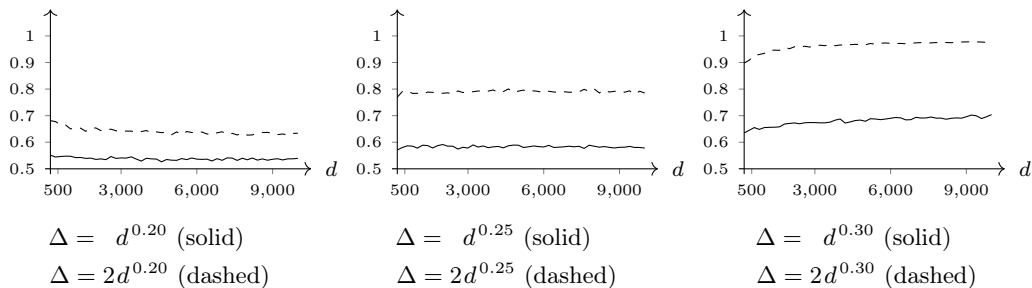
(ii)  $\implies$  (i) Multiplikation mit  $\sqrt{\Delta(d)^2 - 2d} + \sqrt{2d}$  in Bedingung (ii) liefert

$$\forall R > 0 \exists d_0 \in \mathbb{N} \forall d \geq d_0: \Delta(d)^2 \geq R \cdot (\sqrt{\Delta(d)^2 - 2d} + \sqrt{2d}) \geq R \cdot \sqrt{d},$$

woraus per Division durch  $\sqrt{d}$  und nachfolgendes Wurzelziehen (i) gezeigt ist.  $\square$

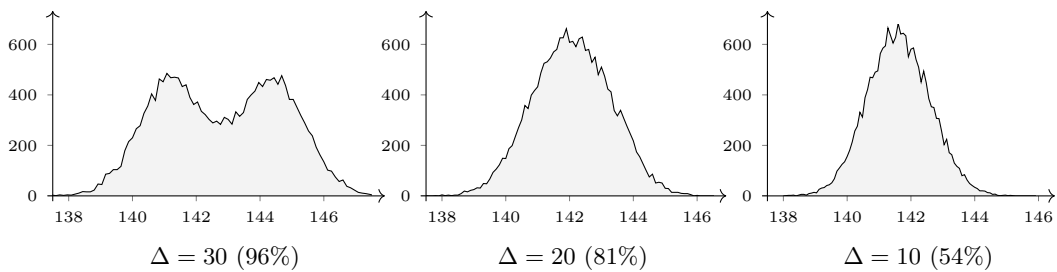
In Kombination besagen Satz 12.10 und Proposition 12.11, dass für  $d \rightarrow \infty$  eine Trennung zweier Gaußianen möglich ist, falls der Abstand  $\Delta$  der Mittelpunkte echt schneller wächst  $\sqrt[4]{d}$ . Im folgenden Beispiel werden wir letzteres illustrieren, und dabei aber auch die Schwächen dieser lediglich asymptotischen Aussage aufzeigen.

**Beispiel 12.12.** (i) Im folgenden betrachten wir den Abstand  $\Delta = \Delta(d)$  als Funktion der Dimension. Wir betrachten also eine Folge von Datenmengen  $D_d \subseteq \mathbb{R}^d$  die jeweils aus der Vereinigung von zwei Gaußianen bestehen, deren Mittelpunkte den Abstand  $\Delta = cd^\alpha$  haben für verschiedene  $\alpha > 0$  und  $c > 0$ . Konkret verwenden wir unten die Dimensionen  $d = 200, 400, 600, \dots, 10000$  und Gaußiane  $G_1, G_2$  mit jeweils 100 Punkten um 0 bzw. um  $(\Delta, 0, \dots, 0)$ . Nach Wahl eines zufälligen Punktes  $z \in G_1 \cup G_2$  weisen wir den 100-nächsten Nachbarn von  $z$  das Label 1 und den restlichen Punkten das Label 2 zu. Danach berechnen wir den Anteil der korrekt klassifizierten Punkte. Um die Genauigkeit zu erhöhen führen wir die Simulation pro Dimension 100-mal durch und plotten dann die Mittelwerte der korrekt klassifizierten Anteile. Dies liefert die folgenden Bilder.



Wir machen insbesondere darauf aufmerksam, dass die für den Asymptotischen Trennungssatz 12.10 irrelevante Konstante 2 im rechten Bild für die betrachteten Dimensionen  $d$  sehr wohl die Güte der Trennung deutlich verbessert. Um ohne die Konstante eine annähernd gute Trennung zu erreichen, muss man beim Exponent 0.30 die Dimension noch sehr viel weiter erhöhen, vergleiche Aufgabe 12.5. Gleichzeitig sieht man aber auch, dass bei  $\Delta = cd^{0.25}$ , sowohl für  $c = 1$  und  $c = 2$ , die Grenze der Trennbarkeit liegt.

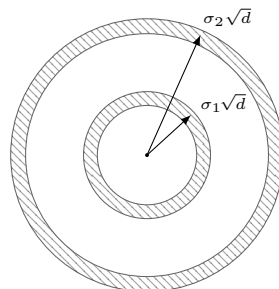
(ii) Jetzt fixieren wir die Dimension  $d = 10000$  und führen hier das Trennungsexperiment dreimal durch mit  $\Delta = 30, 20, 10$ . Die folgenden Bilder zeigen die Verteilung der paarweisen Abstände der 200 Datenpunkte. In Klammern dahinter ist angegeben, wieviele Punkte durch den Trennungsalgorithmus 12.3 korrekt klassifiziert wurden.



Obwohl die zwei ‘Buckel’ in der Verteilung, im Vergleich zur abstrakten Zeichnung auf Seite 154, bereits bei  $\Delta = 30$  teilweise ineinanderübergehen und bei  $\Delta = 20$  bereits zu einem ‘Buckel’ verschmolzen sind, liefert der Trennungsalgorithmus 12.3 ganz links sehr gute und in der Mitte immerhin noch brauchbare Ergebnisse. Rechts, bei  $\Delta = 10 = \sqrt[4]{d}$ , ist der Abstand der Gaußianen dann allerdings zu klein, um ihre Vereinigung wieder erfolgreich zu trennen.

Wir notieren die folgenden Bemerkungen zum Fall, dass die Varianzen der Gaußianen nicht beide Eins sind, und zur Behandlung von mehr als zwei Gaußianen.

**Bemerkung 12.13.** (i) Für  $X \sim \mathcal{N}(\mu, \sigma^2, \mathbb{R}^d)$  kann man leicht nachrechnen, dass  $E(\|X\|^2) = \sigma^2 d$  gilt. Es liegt daher nicht fern anzunehmen, dass unser Trennungsalgorithmus 12.3 auch für eine Vereinigung  $D = G_1 \cup G_2$  von Gaußianen mit Varianzen ungleich Eins gelten wird. Sind die Varianzen von  $G_1$  und  $G_2$  verschieden und ist  $d$  groß genug, so kann man sogar erwarten, dass sich Gaußianen trennen lassen auch wenn ihre Mittelpunkte  $\mu_1 \approx \mu_2$  nah beieinander liegen. Selbst wenn die Mittelpunkte gleich sind, suggeriert der Ringsatz, dass die Punkte sich auf zwei schmale Ringe verteilen, die jetzt zwar konzentrisch sind, aber unterschiedliche Radien haben.



Wir überlassen es dem Leser als Aufgabe 12.6 dies experimentell zu überprüfen.

(ii) Besteht eine gegebene Datenmenge  $D = G_1 \cup \dots \cup G_m$  aus mehreren Gaußianen mit jeweils  $n$  Punkten pro Gaußian, so kann man die Trennung versuchen, indem man zunächst zu einem zufälligen  $z \in D$  die  $n$ -nächsten Nachbarn mit 1 labelt. Dann entfernt man diese Punkte aus  $D$ , wählt einen neuen Punkt  $z$  und labelt dessen  $n$ -nächste Nachbarn mit 2, usw.

## 12.2 Parameteranpassung für Gaußiane

Nach der Trennung von Gaußianen besteht die zweite natürliche Aufgabe darin, für jeden einzelnen Gaußian seinen Mittelpunkt und seine Varianz zu schätzen. Wir betrachten hierzu der Einfachheit halber nur die 1-dimensionale Situation, also einen Gaußian

$$G = \{x_1, \dots, x_n\} \subseteq \mathbb{R}, \quad (12.2)$$

bei dem die  $x_i$  Realisierungen einer gaußverteilten Zufallsvariable  $X: \Omega \rightarrow \mathbb{R}$  auf einem Wahrscheinlichkeitsraum  $(\Omega, \Sigma, P)$  sind, und bei der wir Mittelwert und Varianz nicht kennen. Es ist dann naheliegend, es mit Mittelwert und Varianz der Datenpunkte

$$\mu_s := \frac{1}{n} \sum_{i=1}^n x_i \quad \text{und} \quad \sigma_s^2 := \frac{1}{n} \sum_{i=1}^n (x_i - \mu_s)^2$$

als Schätzung zu versuchen. Das ‘s’ steht hierbei für *Stichprobenmittelwert* bzw. für *Stichprobenvarianz*. Um zu formalisieren, in welchem Sinne dies tatsächlich eine gute Schätzung ist, verwenden wir die Maximum-Likelihood-Methode. Ist eine Menge (12.2) gegeben, so definieren wir hierzu die Likelihood-Funktion

$$L: \mathbb{R} \times (0, \infty) \rightarrow \mathbb{R}, \quad L(\mu, \sigma) := \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}. \quad (12.3)$$

Sind dann  $X_1, \dots, X_n$  unabhängige Kopien von  $X \sim \mathcal{N}(\mu, \sigma^2)$ , so gilt

$$P[X_i \in [x_i - \varepsilon, x_i + \varepsilon] \text{ für alle } i] = \left(\frac{1}{\sqrt{2\pi}\sigma}\right)^n \int_{x_1 - \varepsilon}^{x_1 + \varepsilon} e^{-\frac{(x - \mu)^2}{2\sigma^2}} dx \cdots \int_{x_n - \varepsilon}^{x_n + \varepsilon} e^{-\frac{(x - \mu)^2}{2\sigma^2}} dx.$$

Für kleines festes  $\varepsilon > 0$  wird obige Wahrscheinlichkeit um so größer, je größer die Integranden, jeweils ausgewertet in  $x_i$ , sind. Dies führt zur Aufgabe, die Funktion  $L$  zu maximieren, und zur Heuristik, dass für einen Maximierer  $(\mu^*, \sigma^*)$  von  $L$  die Wahrscheinlichkeit dafür, dass die Datenpunkte in  $G$  von einer Gaußverteilung mit Mittelwert  $\mu^*$  und Varianz  $\sigma^{*2}$  kommen, im Vergleich zu allen anderen Gaußverteilungen maximiert.

**Satz 12.14.** *Sei  $G = \{x_1, \dots, x_n\} \subseteq \mathbb{R}$  gegeben und seien nicht alle  $x_i$  gleich. Dann hat die zugehörige Likelihood-Funktion  $L$ , definiert wie in (12.3), genau einen*

Maximierer und dieser ist gegeben durch

$$\mu^* = \frac{1}{n} \sum_{i=1}^n x_i \quad \text{und} \quad \sigma^{*2} = \frac{1}{n} \sum_{i=1}^n (x_i - \mu^*)^2.$$

*Beweis.* Wie schon in früheren Beweisen betrachten wir anstelle von  $L$  die Log-Likelihood-Funktion  $\ell: \mathbb{R} \times (0, \infty) \rightarrow \mathbb{R}$ ,  $\ell(\mu, \sigma) := \log L$ , für welche sich

$$\begin{aligned} \ell(\mu, \sigma) &= \log \left( \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} \right) = \log \left( \left( \frac{1}{\sqrt{2\pi}\sigma} \right)^n e^{-\frac{(x_1 - \mu)^2}{2\sigma^2} - \dots - \frac{(x_n - \mu)^2}{2\sigma^2}} \right) \\ &= n(0 - \log \sqrt{2\pi} - \log \sigma) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \end{aligned}$$

ergibt. Nullsetzen der partiellen Ableitungen liefert

$$\frac{\partial \ell}{\partial \mu} = \frac{1}{2\sigma^2} \sum_{i=1}^n 2(x_i - \mu) \cdot 1 \stackrel{!}{=} 0 \iff n\mu - \sum_{i=1}^n x_i = 0 \iff \mu = \frac{1}{n} \sum_{i=1}^n x_i$$

und

$$\frac{\partial \ell}{\partial \sigma} = -\frac{n}{\sigma} - \frac{-2}{2\sigma^3} \sum_{i=1}^n (x_i - \mu)^2 \stackrel{!}{=} 0 \iff n = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \iff \sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2.$$

Es gibt also genau einen kritischen Punkt. Um zu sehen, dass dies der eindeutig bestimmte Maximierer ist, zeigen wir, dass die Log-Likelihood-Funktion am Rand und im Unendlichen von  $\mathbb{R} \times (0, \infty)$  gegen  $-\infty$  divergiert. Zur Schreiberleichterung lassen wir den  $(\frac{1}{\sqrt{2\pi}})^n$ -Term weg und betrachten

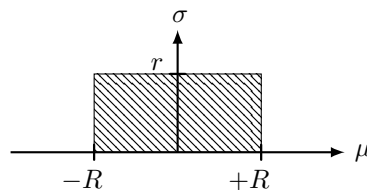
$$f: \mathbb{R} \times (0, \infty) \rightarrow \mathbb{R}, \quad f(\mu, \sigma) := \frac{1}{\sigma^n} e^{-\frac{1}{2\sigma^2} g(\mu)} \quad \text{mit} \quad g: \mathbb{R} \rightarrow \mathbb{R}, \quad g(\mu) := \sum_{i=1}^n (x_i - \mu)^2.$$

① Als erstes schauen wir uns Punkte in  $\{0\} \times (0, \infty)$  an. Gelte  $(\mu, \sigma) \rightarrow (\mu_0, 0)$  mit  $\mu_0 \in \mathbb{R}$ . Da  $\lim_{|\mu| \rightarrow \infty} g(\mu) = \infty$  gilt und  $g$  stetig und ungleich Null ist, ist  $c := \min_{\mu \in \mathbb{R}} g(\mu) > 0$ . Es folgt

$$f(\mu, \sigma) = \frac{1}{\sigma^n} e^{-\frac{1}{2\sigma^2} g(\mu)} \leq \frac{1}{\sigma^n} e^{-\frac{1}{2\sigma^2} c} \xrightarrow{\sigma \rightarrow 0} 0$$

und damit  $\ell(\mu, \sigma) = \log \left( \left( \frac{1}{\sqrt{2\pi}} \right)^n f(\mu, \sigma) \right) \rightarrow -\infty$ .

② Um zu sehen, dass  $\lim_{\|(\mu, \sigma)\| \rightarrow \infty} f(\mu, \sigma) = 0$  gilt, schätzen wir die Werte von  $f$  außerhalb einer Box wie im folgenden Bild ab.





Formal behaupten wir

$$\forall \varepsilon > 0 \exists r, R > 0 \forall (\mu, \sigma) \in \mathbb{R} \times (0, \infty): [|\mu| \geq R \text{ oder } \sigma \geq r] \implies f(\mu, \sigma) < \varepsilon.$$

Die Bedingung in der eckigen Klammer ist äquivalent zu  $\|(\mu, \sigma)\|_\infty \geq \min(R, r)$  und weil  $\|\cdot\|_\infty$  und  $\|\cdot\|_2$  auf  $\mathbb{R}^2$  äquivalent sind, erhalten wir dann genau die gewünschte Grenzwertaussage. Sei also jetzt  $\varepsilon > 0$  gegeben. Als erstes wählen wir  $r > 0$  derart, dass  $\frac{1}{\sigma^n} < \varepsilon$  gilt für alle  $\sigma \geq r$ . Mit  $c = \min_{\mu \in \mathbb{R}} g(\mu) > 0$  wie im ersten Teil wählen wir als nächstes  $\Sigma > 0$ , sodass  $\frac{1}{\sigma^n} e^{-\frac{c}{2\sigma^2}} < \varepsilon$  für alle  $0 < \sigma < \Sigma$  gilt. Weil  $\lim_{\mu \rightarrow \infty} g(\mu) = \infty$  gilt, können wir schließlich  $R > 0$  derart wählen, dass  $\frac{1}{\Sigma^n} e^{-\frac{g(\mu)}{2r^2}} < \varepsilon$  für alle  $\mu$  mit  $|\mu| \geq R$  gilt.

Sei jetzt  $(\mu, \sigma)$  mit  $|\mu| \geq R$  oder  $\sigma \geq r$  gegeben. Ist  $\sigma \geq r$ , so haben wir wegen  $g \geq 0$

$$\frac{1}{\sigma^n} e^{-\frac{g(\mu)}{2\sigma^2}} \leq \frac{1}{\sigma^n} < \varepsilon$$

und nach unserer obigen Wahl von  $r$ . Andernfalls ist  $\sigma < r$  und dann notwendigerweise  $|\mu| \geq R$ . Wir unterscheiden zwei Unterfälle: Ist  $\sigma < \Sigma$ , so ergibt sich

$$f(\mu, \sigma) = \frac{1}{\sigma^n} e^{-\frac{g(\mu)}{2\sigma^2}} \leq \frac{1}{\sigma^n} e^{-\frac{c}{2\sigma^2}} < \varepsilon$$

nach Wahl von  $\Sigma$ . Ist schließlich  $\sigma \geq \Sigma$ , und immer noch  $\sigma < r$  und  $|\mu| \geq R$ , so folgt aus diesen drei Ungleichungen in der genannten Reihenfolge

$$f(\mu, \sigma) = \frac{1}{\sigma^n} e^{-\frac{g(\mu)}{2\sigma^2}} \leq \frac{1}{\Sigma^n} e^{-\frac{g(\mu)}{2\sigma^2}} \leq \frac{1}{\Sigma^n} e^{-\frac{g(\mu)}{2r^2}} < \varepsilon$$

entsprechend der obigen Wahl von  $R$ . Dies zeigt die Behauptung und wir hatten bereits argumentiert dass aus dieser  $\ell(\mu, \sigma) \rightarrow -\infty$  für  $\|(\mu, \sigma)\| \rightarrow \infty$  folgt.  $\square$

Wir notieren die folgenden Bemerkungen zum sogenannten *Maximum-Likelihood-Schätzer*  $(\mu^*, \sigma^*)$ .

**Bemerkung 12.15.** (i) Satz 12.14 induziert eine Abbildung

$$M: \{(x_1, \dots, x_n) \in \mathbb{R}^n \mid \exists i, j: x_i \neq x_j\} \rightarrow \mathbb{R}^2, \quad M(x_1, \dots, x_n) := (\mu^*, \sigma^*)$$

indem man den Schätzer als Funktion der Daten auffasst, die man in Anbetracht der expliziten Formeln für  $\mu^*$  und  $\sigma^*$  auf ganz  $\mathbb{R}^n$  fortsetzen kann. Wir haben in Aufgabe 2.14 bereits gesehen, dass man  $M$  auch als Zufallsvariable auffassen und in diesem Sinne  $E(M)$  und  $V(M)$  berechnen kann. Macht man das, so kann es natürlicher erscheinen, statt des Maximum-Likelihood-Schätzers dieselben Formel zu nehmen, aber bei der Varianz durch  $n - 1$  anstelle von  $n$  zu teilen.

(ii) Ist  $X: \Omega \rightarrow \mathbb{R}^d$ ,  $X \sim \mathcal{N}(\mu, \sigma^2, \mathbb{R}^d)$  mit  $\mu = (\mu_1, \dots, \mu_d) \in \mathbb{R}^d$  und  $\sigma > 0$  ein Zufallsvektor, so zeigt man analog zu Satz 12.14, dass

$$\mu^* = \frac{1}{n} \sum_{i=1}^n x_i \in \mathbb{R}^d \quad \text{und} \quad \sigma^{*2} = \frac{1}{n} \sum_{i=1}^n \|\hat{\mu} - x_i\|^2$$

eine Likelihood-Funktion analog zu (12.3) maximieren.

Wir beenden das Kapitel und auch den Abschnitts dieses Textes zu hochdimensionalen Räumen damit, dass wir darauf aufmerksam machen, dass wir hier stets die sphärische Gaußverteilung betrachtet haben. Allgemein betrachtet man Gaußsche Zufallsvariablen  $X \sim \mathcal{N}(\mu, \Sigma, \mathbb{R}^d)$  gegeben durch die Dichtefunktion

$$\rho(x) = \frac{1}{(2\pi)^{d/2}} \cdot \frac{1}{(\det \Sigma)^{1/2}} e^{-\frac{1}{2} \langle x - \mu, \Sigma^{-1}(x - \mu) \rangle}$$

mit positiv definiter *Kovarianzmatrix*  $\Sigma$ . In zwei Dimensionen sehen Stichproben dann wie folgt aus, wobei das linke Bild den von uns behandelten Fall zeigt, in dem  $\Sigma$  die Einheitsmatrix mal einem Skalar  $\sigma > 0$  ist.



Im mittleren Bild ist  $\Sigma$  eine Diagonalmatrix und rechts dann eine beliebige positiv definite Matrix. Die Behandlung von in diesem Sinne allgemeinen Gaußianen ist verständlicherweise deutlich technischer als das was wir in den vorangegangenen Kapiteln gemacht haben.

## Referenzen

Grundlage für dieses Kapitel ist wieder [BHK20], wobei wir aber auch hier wieder Änderungen und Ergänzungen vorgenommen haben. Wir verweisen außerdem auf die Originalarbeit [DS07]. Beispiele für die in Algorithmus 12.3 erwähnten Sortieralgorithmen findet man etwa in [CLRS22, Part II]. Für mehr Details zur Behandlung von nicht-sphärischen Gaußianen verweisen wir auf [Bis06, Chapter 9].

## Aufgaben

**Aufgabe 12.1.** Die folgende Tabelle enthält die Größen von Amerikanern im Alter von 40-49 Jahren. Die Zahlen geben jeweils an, welcher Prozentsatz kleiner gleich der darüber angegebene Größe und größer als die in der Spalte links daneben angegebene Größe ist.

	147.32	149.86	152.40	154.94	157.48	160.02	162.56	165.10	167.64	170.18
Frauen	1.6	3.4	5.8	9.0	11.0	15.2	12.0	14.2	10.8	8.2
Männer	0	0	0	0	1.9	1.9	1.8	4.2	9.6	10.9
	172.72	175.26	177.80	180.34	182.88	185.42	187.96	190.50	193.04	195.58
Frauen	3.5	3.1	1.6	0.1	0	0	0	0	0.5	0
Männer	10.1	14.0	15.2	9.5	8.3	5.1	5.2	1.3	0.4	0.5

Quelle: <https://www2.census.gov/library/publications/2010/compendia/statab/130ed/tables/11s0205.pdf>

Plotten Sie die zwei (gaußschen) Verteilungen, erstellen Sie dann die Vereinigungsmenge und versuchen Sie, diese mit Algorithmus 12.3 wieder zu trennen.

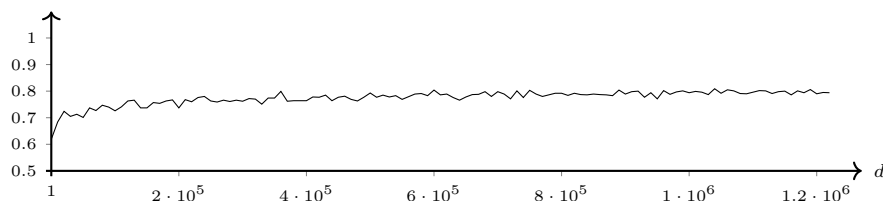
**Aufgabe 12.2.** Erzeugen Sie zwei Gaußiane  $G_1, G_2 \subseteq \mathbb{R}^{100}$ , den einen mit Mittelpunkt Null und den anderen mit Mittelpunkt  $(\Delta, 0, \dots, 0)$  für verschiedene  $\Delta > 0$  und plotten Sie dann die Verteilung der paarweisen Abstände in der Vereinigungsmenge  $D = G_1 \cup G_2$ .

**Aufgabe 12.3.** Ein Schwachpunkt des in diesem Kapitel diskutierten Trennungsalgorithmus 12.3 ist es, dass der Startpunkt  $z$  ein ‘untypischer Punkt’ ist. In zwei Dimensionen könnte er z.B. auf halbem Weg zwischen den Gaußianen liegen, in hohen Dimensionen außerhalb der beiden Ringe. Überlegen Sie sich Strategien zur Wahl von  $z$ , die dem entgegenwirken.

**Aufgabe 12.4.** (i) Wenden Sie den Trennungsalgorithmus 12.3 auf selbst-erzeugte hochdimensionale Gaußiane an und lassen Sie ihr Programm dabei die Rate berechnen mit der die Punkte korrekt getrennt werden.

(ii) Probieren Sie die nach Beispiel 12.8 angedeuteten Variationen des Trennungsalgorithmus aus und vergleichen Sie mit der ursprünglichen Version.

**Aufgabe 12.5.** Beispiel 12.12 suggeriert, dass für  $\Delta = d^{0.30}$  die Güte des Trennungsalgorithmus zwar kleiner ist als mit für  $\Delta = 2d^{0.30}$ , aber dennoch wächst. Simulieren Sie dies, indem Sie möglichst große Dimensionen untersuchen und verifizieren Sie, dass ab  $d = 10^6$  damit zu rechnen ist, dass die Trennung mit  $\geq 80\%$  richtig klassifizierten Punkten gelingt.



*Hinweis:* Auf einem handelsüblichen Computer kann dies eine ganze Weile dauern und je nach Programmiersprache ist irgendwann auch einfach Schluss.

**Aufgabe 12.6.** Erzeugen Sie zwei Gaußiane in  $\mathbb{R}^d$  mit Varianzen  $\sigma_1^2 \neq \sigma_2^2$  und Mittelpunkt Null. Finden Sie durch Ausprobieren heraus, ab welchem  $d$  die Trennung gelingt.

**Aufgabe 12.7.** Erzeugen Sie einen Gaußian  $G$  und schätzen Sie dann dessen Parameter  $\mu$  und  $\sigma$  mithilfe des Maximum-Likelihood-Schätzers.

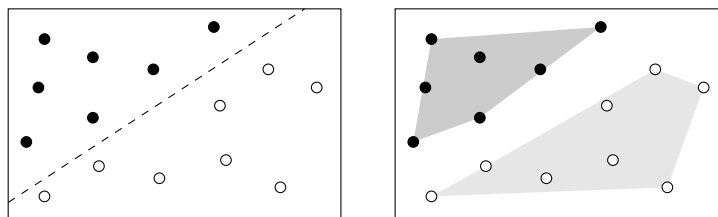
## Kapitel 13

# Das Perzeptron

In diesem Kapitel betrachten wir Klassifikationsprobleme der folgenden Form. Gegeben sei eine gelabelte Datenmenge  $D = \{(x_i, y_i) \mid i = 1, \dots, n\} \subseteq \mathbb{R}^d \times \{-1, 1\}$  mit binären Labels  $-1$  oder  $+1$ . Wir suchen einen *affin-lineare Klassifizierer*, auch *Perzeptron* genannt,  $h: \mathbb{R}^d \rightarrow \mathbb{R}$ , d.h. eine Funktion der Bauart

$$h(x) = \text{sign}(\langle w, x \rangle + b)$$

mit  $0 \neq w \in \mathbb{R}^d$  und  $b \in \mathbb{R}$ , sodass  $h(x_i) = y_i$  für alle  $i = 1, \dots, n$  gilt. Natürlich können wir nicht erwarten, dass dies für beliebige Datenmengen  $D$  überhaupt möglich ist. Als erstes werden wir daher in diesem Kapitel zwei geometrische Charakterisierungen derjenigen Datenmengen angeben, für die ein Klassifizierer  $h$  wie oben existiert. Die erste Charakterisierung besagt, dass die Datenpunkte mit Label  $-1$  von denen mit Label  $+1$  durch eine affine Hyperebene strikt getrennt werden können (linkes Bild). Die zweite charakterisierende Eigenschaft ist, dass die konvexen Hüllen der Punkte mit Label  $+1$ , bzw.  $-1$ , leeren Schnitt haben (rechtes Bild).



Für die zweite Charakterisierung benötigen wir den Satz von Carathéodory über konvexe Hüllen und den Trennungssatz für konvexe Mengen, welche wir beide unten ohne Beweis angeben. Wir verweisen auf das Kapitelende für detaillierte Referenzen.

Wir verwenden nun die Existenz eines Klassifizierers der oben beschrieben Form, um die Klasse von Datenmengen, die wir in diesem Kapitel behandeln, zu definieren. Dies mag etwas tautologisch anmuten, aber erstens werden sich später alle drei Eigenschaften ohnehin als äquivalent herausstellen und zweitens ist die Existenz des Klassifizierers diejenige der drei Eigenschaften, die sich am einfachsten formulieren lässt.

**Definition 13.1.** Eine Datenmenge  $D = \{(x_i, y_i) | i = 1, \dots, n\} \subseteq \mathbb{R}^d \times \{-1, 1\}$  heißt *linear trennbar*, falls eine Funktion  $h: \mathbb{R}^d \rightarrow \mathbb{R}$  der Form  $h(x) = \text{sign}(\langle w, x \rangle + b)$  mit  $0 \neq w \in \mathbb{R}^d$ ,  $b \in \mathbb{R}$  existiert, die alle Datenpunkte korrekt klassifiziert.

Häufig nennt man die  $w_1, \dots, w_d$  *Gewichte* und das  $b$  *Bias* des Klassifizierers, weil  $b$  festlegt, welchen Wert die gewichtete Summe  $w_1x_1 + \dots + w_dx_d$  überschreiten muss, damit der Klassifizier vom Wert  $-1$  auf den Wert  $+1$  wechselt.

Wir weisen darauf hin, dass wir hier auch  $w = 0$  zulassen könnten. In diesem Fall ist  $h \equiv 0$ ,  $h \equiv -1$  oder  $h \equiv +1$ . Ersteres würde bedeuten, dass  $D = \emptyset$  ist, die beiden anderen Fälle korrespondieren mit Datenmengen, in denen alle Datenpunkte das gleiche Label haben. Abgesehen davon, dass dies uninteressante Fälle sind, kann in allen drei Situationen selbstverständlich auch ein korrekter affin-linearer Klassifizierer mit  $w \neq 0$  gefunden werden. Wir schließen also  $w = 0$  hier ohne Einschränkung aus. Der Grund hierfür wird mit der ersten Charakterisierung deutlich werden und mit dieser wollen wir jetzt mit der nachfolgenden Definition beginnen.

**Definition 13.2.** Eine Teilmenge  $\mathcal{H} \subseteq \mathbb{R}^d$  heißt *affine Hyperebene*, falls  $\mathcal{H}$  von der Form  $\mathcal{H} = X + x_0$ , wobei  $X \subseteq \mathbb{R}^d$  ein linearer Unterraum der Dimension  $d - 1$  ist und  $x_0 \in \mathbb{R}^d$ .

Als nächstes zeigen wir, dass affine Hyperebenen genau die Nullstellenmengen affin-linearer Klassifizier sind.

**Lemma 13.3.** Eine Teilmenge  $\mathcal{H} \subseteq \mathbb{R}^d$  ist genau dann eine affine Hyperebene, wenn  $0 \neq w \in \mathbb{R}^d$  und  $b \in \mathbb{R}$  existieren mit  $\mathcal{H} = \{x \in \mathbb{R}^d | \langle w, x \rangle + b = 0\}$ . Dabei sind  $w$  und  $b$  eindeutig bis auf einen von Null verschiedenen Faktor.

*Beweis.* “ $\implies$ ” Sei  $\mathcal{H} = X + x_0$  wie in Definition 13.2. Wir wählen  $0 \neq w \in X^\perp$  und beachten, dass  $X^\perp$  Dimension Eins hat. Dann setzen wir  $b := -\langle w, x_0 \rangle$ . Sei jetzt  $x = x_1 + x_0 \in X + x_0$ . Es gilt dann

$$\langle w, x \rangle + b = \langle w, x_1 + x_0 \rangle - \langle w, x_0 \rangle = \langle w, x_1 \rangle$$

und letzteres ist Null wegen  $w \perp x_1$ . Sei andersherum  $x \in \mathbb{R}^d$  gegeben mit  $\langle w, x \rangle + b = 0$ . Wir setzen  $x_1 := x - x_0$ , also  $x = x_1 + x_0$  und sehen mit derselben Rechnung nun, dass  $w \perp x_1$  folgt.

“ $\impliedby$ ” Sind  $0 \neq w \in \mathbb{R}^d$  und  $b \in \mathbb{R}$  gegeben, so setzen wir  $X := \{w\}^\perp$  und wählen  $x_0$  derart, dass  $\langle w, x_0 \rangle = -b$  gilt.

Für die Eindeutigkeitsaussage seien  $(w_1, b_1), (w_2, b_2) \in \mathbb{R}^{d+1}$  gegeben mit  $w_1 \neq 0 \neq w_2$ . Wir zeigen, dass  $\mathcal{H} := \{x \in \mathbb{R}^d | \langle w_1, x \rangle + b_1 = 0\} = \{x \in \mathbb{R}^d | \langle w_2, x \rangle + b_2 = 0\}$  genau dann gilt, wenn ein  $\alpha \neq 0$  existiert mit  $w_1 = \alpha w_2$  und  $b_1 = \alpha b_2$ .

“ $\implies$ ” Nach Lemma 13.2 existiert ein Unterraum  $X \subseteq \mathbb{R}^d$  mit  $\dim X = d - 1$  und  $x_0 \in \mathbb{R}^d$  mit  $\mathcal{H} = X + x_0$ . Der Beweis des ersten Teils zeigt überdies, dass  $w_1, w_2 \perp X$  gilt. Anders ausgedrückt, gehören  $w_1$  und  $w_2$  beide zum eindimensionalen Unterraum

$X^\perp$ . Da weiterhin beide Vektoren ungleich Null sind, existiert  $\alpha \neq 0$  mit  $w_1 = \alpha w_2$ . Da  $x_0 \in X + x_0$  gilt, sehen wir

$$0 = \langle w_1, x_0 \rangle + b_1 = \langle \alpha w_2, x_0 \rangle + b_1 = \alpha \langle w_2, x_0 \rangle + b_1 = -\alpha b_2 + b_1,$$

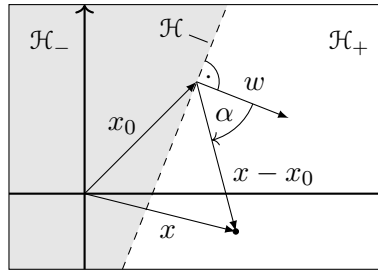
also  $b_1 = \alpha b_2$  wie behauptet.

“ $\Leftarrow$ ” Folgt aus der Bilinearität des Skalarproduktes.  $\square$

**Bemerkung 13.4.** Ist  $(w, b) \in \mathbb{R}^{d+1}$  mit  $w \neq 0$  gegeben, so teilt die zugehörige Hyperebene  $\mathcal{H}$  den ganzen Raum in zwei *affine Halbräume*

$$\mathcal{H}_+ := \{x \in \mathbb{R}^d \mid \langle w, x \rangle + b > 0\} \quad \text{und} \quad \mathcal{H}_- := \{x \in \mathbb{R}^d \mid \langle w, x \rangle + b < 0\}$$

auf. Beachte hierbei, dass  $\mathcal{H}_\pm$ , im Gegensatz zu  $\mathcal{H}$ , insofern von der Wahl  $(w, b)$  abhängen, als dass beim Übergang zu  $(\alpha w, \alpha b)$  mit  $\alpha < 0$  die Halbräume die Plätze tauschen. Fixieren wir  $x_0 \in \mathcal{H}$ , so können wir  $\mathcal{H}$ ,  $\mathcal{H}_+$  und  $\mathcal{H}_-$  durch das folgende Bild veranschaulichen.



Hier sehen wir, dass  $\langle w, x \rangle + b = \langle w, x - x_0 \rangle > 0$  genau dann gilt, wenn der Winkel  $\alpha = \arccos(\langle \frac{w}{\|w\|}, \frac{x-x_0}{\|x-x_0\|} \rangle)$  zwischen  $w$  und  $x - x_0$  im Intervall  $(-\pi/2, \pi/2)$  liegt. Dies erklärt nochmal den Begriff ‘Halbraum’.

**Satz 13.5.** Eine Datenmenge  $D = \{(x_i, y_i) \mid i = 1, \dots, n\} \subseteq \mathbb{R}^d \times \{-1, 1\}$  ist genau dann linear trennbar, wenn eine Hyperebene  $\mathcal{H}$  existiert, sodass alle Punkte mit Label  $+1$  in einem der beiden Halbräume liegen und alle Punkte mit Label  $-1$  im anderen.

*Beweis.* “ $\Rightarrow$ ” Wenn  $S$  linear trennbar ist, existieren per Definition  $(w, b) \in \mathbb{R}^{d+1}$  mit  $w \neq 0$  sodass  $y_i = h(x_i) = \text{sign}(\langle w, x_i \rangle + b)$  für alle  $i = 1, \dots, n$  gilt. Dass heißt aber gerade, dass immer wenn  $y_i = +1$  ist,  $\langle w, x_i \rangle + b > 0$  sein muss, oder in anderen Worten,  $x_i \in \mathcal{H}_+$ . Analog gilt immer wenn  $y_i = -1$  ist,  $x_i \in \mathcal{H}_-$ .

“ $\Leftarrow$ ” Ist die Bedingung im Satz gegeben und die Hyperebene durch  $(w, b) \in \mathbb{R}^{d+1}$  gegeben, so liegen entweder alle Punkte mit Label  $+1$  und  $\mathcal{H}_+$  und alle Punkte mit Label  $-1$  in  $\mathcal{H}_-$  oder es ist genau umgekehrt. Im ersten Fall setzen wir  $h = \text{sign}(\langle w, - \rangle + b)$  und im zweiten Fall  $h = \text{sign}(\langle -w, - \rangle - b)$ .  $\square$

Bevor wir zur zweiten Charakterisierung kommen, bemerken wir noch, dass unsere Definition gegenüber der Bedingung in Satz 13.5 den Vorteil hat, dass wir das

Vorzeichen von  $w$  und  $b$  festlegen: Anschaulich gesprochen zeigt  $w$  in die Richtung der positiv gelabelten Daten.

**Satz 13.6.** *Sei  $D = \{(x_i, y_i) \mid i = 1, \dots, n\} \subseteq \mathbb{R}^d \times \{-1, 1\}$  und seien  $S_{\pm} := \{x_i \mid y_i = \pm 1\}$ . Dann ist  $D$  genau dann linear trennbar, wenn  $\text{conv } D_+ \cap \text{conv } D_- = \emptyset$  gilt, d.h. wenn die konvexen Hüllen derjenigen Datenpunkte mit Label  $+1$ , bzw. derjenigen mit  $-1$ , disjunkt sind.*

*Beweis.* “ $\implies$ ” Wenn  $D$  linear trennbar ist, dann gibt es  $(w, b) \in \mathbb{R}^{d+1}$  mit

$$D_+ \subseteq \{x \in \mathbb{R}^d \mid \langle w, x \rangle + b > 0\} =: \mathcal{H}_+$$

und die Menge rechts ist konvex. Daher folgt  $\text{conv } D_+ \subseteq \mathcal{H}_+$  und analog  $\text{conv } D_- \subseteq \mathcal{H}_-$ . Der Schnitt  $\mathcal{H}_+ \cap \mathcal{H}_-$  ist aber leer per Konstruktion.

“ $\impliedby$ ” Wir zeigen, dass  $K_+ := \text{conv } D_+$  und  $K_- := \text{conv } D_-$  im Sinne von Satz 13.5 durch eine affine Hyperebene getrennt werden können und benutzen dafür den folgenden Satz.

**Satz 13.7.** (Trennungssatz für konvexe Mengen) *Seien  $A, B \subseteq \mathbb{R}^d$  beide nichtleer, konvex und abgeschlossen. Sei eine der Mengen kompakt. Dann existiert eine Hyperebene  $\mathcal{H}$  derart dass  $A$  in einem der Halbräume liegt und  $B$  im anderen.*  $\diamond$

In unserem Fall sind  $K_-$  und  $K_+$  konvex per Definition und wir zeigen jetzt, dass auch beide Mengen kompakt sind. Wir behaupten, etwas allgemeiner, dass für jede kompakte Menge  $X \subseteq \mathbb{R}^d$  deren konvexe Hülle  $K := \text{conv } X$  wieder kompakt ist. Zur Darstellung der konvexen Hülle verwenden wir den folgenden Satz.

**Satz 13.8.** (von Carathéodory über konvexe Mengen) *Sei  $X \subseteq \mathbb{R}^d$  und  $z \in K := \text{conv } X$ . Dann existieren  $z_1, \dots, z_{d+1} \in X$  und  $\alpha_1, \dots, \alpha_{d+1} \in [0, 1]$  mit  $\alpha_1 + \dots + \alpha_{d+1} = 1$  und  $z = \alpha_1 z_1 + \dots + \alpha_{d+1} z_{d+1}$ .*  $\diamond$

Anders ausgedrückt haben wir

$$K = \left\{ \sum_{k=1}^{d+1} \alpha_k z_k \mid z = (z_1, \dots, z_{d+1}) \in X^{d+1}, \alpha = (\alpha_1, \dots, \alpha_{d+1}) \in \Sigma_d \right\}$$

wobei man der Menge

$$\Sigma_d := \left\{ \alpha = (\alpha_1, \dots, \alpha_d) \in \mathbb{R}_{\geq 0}^{d+1} \mid \sum_{k=1}^{d+1} \alpha_k = 1 \right\}$$

sofort ansieht, dass sie kompakt ist. Nun betrachten wir die Abbildung

$$\varphi: X^{d+1} \times \Sigma_d \rightarrow \mathbb{R}^d, (z, \alpha) \mapsto \sum_{k=1}^{d+1} \alpha_k z_k,$$

welche offenbar stetig ist. Weil  $X^{d+1} \times \Sigma_d$  kompakt ist folgt dann, dass  $K = \text{ran } \varphi$  auch kompakt sein muss. Dies zeigt die Behauptung und Anwendung derselben mit  $X = D_{\pm}$  beendet den Beweis.  $\square$

Nachdem wir nun die zwei geometrischen Charakterisierungen von linear trennbaren Mengen bewiesen haben, wenden wir uns der Frage zu, wie wir, vorausgesetzt eine gegebene Datenmenge  $D$  ist linear trennbar, einen möglichen affin-linearen Klassifizierer finden können. Dies ist möglich mithilfe des Perzeptronalgorithmus 13.11. Um diesen effizient zu formulieren, ist es zweckmäßig, Gewichte und Bias zu einem  $(d+1)$ -dimensionalen Vektor zusammenzufassen und die Daten in geeigneter Weise in  $\mathbb{R}^{d+1}$  einzubetten.

**Definition 13.9.** Für Datenpunkte  $z = (z_1, \dots, z_d) \in \mathbb{R}^d$  bezeichnen wir mit  $\hat{z} := (z, 1) = (z_1, \dots, z_d, 1) \in \mathbb{R}^{d+1}$  den (um Eins) erweiterten Datenpunkt. Für Gewichte  $w = (w_1, \dots, w_d) \in \mathbb{R}^d$  und Bias  $b \in \mathbb{R}$  bezeichnen wir mit  $(w, b) = (w_1, \dots, w_d, b) \in \mathbb{R}^{d+1}$  den zusammengefassten Gewichtsvektor.

**Bemerkung 13.10.** Seien  $w \in \mathbb{R}^d$  und  $b \in \mathbb{R}$  gegeben. Dann gilt für  $z = (z_1, \dots, z_d) \in \mathbb{R}^d$  mit der oben eingeführten Notation

$$\langle w, z \rangle + b = \sum_{i=1}^d w_i z_i + b \cdot 1 = \left\langle \begin{bmatrix} w_1 \\ \vdots \\ w_d \\ b \end{bmatrix}, \begin{bmatrix} z_1 \\ \vdots \\ z_d \\ 1 \end{bmatrix} \right\rangle = \langle w, \hat{z} \rangle$$

wobei wir (etwas missbräuchlich!) links mit  $w$  den normalen Gewichtsvektor und rechts, ebenfalls mit  $w$ , den zusammengefassten Gewichtsvektor bezeichnen.

**Algorithmus 13.11.** Der folgende Pseudocode gibt den sogenannten Perzeptronalgorithmus wieder. Diesem wird die Datenmenge  $S = \{(x_i, y_i) \mid i = 1, \dots, n\}$  übergeben, sowie ein initialer Vektor  $w \in \mathbb{R}^{d+1}$ . Als Ausgabe erhalten wir den zusammengefassten (!) Gewichtsvektor  $w^{(j)} \in \mathbb{R}^{d+1}$ . D.h. die ersten  $d$ -vielen Einträge von  $w^{(j)}$  sind die Gewichte und der letzte ist das Bias des Klassifizierers.

```

1: function PERZEPTRON ( $D, w$ )
2:    $w^{(0)} \leftarrow w$ 
3:   for  $j \leftarrow 0$  to  $\infty$  do
4:     if  $\exists i \in \{1, \dots, n\} : y_i \langle w^{(j)}, \hat{x}_i \rangle \leq 0$  then
5:        $w^{(j+1)} \leftarrow w^{(j)} + y_i \hat{x}_i$ 
6:     else
7:       break
8: return  $w^{(j)}$ 

```

Wir werden in Satz 13.13 zeigen, dass der Perzeptronalgorithmus unter der Voraussetzung, dass  $D$  linear trennbar ist, immer terminiert und einen korrekten Klassifizierer liefert. Zunächst diskutieren wir jedoch die Heuristik hinter dem Algorithmus.

**Bemerkung 13.12.** Wir betrachten zuerst die **if**-Abfrage, d.h. die Bedingung die zu einem Update des zusammengefassten Gewichtsvektors führt. Hier gilt

$$\begin{aligned} y_i \langle w^{(j)}, \hat{x}_i \rangle > 0 &\Leftrightarrow [(y_i = +1 \wedge \langle w^{(j)}, \hat{x}_i \rangle > 0) \vee (y_i = -1 \wedge \langle w^{(j)}, \hat{x}_i \rangle < 0)] \\ &\Leftrightarrow h(x_i) = \text{sign}(\langle w^{(j)}, \hat{x}_i \rangle) = y_i \end{aligned}$$

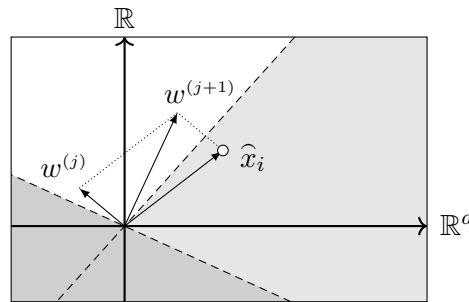
$\uparrow$   
 Beachte  
 Bem. 13.10



und es folgt durch Verneinung, dass  $y_i \langle w^{(j)}, \hat{x}_i \rangle \leq 0$  genau dann erfüllt ist, wenn  $x_i$  falsch klassifiziert wird. Zeile 4 findet also falsch klassifizierte Datenpunkte. Das Update in Zeile 5 ‘versucht’ dann gewissermaßen die Fehlklassifikation zu beheben

$$y_i \langle w^{(j+1)}, x_i \rangle = y_i \langle w^{(j)} + y_i x_i, x_i \rangle = \underbrace{y_i \langle w^{(j)}, x_i \rangle}_{\text{fälschlicherweise kleiner Null}} + \underbrace{y_i^2 \langle x_i, x_i \rangle}_{\text{Addition eines positiven Korrekturterms}}$$

indem ein Korrekturterm addiert wird. Hierbei besteht jedoch keine Garantie, dass  $\hat{x}_i$  danach korrekt klassifiziert wird, vgl. Aufgabe 13.2. Skizzieren wir die Hyper Ebenen die durch  $w^{(j)}$  und  $w^{(j+1)}$  gegeben werden, so sieht man im Beispiel unten, dass das Update diese derart dreht, sodass der Punkt  $\hat{x}_i$ , welcher durch  $w^{(j)}$  nicht korrekt klassifiziert wurde, nach dem Update, also durch  $w^{(j+1)}$  tatsächlich korrekt klassifiziert wird.



Wir weisen darauf hin, dass wir im Bild mit den erweiterten Daten und zusammengefassten Gewichtsvektoren arbeiten. Aus diesem Grund sind die affinen Hyper Ebenen in der Tat Unterräume.

Jetzt liefern wir den formalen Beweis dafür, dass der Perzeptronalgorithmus immer terminiert und geben dabei auch eine Schranke für seine Laufzeit an.

**Satz 13.13.** Sei  $\emptyset \neq D = \{(x_i, y_i) \mid i = 1, \dots, n\} \subseteq \mathbb{R}^d \times \{-1, 1\}$  linear trennbar. Dann existiert  $w^* \in \mathbb{R}^{d+1}$  sodass  $y_i \langle w^*, \hat{x}_i \rangle > 0$  für alle  $i = 1, \dots, n$  und  $\|w^*\| = 1$  gilt. Sei  $\gamma := \min_{i=1, \dots, n} y_i \langle w^*, \hat{x}_i \rangle$  ( $> 0$  nach dem Vorhergehenden!) und sei  $R := \max_{i=1, \dots, n} \|x_i\|$ . Dann gilt:

- (i) Der mit  $w = 0$  initialisierte Perzeptronalgorithmus stoppt nach spätestens  $\lfloor (R/\gamma)^2 \rfloor$ -vielen Iterationen.
- (ii) Stoppt der Perzeptronalgorithmus und gibt den zusammengefassten Gewichtsvektor  $w^{(j)}$  aus, dann klassifiziert der durch diesen definierte Klassifizierer  $h: \mathbb{R}^d \rightarrow \mathbb{R}$  alle Punkte aus  $D$  korrekt.

*Beweis.* Da  $D$  linear trennbar ist, existiert irgendein  $w \in \mathbb{R}^{d+1}$ , sodass alle Punkte korrekt klassifiziert werden. Dieser Vektor  $w$  kann nicht der Nullvektor sein, denn dann wäre  $h \equiv 0$  und alle Datenpunkte würden falsch klassifiziert. Setze  $w^* :=$

$w/\|w\|$ . Dann ist  $\|w^*\| = 1$  und  $y_i\langle w^*, x_i \rangle = \frac{1}{\|w\|} y_i\langle w, x_i \rangle > 0$  für alle  $i = 1, \dots, n$ . Mit  $\gamma$  wie im Satz gilt also jetzt

$$\forall i = 1, \dots, n: y_i\langle w^*, \hat{x}_i \rangle \geq \gamma > 0.$$

Sei nun  $0 = w^{(0)}, w^{(1)}, w^{(2)}, \dots$  die Folge von Gewichtsvektoren, die der mit dem Nullvektor initialisierte Perzeptronalgorithmus berechnet. Sei  $w^{(j)}$  ein Gewichtsvektor, der keine korrekte Klassifikation aller Punkte liefert. Dann führt der Algorithmus das Update  $w^{(j+1)} = w^{(j)} + y_i\hat{x}_i$  durch, wobei  $i \in \{1, \dots, n\}$  so gewählt ist, dass  $\hat{x}_i$  mit  $w^{(j)}$  als Gewichtsvektor falsch klassifiziert wurde. Damit folgt

$$\langle w^*, w^{(j+1)} \rangle = \langle w^*, w^{(j)} + y_i\hat{x}_i \rangle = \langle w^*, w^{(j)} \rangle + y_i\langle w^*, \hat{x}_i \rangle \geq \langle w^*, w^{(j)} \rangle + \gamma.$$

Per Iteration, und weil  $w^{(0)} = 0$  ist, ergibt sich

$$\langle w^*, w^{(j+1)} \rangle \geq \langle w^*, w^{(0)} \rangle + j \cdot \gamma = j \cdot \gamma.$$

Andererseits liefert das Update, mit  $R$  wie im Satz definiert, die Abschätzung

$$\begin{aligned} \|w^{(j+1)}\|^2 &= \|w^{(j)} + y_i\hat{x}_i\|^2 = \langle w^{(j)} + y_i\hat{x}_i, w^{(j)} + y_i\hat{x}_i \rangle \\ &= \|w^{(j)}\|^2 + \underbrace{2y_i\langle w^{(j)}, \hat{x}_i \rangle}_{<0, \text{ weil } \hat{x}_i \text{ falsch klassifiziert wird}} + \underbrace{\|x_i\|^2}_{\leq R} < \|w^{(j)}\|^2 + R^2. \end{aligned}$$

Wieder folgt per Iteration, und weil  $w^{(0)} = 0$  ist, dass

$$\|w^{(j+1)}\|^2 \leq \|w^{(0)}\|^2 + j \cdot R^2 = j \cdot R^2$$

gilt. Zusammensetzen beider Abschätzungen liefert

$$j \cdot \gamma \leq \langle w^*, w^{(j+1)} \rangle \leq \|w^*\|^2 \|w^{(j+1)}\|^2 = \|w^{(j+1)}\|^2 \leq \sqrt{j} \cdot R$$

und Division durch  $\gamma$  und  $\sqrt{j}$  impliziert  $\sqrt{j} = j/\sqrt{j} \leq R/\gamma$ . Da  $j \in \mathbb{N}$ , geht dies nur, wenn  $j \leq \lfloor R/\gamma \rfloor^2$  gilt und wir sehen, dass der Perzeptronalgorithmus nach höchstens  $\lfloor R/\gamma \rfloor^2$ -vielen Updates gestoppt haben muss. Dies zeigt (i). Wenn der Perzeptronalgorithmus in der  $J$ -ten Runde stoppt, dann ist per Definition die Bedingung

$$\exists i \in \{1, \dots, n\}: y_i\langle w^{(J)}, x_i \rangle \leq 0$$

nicht erfüllt. Dies bedeutet aber, dass alle Datenpunkte aus  $D$  via dem Gewichtsvektor  $w^{(J)}$  korrekt klassifiziert werden. Damit ist auch (ii) gezeigt.  $\square$

**Bemerkung 13.14.** (i) In Zeile 4 des Perzeptronalgorithmus 13.11 haben wir nicht spezifiziert, wie der falsch klassifizierte Punkt  $\hat{x}_i$  ausgewählt wird und es hat sich gezeigt, dass die Schranke für die Laufzeit unabhängig davon ist, wie dies in der Praxis gehandhabt wird. Die tatsächliche Laufzeit hängt allerdings sehr wohl davon ab, wie wir hier vorgehen, vergleiche Aufgabe 13.3.

(ii) Starten wir den Perzeptronalgorithmus mit  $w^{(0)} \neq 0$ , so kann man obigen Beweis anpassen und bekommt ein analoges Ergebnis, siehe Aufgabe 13.5.

(iii) Den Spam-Klassifizierer aus Aufgabe 13.3 kann man per

$$h(x_1, \dots, x_5) = \text{sign}(0 \cdot x_1 + 2 \cdot x_2 + 0 \cdot x_3 - 1 \cdot x_4 + 1 \cdot x_5 + 0)$$

notieren und interpretieren, dass die Worte ‘und’ sowie ‘das’ bezüglich Spam neutral sind, wohingegen ‘Bonus’ und ‘Vertrag’ entsprechend ihren unterschiedlichen Gewichten mehr oder weniger indikativ für Spam sind. Schließlich ist ‘Mensa’ indikativ für No-Spam. Man kann argumentieren, dass der Perzeptronalgorithmus *nachvollziehbare* Ergebnisse liefert.

(iv) Anstatt mit Labeln  $-1$  und  $1$  kann man auch z.B. mit  $0$  und  $1$  arbeiten und entsprechend Klassifizierer der Bauart  $h = \mathbb{1}_{[0, \infty)}(\langle w, \cdot \rangle + b)$  verwenden. Der Perzeptronalgorithmus benötigt dann zwei `if`-Abfragen, weil im Fall  $y_i = 0$  das Produkt  $y_i \langle w, \hat{x}_i \rangle$  nichts über das Vorzeichen von  $\langle w, \hat{x}_i \rangle$  verrät. Außerdem kann der so modifizierte Perzeptronalgorithmus einen Klassifizierer liefern, bei dem Datenpunkte genau auf der Entscheidungsgrenze liegen — das wird beim obigen Zugang automatisch verhindert.

(v) Das Perzeptron, damals kein abstrakter Algorithmus sondern eine physische Maschine, wurde in den 1950er Jahren von Frank Rosenblatt entwickelt, der damals überzeugt davon war, dass es der Grundstein sei für “Maschinen, die wahrnehmen, erkennen, sich erinnern und reagieren können wie das menschliche Gehirn”. Ende der 1960er Jahre veröffentlichten Martin Minsky und Seymour Papert dann ein Buch, in dem sie argumentierten, dass das Perzeptrons nicht geeignet sei, um Algorithmen zur komplexen Mustererkennung zu entwickeln, insbesondere weil ein einzelnes Perzeptron die logische Entweder-Oder-Funktion nicht darstellen kann, vgl. Aufgabe 13.3. Dies ist zwar wahr, kann aber durch die Verwendung mehrerer ‘gekoppelter’ Perzeptrons überwunden werden, was auf das Konzept der künstlichen neuronalen Netze führt, siehe Kapitel 16. Letzteres war bereits seit Anfang der 1960er bekannt. Aufgrund des Buches von Minsky und Papert kippte jedoch die öffentliche Meinung, Forschungsmittel wurden gestrichen und erst in den 1980er Jahren erlebte die Forschung zu neuronalen Netzen eine Wiederauferstehung.

## Referenzen

Wir notieren zuerst einmal, dass man den Trennungssatz in der Fassung 13.7 elementar beweisen kann, siehe z.B. [GK02, Satz 2.24]. Alternativ folgt derselbe aber auch aus einer Trennungsversionen des Satzes von Hahn-Banach [Rud91, Theorem 3.4(b)]. Für den elementaren Beweis des Satzes von Carathéodory verweisen wir auf [LL16, Theorem 3.3.10]. Der im Kapitel vorgeführte Beweis für Satz 13.13 basiert auf [Col12]. Wir vermerken weiter, dass das Zitat in Bemerkung 13.14(iv) aus [Ros59] entnommen ist. Weiter verweisen wir auf [Lef19] und für mehr historische Details und eine Analyse der ‘Perzeptron-Kontroverse’ auf [Ola96].

## Aufgaben

**Aufgabe 13.1.** Angenommen, der Perzeptronalgorithmus gibt den zusammengefassten Gewichtsvektor  $w = (2, 1, 1)$  aus. Skizzieren Sie, welche Punkte in  $\mathbb{R}^2$  als  $-1$  und welche als  $+1$  klassifiziert werden.

**Aufgabe 13.2.** Finden Sie ein Dataset  $S$  und einen zusammengefassten Gewichtsvektor  $w$ , sodass ein Datenpunkt  $(x, y) \in S$  nicht korrekt klassifiziert wird und ein (einzelnes) Update an diesem Datenpunkt *nicht* dazu führt, dass  $x$  nach dem Update korrekt klassifiziert wird.

**Aufgabe 13.3.** Die folgende Tabelle enthält für sechs gegebene e-Mails jeweils die Angabe, ob die in der Kopfzeile angegebenen Worte in der e-Mail vorkommen, sowie die Information, ob es sich bei der e-Mail um Spam handelt.

	‘und’	‘Bonus’	‘das’	‘Mensa’	‘Vertrag’	Spam
e-Mail 1	✓	✓	✗	✓	✓	Ja
e-Mail 2	✗	✗	✓	✓	✗	Nein
e-Mail 3	✗	✓	✓	✗	✗	Ja
e-Mail 4	✓	✗	✗	✓	✗	Nein
e-Mail 5	✓	✗	✓	✗	✓	Ja
e-Mail 6	✓	✗	✓	✓	✗	Nein

- Erstellen Sie anhand der Tabelle ein numerisches Dataset aus Punkten  $x_1, \dots, x_6 \in \mathbb{R}^5$  mit Labeln  $y_1, \dots, y_6 \in \{-1, 1\}$ .
- Lassen Sie den Perzeptronalgorithmus zyklisch auf den Datenpunkten aus (i) laufen, bis dieser terminiert. Überzeugen Sie sich davon, dass die Laufzeit von der Reihenfolge abhängt, in der Sie die Datenpunkte eingeben.
- Klassifizieren Sie eine neue e-Mail, welche die Worte ‘Bonus’, ‘Vertrag’, ‘das’ und ‘ohne’ enthält als Spam oder No-Spam.

**Aufgabe 13.4.** Wir beschäftigen uns jetzt mit der Frage ob der Perzeptronalgorithmus die logischen Funktionen *und* sowie *entweder-oder*

		AND			XOR
f	f	f	f	f	f
f	w	f	f	w	w
w	f	f	w	f	w
w	w	w	w	w	f

‘erlernen’ kann.

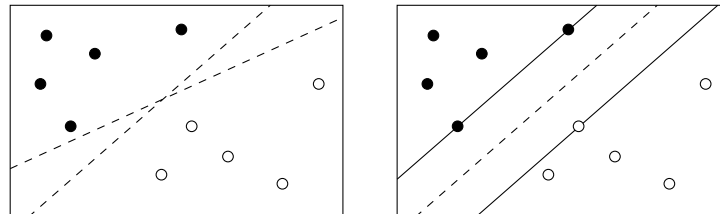
- Ersetzen Sie in der Wertetabelle für AND jeweils ‘w’ durch 1 und ‘f’ durch  $-1$ . Führen Sie dann den Perzeptronalgorithmus auf diesem Trainingsset zyklisch von Hand aus und beginnen Sie mit dem zusammengefassten Gewichtsvektor  $w^{(0)} = (1, 2, 3)$ .
- Machen Sie das gleiche für XOR, beachten Sie Bemerkung 13.14(iv).

**Aufgabe 13.5.** Sei  $S = \{(x_i, y_i) \mid i = 1, \dots, n\} \subseteq \mathbb{R}^d \times \{-1, 1\}$  ein linear trennbares Dataset bei dem nicht alle Label gleich sind. Zeigen Sie, dass der Perzeptronalgorithmus mit beliebigem Startwert  $w^{(0)} \in \mathbb{R}^{d+1}$  in endlich vielen Schritten terminiert und die Laufzeit durch  $\lfloor (C/\gamma)^2 \rfloor$  beschränkt ist, wobei  $C > 0$  eine von  $w^{(0)} \in \mathbb{R}^{d+1}$  abhängige Konstante ist.

## Kapitel 14

# Support Vector Machines

In Kapitel 13 haben wir gesehen, wie man für eine linear trennbare Datenmenge  $D = \{(x_i, y_i) \mid i = 1, \dots, n\} \subseteq \mathbb{R}^d \times \{-1, 1\}$  einen affin-linearen Klassifizierer  $h: \mathbb{R}^d \rightarrow \mathbb{R}$ ,  $h = \text{sign}(\langle w, - \rangle + b)$ , finden kann (Algorithmus 13.11) und wir haben bemerkt, dass per Konstruktion dann immer alle Datenpunkte echt positiven Abstand zur Entscheidungsgrenze  $\mathcal{H} = \{x \in \mathbb{R}^d \mid h(x) = 0\}$  haben (Bemerkung 13.14(iii)). Da es stets mehrere korrekte Klassifizierer gibt (siehe linkes Bild unten) ist es natürlich, nach Klassifizierern zu suchen, bei denen der Abstand zwischen  $\mathcal{H}$  und den am nächsten an  $\mathcal{H}$  liegenden Datenpunkten möglichst groß ist. In der unten gezeichneten Datenmenge gibt es genau einen Klassifizierer, der diesen Abstand maximiert (rechtes Bild) und man kann diesen zeichnerisch leicht ermitteln, vgl. Aufgabe 14.1.



In diesem Kapitel zeigen wir zunächst, dass bei linear trennbaren Daten, die nicht alle das gleiche Label haben, stets genau ein affin-linearer Klassifizierer existiert, der alle Daten korrekt klassifiziert, und der den minimalen Abstand der Datenpunkte zur Entscheidungsgrenze maximiert. Letzteren Klassifizierer nennt man die *Support Vector Machine (SVM)*. Im zweiten Teil des Kapitels geben wir eine Methode an, mithilfe derer die Parameter der SVM durch ein quadratisches Optimierungsproblem berechnet werden können. Hierfür greifen wir auf den Satz von Karush-Kuhn-Tucker zurück, den wir ohne Beweis notieren werden. Außerdem benutzen wir Kapitel 17 über die Existenz und Eindeutigkeit von Minimierern konvexer Funktionen.

**Definition 14.1.** Sei  $h: \mathbb{R}^d \rightarrow \mathbb{R}$ ,  $h(x) = \text{sign}(\langle w, x \rangle + b)$  mit  $0 \neq w \in \mathbb{R}^d$  und  $b \in \mathbb{R}$  ein affin-linearer Klassifizierer für eine Datenmenge  $D = \{(x_i, y_i) \mid i = 1, \dots, n\} \subseteq \mathbb{R}^d \times \{-1, 1\}$  und sei  $\mathcal{H} = \{x \in \mathbb{R}^d \mid \langle w, x \rangle + b = 0\}$  die zugehörige affine Hyperebene. Wir nennen

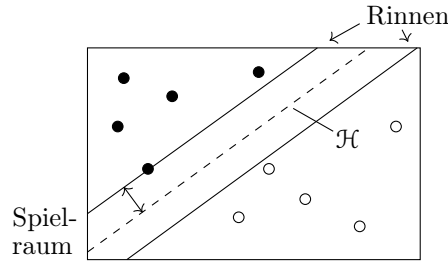
$$\gamma(h) := \min_{i=1, \dots, n} \text{dist}(\mathcal{H}, x_i)$$

den *Spielraum* des Klassifizierers  $h$  und die affinen Hyperebenen

$$\mathcal{R}_{\pm} = \{x \in \mathbb{R}^d \mid \text{dist}(\mathcal{H}, x) = \gamma \text{ und } h(x) = \pm 1\}$$

dessen *Rinnen*.

Das folgende Bild veranschaulicht die eben eingeführten Begriffe. Assoziiert man das Bild mit einer Straße mit Mittellinie  $\mathcal{H}$ , so wird klar woher die Bezeichnung ‘Rinnen’ kommt.



Aus Kapitel 13 wissen wir, dass die Parameter  $(w, b)$  durch den Klassifizierer bis auf einen skalaren Faktor eindeutig bestimmt sind. Wir präzisieren dies nochmal wie folgt und halten es für die nachfolgende Nutzung als Lemma fest.

**Lemma 14.2.** Sei  $h: \mathbb{R}^d \rightarrow \mathbb{R}$ ,  $h = \text{sign}(\langle w, - \rangle + b)$  mit  $0 \neq w \in \mathbb{R}^d$ ,  $b \in \mathbb{R}$  gegeben. Dann gilt  $h = \text{sign}(\langle w', - \rangle + b')$  genau dann wenn ein  $\alpha > 0$  existiert mit  $(w', b') = \alpha(w, b)$ .

*Beweis.* Nach Lemma 13.3 ist  $\mathcal{H} = \{x \in \mathbb{R}^d \mid \langle w, x \rangle + b = 0\}$  durch  $(w, b)$  bis auf einen Faktor  $\alpha \neq 0$  bestimmt. Ist  $(w', b') = \alpha(w, b)$  mit  $\alpha < 0$ , so folgt  $\text{sign}(\langle w', - \rangle + b') = -h$ . Für  $\alpha > 0$  gilt hingegen die Gleichheit.  $\square$

Wir benötigen nun einen Weg, um den Spielraum  $\gamma(h)$  anhand der Parameter  $(w, b)$ , durch die  $h$  gegeben ist, zu berechnen. Das folgende Lemma gibt hierfür vorbereitend eine Formel an, mit der der Abstand eines Punktes von der zu  $h$  gehörenden affinen Hyperebene bestimmt werden kann.

**Lemma 14.3.** Seien  $0 \neq w \in \mathbb{R}^d$ ,  $b \in \mathbb{R}$  und sei  $\mathcal{H} = \{x \in \mathbb{R}^d \mid \langle w, x \rangle + b = 0\}$ . Für  $x_0 \in \mathbb{R}^d$  gilt dann  $\text{dist}(\mathcal{H}, x_0) = \frac{1}{\|w\|} |\langle w, x_0 \rangle + b|$ .

*Beweis.* Wir setzen zunächst  $w' := \frac{w}{\|w\|}$  und  $b' := \frac{b}{\|w\|}$ . Dann liefert  $(w', b')$  nach Lemma 14.2 den gleichen Klassifizierer, sowie die gleiche affine Hyperebene, und die Behauptung des Lemmas lautet

$$\min_{x \in \mathcal{H}} \|x_0 - x\| = |\langle w', x_0 \rangle + b'|.$$

Wir definieren  $x_1 := x_0 - (\langle w', x_0 \rangle + b')w'$  und zeigen die folgenden zwei Gleichungen

$$\min_{x \in \mathcal{H}} \|x_0 - x\| \stackrel{(1)}{=} \|x_0 - x_1\| \stackrel{(2)}{=} |\langle w', x_0 \rangle + b'|.$$

Durch Einsetzen sehen wir zunächst

$$\begin{aligned}\langle w', x_1 \rangle + b' &= \langle w', x_0 - (\langle w', x_0 \rangle + b')w' \rangle + b' \\ &= \langle w', x_0 \rangle - (\langle w', x_0 \rangle + b')\langle w', w' \rangle + b' = 0,\end{aligned}$$

also  $x_1 \in \mathcal{H}$  und wir haben schonmal eine Ungleichung in (1) erledigt. Für die andere sei  $x \in \mathcal{H}$  beliebig. Dann ist  $\langle w', x_1 - x \rangle = 0$  und daher haben wir

$$\begin{aligned}\|x_0 - x\|^2 &= \|x_0 - x_1 + x_1 - x\|^2 \\ &= \|x_0 - x_1\|^2 + 2\langle x_0 - x_1, x_1 - x \rangle + \|x_1 - x\|^2 \\ &\geq \|x_0 - x_1\|^2 + 2\langle (\langle w', x_0 \rangle + b')w', x_1 - x \rangle \\ &= \|x_0 - x_1\|^2 + 2(\langle w', x_0 \rangle + b')\langle w', x_1 - x \rangle \\ &= \|x_0 - x_1\|^2\end{aligned}$$

was den Beweis von (1) durch Wurzelziehen abschließt. Durch Einsetzen sehen wir

$$\text{dist}(\mathcal{H}, x_0) = \|x_0 - x_1\| = \left\| \left( \left\langle \frac{w}{\|w\|}, x_0 \right\rangle + b \right) \frac{w}{\|w\|} \right\| = \frac{1}{\|w\|} |\langle w, x_0 \rangle + b|$$

und damit ist auch Gleichung (2) bewiesen.  $\square$

Betrachten wir, statt eines einzelnen Punktes, jetzt wieder eine ganze Datenmenge  $D = \{(x_i, y_i) \mid i = 1, \dots, n\} \subseteq \mathbb{R}^d \times \{-1, 1\}$ , so folgt aus obigem, dass der Spielraum eines korrekten Klassifizierers  $h_{w,b} = \text{sign}(\langle w, \cdot \rangle + b)$  mit  $0 \neq w \in \mathbb{R}^d$  und  $b \in \mathbb{R}$  durch

$$\gamma(h_{w,b}) = \min_{i=1,\dots,n} \text{dist}(\mathcal{H}, x_i) = \frac{1}{\|w\|} \min_{i=1,\dots,n} |\langle w, x_i \rangle + b|$$

gegeben ist und der rechte Ausdruck muss nun über eine Menge von Parametern maximiert werden, sodass durch deren Elemente alle korrekten Klassifizierer beschrieben werden. Wegen Lemma 14.2 muss diese Menge nicht alle  $(w, b)$ , die zu korrekten Klassifizierern führen, enthalten und man könnte meinen, dass die Reduktion auf  $(w, b)$  mit  $\|w\| = 1$  sinnvoll ist, wäre doch  $w$  dann der Normaleneinheitsvektor der zu  $h$  gehörenden affinen Hyperebene. Dieser Zugang, obwohl machbar, führt allerdings auf ein eher unersprißliches Optimierungsproblem, siehe Aufgabe 14.3. Wir schlagen daher einen anderen Weg ein und notieren zunächst das Folgende, welches uns im Wesentlichen auch schon in Bemerkung 13.12 begegnet ist.

**Lemma 14.4.** *Sei  $D = \{(x_i, y_i) \mid i = 1, \dots, n\} \subseteq \mathbb{R}^d \times \{-1, 1\}$  eine Datenmenge. Für  $(w, b) \in \mathbb{R}^{d+1}$  mit  $w \neq 0$  sind die folgenden Aussagen äquivalent.*

- (i) *Die Funktion  $h: \mathbb{R}^d \rightarrow \mathbb{R}$ ,  $h = \text{sign}(\langle w, \cdot \rangle + b)$  klassifiziert korrekt.*
- (ii) *Es gilt  $y_i(\langle w, x_i \rangle + b) > 0$  für alle  $i = 1, \dots, n$ .*
- (iii) *Es gilt  $\min_{i=1,\dots,n} y_i(\langle w, x_i \rangle + b) > 0$ .*

*Beweis.* Es genügt, sich zu überlegen, dass  $h(x_i) = y_i$  genau dann gilt, wenn das

Produkt  $y_i \cdot h(x_i) > 0$  ist und dann weiter zu benutzen, dass  $h(x_i) > 0$  genau dann gilt, wenn  $\langle w, x_i \rangle + b > 0$  ausfällt. Dies zeigt die Äquivalenz von (i) und (ii). Die Äquivalenz von (ii) und (iii) ist trivial.  $\square$

Datenmengen sind per Definition genau dann linear trennbar, wenn ein affin-linearer Klassifizierer existiert der alle Daten korrekt klassifiziert. Dies gilt auch im pathologischen Fall, dass alle gegebenen Daten das gleiche Label haben, eingeschlossen den Fall, dass  $D$  leer ist. In diesem Fall gibt es allerdings Klassifizierer mit  $w = 0$ , wodurch  $\mathcal{H}$  entweder leer ausfällt oder der ganze Raum ist, auf jeden Fall aber keine affine Hyperebene. Da die vorgenannten Fälle vom Anwendungsstandpunkt her uninteressant sind, schließen wir sie in den folgenden Sätzen aus.

**Bemerkung 14.5.** Um uns im Folgenden nicht in Pathologien zu verlieren, setzen wir ab sofort immer voraus, dass unsere Datenmenge  $D = \{(x_i, y_i) \mid i = 1, \dots, n\} \subseteq \mathbb{R}^d \times \{-1, 1\}$  linear trennbar ist und nicht alle Label gleich sind. Letzteres impliziert, dass  $S \neq \emptyset$  gilt und, dass jeder korrekte affine-lineare Klassifizierer von der Form  $h = \text{sign}(\langle w, \cdot \rangle + b)$  mit  $w \neq 0$  ist.

Um nun die in Lemma 14.2 beschriebene Mehrdeutigkeit loszuwerden, ohne korrekte Klassifizierer auszuschließen, und um am Ende auf ein praktikables Optimierungsproblem zu kommen, schränken wir uns wie folgt auf eine minimale Parametermenge ein.

**Proposition 14.6.** Sei  $D = \{(x_i, y_i) \mid i = 1, \dots, n\} \subseteq \mathbb{R}^d \times \{-1, 1\}$  eine linear trennbare Datenmenge und seien nicht alle Label gleich. Mit

$$\mathcal{R}(S) := \{(w, b) \in \mathbb{R}^{d+1} \mid \min_{i=1, \dots, n} y_i(\langle w, x_i \rangle + b) = 1\}$$

$$\mathcal{K}(S) := \{h: \mathbb{R}^d \rightarrow \mathbb{R} \mid h \text{ ist korrekter affin-linearer Klassifizierer für } S\}$$

gelten die folgenden Aussagen.

- (i) Die Abbildung  $\phi: \mathcal{R}(S) \rightarrow \mathcal{K}(S)$ ,  $(w, b) \mapsto h_{w,b} := \text{sign}(\langle w, \cdot \rangle + b)$  ist bijektiv.
- (ii) Für  $(w, b) \in \mathcal{R}(S)$  gilt  $\gamma(h_{w,b}) = \frac{1}{\|w\|}$ .

*Beweis.* (i) Nach Lemma 14.4 führt jedes  $(w, b) \in \mathcal{R}(S)$  auf einen korrekten Klassifizierer,  $\phi$  ist also wohldefiniert. Ist  $h = \text{sign}(\langle w, \cdot \rangle + b)$  ein beliebiger korrekter Klassifizierer, so ist, wieder nach Lemma 14.4,

$$m := \min_{i=1, \dots, n} y_i(\langle w, x_i \rangle + b) > 0$$

und  $(w/m, b/m)$  gehört zu  $\mathcal{R}(S)$  und liefert nach Lemma 14.2 den gleichen Klassifizierer. Somit ist  $\phi$  surjektiv. Sind schließlich  $(w, b), (w', b') \in \mathcal{R}(S)$  gegeben mit  $h_{w,b} = h_{w',b'}$ , so existiert nach Lemma 14.2 ein  $\alpha > 0$  mit  $(w', b') = (\alpha w, \alpha b)$ . Per Definition von  $\mathcal{R}(S)$  folgt

$$\alpha = \alpha \cdot 1 = \alpha \min_{i=1, \dots, n} y_i(\langle w, x_i \rangle + b) = \min_{i=1, \dots, n} y_i(\langle w', x_i \rangle + b') = 1$$



und  $\phi$  ist injektiv.

(ii) Sei  $(w, b) \in \mathcal{R}(S)$  beliebig. Nach Lemma 14.4 ist dann  $y_i(\langle w, x_i \rangle + b) > 0$  für alle  $i$ , also  $|\langle w, x_i \rangle + b| = y_i(\langle w, x_i \rangle + b)$  für alle  $i$ . Mit Lemma 14.3 folgt

$$\gamma(h_{w,b}) = \min_{i=1,\dots,n} \text{dist}(\mathcal{H}, x_i) = \frac{1}{\|w\|} \min_{i=1,\dots,n} |\langle w, x_i \rangle + b| = \frac{1}{\|w\|} \cdot 1$$

$\uparrow$   
 $(w,b) \in \mathcal{R}(S)$

was den Beweis beendet.  $\square$

**Bemerkung 14.7.** Die Forderung

$$\min_{i=1,\dots,n} y_i(\langle w, x_i \rangle + b) \stackrel{!}{=} 1$$

mit der  $\mathcal{R}(S)$  definiert wurde, nennt man die *Rinnenbedingung*. Erfüllt nämlich ein Parameterpaar  $(w, b)$  dieselbe, so ist  $h_{w,b}$  ein korrekter Klassifizierer und die Rinnen des Klassifizierers können per

$$\mathcal{R}_{\pm} = \{x \in \mathbb{R}^d \mid \langle w, x \rangle + b = \pm 1\}$$

geschrieben werden. Es gilt nämlich  $\text{dist}(x, \mathcal{H}) = \frac{1}{\|w\|} |\langle w, x \rangle + b| = \gamma |\langle w, x \rangle + b|$  und letzteres ist gleich  $\gamma$  genau dann, wenn  $\langle w, x \rangle + b = \pm 1$  gilt.

Die Rinnenbedingung ermöglicht die folgende Charakterisierung der affin-linearen Klassifizierer mit maximalem Spielraum durch ein Optimierungsproblem.

**Satz 14.8.** Sei  $D = \{(x_i, y_i) \mid i = 1, \dots, n\} \subseteq \mathbb{R}^d \times \{-1, 1\}$  eine linear trennbare Datenmenge bei der nicht alle Label gleich sind. Dann ist  $h^*: \mathbb{R}^d \rightarrow \mathbb{R}$  ein korrekter affin-linear Klassifizierer mit maximalem Spielraum genau dann, wenn

$$(w^*, b^*) \in \underset{(w,b) \in \mathcal{R}(S)}{\text{argmax}} \frac{1}{\|w\|}$$

existiert mit  $h^* = h_{w^*, b^*}$ .

*Beweis.* “ $\implies$ ” Sei  $h^*$  wie angegeben, also insbesondere  $h^* \in \mathcal{K}(S)$ . Nach Proposition 14.6(i) gibt es  $(w^*, b^*) \in \mathcal{R}(S)$  mit  $h^* = h_{w^*, b^*}$ . Seien nun  $(w, b) \in \mathcal{R}(S)$  beliebig und sei  $h := h_{w,b} \in \mathcal{K}(S)$  der zugehörige Klassifizierer. Dann gilt mit Proposition 14.6(ii)

$$1/\|w\| = \gamma(h) \leq \gamma(h^*) = 1/\|w^*\|$$

und  $(w^*, b^*)$  ist Maximierer.

“ $\impliedby$ ” Sei  $(w^*, b^*)$  wie angegeben. Dann ist  $h^* := h_{w^*, b^*} \in \mathcal{K}(S)$  und für jedes andere  $h \in \mathcal{K}(S)$  können wir  $(w, b) \in \mathcal{R}(S)$  wählen sodass  $h = h_{w,b}$  gilt. Analog zu oben folgt

$$\gamma(h) = 1/\|w\| \leq 1/\|w^*\| = \gamma(h^*)$$

durch Anwendung von Proposition 14.6(ii).  $\square$

Der obige Satz sagt weder etwas über die Existenz, noch über die Eindeutigkeit, von  $h^*$  aus. Dies erreicht aber der folgende Satz und zwar durch eine Umformulierung des obigen Optimierungsproblems in ein (strikt) konvexes Problem.

**Satz 14.9.** *Sei  $D = \{(x_i, y_i) \mid i = 1, \dots, n\} \subseteq \mathbb{R}^d \times \{-1, 1\}$  eine Datenmenge, die linear trennbar ist und bei der nicht alle Label gleich sind. Dann existiert genau ein affin-linearer Klassifizierer  $h^*: \mathbb{R}^d \rightarrow \mathbb{R}$  mit maximalem Spielraum. Dieser Klassifizierer wird genau durch die Parameter  $\alpha \cdot (w^*, b^*)$  mit  $\alpha > 0$  gegeben, wobei*

$$(w^*, b^*) = \underset{\substack{(w,b) \in \mathbb{R}^{d+1} \\ \forall i: y_i(\langle w, x_i \rangle + b) \geq 1}}{\operatorname{argmin}} \|w\|^2$$

die einzige Lösung des angegebenen konvexen Minimierungsproblems ist.

*Beweis.* ① Wir zeigen die Existenz von  $(w^*, b^*)$ . Da  $D$  linear trennbar ist, existieren  $(w_L, b_L) \in \mathbb{R}^{d+1}$  und  $\varepsilon > 0$  sodass  $y_i(\langle w_L, x_i \rangle + b_L) \geq \varepsilon$  für alle  $i = 1, \dots, n$  gilt und  $w_L \neq 0$  ist. Es folgt dann

$$y_i\left(\left\langle \frac{w_L}{\varepsilon}, x_i \right\rangle + \frac{b_L}{\varepsilon}\right) \geq 1$$

und

$$\begin{aligned} M &:= \{(w, b) \in \mathbb{R}^{d+1} \mid \forall i = 1, \dots, n: y_i(\langle w, x_i \rangle + b) \geq 1\} \\ &= \bigcap_{i=1}^n \{(w, b) \mid y_i(\langle w, x_i \rangle + b) \geq 1\} \end{aligned}$$

ist als Schnitt endlich vieler abgeschlossener affiner Halbräume (Division durch  $y_i$  führt entweder auf  $\langle w, x_i \rangle + b \geq 1$  oder auf  $\langle w, x_i \rangle + b \leq -1$ ) eine nichtleere, konvexe und abgeschlossene Menge. Wir definieren  $f: M \rightarrow \mathbb{R}$ ,  $f(w, b) := \|w\|^2$  und erhalten

$$\inf_{(w,b) \in M} f(w, b) \leq f\left(\frac{w_L}{\varepsilon}, \frac{b_L}{\varepsilon}\right) = \frac{\|w_L\|^2}{\varepsilon^2}.$$

Setzen wir  $C := \|w_L\|/\varepsilon$ , so genügt es also, über  $(w, b) \in M$  mit  $\|w\| \leq C$  zu minimieren. Für solche  $(w, b)$  gilt  $y_i \langle w, x_i \rangle + y_i b = y_i(\langle w, x_i \rangle + b) \geq 1$  woraus mit der Cauchy-Schwarz-Bunjakowski-Ungleichung

$$-y_i b \leq y_i \langle w, x_i \rangle - 1 \leq |\langle w, x_i \rangle| + 1 \leq \|w\| \|x_i\| + 1 \leq C \|x_i\| + 1 \leq R$$

für alle  $i$  folgt, wenn wir  $R := 1 + C \max_{j=1, \dots, n} \|x_j\|$  definieren. Per Voraussetzung gibt es  $i, j \in \{1, \dots, n\}$  mit  $y_i = 1$  und  $y_j = -1$ . Mit diesen ergibt sich  $-b \leq R$  und  $b \leq R$  also  $|b| \leq R$  und es folgt

$$\underset{(w,b) \in M}{\operatorname{argmin}} f(w, b) = \underset{\substack{(w,b) \in M \\ \|w\| \leq C, |b| \leq R}}{\operatorname{argmin}} f(w, b) \neq \emptyset$$

weil wir eine stetige Funktion über ein Kompaktum minimieren.

② Als Vorbereitung des Eindeutigkeitsbeweises für  $(w^*, b^*)$  zeigen wir zuerst, dass

für jeden Minimierer  $(w^*, b^*) \in \operatorname{argmin}_{(w,b) \in M} f(w, b)$  Indizes  $i, j \in \{1, \dots, n\}$  existieren mit

$$y_i = 1 \text{ und } \langle w^*, x_i \rangle + b^* = 1 \text{ sowie } y_j = -1 \text{ und } \langle w^*, x_j \rangle + b^* = -1.$$

Angenommen, dies ist nicht so, dann es gibt einen Minimierer  $(w^*, b^*)$ , sodass für alle  $i \in \{1, \dots, n\}$  mit  $y_i = 1$  stets  $\langle w^*, x_i \rangle + b^* > 1$  gilt, und es ist

$$m := \min_{\substack{i=1, \dots, n \\ y_i=1}} \langle w^*, x_i \rangle + b^* - 1 > 0.$$

Wir setzen  $\hat{w} := \frac{w^*}{1+m/2}$  und  $\hat{b} := \frac{b^* - m/2}{1+m/2}$  und erhalten für  $i$  mit  $y_i = 1$

$$\langle \hat{w}, x_i \rangle + \hat{b} = \frac{\langle w^*, x_i \rangle + b^* - m/2}{1 + m/2} \geq \frac{1 + m/2}{1 + m/2} = 1$$

sowie für  $i$  mit  $y_i = -1$

$$\langle \hat{w}, x_i \rangle + \hat{b} = \frac{\langle w^*, x_i \rangle + b^* - m/2}{1 + m/2} \leq \frac{-1 - m/2}{1 + m/2} = -1.$$

Zusammengenommen folgt  $(\hat{w}, \hat{b}) \in M$  aber  $\|\hat{w}\| < \|w^*\|$  im Widerspruch dazu, dass  $(w^*, b^*) \in D$  ein Minimierer von  $f$  ist. Die Existenz von  $j$  zeigt man analog.

③ Wir zeigen jetzt die Eindeutigkeit von  $(w^*, b^*)$ . Angenommen, wir haben zwei Minimierer  $(w_1^*, b_1^*)$  und  $(w_2^*, b_2^*)$ . Da  $M$  konvex ist, ist  $M_1 := \{w \in \mathbb{R}^d \mid (w, b) \in M\}$  ebenfalls konvex und  $f|_{M_1} = \|\cdot\|^2$  ist strikt konvex nach Beispiel 17.15. Es folgt also nach 17.14, dass  $w_1^* = w_2^* =: w^*$  gilt. Angenommen, es gilt  $b_1^* < b_2^*$ . Dann wählen wir für  $(w^*, b_1^*)$  und  $(w^*, b_2^*)$  jeweils  $i_1$  und  $i_2$  wie in Teil ② und erhalten

$$y_{i_2}(\langle w^*, x_{i_2} \rangle + b_1^*) = \langle w^*, x_{i_2} \rangle + b_1^* < \langle w^*, x_{i_2} \rangle + b_2^* = 1$$

im Widerspruch dazu, dass  $(w^*, b_2^*)$  in  $M$  liegt. Ist  $b_1^* > b_2^*$ , so vertauscht man die Rollen von  $i_1$  und  $i_2$ .

④ Wir müssen nun den Bogen schlagen von den Parametern zum Klassifizierer. Die bisherigen Teile des Beweises zeigen in der folgenden Gleichungskette

$$(w^*, b^*) = \operatorname{argmin}_{\substack{(w,b) \in \mathbb{R}^{d+1} \\ \forall i: y_i(\langle w, x_i \rangle + b) \geq 1}} \|w\|^2 = \operatorname{argmin}_{\substack{(w,b) \in \mathbb{R}^{d+1} \\ \min_{i=1, \dots, n} y_i(\langle w, x_i \rangle + b) = 1}} \|w\|^2 = \operatorname{argmin}_{(w,b) \in \mathcal{R}(S)} \|w\| = \operatorname{argmax}_{(w,b) \in \mathcal{R}(S)} \frac{1}{\|w\|}$$

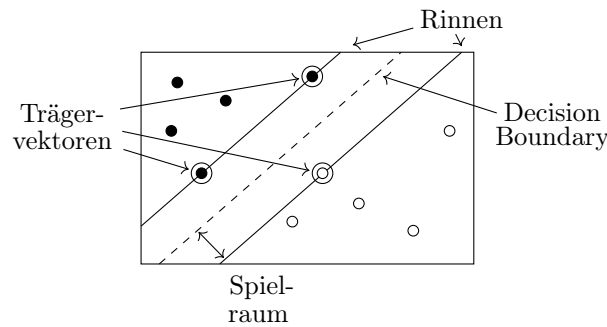
offenbar die erste Gleichung, aber auch die zweite, denn Teil ② impliziert insbesondere, dass es ein  $i$  gibt, sodass  $y_i(\langle w, x_i \rangle + b) = 1$  gilt. Für die dritte Gleichung benutzen wir die Abkürzung  $\mathcal{R}(S)$  aus 14.6 und lassen das Quadrat weg, weil  $(\cdot)^2: (0, \infty) \rightarrow \mathbb{R}$  streng monoton wachsend ist. Die letzte Gleichung gilt schließlich, weil wir bereits wissen, dass das Minimum  $\|w^*\| > 0$  ist und  $(\cdot)^{-1}: (0, \infty) \rightarrow \mathbb{R}$  streng monoton fallend ist.

⑤ Nach 14.8 ist nun  $h_{w^*,b^*} = \operatorname{argmax}_{h \in \mathcal{H}(S)} \gamma(h)$ . Die Aussage über weitere Parameter, die ebenfalls  $h^* := h_{w^*,b^*}$  liefern, folgt aus Lemma 14.2.  $\square$

Beweisteil ③ benutzt nur die erste der in Teil ② gezeigten Aussagen, nämlich dass mindestens einer der positiv gelabelten Punkte in der positiven Rinne  $\mathcal{R}_+$  liegt. Wie im Beweis bemerkt, gibt es aber auch immer mindestens einen negativ gelabelten Punkt in der negativen Rinne  $\mathcal{R}_-$ .

**Definition 14.10.** Sei  $D = \{(x_i, y_i) \mid i = 1, \dots, n\} \subseteq \mathbb{R}^d \times \{-1, 1\}$  eine linear trennbare Datenmenge bei der nicht alle Label gleich sind. Sei  $h^*: \mathbb{R}^d \rightarrow \mathbb{R}$  der nach 14.9 eindeutige Klassifizierer mit maximalem Spielraum  $\gamma^* > 0$ . Seien  $\mathcal{H}$  und  $\mathcal{R}_\pm$  die Entscheidungsgrenze sowie die Rinnen von  $h^*$ . Dann heißen diejenigen Datenpunkte  $x_i$ , die  $\operatorname{dist}(\mathcal{H}, x_i) = \gamma^*$  erfüllen, die *Trägervektoren* und  $h^*$  heißt die zu  $S$  gehörende *Support Vector Machine* oder kurz *SVM*.

Das folgende Bild illustriert nochmal alle neuen Begriffe an unserem anfänglichen Beispiel.



Man könnte jetzt versuchen, die Funktion  $\|\cdot\|^2$  über den Bereich  $D \subseteq \mathbb{R}^{d+1}$  zu minimieren, vgl. Aufgabe 14.4. Stattdessen werden wir die Maximierungsaufgabe in Satz 14.9 auf ein quadratisches Optimierungsproblem zurückführen, für das effiziente Lösungsmethoden existieren. Der zentrale Satz lautet wie folgt.

**Satz 14.11.** Sei  $D = \{(x_i, y_i) \mid i = 1, \dots, n\} \subseteq \mathbb{R}^d \times \{-1, 1\}$  eine Datenmenge bei der nicht alle Label gleich sind und bezeichne mit  $y := (y_1, \dots, y_n) \in \mathbb{R}^n$  den Vektor der Label. Dann ist  $S$  genau dann linear trennbar, wenn mindestens eine Lösung des assoziierten quadratischen Optimierungsproblems

$$\lambda^* \in \operatorname{argmin}_{\substack{\lambda \in \mathbb{R}_{\geq 0}^n \\ \langle \lambda, y \rangle = 0}} \frac{1}{2} \sum_{i,j=1}^n \lambda_i \lambda_j y_i y_j \langle x_i, x_j \rangle - \sum_{i=1}^n \lambda_i$$

existiert. In diesem Fall ist jede Lösung  $\lambda^*$  verschieden von Null und führt durch die folgenden Schritte auf die Support Vector Machine  $h^* = h_{w^*,b^*}$  für  $S$ :

1. Definiere  $w^* := \lambda_1^* y_1 x_1 + \dots + \lambda_n^* y_n x_n$ .
2. Wähle  $i_0 \in \{1, \dots, n\}$  derart, dass  $\lambda_{i_0} \neq 0$  gilt.

3. Definiere  $b^* := y_{i_0} - \langle w^*, x_{i_0} \rangle$ .

*Beweis.* Wir benötigen hierfür eine Version des sogenannten Karush-Kuhn-Tucker-Theorems, das genau auf unsere Anwendungssituationen passt. Am Ende des Kapitels finden sich detaillierte Referenzen zu den einzelnen Teilen.

Satz 14.12. (Karush-Kuhn-Tucker) Sei  $f: \mathbb{R}^p \rightarrow \mathbb{R}$  konvex und stetig differenzierbar und seien  $g_1, \dots, g_q, h_1, \dots, h_r: \mathbb{R}^p \rightarrow \mathbb{R}$  affin-linear. Sei

$$L: \mathbb{R}^p \times \mathbb{R}^q \times \mathbb{R} \rightarrow \mathbb{R}, \quad L(x, \theta, \mu) := f(x) - \sum_{i=1}^q \theta_i g_i(x) + \mu \sum_{j=1}^r \mu_j h_j(x)$$

die zugehörige Lagrangefunktion. Dann gilt:

(i) Für jeden Minimierer

$$x^* \in \underset{\substack{x \in \mathbb{R}^p \\ g(x) \geq 0, h(x) = 0}}{\operatorname{argmin}} f(x)$$

gibt es Lagrangemultiplikatoren  $\theta^* = (\theta_1^*, \dots, \theta_q^*) \in \mathbb{R}^q$  und  $\mu^* = (\mu_1^*, \dots, \mu_r^*) \in \mathbb{R}^r$ , sodass  $(x^*, \theta^*, \mu^*)$  die KKT-Bedingungen

$$(KKT-1) \quad \nabla_x L(x^*) = 0$$

$$(KKT-2) \quad h(x) = 0$$

$$(KKT-3) \quad \theta^* \geq 0, g(x^*) \geq 0, \langle \theta^*, g(x^*) \rangle = 0$$

erfüllt, wobei Ungleichungen bei Vektoren eintragsweise gemeint sind.

(ii) Für jedes Tripel  $(x^*, \theta^*, \mu^*) \in \mathbb{R}^p \times \mathbb{R}^q \times \mathbb{R}^r$ , das die KKT-Bedingungen erfüllt, ist  $x^*$  ein Maximierer des oben angegebenen Optimierungsproblems.  $\diamond$

Im Folgenden werden wir das konvexe Problem aus Satz 14.9 wie auch das quadratische Problem aus diesem Satz jeweils in die obige Form bringen und für jedes die KKT-Bedingungen verwenden. In der Tat benutzen wir jeden der zwei Teile von Satz 14.12 unten jeweils zweimal.

Wir beginnen mit der Existenzaussage.

① Sei dazu  $D$  linear trennbar und  $(w^*, b^*)$  die eindeutige Lösung des konvexen Optimierungsproblems in Satz 14.9. Durch Hinzufügen eines rein kosmetischen Faktors  $1/2$  erhalten wir

$$(w^*, b^*) = \underset{\substack{(w,b) \in \mathbb{R}^{d+1} \\ \forall i: y_i(\langle w, x_i \rangle + b) \geq 1}}{\operatorname{argmin}} \frac{1}{2} \|w\|^2 \quad (14.1)$$

was von der in Satz 14.12 angegebenen Form ist, wenn wir  $p = d + 1$ ,  $q = n$ ,  $r = 0$ , sowie  $f(w, b) = \frac{1}{2} \|w\|^2$  und  $g_i(w, b) = y_i(\langle w, x_i \rangle + b) - 1$  setzen. Die Lagrangefunktion lautet (mit  $\lambda$  statt  $\theta$  aus guten Gründen!)

$$L(w, b, \lambda) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \lambda_i [y_i(\langle w, x_i \rangle + b) - 1]$$

und nach Satz 14.12(i) existiert  $\lambda^* \in \mathbb{R}^n$  sodass  $(w^*, b^*, \lambda^*)$  die KKT-Bedingungen erfüllt. Nun liefert (KKT-1)

$$0 = \nabla_{(w,b)} L(w^*, b^*, \lambda^*) = \left( \frac{\partial L}{\partial w}, \frac{\partial L}{\partial b} \right) (w^*, b^*, \lambda^*) = \left( w^* - \sum_{i=1}^n \lambda_i^* y_i x_i, - \sum_{i=1}^n \lambda_i^* y_i \right)$$

woraus durch Auflösen  $w^* = \lambda_1^* y_1 x_1 + \dots + \lambda_n^* y_n x_n$  und  $\langle \lambda^*, y \rangle = 0$  folgen. Nach (KKT-3) gelten ferner  $\lambda^* \geq 0$ ,  $(y_i(\langle w^*, x_i \rangle + b^*) - 1) \geq 0$  für alle  $i = 1, \dots, n$ , sowie  $\langle \lambda^*, (y_i(\langle w^*, x_i \rangle + b^*) - 1)_{i=1, \dots, n} \rangle = 0$ .

② Wir behaupten nun, dass

$$\lambda^* \in \underset{\substack{\lambda \in \mathbb{R}_{\geq 0}^n \\ \langle \lambda, y \rangle = 0}}{\operatorname{argmin}} \frac{1}{2} \sum_{i,j=1}^n \lambda_i \lambda_j y_i y_j \langle x_i, x_j \rangle - \sum_{i=1}^n \lambda_i \quad (14.2)$$

gilt und wir zeigen dies mithilfe von Satz 14.12(ii). In der Tat ist (14.2) auch von der in Satz 14.12 betrachteten Form, aber diesmal ist  $p = q = n$ ,  $r = 1$ ,  $f(\lambda)$  ist der in (14.2) zu minimierende Ausdruck,  $g_i(\lambda) = \lambda_i$  und  $h_1(\lambda) = \langle \lambda, y \rangle$ . Wir weisen insbesondere darauf hin, dass  $\lambda$  jetzt die Rolle der Variable einnimmt und nicht mehr Lagrangemultiplikator ist! Die Lagrangefunktion lautet

$$\begin{aligned} L(\lambda, \theta, \mu) &= f(\lambda) - \sum_{i=1}^n \theta_i \lambda_i + \mu_1 \langle \lambda, y \rangle \\ &= \frac{1}{2} \sum_{i,j=1}^n \lambda_i \lambda_j y_i y_j \langle x_i, x_j \rangle - \sum_{i=1}^n \lambda_i - \sum_{i=1}^n \theta_i \lambda_i + \mu_1 \sum_{i=1}^n \lambda_i y_i \\ &= \frac{1}{2} \left\langle \begin{bmatrix} \lambda_1 \\ \vdots \\ \lambda_n \end{bmatrix}, \underbrace{\begin{bmatrix} y_1 y_1 \langle x_1, x_1 \rangle & \dots & y_1 y_n \langle x_1, x_n \rangle \\ \vdots & & \vdots \\ y_n y_1 \langle x_n, x_1 \rangle & \dots & y_n y_n \langle x_n, x_n \rangle \end{bmatrix}}_{=:P} \begin{bmatrix} \lambda_1 \\ \vdots \\ \lambda_n \end{bmatrix} \right\rangle + \left\langle \underbrace{\begin{bmatrix} -1 - \theta_1 + \mu_1 y_1 \\ \vdots \\ -1 - \theta_n + \mu_1 y_n \end{bmatrix}}_{=:q}, \begin{bmatrix} \lambda_1 \\ \vdots \\ \lambda_n \end{bmatrix} \right\rangle \end{aligned}$$

Schauen wir jetzt nochmal auf die letzten zwei Zeilen in Beweisteil ①, so sehen wir, dass  $\lambda^*$  die Bedingung (KKT-2) und die beiden letzten Bedingungen aus (KKT-3) für das Problem (14.2) erfüllt. Es bleibt also zu zeigen, dass  $\theta^* \geq 0$  und  $\mu^* \in \mathbb{R}$  existieren, sodass (KKT-1) gilt. Um dies zu sehen, berechnen wir

$$\nabla_{\lambda} L(\lambda, \theta, \mu) = \nabla_{\lambda} \left( \frac{1}{2} \langle \lambda, P \lambda \rangle + \langle q, \lambda \rangle \right) = P \lambda + q$$

und notieren hiervon den  $i$ -ten Eintrag, ausgewertet in  $\lambda = \lambda^*$ , und unter Ausnutzung der Gleichung  $w^* = \lambda_1^* y_1 x_1 + \dots + \lambda_n^* y_n x_n$  die wir am Ende von Teil ① gezeigt haben. Dies liefert

$$\begin{aligned} (\nabla_{\lambda} L(\lambda^*, \theta, \mu))_i &= \sum_{j=1}^n y_i y_j \langle x_i, x_j \rangle \lambda_j^* - 1 - \theta_i + \mu_1 y_i \\ &= y_i \left\langle x_i, \sum_{j=1}^n \lambda_j^* y_j x_j \right\rangle - 1 - \theta_i + \mu_1 y_i \end{aligned}$$

$$\begin{aligned}
&= y_i \langle x_i, w^* \rangle - 1 - \theta_i + \mu_1 y_i \\
&= y_i (\langle x_i, w^* \rangle + \mu_1) - 1 - \theta_i.
\end{aligned}$$

Wählen wir  $\mu_1^* := b^*$  und  $\theta_i^* := y_i(\langle x_i, w^* \rangle + \mu_1^*) - 1$  so erhalten wir  $\theta^* \geq 0$ , wegen der am Ende von Teil ③ gezeigten (Un-)gleichungen, und  $\nabla_\lambda L(\lambda^*, \theta^*, \mu^*) = 0$  per Konstruktion. Nach Satz 14.12(ii) löst  $\lambda^*$  das quadratische Optimierungsproblem wie behauptet.

Jetzt kommen wir zur Allaussage.

③ Sei dazu  $\lambda^*$  ein beliebiger Minimierer des quadratischen Problems im Satz und damit auch des Problems (14.2). Wir haben Teil ② das Letztere bereits auf die entsprechende Form gebracht und wissen daher nun, dass aufgrund von Satz 14.12(i) Multiplikatoren  $\theta^* \geq 0$  und  $\mu \geq 0$  existieren, sodass  $(\lambda^*, \theta^*, \mu^*)$  die KKT-Bedingungen erfüllt. Den Gradient der Lagrangefunktion haben wir ebenfalls in Teil ② schon berechnet, erhalten jetzt also

$$0 = (\nabla_\lambda L(\lambda^*, \theta^*, \mu^*))_i = y_i \left\langle x_i, \sum_{j=1}^n \lambda_j^* y_j x_j \right\rangle - 1 - \theta_i^* + \mu_1^* y_i$$

für alle  $i = 1, \dots, n$  wegen (KKT-1) und überdies  $\langle \lambda^*, y \rangle = 0$  wegen (KKT-2) sowie  $\theta^* \geq 0$ ,  $\lambda^* \geq 0$  und  $\langle \theta^*, \lambda^* \rangle = 0$  wegen (KKT-3).

④ Wir definieren nun  $w^* := \lambda_1^* y_1 x_1 + \dots + \lambda_n^* y_n x_n$ ,  $b^* := \mu_1^*$  und behaupten, dass  $(w^*, b^*, \lambda^*)$  die KKT-Bedingungen des Problems (14.1) erfüllt. In Teil ① hatten wir bereits erklärt, wie Letzteres auf die in Satz 14.12 angegebene Form gebracht werden kann. Beachte, dass jetzt wieder  $(w, b)$  die Variable ist und  $\lambda$  die Lagrangemultiplikatoren für die Ungleichungsbedingungen enthält! Wir prüfen zuerst (KKT-1) und hier gilt

$$\nabla_{(w,b)} L(w^*, b^*, \lambda^*) = \left( \frac{\partial L}{\partial w}, \frac{\partial L}{\partial b} \right)(w^*, b^*, \lambda^*) = \left( w^* - \sum_{i=1}^n \lambda_i^* y_i x_i, - \sum_{i=1}^n \lambda_i^* y_i \right) = 0$$

wenn wir die Definition von  $w^*$  und die Gleichung  $\langle \lambda^*, y \rangle = 0$  aus Teil ③ nutzen. Da wir in (14.1) keine Gleichheitsnebenbedingungen haben, entfällt (KKT-2). Für (KKT-3) haben wir  $\lambda^* \geq 0$  und

$$\begin{aligned}
g_i(w^*, b^*) &= y_i (\langle w^*, x_i \rangle + b^*) - 1 \\
&\stackrel{\substack{\uparrow \\ w^* \text{ und } b^* \\ \text{einsetzen}}}{=} y_i \left( \left\langle \sum_{j=1}^n \lambda_j^* y_j x_j, x_i \right\rangle + \mu_1^* \right) - 1 \\
&\stackrel{\substack{\uparrow \\ \text{nährhafte} \\ \text{Null}}}{=} y_i \underbrace{\left\langle \sum_{j=1}^n \lambda_j^* y_j x_j, x_i \right\rangle - 1 - \theta_i^* + y_i \mu_1^* + \theta_i^*}_{=0 \text{ nach Teil ③}} = \theta_i^*
\end{aligned}$$

für  $i = 1, \dots, n$ , also  $g(w^*, b^*) = \theta^* \geq 0$  und  $\langle \lambda^*, g(w^*, b^*) \rangle = \langle \lambda^*, \theta^* \rangle = 0$  jeweils

durch Ausnutzung der (Un-)gleichungen die wir in Teil ③ gezeigt haben. Nach Satz 14.12(ii) ist  $(w^*, b^*)$  ein Minimierer des Problems (14.1), der insbesondere alle Daten korrekt klassifiziert. Damit ist  $D$  linear trennbar und nach Satz 14.9 ist  $(w^*, b^*)$  der einzige Minimierer für (14.1) und  $h_{w^*, b^*}$  ist die SVM für die Datenmenge  $D$ . Da  $h_{w^*, b^*}$  alle Daten korrekt klassifiziert und nicht alle Daten das gleiche Label haben, kann  $w^*$  nicht der Nullvektor sein. Mit der Gleichung  $w^* = \lambda_1^* y_1 x_1 + \cdots + \lambda_n^* y_n x_n$  folgt, dass  $\lambda^*$  ebenfalls nicht der Nullvektor sein kann.

⑤ Es bleibt noch zu verifizieren, dass, bei gegebenem  $\lambda^*$ , die Schritte 1–3 tatsächlich auf  $(w^*, b^*)$  führen: In der Tat haben wir  $w^* = \lambda_1^* y_1 x_1 + \cdots + \lambda_n^* y_n x_n$  genau wie in Schritt 1 definiert. Wegen  $\lambda^* \neq 0$  ist es möglich  $i_0$  mit  $\lambda_{i_0} \neq 0$  wie in Schritt 2 gefordert zu wählen. Für jedes beliebige  $i$  gilt, unter Beachtung von  $y_i \in \{-1, 1\}$ , also  $1/y_i = y_i$ ,

$$b^* = \underset{\substack{\uparrow \\ \text{Dfn} \\ \text{von } b^*}}{\mu_1^*} = \underset{\substack{\uparrow \\ \text{1. Gleich-} \\ \text{ung in ③}}}{\frac{1}{y_i} \left( -y_i \left\langle x_i, \sum_{j=1}^n \lambda_j^* y_j x_j \right\rangle + 1 + \theta_i^* \right)} = \underset{\substack{\uparrow \\ \text{Dfn} \\ \text{von } w^*}}{-\langle x_i, w^* \rangle + y_i + y_i \theta_i^*}.$$

Wir behaupten, dass  $\theta_i^* = 0$  gilt, falls  $\lambda_i^* \neq 0$  ist. Hierfür beachten wir, dass nach Teil ③ alle  $\theta_i^*$  und alle  $\lambda_i^*$  größer gleich Null sind, während ihr Skalarprodukt  $\langle \theta^*, \lambda^* \rangle$  verschwindet. In der folgenden Summe muss daher

$$\sum_{i=1}^n \theta_i^* \lambda_i^* = 0$$

jeder einzelne Summand gleich Null sein. Dann impliziert aber  $\lambda_i^* \neq 0$  stets  $\theta_i = 0$  wie behauptet.  $\square$

**Korollar 14.13.** *Sei  $D$  wie Satz 14.11,  $(w^*, b^*)$  sei die Lösung des konvexen Optimierungsproblems und  $\lambda^*$  eine Lösung des quadratischen Optimierungsproblems. Dann gilt:*

- (i) *Alle  $x_i$  mit  $\lambda_i \neq 0$  sind Trägervektoren.*
- (ii) *Die Summe  $\lambda_1^* + \cdots + \lambda_n^* = \|w^*\|^2 = (1/\gamma^*)^2$  ist unabhängig von  $\lambda^*$  und gleich dem Kehrwert des Spielraums  $\gamma^*$  der SVM zum Quadrat.*

*Beweis.* (i) Nach Satz 14.11 gilt  $b^* = y_i + \langle w^*, x_i \rangle$  für jedes  $i \in \{1, \dots, n\}$  mit  $\lambda_i \neq 0$ . Umstellen liefert  $\langle w^*, x_i \rangle - b^* = y_i \in \{-1, 1\}$ . Da  $h_{w^*, b^*}$  die SVM für  $S$  ist, erfüllt  $(w^*, b^*)$  die Rinnenbedingung und es folgt mit Bemerkung 14.7, dass  $x_i \in \mathcal{R}_+ \cup \mathcal{R}_-$  ein Trägervektor ist.

- (ii) Projektion der Gleichung  $w^* = \lambda_1^* y_1 x_1 + \cdots + \lambda_n^* y_n x_n$  auf  $\text{span}\{w^*\}$  liefert

$$\langle w^*, w^* \rangle = \left\langle w^*, \sum_{i=1}^n \lambda_i^* y_i x_i \right\rangle = \sum_{i=1}^n \lambda_i^* y_i \langle w^*, x_i \rangle$$



woraus mit  $y_i(\langle w^*, x_i \rangle + b^*) = 1$  und  $\langle \lambda^*, y \rangle = 0$  die Identität

$$\|w^*\|^2 = \sum_{i=1}^n \lambda_i^* y_i \langle w^*, x_i \rangle = \sum_{i=1}^n \lambda_i^* (1 - y_i b^*) = \sum_{i=1}^n \lambda_i^* + b^* \langle \lambda^*, y \rangle = \sum_{i=1}^n \lambda_i^*$$

folgt. Dass  $\|w^*\| = 1/\gamma^*$  ist, folgt aus Proposition 14.6.  $\square$

**Bemerkung 14.14.** (i) Lösungen des quadratischen Optimierungsproblems in Satz 14.11 kann man mit fertigen Paketen zur quadratischen Optimierung berechnen lassen. Dazu muss man in der Regel das Problem in der Form

$$\underset{\substack{\lambda \in \mathbb{R}^n \\ G\lambda \leq h, A\lambda = b}}{\operatorname{argmin}} \quad \frac{1}{2} \lambda^\top P \lambda + q^\top \lambda$$

eingeben und dann Matrizen  $P, G, A$  und (evtl. 1-dimensionale) Vektoren  $q, h, b$  übergeben. In unserem Fall haben wir  $P$  und  $q$  im Beweis von Satz 14.11 bereits bestimmt und können  $G = -\operatorname{id}_{\mathbb{R}^n}$ ,  $h = 0_{\mathbb{R}^n}$ ,  $A = [y_1 \cdots y_n]$  und  $b = 0_{\mathbb{R}^1}$  wählen, vgl. Aufgabe 14.7. Auf die interessante Frage, wie solche Probleme numerisch gelöst werden, können wir in diesem Buch nicht eingehen, sondern verweisen auf die klassische Optimierungsliteratur, siehe u.a. die Referenzen am Ende des Kapitels.

(ii) Im Allgemeinen hat das quadratischen Optimierungsproblems in Satz 14.11 mehrere Lösungen, siehe Beispiel 14.15 und Aufgabe 14.7. Korollar 14.13(ii) besagt, dass die Summe der  $\lambda_i$  auf Lösungen konstant ist, also nur von der Datenmenge abhängt. Je nach Datenmenge, kann letztere aber beliebig groß oder beliebig klein ausfallen; man denke z.B. an ein Datenmenge mit nur zwei Datenpunkten. Korollar 14.13(i) ist keine Äquivalenz, es kann also durchaus Trägervektoren  $x_i$  mit  $\lambda_i = 0$  geben, siehe ebenfalls Beispiel 14.15 und Aufgabe 14.9.

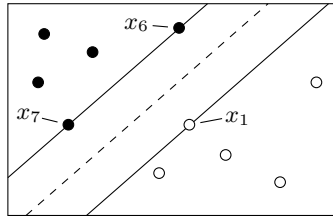
**Beispiel 14.15.** (i) Wir betrachten die folgende Datenmenge

Datenpunkt	1	2	3	4	5	6	7	8	9	10
Abszisse	3.00	4.50	5.10	2.50	3.60	2.83	1.00	0.50	1.40	0.60
Ordinate	1.50	0.55	2.20	0.70	1.00	3.10	1.50	2.20	2.70	3.00
Label	+1	+1	+1	+1	+1	-1	-1	-1	-1	-1

welche zu dem in diesem Kapitel bereits mehrfach gezeigten Bild gehört. Wir lösen mit einem geeigneten Paket das entsprechende quadratische Optimierungsproblem und erhalten

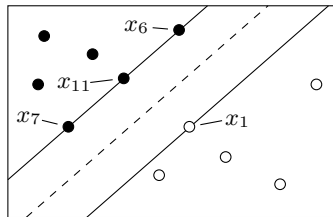
$$\lambda^* = \overset{1}{(1.154, 0.000, 0.000, 0.000, 0.000, 0.714, 0.439, 0.000, 0.000, 0.000)}$$

was auf  $w^* = 1.154 \cdot \begin{bmatrix} 3.00 \\ 1.50 \end{bmatrix} - 0.714 \cdot \begin{bmatrix} 2.83 \\ 3.10 \end{bmatrix} - 0.439 \cdot \begin{bmatrix} 1.00 \\ 1.50 \end{bmatrix} = \begin{bmatrix} 1.00 \\ -0.75 \end{bmatrix}$  und  $b^* = 1 - \langle \begin{bmatrix} 1.00 \\ -0.75 \end{bmatrix}, \begin{bmatrix} 3.00 \\ 1.50 \end{bmatrix} \rangle = 1.87$  führt. Im folgenden Bild



sieht man, dass in diesem Beispiel genau die Datenpunkte  $x_i$  Trägervektoren sind für die  $\lambda_i \neq 0$  sind. In der Tat ist  $\lambda^*$  die eindeutige Lösung des quadratischen Optimierungsproblems: Da nur  $x_1, x_6$  und  $x_7$  in den Rinnen liegen, können höchstens  $\lambda_1, \lambda_6$  und  $\lambda_7$  ungleich Null sein. Da  $w^*$  durch die Datenmenge eindeutig bestimmt ist, liefern Satz 14.11(i) und Korollar 14.13(ii) drei lineare Gleichungen für die drei zu bestimmenden Variablen.

(ii) Wir fügen jetzt der Datenmenge aus (i) einen 11-ten Datenpunkt mit Koordinaten (1.915, 2.300) und Label  $-1$  hinzu. Dieser Punkt liegt genau zwischen  $x_6$  und  $x_7$  und damit in der Rinne  $\mathcal{R}_-$ :



Die SVM bleibt also unverändert, aber es gibt jetzt vier Trägervektoren. Lösen wir per Computer das Optimierungsproblem für die erweiterte Datenmenge, so erhalten wir die Lösung

$$\mu^* = (\overset{1}{1.154}, \overset{6}{0.000}, \overset{7}{0.000}, \overset{11}{0.000}, 0.000, 0.000, 0.532, 0.256, 0.000, 0.000, 0.000, 0.364)$$

in der nun vier Einträge ungleich Null sind. Andererseits ist die Lösung

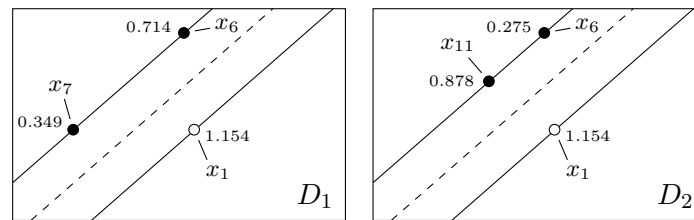
$$\lambda^* = (\overset{1}{1.154}, \overset{6}{0.000}, \overset{7}{0.000}, \overset{11}{0.000}, 0.000, 0.000, 0.714, 0.439, 0.000, 0.000, 0.000, 0.000),$$

aus (i), durch eine Null an der 11-ten Stelle erweitert, eine zweite Lösung. Bei dieser Lösung ist also  $\lambda_{11} = 0$  und  $x_{11}$  Trägervektor.

**Beispiel 14.16.** (i) Wir bleiben bei den in Beispiel 14.15 angegebenen Daten und betrachten jetzt

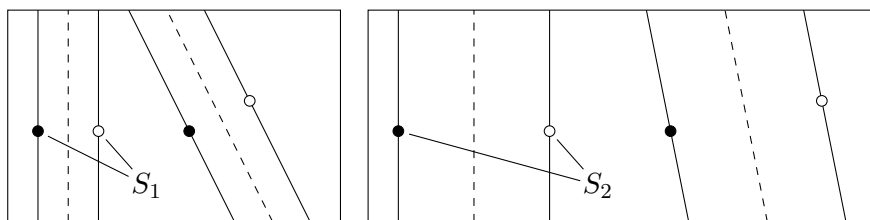
$$D_1 = \{(x_i, y_i) \mid i = 1, 6, 11\} \quad \text{und} \quad D_2 = \{(x_i, y_i) \mid i = 1, 6, 7\}.$$

Durch Lösung des Optimierungsproblems für  $D_1$  und Streichen von Nullen in Beispiel 14.15(i) erhalten wir die im Bild angegebenen Werte  $\lambda_i^*$  für die jeweiligen Datenmengen, die mit dem Argument aus Beispiel 14.15(i) eindeutig sind. Es ergibt sich das folgende Bild



in welchem die Werte der  $\lambda_i^*$  quantitativ wie folgt interpretiert werden können. Zuerst sehen wir, dass die Summe der  $\lambda_i^*$  links und rechts jeweils gleich 2.307 ist. Da beide Datenmengen durch Weglassen von Punkten aus Beispiel 14.15 entstanden sind, und sich dabei die SVM nicht geändert hat, ist dies kein Zufall, sondern folgt aus Korollar 14.13. Weiter hat  $\lambda_1^*$  links wie rechts denselben Wert und dieser ist modulo Rundungsfehlern gleich  $0.5 \cdot 2.307$ , also  $1/2$  mal der Summe aller Einträge. Dies spiegelt wider, dass  $x_1$  alleine auf einer Seite der Entscheidungsgrenze liegt. Links sind die Werte von  $\lambda_6^*$  und  $\lambda_7^*$  nicht gleich, aber ähnlich, was dazu passt, dass die Punkte  $x_6$  und  $x_7$  beinahe symmetrisch zu  $x_1$  angeordnet sind. Lässt man einen der beiden weg, so ändert sich die SVM in beiden Fällen ähnlich stark. Auf der rechten Seite ist dies ganz anders: Dort ist  $\lambda_{11}^*$  deutlich größer als  $\lambda_6^*$  und man sieht aufgrund der Anordnung der Punkte, dass sich die SVM fast gar nicht ändert, wenn man  $x_6$  weglässt, aber dass sie sich deutlich ändert, wenn man  $x_{11}$  weglässt.

(ii) Als nächstes betrachten wir zwei Datenmengen  $S_1$  und  $S_2$ , die jeweils nur aus zwei Punkten bestehen, von denen einer Label  $+1$  und der andere Label  $-1$  hat. Der Abstand der zwei Punkte sei in den Datenmenge  $S_1$  kleiner als in der Datenmenge  $S_2$ . Die Lösungsvektoren  $\lambda_{S_i}^* = (\lambda_{S_i}, \lambda_{S_i})$  haben dann jeweils zwei gleiche Einträge, und weil der Spielraum der SVM für  $S_1$  kleiner als der Spielraum der SVM für  $S_2$  ist, gilt  $\lambda_{S_1} > \lambda_{S_2}$  nach Korollar 14.13(ii). Ändert man nun jeweils einen der zwei Datenpunkte in gleicher Weise (im folgenden Bild wird z.B. jeweils die Ordinate des mit  $+1$  gelabelte Punkt um denselben Wert erhöht)



so ändert sich die SVM für die Datenmenge  $S_1$ , also für diejenige mit größerem  $\lambda_{S_1}^*$ , stärker als für die Datemenge  $S_2$  mit dem kleineren  $\lambda_{S_2}^*$ .

**Bemerkung 14.17.** (i) Die Beispiele suggerieren, dass die Einträge  $\lambda_i^*$  einer jeden Lösung  $\lambda^*$  des quadratischen Optimierungsproblems als ‘Wichtigkeit’ des Datenpunktes  $x_i$  für die SVM interpretiert werden können. Dies kann man qualitativ verstehen, im Sinne dass das Weglassen von  $x_i$  mit  $\lambda_i^* = 0$  auf die gleiche SVM führt, siehe Beispiel 14.15. Man kann ‘Wichtigkeit’ aber auch quantitativ auffassen: Bei großem  $\lambda_i$  führt ein Ändern oder Weglassen des Datenpunktes  $x_i$  zu einer größeren Veränderung der SVM, als bei kleinem  $\lambda_i$ . Hier kann man ‘groß’ und ‘klein’ relativ

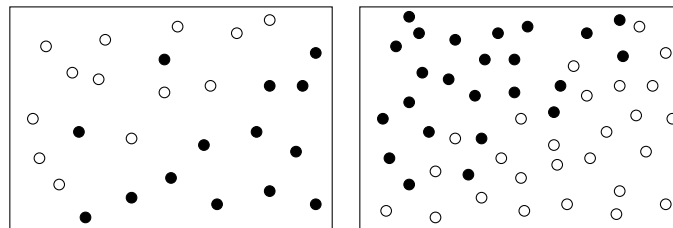
zur Summe der  $\lambda_i^*$  verstehen, siehe Beispiel 14.16(i), aber auch absolut wie in Beispiel 14.16(ii).

(ii) Formal bewiesen haben wir, dass das Weglassen eines Nicht-Trägervektors die SVM nicht verändert (dies folgt aus der zweiten Gleichung in Teil ④ des Beweises von Satz 14.9) und dass die Summe der  $\lambda_i^*$  bei einer festen Datenmenge stets gleich dem Quadrat des Kehrwertes des Spielraums der SVM ist, siehe Korollar 14.13. Die lediglich heuristischen Ausführungen in (i) sind daher mit Vorsicht zu genießen, vgl. Aufgabe 14.6.

Zum Abschluß des Kapitels weisen wir darauf hin, dass wir hier sogenannte *harte* SVMs betrachtet haben. Im Gegensatz dazu werden bei *weichen* SVMs ein paar Fehlklassifizierungen erlaubt. Dies erreicht man durch sogenannte *Schlupfvariablen*  $\xi_1, \dots, \xi_n \geq 0$  und Ersetzen der Nebenbedingung  $y_i(\langle w, x_i \rangle + b) \geq 1$  durch  $y_i(\langle w, x_i \rangle + b) \geq 1 - \xi_i$ . D.h. man erlaubt eine Verletzung der ursprünglichen Nebenbedingung um  $\xi_i$ , wobei man natürlich will, dass die  $\xi_i$  möglichst klein sind. Dies kann man wiederum erreichen, indem man per

$$\underset{\substack{(w,b,\xi) \in \mathbb{R}^{d+1} \times \mathbb{R}_{\geq 0}^n \\ \forall i: y_i(\langle w, x_i \rangle + b) \geq 1 - \xi_i}}{\operatorname{argmin}} \quad \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \xi_i$$

den Durchschnitt der  $\xi_i$  als Summand in das Minimierungsproblems einführt und auf diese Weise Verletzungen der ursprünglichen Nebenbedingung ‘bestraft’. In diesem Fall verliert man im Allgemeinen die Eindeutigkeitsaussagen des harten Falls, kann dann aber Datenmengen behandeln, die ‘fast’ linear trennbar sind, wie z.B. in den folgenden zwei Beispielen



in denen (linkes Bild) nur eine wenige Datenpunkte die lineare Trennbarkeit stören, oder in denen (rechtes Bild) die mit +1 gelabelten Daten entlang einer Hyperebene in die mit -1 gelabelten übergehen.

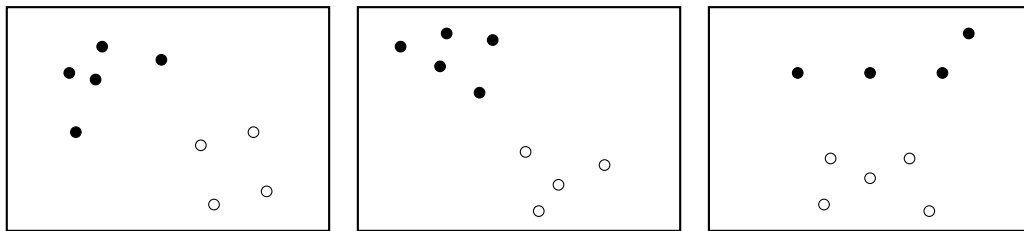
## Referenzen

Wir notieren zunächst, dass das Minimierungsproblem in Satz 14.12 gerade [GK02, (2.21)] ist, jedoch mit  $g_i$  ersetzt durch  $-g_i$ . Entsprechend ist bei der Lagrangefunktion [GK02, Definition 2.34] ein Minus vor der Summe über die  $g_i$  zu ergänzen und bei den KKT-Bedingungen [GK02, Definition 2.35] dreht sich die Ungleichung für  $g$  um. Satz 14.12(i) ist dann [GK02, Satz 2.42] und Satz 14.12(i) ist [GK02, Satz 2.46]. Vieles in diesem Kapitel basiert auf dem ausgezeichneten Video [AM12], wurde aber detailliert ausgearbeitet und

mit vollständigen Beweisen versehen. Der Eindeigkeitssteil ② des Beweises von Satz 14.9 folgt dabei [Ber, BC99]. Aufgabe 14.4 stammt aus [Nos16]. Für mehr Informationen zu Soft-SVM verweisen wir auf [SSBD14, Chapter 15.2]. Eine gute Anleitung zur Benutzung des Python-Paketes `CVXOPT` findet sich in [MC].

## Aufgaben

**Aufgabe 14.1.** Zeichnen Sie für die folgenden Datenmengen denjenigen Klassifizierer, der maximalen Spielraum hat. Skizzieren Sie die Rinnen und markieren Sie die Trägervektoren.



Wenn man aus den Trägervektoren eine minimale Menge wählen will, welche die SVM eindeutig bestimmt, wieviele Trägervektoren benötigt man dann in den Beispielen?

**Aufgabe 14.2.** Berechnen Sie den Abstand des Punktes  $(1, 4, 0)$  von der durch die Gleichung  $3x + 4y = 4$  gegebenen Ebene in  $\mathbb{R}^3$  mithilfe von Lemma 14.3. Kommt Ihnen die Formel aus Lemma 14.3, vielleicht mit der zusätzlichen Bedingung  $\|w\| = 1$ , bekannt vor? Vielleicht aus dem Schulunterricht?

**Aufgabe 14.3.** Sei  $D = \{(x_i, y_i) \mid i = 1, \dots, n\} \subseteq \mathbb{R}^d \times \{-1, 1\}$  eine Datenmenge. Zeigen Sie, dass die Normierungsbedingung  $\|w\| = 1$ , anstelle der Rinnenbedingung, auf das Optimierungsproblem

$$\begin{aligned} \operatorname{argmax}_{\substack{(w,b) \in \mathbb{S}^1 \times \mathbb{R} \\ \forall i: y_i(\langle w, x_i \rangle + b) > 0}} \min_{i=1, \dots, n} |\langle w, x_i \rangle + b| \end{aligned}$$

führt, wobei  $\mathbb{S}^1 = \{w \in \mathbb{R}^d \mid \|w\| = 1\}$  die  $d$ -dimensionale Einheitssphäre bezeichnet.

**Aufgabe 14.4.** Wir betrachten die Datenmenge

$$D = \{(x_1, y_1) := ([1], 1), (x_2, y_2) := ([1/3], 1), (x_3, y_3) := ([2], -1), (x_4, y_4) := ([3/2], -1)\}.$$

- (i) Notieren Sie die Optimierungsaufgabe aus Satz 14.9, aber lassen Sie das Quadrat weg. Veranschaulichen Sie sich den Bereich  $M \subseteq \mathbb{R}^3$  über den minimiert wird per 3D-Plot in Geogebra.
- (ii) Lesen Sie an Ihrem Plot  $w^* = \operatorname{argmin}_{w \in M} \|w\|$  (=derjenige Punkt aus  $M$  der am nächsten am Ursprung liegt) ab und finden Sie eine Ungleichungsnebenbedingung, bei der mit  $w^*$  die Gleichheit gilt. Bestimmen Sie daraus  $b^*$ .
- (iii) Skizzieren Sie nun die Datenmenge, sowie Entscheidungsgrenze und Rinnen des durch  $(w^*, b^*)$  gegebenen Klassifizierers mit maximalem Spielraum. Ärgern Sie sich nicht, wenn Sie bemerken, dass man in diesem einfachen Beispiel die Parameter  $(w^*, b^*)$  auch leicht hätte erraten können.

**Aufgabe 14.5.** Wir betrachten die Datenmenge

$$D := \{(x_1, y_1) := ([0], -1), (x_2, y_2) := ([1], 1), (x_3, y_3) := ([1], 1)\}$$

aus drei Punkten in  $\mathbb{R}^2$ . Notieren Sie die quadratische Optimierungsaufgabe aus Satz 14.11 und berechnen Sie per Hand eine Lösung.

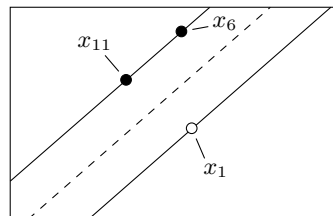
**Aufgabe 14.6.** Testen Sie die Heuristik aus Bemerkung 14.17 wie folgt.

- (i) Betrachten Sie in der Situation von Beispiel 14.15(ii) die nur aus  $x_1$ ,  $x_6$ ,  $x_7$  und  $x_{11}$  bestehende Datmenge. Es ist dann

$$\mu^* = (\overset{1}{1.154}, \overset{6}{0.532}, \overset{7}{0.256}, \overset{11}{0.364})$$

eine Lösung des zugehörigen quadratischen Optimierungsproblems. Wie stark ändert sich die entsprechende SVM, wenn man einen der Punkte  $\mu_6^*$ ,  $\mu_7^*$ , oder  $\mu_{11}^*$  weglässt?

- (ii) Betrachten Sie nun die Datenmenge  $D_2$  aus Beispiel 14.16(i):



Können Sie der Datenmenge einen vierten Punkt hinzufügen, sodass (a) das quadratische Optimierungsproblem weiterhin die eindeutige Lösung  $\lambda^* = (1.154, 0.275, 0.878)$  hat und sich (b) bei Entfernung des Punktes  $x_{11}$  die SVM so gut wie gar nicht ändert?

**Aufgabe 14.7.** Gegeben sei die Datenmenge  $D = \{(x_i, y_i) \mid i = 1, \dots, 7\} \subseteq \mathbb{R}^2 \times \{-1, 1\}$ , wobei die Koordinaten der  $x_i$  und die Werte der Label  $y_i$  in der folgenden Tabelle aufgelistet sind.

Datenpunkt	1	2	3	4	5	6	7
Abszisse	0.00	-1.0	2.5	1.0	-1.5	-3.0	1.0
Ordinate	3.0	-3.0	3.0	1.0	0.0	1.0	-1.0
Label	+1	-1	-1	-1	+1	+1	-1

- (i) Notieren Sie die quadratische Optimierungsaufgabe aus Satz 14.11 und berechnen Sie, z.B. mithilfe des Pythonpakets `CVXOPT`, eine Lösung.
- (ii) Bestimmen Sie, basierend auf obigem, den affin-linearen Klassifizierer  $h: \mathbb{R}^2 \rightarrow \mathbb{R}$  mit maximalem Spielraum und geben Sie den Spielraum an.
- (iii) Skizzieren Sie die Datenmenge, die Entscheidungsgrenze und die Rinnen. Markieren Sie die Trägervektoren und tragen Sie ein, wo man den Spielraum im Bild sehen kann. Finden Sie in (i) mindestens eine weitere Lösung.

**Aufgabe 14.8.** Recherchieren Sie, wie man z.B. mit Python die Support Vector Machine berechnen kann, ohne sich die Finger schmutzig zu machen. Bestimmen Sie eine solche für die Datenmenge aus Aufgabe 3.2 und klassifizieren Sie das dort angegebene ungelabelte Bild als ‘Eins’ bzw. ‘Sieben’.

**Aufgabe 14.9.** Zeigen Sie, dass das  $\gamma$  im Beweis von Satz 13.13 durch den Spielraum  $\gamma^*$  der durch die Daten eindeutig bestimmten Support Vector Machine nach unten abgeschätzt werden kann. Folgern Sie daraus eine Abschätzung nach oben für die Laufzeit des Perzeptronalgorithmus. Überlegen Sie sich schließlich, wie man den Spielraum ungefähr mittels der Daten schätzen kann *ohne* die SVM konkret auszurechnen.

---

*Hinweis:* Beachten Sie, dass wir in Satz 13.13 die um Eins erweiterten Datenpunkte  $\hat{x}_i$  benutzt haben, aber in Kapitel 14 nicht.

## Kapitel 15

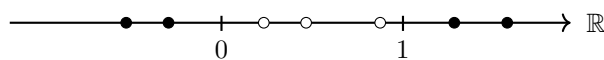
# Die Kernmethode

Als nächstes betrachten wir Datenmengen, die nicht linear trennbar sind. Wir werden dazu unsere gegebene Datenmenge in einen höherdimensionalen Raum (manchmal sogar unendlichdimensional!) abbilden, um für das Bild die lineare Trennbarkeit zu erreichen. Bemerkenswert ist insbesondere, dass wir jetzt in gewissem Sinne das Gegenteil zur Dimensionalitätsreduktion, vgl. Kapitel 11 und 7.1, machen.

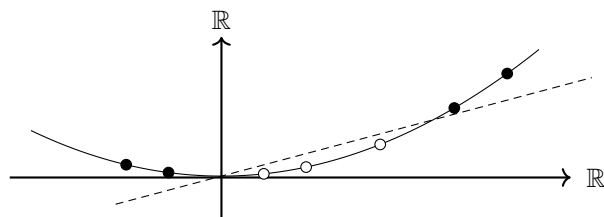
In diesem Kapitel benutzen wir mehrere Konzepte aus der Hilbertraumtheorie; den Satz von der Orthogonalprojektion und die Existenz der Vervollständigung eines Prähilbertraumes werden wir unten ohne Beweis zitieren.

Wir beginnen jetzt erstmal mit einem Beispiel für oben skizzierte Idee der ‘Einbettung’ von nicht linear trennbaren Daten in einen höherdimensionalen Raum.

**Beispiel 15.1.** Sei  $D = \{(x_i, y_i) \mid i = 1, \dots, n\} \subseteq \mathbb{R} \times \{-1, 1\}$  die im folgenden Bild dargestellte 1-dimensionale Datenmenge, welche offenbar nicht linear trennbar ist.



Wir bilden die Datenpunkte via der Abbildung  $\psi: \mathbb{R} \rightarrow \mathbb{R}^2$ ,  $\psi(x) = (x, x^2)$  nach  $\mathbb{R}^2$  ab und weisen den Bildern jeweils dasselbe Label wie den Urbildern zu. D.h. wir betrachte die 2-dimensionale Datenmenge  $\hat{D} := \{(\psi(x_i), y_i) \mid i = 1, \dots, n\} \subseteq \mathbb{R}^2 \times \{-1, 1\}$ , welches wir die *abgebildete Datenmenge* nennen. Dieses wird durch das folgende Bild dargestellt und wir sehen, dass  $\hat{D}$  linear trennbar ist.



Für  $\hat{D}$  können wir also die Support Vector Machine  $h: \mathbb{R}^2 \rightarrow \mathbb{R}$  mit den in Kapitel 14 diskutierten Methoden bestimmen. Verketteten wir diese mit der Abbildung  $\psi$ , so



liefert dies einen Klassifizierer  $h \circ \psi: \mathbb{R} \rightarrow \mathbb{R}$  für die Datenmenge  $D$ , den wir den *zurückgezogenen Klassifizierer* nennen. Dieser ist natürlich nicht mehr affin-linear.

Im obigen Beispiel erscheint die Wahl der Abbildung  $\psi$  einerseits etwas willkürlich, andererseits suggeriert das erste Bild in Beispiel 15.1, dass lineare Trennbarkeit mit einer quadratischen Funktion erreicht werden kann. Unsere Wahl von  $\psi$  ist dann der einfachste Kandidat, und wie sich gezeigt hat, erreichen wir mit diesem das Gewünschte. Alternativ kann man  $\psi$  auch noch weiter an die Daten anpassen, z.B., erreicht man mit  $\psi(x) = (x, p(x))$  und  $p(x) = -(x-0)(x-1)$ , dass bei der abgebildeten Datenmenge  $\hat{D}$  alle Punkte mit Label +1 in der oberen Halbebene und alle Punkte mit Label -1 in der unteren Halbebene liegen. Man kann einen affin-linearen Klassifizierer, nämlich  $h = \text{sign}(\langle \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \cdot \rangle + 0)$ , dann sofort ablesen. Man sieht sofort, dass die zugrundeliegende Vorgehensweise, nämlich  $\psi$  mittels eines Polynoms zu konstruieren, dessen Linearfaktoren man an den Daten abliest, für beliebige Datenmengen in  $\mathbb{R}$  funktioniert, siehe Aufgabe 15.1 für ein weiteres Beispiel.

Ist man, für eine gegebene Datenmenge, auf der Suche nach einer Funktion  $\psi$  mit Eigenschaften wie oben, so steht zu vermuten, dass man nicht immer so einfach davon kommt wie in Beispiel 15.1, sondern dass komplizierte Funktionen mit hochdimensionalen Zielbereichen in Betracht gezogen werden müssen. Wir definieren zuerst, was für Zielbereiche in Frage kommen. Da wir dort affin-linear klassifizieren wollen, benötigen wir auf jeden Fall ein Skalarprodukt, welches allerdings mitnichten das Standardskalarprodukt sein muss.

**Definition 15.2.** Sei  $H$  ein (möglicherweise unendlichdimensionaler aber stets reeller) Vektorraum. Eine Abbildung  $\langle \cdot, \cdot \rangle: H \times H \rightarrow \mathbb{R}$  heißt *Skalarprodukt* auf  $H$ , falls für alle  $u, v, w \in H$  und  $\alpha \in \mathbb{R}$  gilt

(SP1)  $\langle u, u \rangle \geq 0$ , (Positivität)

(SP2)  $\langle u, u \rangle = 0 \implies u = 0$ , (Definitheit)

(SP3)  $\langle u, v \rangle = \langle v, u \rangle$ , (Symmetrie)

(SP4)  $\langle \alpha u + v, w \rangle = \alpha \langle u, w \rangle + \langle v, w \rangle$ , (Linearität im 1. Argument)

wobei man sofort sieht, dass (SP3) und (SP4) implizieren, dass  $\langle \cdot, \cdot \rangle$  auch im 2. Argument linear, also insgesamt bilinear, ist. Aus Letzterem folgt weiter  $\langle \alpha u, \alpha u \rangle = \alpha^2 \langle u, u \rangle$  und damit, dass in (SP2) die umgekehrte Implikation gilt.

Wir zeigen nun die folgende fundamentale Ungleichung, wobei wir aber ganz genau notieren, welche Eigenschaften eines Skalarproduktes wir tatsächlich benutzen. Dies wird in Satz 15.16 von Nutzen sein.

**Lemma 15.3.** Sei  $H$  ein Vektorraum und sei  $\langle \cdot, \cdot \rangle: H \times H \rightarrow \mathbb{R}$  eine Abbildung mit den Eigenschaften (SP1), (SP3) und (SP4). Dann gilt für alle  $u, v \in H$  die Cauchy-Schwarz-Bunjakowski-Ungleichung

$$\langle u, v \rangle^2 \leq \langle u, u \rangle \langle v, v \rangle.$$

Erfüllt  $\langle \cdot, \cdot \rangle$  auch noch (SP2), so gilt Gleichheit genau dann wenn  $u$  und  $v$  linear abhängig sind.

*Beweis.* Für  $\alpha \in \mathbb{R}$  haben wir

$$0 \underset{\substack{\uparrow \\ \text{(SP1)}}}{\leq} \langle u + \alpha v, u + \alpha v \rangle \underset{\substack{\uparrow \\ \text{(SP3) \& (SP4)}}}{=} \langle u, u \rangle + 2\alpha \langle u, v \rangle + \alpha^2 \langle v, v \rangle$$

und wir betrachten jetzt zwei Fälle.

① Gilt  $\langle u, u \rangle = \langle v, v \rangle = 0$ , so liefert die Wahl von  $\alpha = \pm 1$ , nach Division durch 2 in der obigen Ungleichung, einerseits  $0 \leq \langle u, v \rangle$  und andererseits  $0 \leq -\langle u, v \rangle$ , was nur möglich ist, wenn  $\langle u, v \rangle = 0$  gilt.

② Sei nun ohne Einschränkung  $\langle v, v \rangle \neq 0$ . Dann setzen wir  $\alpha = -\langle u, v \rangle / \langle v, v \rangle$  in die anfängliche Ungleichung ein und erhalten

$$0 \leq \langle u, u \rangle - 2 \frac{\langle u, v \rangle}{\langle v, v \rangle} \langle u, v \rangle + \frac{\langle u, v \rangle^2}{\langle v, v \rangle^2} \langle v, v \rangle = \langle u, u \rangle - \frac{\langle u, v \rangle^2}{\langle v, v \rangle},$$

woraus durch Addition von  $-\langle u, v \rangle^2 / \langle v, v \rangle$  und Multiplikation mit dem wegen (SP2) positiven Term  $\langle v, v \rangle$  die behauptete Ungleichung folgt.

Nun zur Charakterisierung der Gleichheit: Ist z.B.  $u = \alpha v$ , so sieht man durch Einsetzen  $\langle u, v \rangle^2 = \langle \alpha v, v \rangle^2 = \alpha^2 \langle v, v \rangle \langle v, v \rangle = \langle u, u \rangle \langle v, v \rangle$ . Für die andere Richtung bemerken wir, dass jetzt im Fall ① nur  $u = v = 0$  möglich ist. Für Fall ② bemerken wir, dass wir nur an einer einzigen Stelle abgeschätzt haben. Gilt dort Gleichheit, so folgt  $u + \alpha v = 0$  mit (SP2) wobei  $\alpha = -\langle u, v \rangle^2 / \langle v, v \rangle$  ist.  $\square$

**Lemma 15.4.** Sei  $\langle \cdot, \cdot \rangle: H \times H \rightarrow \mathbb{R}$  ein Skalarprodukt. Dann ist  $\| \cdot \|: H \rightarrow \mathbb{R}$ ,  $\|u\| := \sqrt{\langle u, u \rangle}$  eine Norm auf  $H$ , d.h. für  $u, v \in H$  und  $\alpha \in \mathbb{R}$  gelten

(N1)  $\|u\| \geq 0$  und  $\|u\| = 0 \iff u = 0$ , (Positive Definitheit)

(N2)  $\|\alpha u\| = |\alpha| \|u\|$ , (Homogenität)

(N3)  $\|u + v\| \leq \|u\| + \|v\|$ . (Dreiecksungleichung)

*Beweis.* Die Eigenschaften (N1) und (N2) folgen sofort aus der Definition und aus (SP1)–(SP4). Eigenschaft (N3) sieht man per

$$\begin{aligned} \|u + v\|^2 &= \langle u + v, u + v \rangle = \|u\|^2 + 2\langle u, v \rangle + \|v\|^2 \\ &\leq \|u\|^2 + 2\|u\|\|v\| + \|v\|^2 = (\|u\| + \|v\|)^2 \end{aligned}$$

wobei die Abschätzung gerade die Cauchy-Schwarz-Bunjakowski-Ungleichung aus Lemma 15.3 ist.  $\square$

**Definition 15.5.** Sei  $H$  ein Vektorraum und sei  $\langle \cdot, \cdot \rangle$  ein Skalarprodukt auf  $H$ . Dann nennen wir  $(H, \langle \cdot, \cdot \rangle)$  einen *Prähilbertraum*. Ist  $(H, \| \cdot \|)$ , wobei  $\| \cdot \| = \sqrt{\langle \cdot, \cdot \rangle}$  die durch das Skalarprodukt induzierte Norm bezeichnet, vollständig im Sinne dass jede  $\| \cdot \|$ -Cauchyfolge in  $H$  konvergiert, dann nennen wir  $(H, \langle \cdot, \cdot \rangle)$  einen *Hilbertraum*.

Das Studium von Hilberträumen ist ein Gegenstand des Gebietes der Funktionalanalysis und es gibt dazu eine fruchtbare Theorie, aus der wir hier allerdings nur einige, für unsere Zwecke in diesem Kapitel relevante, Einzelheiten präsentieren wollen. Wir beginnen mit zwei Beispielen.

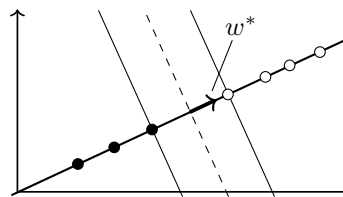
**Beispiel 15.6.** (i) Der euklidische Raum  $(\mathbb{R}^d, \langle \cdot, \cdot \rangle)$  mit dem Standardskalarprodukt ist ein Hilbertraum. In der Tat sieht man, dass jeder endlichdimensionale Hilbertraum  $(H, \langle \cdot, \cdot \rangle)$  isometrisch isomorph zu  $(\mathbb{R}^d, \langle \cdot, \cdot \rangle)$  mit  $d = \dim H$  ist: Wählt man eine Orthonormalbasis  $\mathcal{U} = \{u_1, \dots, u_d\}$  in  $H$  und definiert  $A: H \rightarrow \mathbb{R}^d$  als lineare Erweiterung von  $Au_i = e_i$ , so liefert dies einen Isomorphismus von Vektorräumen und es folgt außerdem  $\langle Au, Av \rangle = \langle u, v \rangle$  durch Einsetzen der  $\mathcal{U}$ -Entwicklungen von  $u$  und  $v$ .

(ii) Der Raum der quadratsummierbaren reellen Folgen

$$\ell^2 = \left\{ (u_i)_{i \in \mathbb{N}} \subseteq \mathbb{R} \mid \sum_{i=1}^{\infty} |u_i|^2 < \infty \right\} \quad \text{mit} \quad \langle (u_i)_{i \in \mathbb{N}}, (v_i)_{i \in \mathbb{N}} \rangle := \sum_{i=1}^{\infty} x_i y_i$$

ist ein unendlichdimensionaler Hilbertraum: Dass das Skalarprodukt wohldefiniert ist, sieht man indem man die Cauchy-Schwarz-Bunjakowski-Ungleichung auf die bei  $d$  abgeschnittene Reihe anwendet. Dann nutzt man (N3) aus, um zu zeigen, dass  $\ell^2$  ein Unterraum des Vektorraumes  $\mathbb{R}^{\mathbb{N}}$  aller Abbildungen von  $\mathbb{N}$  nach  $\mathbb{R}$  ist. Für die Vollständigkeit überlegt man sich, dass eine  $\ell^2$ -Cauchyfolge insbesondere koordinatenweise Cauchy ist, also dank der Vollständigkeit von  $\mathbb{R}$  koordinatenweise konvergiert. So erhält man einen Kandidat für den Grenzwert und zeigt dann mithilfe der Dreiecksungleichung und zwei nahrhaften Nullen, dass (a) dieser Kandidat tatsächlich in  $\ell^2$  liegt und (b) die gegebene Cauchyfolge in der vom Skalarprodukt induzierten  $\ell^2$ -Norm tatsächlich gegen den Kandidat konvergiert.

Nach diesem Exkurs in die Funktionalanalysis ist klar, dass wir Abbildungen  $\psi$  betrachten wollen, die unsere Daten in einen Hilbertraum abbilden. Handelt es sich hierbei um einen endlichdimensionalen Raum, so können wir dort die Support Vector Maschine berechnen, indem wir das quadratische Optimierungsproblem aus Satz 14.11 lösen. Ist  $H$  unendlichdimensional, so zeigt die folgende Proposition, dass man sich auf einen endlichdimensionalen, die Datenpunkte enthaltenden, Unterraum zurückziehen kann. Anschaulich ist dies sehr einleuchtend. Ist z.B. eine zweidimensionale Datenmenge gegeben, bei dem alle Punkte in einem eindimensionalen Unterraum liegen, so sieht man am folgenden Bild



sofort, dass der Vektor  $w^*$ , der zur SVM führt, ebenfalls in diesem eindimensionalen Unterraum liegen muss.

Für abstrakte und möglicherweise unendlichdimensionale Hilberträume  $(H, \langle \cdot, \cdot \rangle)$  definieren wir zunächst lineare Trennbarkeit genauso wie in Definition 13.1:  $S \subseteq H \times \{-1, 1\}$  ist *linear trennbar*, wenn für  $D$  ein korrekter Klassifizierer der Form  $h = \text{sign}(\langle w, \cdot \rangle + b)$  existiert, wobei nun natürlich  $(w, b) \in H \times \mathbb{R}$  sein muss.

**Proposition 15.7.** *Sei  $(H, \langle \cdot, \cdot \rangle)$  ein Hilbertraum. Sei  $D = \{(x_i, y_i) \mid i = 1, \dots, n\} \subseteq H \times \{-1, 1\}$  eine linear trennbare Datenmenge bei der nicht alle Label gleich sind. Dann existiert genau ein affin-linearer Klassifizierer mit maximalem Spielraum  $h^*: H \rightarrow \mathbb{R}$  und dieser ist gegeben durch  $h^* = \text{sign}(\langle w^*, \cdot \rangle + b^*)$  mit*

$$(w^*, b^*) = \underset{\substack{(w,b) \in \text{span } D_1 \times \mathbb{R} \\ \forall i: y_i(\langle w, x_i \rangle + b) \geq 1}}{\text{argmin}} \|w\|^2 = \underset{\substack{(w,b) \in H \times \mathbb{R} \\ \forall i: y_i(\langle w, x_i \rangle + b) \geq 1}}{\text{argmin}} \|w\|^2$$

wobei  $D_1 = \{x_1, \dots, x_n\}$  die Menge der Features ist, also  $\text{span } D_1 \subseteq H$  der von den Datenpunkten aufgespannte endlichdimensionale Unterraum.

*Beweis.* Für den Beweis benötigen wir die folgende Verallgemeinerung des aus der Linearen Algebra bekannten Projektionssatzes. Wir verweisen auf das Kapitelende für Referenzen und zusätzliche Bemerkungen.

**Satz 15.8.** (von der Orthogonalprojektion) *Sei  $(H, \langle \cdot, \cdot \rangle)$  ein Hilbertraum und  $\|\cdot\|$  die vom Skalarprodukt induzierte Norm. Sei  $U \subseteq H$  ein im Sinne dieser Norm abgeschlossener Unterraum. Dann kann jedes Element  $w \in H$  in eindeutiger Weise als Summe  $w = u + v$  geschrieben werden, sodass  $u \in U$  ist und  $v \in U^\perp := \{x \in H \mid \forall z \in U: \langle x, z \rangle = 0\}$  im orthogonalen Komplement von  $U$  liegt. Die lineare Projektion  $\pi_U: H \rightarrow H$  erfüllt  $\|\pi_U w\| \leq \|w\|$  für alle  $w \in H$ .  $\diamond$*

Wir verwenden jetzt die Abkürzung  $U := \text{span } D_1$ . Der Raum  $(U, \langle \cdot, \cdot \rangle)$  ist dann ein endlichdimensionaler Hilbertraum und insbesondere ein abgeschlossener Unterraum von  $H$ . Wir bezeichnen mit  $\pi_U: H \rightarrow H$  die orthogonale Projektion entsprechend Satz 15.8.

① Da  $D \subseteq H \times \{-1, 1\}$  linear trennbar ist, existieren per Definition  $(w, b) \in H \times \mathbb{R}$  sodass  $y_i(\langle w, x_i \rangle + b) > 0$  für alle  $i$  gilt. Wir behaupten, dass  $(\pi_U w, b) \in U \times \mathbb{R}$  ebenfalls einen korrekten Klassifizierer liefert. Schreiben wir  $w = \pi_U w + v \in U + U^\perp$  wie in Satz 15.8, so folgt nämlich für alle  $i$

$$y_i(\langle \pi_U w, x_i \rangle + b) = y_i(\langle w - v, x_i \rangle + b) = y_i(\langle w, x_i \rangle - \langle v, x_i \rangle + b) = y_i(\langle w, x_i \rangle + b),$$

weil  $v \in U^\perp$  und  $x_i \in U$ . Damit ist insbesondere  $D \subseteq U \times \{-1, 1\}$  linear trennbar.

② Der endlichdimensionale Raum  $(U, \langle \cdot, \cdot \rangle)$  ist nach Beispiel 15.6 isometrisch isomorph zu  $\mathbb{R}^d$  mit geeignetem  $d$ . Hierbei ist  $\mathbb{R}^d$  mit dem Standardskalarprodukt ausgestattet. Nach Satz 14.9 existiert also genau ein Minimierer

$$(w^*, b^*) = \underset{\substack{(w,b) \in U \times \mathbb{R} \\ \forall i: y_i(\langle w, x_i \rangle + b) \geq 1}}{\text{argmin}} \|w\|^2.$$

Sei nun  $(w, b) \in H \times \mathbb{R}$  beliebig sodass  $y_i(\langle w, x_i \rangle + b) \geq 1$  für alle  $i$  gilt. Mit derselben Rechnung wie in Teil ① sieht man, dass  $(\pi_U w, b) \in U \times \mathbb{R}$  dann ebenfalls  $y_i(\langle \pi_U w, x_i \rangle + b) \geq 1$  für alle  $i$  erfüllt. Wir haben also  $\|w^*\|^2 \leq \|\pi_U w\|^2$  nach Wahl von  $(w^*, b^*)$ . Jetzt schreiben wir wieder  $w = \pi_U w + v \in U + U^\perp$ . Damit erhalten wir mit Satz 15.8 die Abschätzung  $\|w\|^2 \geq \|\pi_U w\|^2$ . Beides zusammen zeigt

$$(w^*, b^*) \in \underset{\substack{(w,b) \in H \times \mathbb{R} \\ \forall i: y_i(\langle w, x_i \rangle + b) \geq 1}}{\operatorname{argmin}} \|w\|^2.$$

③ Als nächstes behaupten wir, unabhängig von der bereits gezeigten Existenz, dass auch auf einem möglicherweise unendlichdimensionalen Raum höchstens ein Minimierer für die oben angegebene Aufgabe existieren kann. Um dies zu sehen, gehen wir durch die Beweisteile ② und ③ von Satz 14.9 nochmal durch: In Teil ② haben wir nirgends benutzt, dass der zugrundeliegende Raum endlichdimensional ist. In Teil ③ haben wir Proposition 17.14 verwendet. Dort betrachten wir Funktionen, die auf konvexen Teilmengen von  $\mathbb{R}^d$  definiert sind. Man sieht aber, dass die Implikation

$$f: A \rightarrow \mathbb{R} \text{ strikt konvex} \implies \exists \text{ höchstens ein Minimierer}$$

genauso formuliert und bewiesen werden kann, wenn  $A$  nur Teilmenge eines beliebigen Vektorraumes ist. Es bleibt, sich zu überlegen, dass das Normquadrat  $\|\cdot\|^2: H \rightarrow \mathbb{R}$  in unserem abstrakten Hilbertraum strikt konvex ist. Aber auch hier sieht man, dass Beispiel 17.15 nicht die Endlichdimensionalität benutzt.

④ Aus dem Bisherigen folgt die Gleichungskette in der Proposition und die Existenz des Klassifizierers  $h^*$ . Für die Eindeutigkeit von  $h^*$  können wir ebenfalls unsere Argumente aus Kapitel 13 und 14 recyceln. Eine Inspektion der Beweise von Lemmas 13.3 und 14.2 zeigt nämlich, dass diese auch unendlichdimensional korrekt sind: Man muss hier lediglich in Definition 13.2 aufpassen, und Hyperebenen per  $\mathcal{H} = X + x_0$  mit  $x_0 \in H$  und einem abgeschlossenen Unterraum  $X \subseteq H$  mit  $\dim(H/X) = 1$  definieren. Alle weiteren Argumente basieren auf dem Projektionssatz aus der Linearen Algebra, und können mit Satz 15.8 erledigt werden.  $\square$

**Bemerkung 15.9.** Haben wir ein Datenmenge  $D = \{(x_i, y_i) \mid i = 1 \dots, n\}$  via einer Abbildung  $\psi$  in einen unendlichdimensionalen Hilbertraum  $H$  abgebildet, derart dass die abgebildete Datenmenge  $\hat{D} = \{(\psi(x_i), y_i) \mid i = 1 \dots, n\} \subseteq H \times \{-1, 1\}$  linear trennbar ist, so garantiert Proposition 15.7 einerseits, dass eine SVM für  $\hat{D}$  existiert und andererseits, dass deren Parameter  $(w^*, b^*)$  durch Lösung eines endlichdimensionalen Optimierungsproblems gefunden werden können. Dass  $w^* \in \operatorname{span} D_1$  gilt, heißt explizit, dass es Koeffizienten  $\alpha_1, \dots, \alpha_n \in \mathbb{R}$  gibt mit

$$w^* = \sum_{i=1}^n \alpha_i \psi(x_i).$$

Die (erste Komponente der) Lösung des a priori unendlichdimensionalen Optimie-

rungsproblems kann also als Linearkombination der Bilder der (endlich vielen) Daten dargestellt werden. Auf Englisch werden Resultate der Art von Proposition 15.7 deswegen als *Representer Theorems* bezeichnet.

Das letzte uns nun noch fehlende Puzzleteil ist der natürliche Definitionsbereich von  $\psi$ . In Beispiel 15.1 war unsere Datenmenge eine Teilmenge von  $\mathbb{R}$  und es schien dort natürlich,  $\psi$  auch auf ganz  $\mathbb{R}$  zu definieren. Es spricht einerseits nichts dagegen, dies auf Daten in  $\mathbb{R}^d$  zu verallgemeinern und nach Abbildungen  $\psi: \mathbb{R}^d \rightarrow (H, \langle \cdot, \cdot \rangle)$  zu suchen. Andererseits ist unser Ziel ja gerade, via  $\psi$  die Suche nach affin-linearen Klassifizierern in den Raum  $(H, \langle \cdot, \cdot \rangle)$  zu verlagern. Und dies bedeutet, dass die Daten selbst gar nicht mehr in einem Raum liegen müssen, in dem man über lineare Trennbarkeit sprechen kann! Es reicht also völlig, zusammen mit den Datenpunkten selbst, eine Obermenge  $X$  derselben zu spezifizieren, die am Ende der Definitionsbereich des Klassifizierers werden soll.

In der folgenden Bemerkung fassen wir unser Setup nochmal zusammen und formulieren dann eine zentrale Beobachtung, die häufig ‘Kernel Trick’ genannt wird.

**Bemerkung 15.10.** (Kernel Trick) Sei  $X$  eine nichtleere Menge und sei weiter  $D = \{(x_i, y_i) \mid i = 1, \dots, n\} \subseteq X \times \{-1, 1\}$  eine nichtleere Datenmenge. Sei  $\psi: X \rightarrow (H, \langle \cdot, \cdot \rangle)$  eine Abbildung von  $X$  in einen Hilbertraum und sei  $\hat{D} = \{(\psi(x_i), y_i) \mid i = 1, \dots, n\} \subseteq X \times \{-1, 1\}$  die abgebildete Datenmenge. Angenommen,  $\hat{D}$  ist linear trennbar, dann besagen Proposition 15.7 und Satz 14.11, dass die SVM für  $\hat{D}$  per

$$\lambda^* \in \operatorname{argmin}_{\substack{\lambda \in \mathbb{R}_{\geq 0}^n \\ \langle \lambda, y \rangle = 0}} \frac{1}{2} \sum_{i,j=1}^n \lambda_i \lambda_j y_i y_j \langle \psi(x_i), \psi(x_j) \rangle - \sum_{i=1}^n \lambda_i$$

gefunden werden kann. Dabei muss allerdings erwartet werden, dass  $\psi$  nicht-linear und evtl. kompliziert zu berechnen ist. Ferner erinnern wir nochmal daran, dass  $\langle \cdot, \cdot \rangle$  nicht das Standardskalarprodukt ist, sondern irgendein Skalarprodukt in einem möglicherweise unendlichdimensionalen Raum.

Die entscheidende Beobachtung ist jetzt, dass wir, um  $\lambda^*$  zu finden, die  $\psi(x_i)$  selbst, und sogar den Hilbertraum samt seines Skalarproduktes, eigentlich gar nicht zu kennen brauchen: Um  $\lambda^*$  zu finden, genügt es völlig, wenn wir für die endlich vielen Datenpunkte  $x_1, \dots, x_n$  jeweils die Zahlen  $\langle \psi(x_i), \psi(x_j) \rangle$  kennen! Falls wir also eine Funktion  $k: X \times X \rightarrow \mathbb{R}$  hätten, sodass

$$\forall x, y \in X: k(x, y) = \langle \psi(x), \psi(y) \rangle \quad (15.1)$$

gilt, könnten wir, nur anhand dieser, zuerst  $\lambda^*$  berechnen und damit dann die SVM.

Da obiges eine gewisse gedankliche Wendung enthält, weisen wir nochmal darauf hin, dass wir in der Gleichung (15.1) links Punkte aus unserer Menge  $X$  einsetzen und dann ‘direkt’ deren Skalarprodukt nach Abbilden herausbekommen, während wir rechts erst die Vektoren  $\psi(x)$  und  $\psi(y)$  explizit berechnen müssten, um dann deren Skalarprodukt zu bestimmen.

Die Beobachtung in Bemerkung 15.10 führt auf den folgenden Begriff.

**Definition 15.11.** Sei  $X$  eine nichtleere Menge. Eine Funktion  $k: X \times X \rightarrow \mathbb{R}$  heißt *Kernfunktion*, falls ein Hilbertraum  $(H, \langle -, - \rangle)$  und eine Abbildung  $\psi: X \rightarrow H$  existieren, sodass  $k(x, y) = \langle \psi(x), \psi(y) \rangle$  für alle  $x, y \in X$  gilt.

Wir geben nun zwei typische Beispiele an, in denen man die Kerneigenschaft überprüfen kann, indem man die Funktion  $\psi$  explizit angibt.

**Beispiel 15.12.** (i) Sei  $\emptyset \neq X \subseteq \mathbb{R}^d$ . Die Funktion  $k: X \times X \rightarrow \mathbb{R}$ ,  $k(x, y) = (1 + \langle x, y \rangle)^m$  mit  $m \in \mathbb{N}$  heißt *Polynomkern vom Grad  $m$* . Wir überprüfen Definition 15.11 nur im Spezialfall  $d = 3$ ,  $m = 2$ . Hier ergibt sich

$$\begin{aligned} k(x, y) &= \left(1 + \sum_{i=1}^3 x_i y_i\right)^2 = 1 + 2 \sum_{i=1}^3 x_i y_i + \left(\sum_{i=1}^3 x_i y_i\right)^2 \\ &= 1 + 2x_1 y_1 + 2x_2 y_2 + 2x_3 y_3 + x_1^2 y_1^2 + x_2^2 y_2^2 + x_3^2 y_3^2 \\ &\quad + 2x_1 x_2 y_2 y_1 + 2x_1 x_3 y_3 y_1 + 2x_2 x_3 y_2 y_3 \\ &= \left\langle \begin{bmatrix} 1 \\ \sqrt{2}x_1 \\ \sqrt{2}x_2 \\ \sqrt{2}x_3 \\ x_1^2 \\ x_2^2 \\ x_3^2 \\ \sqrt{2}x_1 x_2 \\ \sqrt{2}x_1 x_3 \\ \sqrt{2}x_2 x_3 \end{bmatrix}, \begin{bmatrix} 1 \\ \sqrt{2}y_1 \\ \sqrt{2}y_2 \\ \sqrt{2}y_3 \\ y_1^2 \\ y_2^2 \\ y_3^2 \\ \sqrt{2}y_1 y_2 \\ \sqrt{2}y_1 y_3 \\ \sqrt{2}y_2 y_3 \end{bmatrix} \right\rangle = \langle \psi(x), \psi(y) \rangle \end{aligned}$$

mit  $\psi: X \rightarrow \mathbb{R}^{10}$ , wobei  $\mathbb{R}^{10}$  mit dem Standardskalarprodukt ausgestattet ist. Für beliebige Dimension  $d$  und beliebigen Grad  $m$  rechnet man  $(1 + \langle x, y \rangle)^m$  per Multinomialformel aus und erhält dann  $\psi: X \rightarrow \mathbb{R}^q$  mit  $q = \binom{d+m}{m} = \frac{(d+m)!}{m!d!}$ , siehe Aufgabe 15.3. Um an diesem Beispiel nochmal den ‘kernel trick’ zu illustrieren, vermerken wir, dass sich  $q$  mit wachsendem  $d$  und  $m$  schnell erhöht und die Auswertung von  $\psi$  viele Rechenschritte erfordert, wenn man implementiert. Jeweils zuerst  $\psi(x_i)$  und  $\psi(x_j)$ , und dann deren  $q$ -dimensionales Skalarprodukt zu berechnen, wäre daher aufwendig. Wir müssen aber in unserem Setup nur  $k(x_i, x_j) = (1 + \langle x_i, x_j \rangle)^m$  berechnen, und das erfordert lediglich ein  $d$ -dimensionales Skalarprodukt, eine Addition und dann das Bilden der  $m$ -ten Potenz.

(ii) Sei  $\emptyset \neq X \subseteq \mathbb{R}^d$ . Die Funktion  $k: X \times X \rightarrow \mathbb{R}$ ,  $k(x, y) = e^{-\frac{\|x-y\|^2}{2\sigma^2}}$  mit  $\sigma > 0$  heißt *RBF-Kern*, wobei RBF für *radiale Basisfunktion* steht. Wir überprüfen Definition 15.11 wieder nur im Spezialfall  $\sigma = 1$ ,  $d = 2$ , in welchem sich

$$\begin{aligned} k(x, y) &= e^{-\|x-y\|^2} = e^{-\langle x-y, x-y \rangle} = e^{\langle x, y \rangle} e^{-\|x\|^2/2} e^{-\|y\|^2/2} \\ &= \sum_{j=0}^{\infty} \frac{\langle x, y \rangle^j}{j!} e^{-\|x\|^2/2} e^{-\|y\|^2/2} = \sum_{j=0}^{\infty} \frac{(x_1 y_1 + x_2 y_2)^j}{j!} e^{-\|x\|^2/2} e^{-\|y\|^2/2} \\ &= \sum_{j=0}^{\infty} \sum_{i=0}^j \frac{1}{j!} \binom{j}{i} (x_1 y_1)^i (x_2 y_2)^{j-i} e^{-\|x\|^2/2} e^{-\|y\|^2/2} \end{aligned}$$

$$= \sum_{j=0}^{\infty} \sum_{i=0}^j \left( e^{-\|x\|^2/2} \frac{x_1^i x_2^{j-i}}{(i!(j-i)!)^{1/2}} \right) \left( e^{-\|y\|^2/2} \frac{y_1^i y_2^{j-i}}{(i!(j-i)!)^{1/2}} \right) = \langle \psi(x), \psi(y) \rangle$$

ergibt, wobei

$$\psi: X \rightarrow \ell^2, \quad \psi(x) = e^{\|x\|^2/2} \left( \underbrace{1}_{j=0}, \underbrace{x_2^1, x_1^1}_{j=1}, \underbrace{\frac{x_1^2}{2}, \frac{x_1 x_2}{2}, \frac{x_2^2}{2}}_{j=2}, \dots \right).$$

Dass  $\psi$  wohldefiniert ist, folgt aus  $\|\psi(x)\|^2 = e^{\|x-y\|^2/2} < \infty$ , was nach der obigen Rechnung gilt. Für beliebige  $d$  kann wieder der Multinomialsatz angewandt werden und für beliebige  $\sigma > 0$  kommen entsprechende Konstanten hinzu. Analog zu (i) kann  $k(x_i, x_j)$  im Wesentlichen durch eine Normberechnung und Auswertung der Exponentialfunktion bestimmt werden — da der Zielbereich von  $\psi$  unendlichdimensional ist, wäre es hier sogar unmöglich, auf numerische Weise zuerst  $\psi(x_i)$  und  $\psi(x_j)$  exakt zu berechnen.

Wir formulieren nun die sich aus dem ‘Kernel Trick’ ergebende Methode als Satz.

**Satz 15.13.** *Sei  $X \subseteq \mathbb{R}^d$  und  $D = \{(x_i, y_i) \mid i = 1, \dots, n\} \subseteq X \times \{-1, 1\}$  eine Datenmenge, bei der nicht alle Label gleich sind. Sei  $k: X \times X \rightarrow \mathbb{R}$  eine Kernfunktion und sei*

$$\lambda^* \in \underset{\substack{\lambda \in \mathbb{R}_{\geq 0}^n \\ \langle \lambda, y \rangle = 0}}{\operatorname{argmin}} \frac{1}{2} \sum_{i,j=1}^n \lambda_i \lambda_j y_i y_j k(x_i, x_j) - \sum_{i=1}^n \lambda_i$$

eine beliebige Lösung des zugehörigen Optimierungsproblems. Dann ist  $\lambda^* \neq 0$  und für jedes  $i_0 \in \{1, \dots, n\}$  mit  $\lambda_{i_0} \neq 0$  ist

$$h_k: X \rightarrow \mathbb{R}, \quad h_k(x) = \operatorname{sign} \left( \sum_{i=1}^n \lambda_i^* k(x_i, x) + y_{i_0} - \sum_{i=1}^n \lambda_i^* y_i k(x_i, x_{i_0}) \right)$$

ein korrekter Klassifizierer für  $D$ . Dieser ist von der Wahl von  $\lambda$  und  $i_0$  unabhängig und kann überdies nur anhand der Kernfunktion ausgewertet werden.

*Beweis.* Per Definition 15.11 wählen wir  $\psi: X \rightarrow (H, \langle \cdot, \cdot \rangle)$  derart dass  $k(x, y) = \langle \psi(x), \psi(y) \rangle$  für alle  $x, y \in X$  gilt. Wir bezeichnen mit  $\hat{D}$  die abgebildete Datenmenge, welche im endlichdimensionalen Unterraum  $U := \operatorname{span}\{\psi(x_i) \mid i = 1, \dots, n\} \subseteq H$  liegt. Letzterer Raum ist isometrisch isomorph zu  $\mathbb{R}^d$  mit  $d = \dim U$ . Nach Satz 14.11 impliziert die Existenz von  $\lambda^*$ , die wir oben vorausgesetzt haben, dass  $\hat{D}$  linear trennbar ist. In Kombination mit Proposition 15.7 impliziert Satz 14.11 weiter, dass

$$w^* = \sum_{i=1}^n \lambda_i^* y_i \psi(x_i) \quad \text{und} \quad b^* = y_{i_0} - \langle w^*, x_{i_0} \rangle$$

den eindeutig bestimmten affin-linearen Klassifizierer  $h^*: H \rightarrow \mathbb{R}$  für  $\hat{D}$  liefern, wobei  $w^*$  unabhängig von der Wahl von  $\lambda^*$  ist, und  $b^*$  unabhängig von der Wahl von  $i_0$ , solange  $\lambda_{i_0}^* \neq 0$  gilt. Die Verkettung  $h^* \circ \psi: X \rightarrow \mathbb{R}$  liefert per Konstruktion einen



korrekten Klassifizierer der Originaldaten  $D$ . Durch Einsetzen erhalten wir

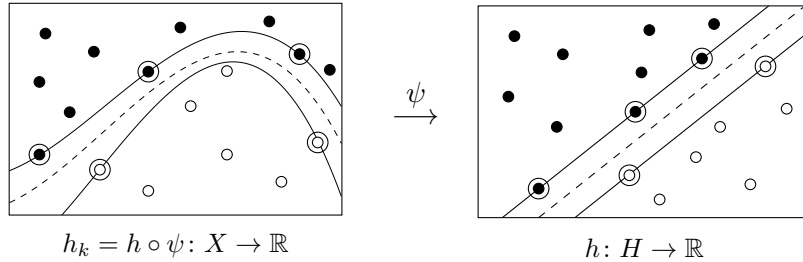
$$h^*(\psi(x)) = \text{sign}(\langle w^*, x \rangle + b^*) = \text{sign}\left(\sum_{i=1}^n \lambda_i^* k(x_i, x) + y_{i_0} - \sum_{i=1}^n \lambda_i^* y_i k(x_i, x_{i_0})\right)$$

also genau den im Satz für  $h_k(x)$  angegebenen Ausdruck.  $\square$

**Definition 15.14.** In der Situation von Satz 15.13 nennen wir  $h_k: X \rightarrow \mathbb{R}$  den *von  $k$  induzierten Klassifizierer*.

**Bemerkung 15.15.** (i) Der Beweis von Satz 15.13 hat  $h_k = h^* \circ \psi$  gezeigt, d.h. der durch  $k$  induzierte Klassifizierer stimmt mit dem unter  $\psi$  zurückgezogenen Klassifizierer überein. Bemerkenswert ist hierbei, dass  $\psi$  durch  $k$  nicht eindeutig bestimmt ist, vgl. Aufgabe 15.6, der induzierte Klassifizierer aber nur von  $k$  abhängt.

(ii) In gutartigen Fällen kann man  $h_k$ ,  $h$  und  $\psi$  wie folgt illustrieren.



Dabei sind im rechten Bild die Trägervektoren von  $h$  durch Kreise markiert, während Kreise im linken Bild *Urbilder von Trägervektoren* markieren. Beachte, dass letztere Punkte im Allgemeinen nicht alle dieselben Abstände von der Entscheidungsgrenze haben. Es kann sogar Datenpunkte geben, deren Bilder keine Trägervektoren sind, die aber näher an der Entscheidungsgrenze liegen, als alle Urbilder von Trägervektoren.

(iii) Das Bild oben wirft die Frage auf, welche Eigenschaften die Funktion  $\psi: X \rightarrow H$  eigentlich haben sollte. In unseren Beispielen 15.1 und 15.12 waren die Abbildungen  $\psi$  notwendigerweise nichtlinear, aber stets injektiv und stetig. In Beispiel 15.1 sieht man sogar ohne Mühe, dass die Funktion  $\psi$  offen auf ihr Bild ist. Andererseits benötigt der ‘kernel trick’ weder eine Topologie auf  $X$  und beinhaltet auch nur eine sehr schwache notwendige Injektivitätsbedingung: Ist die abgebildete Datenmenge  $\hat{D} \subseteq H \times \{-1, 1\}$  linear trennbar, so gilt sicher

$$\forall i = 1, \dots, n: y_i \neq y_j \implies \psi(x_i) \neq \psi(x_j),$$

aber es spricht z.B. nichts dagegen, dass  $\psi$  Datenpunkte mit demselben Label auf ein und denselben Punkt schickt, vgl. Aufgabe 15.3. Andererseits ist es, im Hinblick auf den induzierten Klassifizierer, natürlich nicht von Nachteil, wenn  $k$  so gewählt wird, dass es eine zugehörige Abbildung  $\psi: X \rightarrow H$  gibt, die injektiv und stetig ist. Wir verweisen hierzu auf Bemerkung 15.18.

Als letzten Punkt diesen Kapitels beschäftigen wir uns mit der Frage, wie man überhaupt Kernfunktionen finden kann. Eine Methode haben wir schon in Beispiel

15.12 gesehen, wo wir den Nachweis der Kerneigenschaft durch Angabe der Funktion  $\psi: X \rightarrow H$  erbracht haben. Eine andere Methode stellen wir im folgenden Satz vor, der ein Kriterium bereitstellt, mit dem die Kerneigenschaft ohne explizite Kenntnis von  $\psi$  und  $H$  geprüft werden kann.

Bedingung (ii) im folgenden Satz wird mitunter als *Mercer's Condition* bezeichnet, vergleiche die Bemerkungen und Referenzen am Ende des Kapitels.

**Satz 15.16.** *Sei  $X$  eine nichtleere Menge und  $k: X \times X \rightarrow \mathbb{R}$  eine Funktion. Dann sind folgende Aussagen äquivalent.*

(i)  $k$  ist eine Kernfunktion.

(ii) Für alle  $x_1, \dots, x_m \in X$  ist die folgende Gram-Matrix symmetrisch und positiv semidefinit:

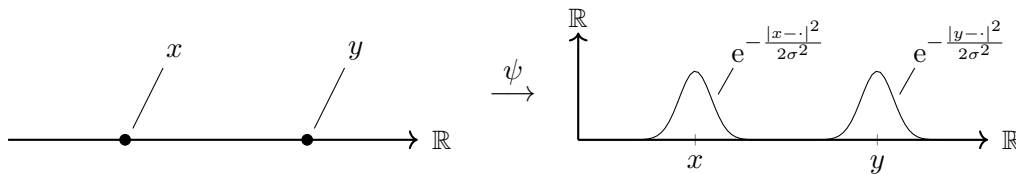
$$G := \begin{bmatrix} k(x_1, x_1) & \cdots & k(x_1, x_m) \\ \vdots & & \vdots \\ k(x_m, x_1) & \cdots & k(x_m, x_m) \end{bmatrix}.$$

*Beweis.* (i)  $\Rightarrow$  (ii): Seien  $x_1, \dots, x_m \in X$  gegeben und  $G$  definiert wie oben. Per Voraussetzung existiert  $\psi: X \rightarrow (H, \langle \cdot, \cdot \rangle)$ , sodass  $k(x, y) = \langle \psi(x), \psi(y) \rangle$  für alle  $x, y \in X$  gilt. Da das Skalarprodukt symmetrisch ist, impliziert dies, dass  $G$  symmetrisch ist. Sei jetzt  $\xi \in \mathbb{R}^m$  gegeben. Dann folgt aus

$$\begin{aligned} \langle \xi, G\xi \rangle &= \left\langle \begin{bmatrix} \xi_1 \\ \vdots \\ \xi_m \end{bmatrix}, \begin{bmatrix} k(x_1, x_1) & \cdots & k(x_1, x_m) \\ \vdots & & \vdots \\ k(x_m, x_1) & \cdots & k(x_m, x_m) \end{bmatrix} \begin{bmatrix} \xi_1 \\ \vdots \\ \xi_m \end{bmatrix} \right\rangle \\ &= \sum_{i,j=1}^m \xi_i \xi_j k(x_i, x_j) = \sum_{i,j=1}^m \xi_i \xi_j \langle \psi(x_i), \psi(x_j) \rangle \\ &= \left\langle \sum_{i=1}^m \xi_i \psi(x_i), \sum_{j=1}^m \xi_j \psi(x_j) \right\rangle = \left\| \sum_{i=1}^m \xi_i \psi(x_i) \right\|^2 \geq 0 \end{aligned}$$

dass  $G$  positiv semidefinit ist.

(ii)  $\Rightarrow$  (i): Wir definieren  $\psi: X \rightarrow \mathbb{R}^X = \{f: X \rightarrow \mathbb{R}\}$  per  $\psi(x) := k(\cdot, x)$  und denken im folgenden zur Veranschaulichung an das Beispiel des 1-dimensionalen RBF-Kerns. In diesem Spezialfall schickt  $\psi$  jeden Punkt  $x$  in  $\mathbb{R}$  auf eine Gaußfunktion die  $x$  als Mittelwert hat:



① Da  $\mathbb{R}^X$  zwar unter punktweise definierten Operationen ein Vektorraum ist, wäre die erste Idee nun, ein geeignetes Skalarprodukt auf diesem Raum zu definieren, und zwar so, dass die Identität  $k(x, y) = \langle \psi(x), \psi(y) \rangle$  erfüllt ist. Da  $\mathbb{R}^X$  sehr groß ist, ist es allerdings besser, letzteres nur auf einem Teilraum zu machen. Der kleinste

Kandidat hierfür ist  $H := \text{span ran } \psi$  und wir probieren es erstmal mit diesem. In der Tat kann jedes  $f \in H$  dann geschrieben werden als

$$f = \sum_{i=1}^n \alpha_i k(\cdot, x_i)$$

mit  $n \in \mathbb{N}$ ,  $\alpha_i \in \mathbb{R}$ ,  $x_i \in X$ .

② Ist  $f$  wie oben und  $g = \sum_{j=1}^m \beta_j k(\cdot, y_j)$ , dann definieren wir

$$\langle f, g \rangle := \sum_{i=1}^n \sum_{j=1}^m \alpha_i \beta_j k(x_i, y_j).$$

Da die Darstellungen von  $f$  und  $g$  nicht eindeutig sind, müssen wir zuerst zeigen, dass oben die rechte Seite unabhängig von der Wahl der  $n, m, \alpha_i, \beta_j, x_i, y_j$  ist. In der Tat gilt

$$\sum_{i=1}^n \sum_{j=1}^m \alpha_i \beta_j k(x_i, y_j) = \sum_{j=1}^m \beta_j \sum_{i=1}^n \alpha_i k(x_i, y_j) = \sum_{j=1}^m \beta_j f(y_j)$$

woraus folgt, dass  $\langle f, g \rangle$  unabhängig von der Wahl der  $m, \alpha_i, x_i$  ist. Die Unabhängigkeit von  $n, \beta_j, y_j$  sieht man analog. Damit ist  $\langle \cdot, \cdot \rangle: H \rightarrow \mathbb{R}$  wohldefiniert und es müssen die Skalarprodukteigenschaften nachgewiesen werden. Hierbei folgt (SP3) aus der Annahme, dass die Matrizen  $G$  symmetrisch sind und (SP4) gilt per Definition. Für  $f$  wie oben gilt per Voraussetzung an die Gram-Matrix

$$\langle f, f \rangle = \sum_{i,j=1}^m \alpha_i \alpha_j k(x_i, x_j) = \langle \alpha, G\alpha \rangle \geq 0$$

womit gezeigt ist, dass  $\langle \cdot, \cdot \rangle$  Bedingung (SP1) erfüllt. Es bleibt (SP2), also die Implikation  $\langle f, f \rangle = 0 \implies f = 0$ , zu zeigen. Dazu beobachten wir, dass für  $f$  wie oben gilt

$$\begin{aligned} \langle k(\cdot, x), f \rangle &= \left\langle k(\cdot, x), \sum_{i=1}^m \alpha_i k(\cdot, x_i) \right\rangle \\ &= \sum_{i=1}^m \alpha_i \langle k(\cdot, x), k(\cdot, x_i) \rangle \\ &= \sum_{i=1}^m \alpha_i k(x, x_i) = f(x). \\ &\quad \uparrow \\ &\quad \text{Dfn des SP!} \end{aligned}$$

Daraus folgt mit der Cauchy-Schwarz-Bunjakowski-Ungleichung, von der wir in Lemma 15.3 gesehen haben, dass sie ohne (SP2) gilt, dass

$$f(x)^2 = \langle k(\cdot, x), f \rangle^2 \leq \langle k(\cdot, x), k(\cdot, x) \rangle \cdot \langle f, f \rangle = k(x, x) \langle f, f \rangle$$

für beliebige  $x \in X$  gilt. Wenn also  $\langle f, f \rangle = 0$  ist für alle  $x \in X$ , so ist  $f(x) = 0$  für alle  $x \in X$ , und dies bedeutet, dass  $f \in H$  der Nullvektor ist.

③ Wir haben soweit jetzt zumindest einen Prähilbertraum  $(H, \langle \cdot, \cdot \rangle)$  und eine Abbildung  $\psi: X \rightarrow H$ , welche die Identität  $k(x, y) = \langle \psi(x), \psi(y) \rangle$  für alle  $x \in X$  erfüllt. Für die Vollständigkeit benötigen wir den folgenden Satz, welcher zum Standardrepertoire der Funktionalanalysis gehört und den man sich als Verallgemeinerung des Übergangs von  $\mathbb{Q}$  nach  $\mathbb{R}$  vorstellen kann.

**Satz 15.17.** (Vervollständigung) *Sei  $(H, \langle \cdot, \cdot \rangle_H)$  ein Prähilbertraum. Dann existiert ein eindeutig bestimmter Hilbertraum  $(\hat{H}, \langle \cdot, \cdot \rangle_{\hat{H}})$  sodass  $H \subseteq \hat{H}$ ,  $\langle x, y \rangle_H = \langle x, y \rangle_{\hat{H}}$  für alle  $x, y \in H$ , und  $\overline{H}^{\|\cdot\|_{\hat{H}}} = \hat{H}$  gelten. Den Raum  $\hat{H}$  nennt man die Vervollständigung von  $H$ .*  $\diamond$

Falls unser Raum  $H$  aus ① also nicht vollständig ist, können wir diesen mithilfe von Satz 15.17 durch seine Vervollständigung ersetzen und  $\psi$  als Abbildung in den eventuell vergrößerten Raum auffassen. Die Identität  $k(x, y) = \langle \psi(x), \psi(y) \rangle$  bleibt davon unberührt.  $\square$

Wir weisen darauf hin, dass in Satz 15.16(ii) die Punkte  $x_1, \dots, x_m$  nicht die Datenpunkte einer gegebenen Datenmenge sind, sondern dass hier beliebige Punkte aus  $X$  getestet werden müssen. Darüber hinaus ist deren Anzahl  $m \in \mathbb{N}$  nicht fest, Symmetrie und positive Definitheit der Gram-Matrix müssen also für unendlich viele Matrizen unbeschränkten Formats gegeben sein.

Wir haben nun zwei Methoden gesehen mithilfe derer man Kernfunktionen finden, bzw. testen kann, ob eine gegebene Funktion eine Kernfunktion ist. Jenseits dessen hat die Kerneigenschaft überaus gute Vererbungseigenschaften. So sind z.B. sind

- (a) Einschränkungen von Kernen wieder Kerne,
- (b) Linearkombinationen von Kernen wieder Kerne,
- (c) Produkte von Kernen wieder Kerne,
- (d) punktweise Grenzwerte von Kernen wieder Kerne,

wobei wir (a), (b) und (d) in Aufgabe 15.7 behandeln und für (d) auf die Referenzen am Kapitelende verweisen.

Basierend auf wohlbekannten Beispielen und den Vererbungseigenschaften gibt es eine Fülle von wohlbekannten Kernfunktion, die man in geeigneten Büchern findet. Nichtsdestotrotz wird der ‘Kernel Trick’ manchmal auch so verstanden, dass es eigentlich gar nicht nötig ist, zu wissen, ob das verwandte  $k$  die Kerneigenschaft wirklich hat: Denn hat man eine Lösung  $\lambda^*$  des Optimierungsproblems in Satz 15.13 gefunden, so kann man den induzierten Klassifizierer  $h_k$  implementieren und durch Einsetzen überprüfen, ob der  $h_k$  korrekt klassifiziert — ob  $k$  nun Kern ist oder nicht.

Wir schließen dieses Kapitel mit einem Hinweis darauf ab, dass eine sehr viel umfassendere Theorie der Kernfunktionen existiert als wir es hier in diesem Kapitel darlegen konnten.

**Bemerkung 15.18.** Sei  $X$  eine nichtleere Menge und  $(H, \langle -, - \rangle)$  ein Hilbertraum dessen Elemente Funktionen von  $X$  nach  $\mathbb{R}$  sind. Eine Funktion  $k: X \times X \rightarrow \mathbb{R}$  heißt

reproduzierender Kern, falls

$$\forall f \in H, x \in X: \langle k(\cdot, x), f \rangle,$$

gilt, also jede Funktion in  $H$  durch das Skalarprodukt und  $k$  ‘reproduziert’ wird. Den Raum  $H$  nennt man dann *Hilbertraum mit reproduzierendem Kern*. Die Spezialisierung auf reproduzierende Kerne erlaubt eine reichhaltigere Theorie als die der Kerne im Sinne von Definition 15.11 für die wir jedoch auf die unten genannten Referenzen verweisen.

## Referenzen

Der Satz über die Orthogonalprojektion 15.8 kann in [Wer18, Theorem V.3.4] nachgelesen werden und der Satz 15.17 über die Vervollständigung von Prähilberträumen in [Wer18, Satz V.1.8]. Wir notieren noch, dass Satz 15.8 oben nur im Spezialfall eines endlichdimensionalen Unterraumes  $U \subseteq X$  bzw. eines endlichdimensionalen Quotienten  $X/U$  gebraucht wird. In diesem Fall könnte man die Projektion auf  $U$ , bzw. auf das algebraische Komplement von  $U$ , per Wahl einer endlichen Basis wie in der Linearen Algebra beweisen, und käme so auch ohne Funktionalanalysis aus.

Die Hauptreferenzen für dieses Kapitel sind [SSBD14, AM12, SC08]. Ausführliche Rechnungen zu den Beispielen 15.12 findet man in [Sha08, Chapter 4.3]. *Mercers Bedingung* ist nach Mercers Theorem von 1909 benannt, siehe [SC08, S. 159]. Es gibt noch sehr viel allgemeinere Versionen des Representer Theorems, siehe [SC08, Chapter 4.5].

Wir weisen noch darauf hin, dass manche Autoren die Abbildung  $\psi$  als *feature map* bezeichnen und den Hilbertraum  $H$  als *feature space*. Außerdem sagt man manchmal, dass ‘die Kernfunktion  $k$  das Skalarprodukt bezüglich  $\psi$  implementiert’. Dies macht besonders Sinn im Hinblick auf Bemerkung 15.10, nach der man zur Lösung des quadratischen Optimierungsproblems nur  $k$  zu kennen braucht, nicht aber  $H$  und  $\psi$ .

## Aufgaben

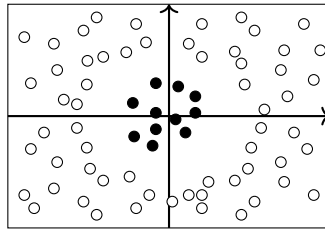
**Aufgabe 15.1.** Wir betrachten die folgende 1-dimensionale Datenmenge

$$D := \{(-3.5, 1), (0.5, 1), (0.75, 1), (-2.5, -1), (-1, -1), (5, -1)\} \subseteq \mathbb{R} \times \{-1, 1\}$$

bei welchem der erste Eintrag das Feature und der zweite Eintrag das Label angibt.

- (i) Skizzieren Sie die Datenmenge  $D$ .
- (ii) Finden Sie eine Abbildung  $\psi: \mathbb{R} \rightarrow \mathbb{R}^2$ , derart dass die abgebildete Datenmenge  $\hat{D} := \{(\psi(x), y) \mid (x, y) \in D\} \subseteq \mathbb{R}^2 \times \{-1, 1\}$  linear trennbar ist und skizzieren Sie  $\hat{D}$ .
- (iii) Bestimmen Sie einen Klassifizierer der Form  $h = \text{sign}(\langle w, - \rangle + b)$  für  $\hat{D}$ .
- (iv) Geben Sie an, welche  $x \in \mathbb{R}$  vom induzierten Klassifizierer  $h \circ \psi$  mit 1 bzw. mit  $-1$  klassifiziert werden.

**Aufgabe 15.2.** Finden Sie eine Funktion  $\psi: \mathbb{R}^2 \rightarrow \mathbb{R}^3$ , die eine Datenmenge der Gestalt wie im folgenden Bild injektiv und stetig auf eine linear trennbare Datenmenge abbilden kann.



**Aufgabe 15.3.** Zeigen Sie, dass die Funktion  $k: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ ,  $k(x, y) = (1 + \langle x, y \rangle)^m$  mit  $m \in \mathbb{N}$  eine Kernfunktion ist. Gehen Sie wie in Beispiel 15.12(i) vor und benutzen Sie den Multinomialssatz um die  $m$ -te Potenz auszurechnen.

**Aufgabe 15.4.** Sei  $\emptyset \neq X \subseteq \mathbb{R}$  und  $D = \{(x_i, y_i) \mid i = 1, \dots, n\} \subseteq X \times \{-1, 1\}$  eine Datenmenge.

- (i) Zeigen Sie, dass für jede Funktion  $\psi: X \rightarrow \mathbb{R}^n$  mit  $\psi(x_i) = e_i$  die abgebildete Datenmenge  $\hat{D} := \{(\psi(x_i), y_i) \mid i = 1, \dots, n\} \subseteq \mathbb{R}^n \times \{-1, 1\}$  linear trennbar ist, indem Sie die Parameter eines Klassifizierers  $h = \text{sign}(\langle w, \cdot \rangle + b)$  für  $\hat{D}$  erraten.
- (ii) Überlegen Sie sich einen Weg, um oben mit  $\mathbb{R}^2$  anstelle von  $\mathbb{R}^n$  auszukommen.
- (iii) Sei  $\psi: X \rightarrow H \subseteq \mathbb{R}^X$ ,  $\psi(x) = e^{-\frac{(x-\cdot)^2}{2\sigma^2}}$  die Abbildung aus dem Beweis von Satz 15.16. Zeigen Sie, dass  $\hat{D} := \{(\psi(x_i), y_i) \mid i = 1, \dots, n\} \subseteq H \times \{-1, 1\}$  für eine geeignete Wahl von  $\sigma$  linear trennbar ist.
- (iv) Welche der drei ‘Einbettungsmethoden’ finden Sie besser und warum?

**Aufgabe 15.5.** Bestimmen Sie für die Datenmenge aus Beispiel 15.1

Datenpunkt	1	2	3	4	5	6	7
Wert	-0.9	-0.3	0.2	0.4	0.9	1.3	1.6
Label	-1	-1	+1	+1	+1	-1	-1

Lösungen des Optimierungsproblems aus Satz 15.13 und zwar

- (i) für die Kernfunktion die sich aus dem in Beispiel 15.1 angegebenen  $\psi$  ergibt,
- (ii) für Polynom- und RBF-Kerne wie in Beispiel 15.12,
- (iii) für andere Funktionen  $k: \mathbb{R}^2 \rightarrow \mathbb{R}$ .

Plotten Sie die jeweiligen Klassifizierer als Funktion  $h: \mathbb{R} \rightarrow \mathbb{R}$  und die Paare (Wert, Label) um zu sehen, ob alle Punkte jeweils korrekt klassifiziert werden oder nicht.

**Aufgabe 15.6.** Zeigen Sie, dass für einen Kern  $k: X \times X \rightarrow \mathbb{R}$  weder der Hilbertraum  $(H, \langle \cdot, \cdot \rangle)$  noch die Funktion  $\psi: X \rightarrow H$  in Definition 15.11 eindeutig sind. Betrachten Sie hierfür  $k: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ ,  $k(x, y) := xy$ , sowie  $\psi_i: \mathbb{R} \rightarrow \mathbb{R}^i$ , definiert durch  $\psi_1(x) = x$  und  $\psi_2(x) = \frac{1}{\sqrt{2}}(x, x)$ .

**Aufgabe 15.7.** Zeigen Sie, dass die Kerneigenschaft stabil unter Einschränkung, unter Linearkombinationen und unter punktweisen Grenzwerten ist.

*Hinweis:* Für die Einschränkung und Linearkombinationen kann man die entsprechenden Hilberträume  $H$  und Abbildungen  $\psi$  erraten. Das der punktweise Grenzwert  $k$  eine Folge von Kernen  $(k_n)_{n \in \mathbb{N}}$  wieder ein Kern ist, kann mithilfe von Satz 15.16 gezeigt werden.

## Kapitel 16

# Neuronale Netze

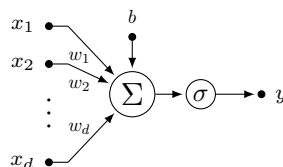
In diesem Kapitel beschäftigen wir uns mit künstlichen neuronalen Netzen. Deren zentralem Baustein, dem künstlichen Neuron, sind wir schon in früheren Kapiteln begegnet, ohne jedoch diesen Namen benutzt zu haben. Bevor wir die formale Definition angeben, weisen wir die Lesenden darauf hin, dass Teile dieses Kapitels Kenntnisse aus der Funktionalanalysis, der Maßtheorie und der Theorie der Fouriertransformation voraussetzen: Zwei Versionen des Satzes von Stone-Weierstraß, den Satz von Hahn-Banach, sowie den Rieszschen Darstellungssatz geben wir daher unten ohne Beweis an. Referenzen finden sich am Kapitelende. Das vorgenannte betrifft nur die zweite Hälfte des Kapitels 16.1 zur sogenannten Expressivität neuronaler Netze. Die dortigen Hauptresultate 16.20, 16.21, 16.22 und 16.25 sind außerdem so gehalten, dass sie vom Leser auch erstmal ohne Beweis zur Kenntnis genommen werden können.

Jetzt beginnen wir mit der Definition eines künstlichen Neurons.

**Definition 16.1.** Eine Funktion  $n: \mathbb{R}^d \rightarrow \mathbb{R}$  der Form

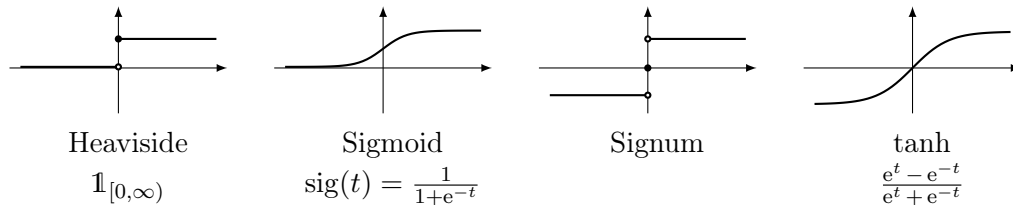
$$n(x) = \sigma(\langle w, x \rangle + b) = \sigma\left(\sum_{i=1}^d w_i x_i + b\right)$$

heißt (*künstliches*) Neuron mit Gewicht  $w = (w_1, \dots, w_d) \in \mathbb{R}^d$ , Bias  $b \in \mathbb{R}$  und Aktivierungsfunktion  $\sigma: \mathbb{R} \rightarrow \mathbb{R}$ . Häufig werden Neuronen durch Bilder wie das folgende veranschaulicht.



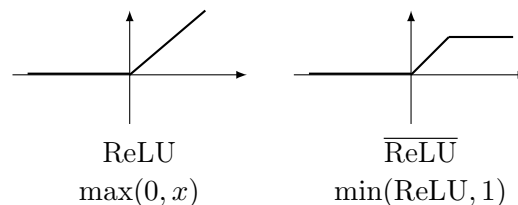
Künstliche Neuronen sind von den biologischen Neuronen des menschlichen Gehirns inspiriert. Bei letzteren sagt man häufig, dass ein solches Neuron *feuert*, wenn die gewichteten eingehenden Signale einen gewissen Schwellwert überschreiten. Diese Analogie passt besonders gut, wenn man bei der Aktivierungsfunktion an die Hea-

visidefunktion denkt, die, zusammen mit anderen typischen Aktivierungsfunktionen, unten dargestellt ist.



Ein Heaviside-aktiviertes Neuron liefert dann den Wert Eins, wenn  $\langle w, x \rangle \geq -b$  ausfällt; ansonsten liefert es den Wert Null. Eine glatte Variante hiervon erreicht man mit der Sigmoidfunktion. In der Tat sind die logistischen Regressoren, die wir in Kapitel 2 betrachtet haben, Beispiele für Sigmoid-aktivierte Neuronen. Das Perzeptron und damit auch jede Support Vector Maschine, siehe Kapitel 13 und 14, ist hingegen ein Beispiel für ein Signum-aktiviertes Neuron. Die Signumsfunktion kann auch wieder durch eine glatte Funktion, wie z.B. den Tangens Hyperbolicus, ersetzt werden.

Neben den oben genannten Aktivierungsfunktionen werden wir in diesem Kapitel insbesondere den *Rectified Linear Unifier (ReLU)* betrachten und dessen beschränkte Version (s.u.), welche wir mit  $\overline{\text{ReLU}}$  bezeichnen werden.



Im Gegensatz zu den vorgenannten Kapiteln, in denen stets einzelne Neuronen untersucht wurden, werden wir uns jetzt mit *neuronalen Netzen* beschäftigen, d.h. Funktionen des Typs in Definition 16.1 in geeigneter Weise verketteten oder linear kombinieren. Diese Idee ist ebenfalls von der Natur inspiriert, in der sich Neuronen durch *Synapsen* verbinden können. Hierbei hat dann jedes Neuron seine eigenen Gewichte und sein eigenes Bias, aber alle Neuronen in einem Netz haben in der Regel die gleiche Aktivierungsfunktion.

Im folgenden Unterkapitel 16.1 konzentrieren wir uns erstmal auf die Frage, welche Funktionen ein neuronales Netz exakt darstellen, bzw. in einem noch zu präzisierenden Sinn, approximieren kann. Dies bezeichnet man häufig als *Expressivität* eines Netzes oder einer ganzen Klasse von Netzen.

## 16.1 Expressivität

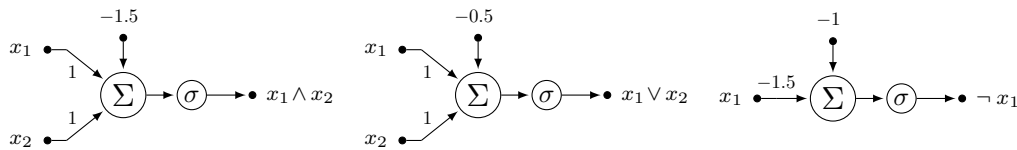
Bei der oben erwähnten Frage, welche Funktionen mit einem gewissen Typ neuronalen Netzes dargestellt oder approximiert werden können, ist es sinnvoll zwischen



Klassifikationsproblemen und Regressionsproblemen zu unterscheiden. Wir beginnen mit den folgenden diskreten, und daher eher zu Klassifikationsproblemen passenden, Resultaten.

**Proposition 16.2.** *Mit  $\sigma = \mathbb{1}_{[0,\infty)}$  können die logischen Funktionen  $\wedge, \vee, \neg$  jeweils durch ein Neuron dargestellt werden.*

*Beweis.* In der Tat haben wir

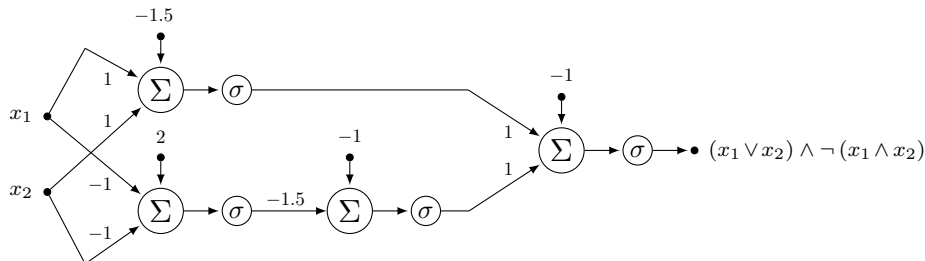


wie behauptet. □

**Bemerkung 16.3.** Im Gegensatz zu Proposition 16.2 lässt sich die Entweder-Oder-Funktion XOR nicht durch ein einzelnes Neuron darstellen, siehe Aufgabe 16.1. Via der Formel

$$\text{XOR}(x_1, x_2) = (x_1 \vee x_2) \wedge \neg(x_1 \wedge x_2)$$

können wir allerdings Neuronen aus Proposition 16.3 wie folgt zusammensetzen



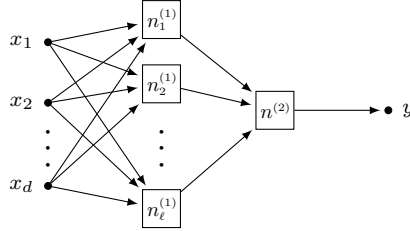
und erhalten XOR mittels vier Neuronen. In Aufgabe 16.1 werden wir zeigen, dass auch eine Darstellung mit nur drei Neuronen möglich ist.

Aus dem obigem folgt direkt, dass jede beliebige Funktion  $f: \{0,1\}^d \rightarrow \{0,1\}$ , bzw. ihre konjunktive Normalform

$$N(z_1, \dots, z_d) = \bigwedge_{f(x_i)=0} y_{i1} \vee \dots \vee y_{id}, \quad \text{wobei } y_{ik} = \begin{cases} \neg z_k, & \text{falls } x_{ik} = 1, \\ z_k, & \text{sonst.} \end{cases}$$

durch ein Heaviside-aktiviertes neuronales Netz dargestellt werden kann. Ein konkretes Beispiel findet sich in Aufgabe 16.2. Wie wir bereits im Kapitel über das Perzeptron beobachtet haben, vergleiche Bemerkung 13.14, ist es manchmal von Vorteil, sowohl bei den Features als auch bei den Labels die Werte  $\{-1, 1\}$  anstelle von  $\{0, 1\}$  zu benutzen und dann entsprechend die Heavisidefunktion durch die Signumsfunktion zu ersetzen. Machen wir dies, so erhalten wir die folgende, sehr elegante, Darstellung der oben diskutierten ‘boolschen Funktionen’.

**Satz 16.4.** Jede Funktion  $f: \{-1, 1\}^d \rightarrow \{-1, 1\}$  kann durch ein neuronales Netz der folgenden Form



dargestellt werden, wobei die  $n_i^{(k)}$  Signum-aktivierte Neuronen sind.

*Beweis.* Seien  $w_i \in \{-1, 1\}^d$ ,  $i = 1, \dots, \ell$ , diejenigen Punkte, für die  $f(w_i) = 1$  gilt. Dann erhalten wir

$$f = n^{(2)} \circ \begin{bmatrix} n_1^{(1)} \\ \vdots \\ n_\ell^{(1)} \end{bmatrix}$$

mit  $n^{(2)} = \text{sign}(\langle \mathbf{1}, \cdot \rangle + \ell - 1.5)$  und  $n_i^{(1)} = \text{sign}(\langle w_i, \cdot \rangle - d + 1)$  für  $i = 1, \dots, \ell$ . In der Tat haben wir für  $x = w_i$

$$n_i^{(1)}(x) = \text{sign}(\langle w_i, w_i \rangle - d + 1) = \text{sign}\left(\sum_{j=1}^d w_{ij}^2 - d + 1\right) = 1.$$

Für  $x = (x_1, \dots, x_d) \neq w_i$  gibt es mindestens ein  $j_0$  derart, dass  $x_{j_0} \neq w_{ij_0}$  ist, und daher gilt  $x_{j_0} \cdot w_{ij_0} = -1$ . Daraus folgt für solche  $x \in \{-1, 1\}^d$

$$n_i^{(1)}(x) = \text{sign}\left(w_{ij_0}x_{j_0} + \sum_{\substack{j=1 \\ j \neq j_0}}^d w_{ij}x_j - d + 1\right) = -1,$$

da wir das Argument der Signumfunktion nach oben durch  $-1$  abschätzen können. Für ein beliebiges  $x \in \{-1, 1\}^d$  liefert also die sogenannte erste Schicht des neuronalen Netzes einen Vektor  $(n_1^{(1)}(x), \dots, n_\ell^{(1)}(x)) \in \{-1, 1\}^\ell$  bei welchem der  $i$ -te Eintrag genau dann Eins ist, wenn  $x = w_i$  ist. Es folgt

$$\left\langle \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}, \begin{bmatrix} n_1^{(1)}(x) \\ \vdots \\ n_\ell^{(1)}(x) \end{bmatrix} \right\rangle = \sum_{i=1}^{\ell} n_i^{(1)}(x) = \begin{cases} -\ell + 2, & \text{falls } \exists i \in \{1, \dots, \ell\}: x = w_i, \\ -\ell, & \text{sonst.} \end{cases}$$

Da die  $w_i$  per Konstruktion genau diejenigen Elemente von  $\{-1, 1\}^d$  sind, für die  $f(w_i) = 1$  gilt, erhalten wir

$$n^{(2)}\left(\begin{bmatrix} n_1^{(1)}(x) \\ \vdots \\ n_\ell^{(1)}(x) \end{bmatrix}\right) = \text{sign}\left(\left\langle \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}, \begin{bmatrix} n_1^{(1)}(x) \\ \vdots \\ n_\ell^{(1)}(x) \end{bmatrix} \right\rangle + \ell - 1.5\right) = f(x)$$

für jedes  $x \in \{-1, 1\}^d$  wie behauptet.  $\square$

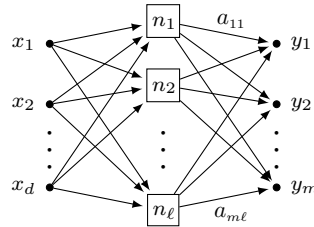
Wir bemerken, dass es bei binären Klassifizierern Sinn macht, ein neuronales Netz mit *Architektur* wie in Satz 16.4 zu verwenden, d.h. nach einer *versteckten Schicht* eine *Ausgangsschicht* mit nur einem Neuron zu verwenden, welches dafür sorgt, dass das Netz  $\{0,1\}$ -wertig wird, wenn man mit Heaviside-Aktivierung arbeitet. Alternativ erhält man mit Signum-Aktivierung, für alle Neuronen oder nur bei  $n^{(2)}$ , ein  $(0,1)$ -wertiges Netz. Im letzteren Fall kann man die Ausgabe als Wahrscheinlichkeit interpretieren, vgl. Aufgabe 16.5. Wir kommen später auch noch auf Klassifizierer mit mehrdimensionalen Labels zu sprechen, siehe Satz 16.34 und die Bemerkungen davor.

Jetzt betrachten wir erstmal die Expressivität in einem analytischen Kontext. Wir wollen neuronale Netze als Regressoren einsetzen und betrachten dazu zunächst sogenannte *flache neuronale Netze*. Deren Architektur ist sehr ähnlich zu der in Satz 16.4 betrachteten, allerdings verzichten wir auf das *Ausgangsneuron*  $n^{(2)}$  und ersetzen dieses durch eine *Ausgangsmatrix*.

**Definition 16.5.** Eine Funktion  $N: \mathbb{R}^d \rightarrow \mathbb{R}^m$  der Form

$$N = \begin{bmatrix} a_{11} & \cdots & a_{1\ell} \\ \vdots & & \vdots \\ a_{m1} & \cdots & a_{m\ell} \end{bmatrix} \begin{bmatrix} n_1 \\ \vdots \\ n_\ell \end{bmatrix}$$

mit Neuronen  $n_1, \dots, n_\ell$  und einer Matrix  $(a_{ij})_{i,j} \in \mathbb{R}^{\ell \times m}$  heißt *flaches, vollständig verbundenes, vorwärtspropagierendes Netz der Breite  $\ell$  mit linearem Ausgang und Aktivierung  $\sigma: \mathbb{R} \rightarrow \mathbb{R}$* . Ein solches Netz kann durch das folgende Bild beschrieben werden.



Wir bezeichnen mit  $\mathcal{S}^{\sigma,\ell}(\mathbb{R}^d, \mathbb{R}^m)$  die Menge aller neuronalen Netze der obigen Form und setzen

$$\mathcal{S}^{\sigma}(\mathbb{R}^d, \mathbb{R}^m) := \bigcup_{\ell=1}^{\infty} \mathcal{S}^{\sigma,\ell}(\mathbb{R}^d, \mathbb{R}^m).$$

Die Elemente von  $\mathcal{S}^{\sigma}(\mathbb{R}^d, \mathbb{R}^m)$  bezeichnen wir im Folgenden kurz als *flache Netze*. Wir setzen  $\mathcal{S}^{\sigma,\ell}(\mathbb{R}^d) := \mathcal{S}^{\sigma,\ell}(\mathbb{R}^d, \mathbb{R})$  und  $\mathcal{S}^{\sigma}(\mathbb{R}^d) := \mathcal{S}^{\sigma}(\mathbb{R}^d, \mathbb{R})$ . Liest man die Anwendung der Aktivierungsfunktion koordinatenweise, so kann man ein flaches Netz  $N \in \mathcal{S}^{\sigma}(\mathbb{R}^d, \mathbb{R}^m)$  durch die Formel

$$N(x) = A\sigma(Wx + b)$$

mit  $W \in \mathbb{R}^{\ell \times d}$ ,  $b \in \mathbb{R}^{\ell}$  und  $A \in \mathbb{R}^{\ell \times m}$  darstellen. Hierbei enthält die  $i$ -te Zeile von  $W$  die Gewichte des Neurons  $n_i$ .

Wir wollen nun untersuchen, für welche Aktivierungsfunktionen  $\sigma$  jede stetige

Funktion  $f: \mathbb{R}^d \rightarrow \mathbb{R}^m$  durch ein flaches neuronales Netz approximiert werden kann, und zwar im Sinne von ‘gleichmäßiger Konvergenz auf Kompakta’. Wir bemerken dabei, dass dies einerseits ein sehr natürlicher Konvergenzbegriff ist, der bereits aus den grundlegenden Analysisvorlesungen bekannt ist; z.B. konvergieren die Reihenentwicklungen von  $\sin$ ,  $\cos$  und  $\exp$  in genau diesem Sinne. Andererseits liegen Datenmengen oft in einem natürlich gegebenen Kompaktum, z.B. weil bei einem Vektor aus Messwerten für jeden Eintrag eine obere und untere Schranke aus der Anwendungssituation heraus bekannt ist. Wir erhalten dann auf diesem Kompaktum gleichmäßige Konvergenz und damit eine sehr starke Approximationsaussage.

**Definition 16.6.** Eine Teilmenge  $\mathcal{S} \subseteq C(\mathbb{R}^d, \mathbb{R}^m)$  des Raumes der stetigen Funktionen von  $\mathbb{R}^d$  nach  $\mathbb{R}^m$  heißt *dicht*, oder genauer *gleichmäßig dicht auf Kompakta*, falls mit einer beliebigen Norm  $\|\cdot\|_{\mathbb{R}^m}$  auf  $\mathbb{R}^m$  gilt

$$\forall F \in C(\mathbb{R}^d, \mathbb{R}^m), \Omega \subseteq \mathbb{R}^d \text{ kompakt}, \varepsilon > 0 \exists S \in \mathcal{S}: \\ \|F - S\|_{\Omega, \infty} := \sup_{x \in \Omega} \|F(x) - S(x)\|_{\mathbb{R}^m} < \varepsilon.$$

Ist  $\sigma: \mathbb{R} \rightarrow \mathbb{R}$  stetig, so gilt offenbar  $\mathcal{S}^\sigma(\mathbb{R}^d, \mathbb{R}^m) \subseteq C(\mathbb{R}^d, \mathbb{R}^m)$ , und wir können fragen, ob dies ein dichter Teilraum im Sinne von Definition 16.6 ist. Wir geben zwei Beispiele.

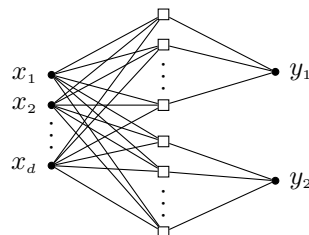
**Beispiel 16.7.** (i) Sei  $P \in \mathbb{R}[X]$  ein Polynom. Schreibt man  $\mathcal{S}^P(\mathbb{R})$  explizit auf, so sieht man sofort, dass dieser Raum endlichdimensional ist. Für kompaktes  $\Omega \subseteq \mathbb{R}$  ist dann  $\mathcal{S}^P(\Omega) \subset (C(\Omega), \|\cdot\|_{\Omega, \infty})$  abgeschlossen und nicht dicht.

(ii) Sei  $\exp: \mathbb{R} \rightarrow \mathbb{R}$  die Exponentialfunktion. Dann ist  $\mathcal{S}^{\exp}(\mathbb{R}) \subseteq C(\mathbb{R})$  gleichmäßig dicht auf Kompakta. Dies folgt durch Anwendung des Satzes 16.10 von Stone-Weierstraß, siehe Aufgabe 16.3.

Die folgende Proposition erlaubt es in der Folge auf den skalarwertigen Fall zu reduzieren.

**Proposition 16.8.** Sei  $\sigma: \mathbb{R} \rightarrow \mathbb{R}$  stetig. Falls  $\mathcal{S}^\sigma(\mathbb{R}^d) \subseteq C(\mathbb{R}^d)$  gleichmäßig dicht auf Kompakta ist, dann ist auch  $\mathcal{S}^\sigma(\mathbb{R}^d, \mathbb{R}^m) \subseteq C(\mathbb{R}^d, \mathbb{R}^m)$  gleichmäßig dicht auf Kompakta.

*Beweis.* Bei gegebenem  $F \in C(\mathbb{R}^d, \mathbb{R}^m)$  approximieren wir jede Koordinatenfunktion  $F_i$  entsprechend der Voraussetzung mit einem  $N_i \in \mathcal{S}^\sigma(\mathbb{R}^d)$  und legen diese, wie im folgenden Bild für  $m = 2$  gezeigt, übereinander. Die unten nicht eingezeichneten Linien entsprechen Nullen in der Ausgangsmatrix  $A$ .



Auf diese Weise erhalten wir ein flaches Netz  $N$ , welches  $F$  im gewünschten Sinn approximiert.  $\square$

Als nächstes wollen wir uns ebenso beim Definitionsbereich auf den 1-dimensionalen Fall zurückziehen. Dies erfordert allerdings etwas Vorbereitung.

**Lemma 16.9.** *Sei  $\mathcal{P} := \{\langle \cdot, a \rangle^r : \mathbb{R}^d \rightarrow \mathbb{R} \mid r \in \mathbb{N}_0, a \in \mathbb{R}^d\}$ . Dann ist  $\text{span } \mathcal{P} \subseteq C(\mathbb{R}^d)$  gleichmäßig dicht auf Kompakta.*

*Beweis.* Wir benötigen für den Beweis die folgende Version des Satzes von Stone-Weierstraß. Am Ende des Kapitels geben wir detaillierte Referenzen zu diesem Satz an.

**Satz 16.10.** (Stone-Weierstraß, Version 1) *Sei  $\Omega \subseteq \mathbb{R}^d$  kompakt und sei  $A \subseteq C(\Omega)$  eine Teilmenge derart, dass die folgenden Eigenschaften erfüllt sind.*

- (i)  *$A$  ist eine Unteralgebra von  $C(\Omega)$ ,*
- (ii) *Es gilt  $1_\Omega \in A$ ,*
- (iii)  *$A$  ist punktetrennend, d.h.  $\forall x \neq y \in \Omega \exists g \in A: g(x) \neq g(y)$ .*

*Dann ist  $A \subseteq C(\Omega)$  dicht bezüglich der Supremumsnorm.*  $\diamond$

Wir fixieren ein Kompaktum  $\Omega \subseteq \mathbb{R}^d$ . Da die Elemente von  $\mathcal{P}$  eindeutig durch die Parameter  $a$  und  $r$  gegeben sind, können wir  $A := \text{span } \mathcal{P} \subseteq C(\Omega)$  als Teilmenge auffassen und müssen nun die Voraussetzungen in Satz 16.10 überprüfen. Da  $A$  per Definition ein Untervektorraum ist, muss für die Unteralgebra-Eigenschaft nur die Abgeschlossenheit unter punktweiser Multiplikation gezeigt werden. Seien dazu  $a, b \in \mathbb{R}^d$  und  $r, s \in \mathbb{N}_0$ . Wir wählen paarweise verschiedene  $\beta_0, \dots, \beta_{r+s} \in \mathbb{R}$  und setzen

$$c = \begin{bmatrix} c_0 \\ \vdots \\ c_{r+s} \end{bmatrix} := \binom{r+s}{s}^{-1} \underbrace{\begin{bmatrix} \beta_0^0 & \cdots & \beta_{r+s}^0 \\ \vdots & & \vdots \\ \beta_0^{r+s} & \cdots & \beta_{r+s}^{r+s} \end{bmatrix}}_{=:B}^{-1} \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

wobei  $B$  als transponierte Vandermondematrix invertierbar ist. Bezeichnen wir mit  $\delta_{ij}$  das Kroneckersymbol, d.h.  $\delta_{ij} = 1$  für  $i = j$  und  $\delta_{ij} = 0$  für  $i \neq j$ , so erhalten wir für  $0 \leq i \leq r+s$

$$\sum_{j=0}^{r+s} \beta_j^i c_j = (Bc)_i = \binom{r+s}{s}^{-1} \delta_{is} = \binom{r+s}{i}^{-1} \delta_{is}$$

also

$$\delta_{is} = \binom{r+s}{i} \sum_{j=0}^{r+s} \beta_j^i c_j.$$

Für  $x \in \mathbb{R}^d$  berechnen wir nun

$$\begin{aligned}
 \langle x, a \rangle^r \langle x, b \rangle^s &= \sum_{i=0}^{r+s} \delta_{is} \langle x, a \rangle^{r+s-i} \langle x, b \rangle^i \\
 &= \sum_{i=0}^{r+s} \left[ \binom{r+s}{i} \sum_{j=0}^{r+s} \beta_j^i c_j \right] \langle x, a \rangle^{r+s-i} \langle x, b \rangle^i \\
 &= \sum_{j=0}^{r+s} c_j \left[ \sum_{i=0}^{r+s} \binom{r+s}{i} \beta_j^i \langle x, a \rangle^{r+s-i} \langle x, b \rangle^i \right] \\
 &= \sum_{j=0}^{r+s} c_j [\langle x, a \rangle + \beta_j \langle x, b \rangle]^{r+s} \\
 &= \sum_{j=0}^{r+s} c_j \langle x, a + \beta_j b \rangle^{r+s}
 \end{aligned}$$

und sehen, dass  $\langle \cdot, a \rangle^r \langle \cdot, b \rangle^s \in A$  gilt und  $A$  damit eine Unteralgebra ist. Durch Wahl von  $r = 0$  folgt  $\mathbf{1}_\Omega \in A$  und für  $x \neq y$  in  $\Omega$  bildet z.B.  $\langle \cdot, x - y \rangle$  diese Punkte auf unterschiedliche Werte ab, andernfalls wäre nämlich  $\|x - y\|^2 = \langle x - y, x - y \rangle = 0$ . Damit ist  $A$  also auch punktetrennend.

Sind jetzt  $F \in C(\mathbb{R}^d)$  und  $\varepsilon > 0$  gegeben, so betrachten wir  $f := F|_\Omega$  und erhalten nach Satz 16.10 ein  $g \in A$  mit  $\|F - g\|_{\Omega, \infty} < \varepsilon$ . Hierbei kann  $g$  in natürlicher Weise als Funktion auf  $\mathbb{R}^d$  betrachtet werden und wir sind fertig.  $\square$

**Lemma 16.11.** *Sei  $\mathcal{F} \subseteq C(\mathbb{R})$  gleichmäßig dicht auf Kompakta. Dann ist auch  $\mathcal{G} := \text{span}\{f(\langle w, \cdot \rangle) \mid w \in \mathbb{R}^d, f \in \mathcal{F}\} \subseteq C(\mathbb{R}^d)$  gleichmäßig dicht auf Kompakta.*

*Beweis.* Sei  $F \in C(\mathbb{R}^d)$ ,  $\Omega \subseteq \mathbb{R}^d$  kompakt, und  $\varepsilon > 0$  gegeben. Sei  $\mathcal{P}$  so definiert wie in Lemma 16.9. Dann existieren  $w_1, \dots, w_m \in \mathbb{R}^d$ ,  $r_1, \dots, r_m \in \mathbb{N} \cup \{0\}$ , sowie  $\alpha_1, \dots, \alpha_m \in \mathbb{R}$  sodass

$$\left\| F - \sum_{j=1}^m \alpha_j \langle w_j, \cdot \rangle^{r_j} \right\|_{\Omega, \infty} < \varepsilon/2$$

gilt. Für fixes  $1 \leq j \leq m$  definieren wir  $K_j := \{\langle w_j, x \rangle \mid x \in \Omega\}$ . Da  $K_j \subseteq \mathbb{R}$  kompakt ist, können wir per Voraussetzung  $f_{j1}, \dots, f_{j\ell} \in \mathcal{F}$  und  $\alpha_{j1}, \dots, \alpha_{j\ell} \in \mathbb{R}$  wählen, sodass

$$\left\| (\cdot)^{r_j} - \sum_{k=1}^{\ell} \alpha_{jk} f_{jk} \right\|_{K_j, \infty} < \frac{\varepsilon}{2\|\alpha\|_1}$$

gilt. Per Konstruktion ist dann

$$\left\| \langle w_j, \cdot \rangle^{r_j} - \sum_{k=1}^{\ell} \alpha_{jk} f_{jk}(\langle w_j, \cdot \rangle) \right\|_{\Omega, \infty} < \frac{\varepsilon}{2\|\alpha\|_1}$$

wobei die Linearkombination der  $f_{jk}(\langle w_j, \cdot \rangle)$  per Definition zu  $\mathcal{G}$  gehört. Folglich ist

auch

$$g := \sum_{j=1}^m \alpha_j \sum_{k=1}^{\ell} \alpha_{jk} f_{jk}(\langle w_j, \cdot \rangle) \in \mathcal{G}$$

und wir erhalten die Abschätzung

$$\begin{aligned} \|F - g\|_{\Omega, \infty} &\leq \|F - \sum_{j=1}^m \alpha_j \langle w_j, \cdot \rangle^{r_j}\|_{\Omega, \infty} + \|\sum_{j=1}^m \alpha_j \langle w_j, \cdot \rangle^{r_j} - g\|_{\Omega, \infty} \\ &< \varepsilon/2 + \sum_{j=1}^m |\alpha_j| \|\langle w_j, \cdot \rangle^{r_j} - \sum_{k=1}^{\ell} \alpha_{jk} f_{jk}(\langle w_j, \cdot \rangle)\|_{\Omega, \infty} \\ &< \varepsilon/2 + (\varepsilon/2) \sum_{j=1}^m |\alpha_j| / \|\alpha\|_1 = \varepsilon \end{aligned}$$

wie gewünscht.  $\square$

Mit den beiden Lemmas 16.9 und 16.11 können wir jetzt das folgende Reduktionsresultat für den Zielbereich beweisen.

**Proposition 16.12.** (Satz von Chui, Li, Lin, Pinkus) *Sei  $\sigma: \mathbb{R} \rightarrow \mathbb{R}$  stetig und derart dass  $\mathcal{S}^\sigma(\mathbb{R}) \subseteq C(\mathbb{R})$  gleichmäßig dicht auf Kompakta ist. Dann ist  $\mathcal{S}^\sigma(\mathbb{R}^d) \subseteq C(\mathbb{R}^d)$  ebenfalls gleichmäßig dicht auf Kompakta.*

*Beweis.* Per Definition sind die Elemente  $f \in \mathcal{S}^\sigma(\mathbb{R})$  von der Form  $f(x) = \sigma(vx + b)$  mit  $v, b \in \mathbb{R}$ . Damit erhalten wir

$$\begin{aligned} \mathcal{S}^\sigma(\mathbb{R}^d) &= \left\{ \sum_{i=1}^{\ell} a_i \sigma(\langle w_i, \cdot \rangle + b_i) \mid a_i \in \mathbb{R}, w_i \in \mathbb{R}^d, b_i \in \mathbb{R}, \ell \geq 1 \right\} \\ &= \text{span}\{\sigma(\langle w, \cdot \rangle + b) \mid w \in \mathbb{R}^d, b \in \mathbb{R}\} \\ &= \text{span}\{\sigma(v \langle w, \cdot \rangle + b) \mid v \in \mathbb{R}, w \in \mathbb{R}^d, b \in \mathbb{R}\} \\ &= \text{span}\{f(\langle w, \cdot \rangle) \mid w \in \mathbb{R}^d, f \in \mathcal{S}^\sigma(\mathbb{R})\}. \\ &\quad \uparrow \\ &\quad f = \sigma(v \cdot + b) \\ &\quad v, b \in \mathbb{R} \end{aligned}$$

Die letztere Menge ist aber gerade gleich der Menge  $\mathcal{G}$  aus Lemma 16.11 und nach diesem dicht in  $C(\mathbb{R}^d)$ .  $\square$

Nach den Propositionen 16.8 und 16.12 können wir uns bei der Frage, ob die flachen Netze  $\mathcal{S}^\sigma(\mathbb{R}^d, \mathbb{R}^m) \subseteq C(\mathbb{R}^d, \mathbb{R}^m)$  in den stetigen Funktionen gleichmäßig dicht auf Kompakta liegen, ohne Einschränkung auf den Fall  $d = m = 1$  zurückziehen. Für diesen führen wir einen weiteren Begriff ein.

**Definition 16.13.** Für kompaktes  $\Omega \subseteq \mathbb{R}$  bezeichnen wir mit  $M(\Omega)$  den Raum der komplexen Borelmaße und nennen eine messbare Funktion  $\sigma: \mathbb{R} \rightarrow \mathbb{R}$  *diskriminativ*, falls für jede kompakte Menge  $\Omega \subseteq \mathbb{R}$  und für  $\mu \in M(\Omega)$  gilt

$$\left( \forall w, b \in \mathbb{R}: \int_{\Omega} \sigma(wx + b) d\mu(x) = 0 \right) \implies \mu = 0.$$

Die nächsten zwei Lemmas sind die finalen Zutaten für die danach folgenden Approximationssätze für stetige Aktivierungsfunktionen.

**Lemma 16.14.** *Ist die Fouriertransformierte  $\hat{\mu}(s) = \int_{\mathbb{R}^d} e^{-ist} d\mu(t)$  eines Maßes  $\mu \in M(\mathbb{R})$  gleich Null für jedes  $s \in \mathbb{R}$ , so ist bereits das Maß  $\mu = 0$ .*

*Beweis.* ① Für  $\varphi \in L^1(\mathbb{R}, \mathbb{C})$  betrachten wir die ‘normale’ Fouriertransformierte

$$[\mathcal{F}\varphi](\xi) \equiv \hat{\varphi}(\xi) := \int_{\mathbb{R}} e^{-ix\xi} \varphi(x) dx$$

für welche nach dem Riemann-Lebesgue-Lemma  $\hat{f} \in C_0(\mathbb{R}, \mathbb{C})$  gilt. Wir betrachten die Menge  $A := \{\hat{f} \mid f \in L^1(\mathbb{R}, \mathbb{C})\} \subseteq C_0(\mathbb{R}, \mathbb{C})$  und zeigen mithilfe der folgenden Version des Satzes von Stone-Weierstraß, dass diese dicht in  $C_0(\mathbb{R}, \mathbb{C})$  liegt.

**Satz 16.15.** (Stone-Weierstraß, Version 2) *Sei  $A$  eine Teilmenge des Raumes  $C_0(\mathbb{R}, \mathbb{C})$  :=  $\{f \in C(\mathbb{R}, \mathbb{C}) \mid \lim_{|x| \rightarrow \infty} f(x) = 0\}$  der stetigen Funktionen die im Unendlichen verschwinden, derart dass die folgenden Eigenschaften erfüllt sind.*

- (i)  *$A$  ist eine Unteralgebra von  $C_0(\mathbb{R}, \mathbb{C})$ .*
- (ii)  *$A$  ist punktetrennend, d.h.  $\forall x \neq y \in \mathbb{R} \exists g \in A: g(x) \neq g(y)$ .*
- (iii)  *$A$  ist nicht verschwindend, d.h.  $\forall x \in \mathbb{R} \exists g \in A: g(x) \neq 0$ .*
- (iv)  *$A$  ist abgeschlossen unter komplexer Konjugation, d.h.  $\forall g \in A: \bar{g} \in A$ .*

*Dann ist  $A \subseteq C_0(\mathbb{R}, \mathbb{C})$  dicht bezüglich Supremumsnorm.* ◇

Da die Fouriertransformation linear ist, folgt mit der leicht zu prüfenden und wohlbekannten Identität  $\mathcal{F}(f_1 \star f_2) = \mathcal{F}f_1 \cdot \mathcal{F}f_2$  zunächst, dass  $A$  eine Unteralgebra ist. Fouriertransformation von  $f(x) := e^{-x} \mathbf{1}_{[0, \infty)}(x)$  liefert

$$\hat{f}(\xi) = \int_0^\infty e^{-ix\xi} e^{-x} dx = e^{-(i\xi+1)x} \frac{1}{i\xi+1} \Big|_0^\infty = \frac{1}{i\xi+1}$$

woran wir sehen, dass  $\hat{f}(\xi) \neq \hat{f}(\zeta)$  für alle  $\xi \neq \zeta$  sowie  $\hat{f}(\xi) \neq 0$  für jedes  $\xi \in \mathbb{R}$  gelten. Damit ist  $A$  nicht verschwindend und punktetrennend. Mit einer weiteren wohlbekannten Identität, nämlich  $[\mathcal{F}\bar{f}](\xi) = \overline{[\mathcal{F}f]}(-\xi)$  für  $\xi \in \mathbb{R}$ , erhalten wir schließlich die Abgeschlossenheit von  $A$  unter komplexer Konjugation.

② Wir bezeichnen jetzt mit  $C_0(\mathbb{R}, \mathbb{C})'$  den Dualraum von  $C_0(\mathbb{R}, \mathbb{C})$ , d.h. den Raum aller linearen und stetigen Funktionalen  $\varphi: C_0(\mathbb{R}, \mathbb{C}) \rightarrow \mathbb{C}$  ausgestattet mit der Operatornorm. Mit  $M(\Omega)$  bezeichnen wir den Raum der komplexen, regulären Borelmaße mit endlicher Variationsnorm. Für das folgende benötigen wir:

**Satz 16.16.** (Rieszscher Darstellungssatz, Version 1) *Die Abbildung  $T: M(\mathbb{R}) \rightarrow C_0(\mathbb{R}, \mathbb{C})'$ ,  $[T\mu](f) = \int_{\mathbb{R}} f d\mu$ , ist ein Isomorphismus von Banachräumen.* ◇

Um das Lemma zu beweisen sei nun  $\mu \in M(\mathbb{R})$  mit  $\hat{\mu} = 0$  gegeben. Sei  $T$  der Isomorphismus aus Satz 16.16. Dann ist  $T\mu \in C_0(\mathbb{R}, \mathbb{C})'$  und wir zeigen, dass dieses



auf  $A$  verschwindet. Sei dazu  $f \in L^1(\mathbb{R}, \mathbb{C})$ . Dann folgt mit dem Satz von Fubini (für komplexe Maße!)

$$\begin{aligned} [T\mu](\hat{f}) &= \int_{\mathbb{R}} \hat{f}(\xi) d\mu(\xi) = \int_{\mathbb{R}} \int_{\mathbb{R}} e^{-ix\xi} f(x) dx d\mu(\xi) \\ &= \int_{\mathbb{R}} f(x) \int_{\mathbb{R}} e^{-ix\xi} d\mu(\xi) dx = \int_{\mathbb{R}} f(x) \hat{\mu}(\xi) dx = 0. \end{aligned}$$

Nach Teil ① ist  $A \subseteq C_0(\mathbb{R}, \mathbb{C})$  dicht, d.h.  $[T\mu](f) = 0$  für jedes  $f \in C_0(\mathbb{R}, \mathbb{C})$ . Damit ist aber dann  $T\mu = 0$  und da  $T$  ein Isomorphismus ist, folgt  $\mu = 0$ .  $\square$

**Lemma 16.17.** *Sei  $\sigma: \mathbb{R} \rightarrow \mathbb{R}$  stetig und diskriminatorisch. Dann ist  $\mathcal{S}^\sigma(\mathbb{R}) \subseteq C(\mathbb{R})$  gleichmäßig dicht auf Kompakta.*

*Beweis.* Für den Beweis benötigen wir die folgende Konsequenz aus dem Satz von Hahn-Banach, die wir ohne Beweis notieren.

**Satz 16.18.** (Korollar zum Satz von Hahn-Banach) *Sei  $X$  ein normierter Raum und  $Y \subseteq X$  ein abgeschlossener Unterraum und  $x \in X \setminus Y$ . Dann existiert eine lineare und stetige Abbildung  $\varphi: X \rightarrow \mathbb{R}$  mit  $\varphi|_Y = 0$  und  $\varphi(x) \neq 0$ .*  $\diamond$

Angenommen,  $\mathcal{S}^\sigma(\mathbb{R}) \subseteq C(\mathbb{R})$  ist nicht gleichmäßig dicht auf Kompakta. Dann gibt es ein Kompaktum  $\Omega \subseteq \mathbb{R}$  sodass  $\overline{\mathcal{S}^\sigma(\Omega)} \subset C(\Omega)$  ein abgeschlossener echter Unterraum ist. Nach Satz 16.18 existiert  $0 \neq \varphi \in C(\Omega)'$ , das auf  $\mathcal{S}^\sigma(\Omega)$  verschwindet. Jetzt wenden wir die folgende, etwas andere, Variante des Rieszscher Darstellungssatzes an.

**Satz 16.19.** (Rieszscher Darstellungssatz, Version 2) *Sei  $\Omega \subseteq \mathbb{R}$  kompakt. Dann ist die Abbildung  $T: M(\Omega) \rightarrow C(\Omega, \mathbb{C})'$ ,  $[T\mu](f) = \int_{\Omega} f d\mu$ , ist ein Isomorphismus von Banachräumen.*  $\diamond$

Mithilfe von dieser wählen wir  $\mu \in M(\Omega)$  mit  $T\mu = \varphi$ . Das heißt, es gilt

$$\forall f \in C(\Omega): \varphi(f) = \int_{\Omega} f(x) d\mu(x).$$

Wir spezialisieren  $f := \sigma(w \cdot + b) \in \mathcal{S}^\sigma(\Omega)$  und erhalten

$$\forall w, b \in \mathbb{R}: \int_{\Omega} \sigma(wx + b) d\mu(x) = \varphi(f) = 0.$$

Da  $\sigma$  diskriminatorisch ist, ist  $\mu = 0$  und damit  $\varphi = 0$ . Widerspruch.  $\square$

Nun kommen wir zum ersten Approximationssatz. Beachte, dass dieser insbesondere die ‘S-förmigen’ Aktivierungsfunktionen tanh, Sigmoid und  $\overline{\text{ReLU}}$  vom Anfang des Kapitels abdeckt, dabei aber natürlich weit über diese drei Beispiele hinausgeht.

**Satz 16.20.** (von Cybenko) *Sei  $\sigma: \mathbb{R} \rightarrow \mathbb{R}$  stetig und habe endliche Grenzwerte  $\ell := \lim_{t \rightarrow -\infty} \sigma(t) \neq \lim_{t \rightarrow +\infty} \sigma(t) =: r$ . Dann ist  $\mathcal{S}^\sigma(\mathbb{R}^d, \mathbb{R}^m) \subseteq C(\mathbb{R}^d, \mathbb{R}^m)$  gleichmäßig dicht auf Kompakta. Genauer gilt das folgende für jede Norm  $\|\cdot\|_{\mathbb{R}^m}$  auf  $\mathbb{R}^m$ :*

$$\forall F \in C(\mathbb{R}^d, \mathbb{R}^m), \Omega \subseteq \mathbb{R}^d \text{ kompakt, } \varepsilon > 0 \exists N \in \mathcal{S}^\sigma(\mathbb{R}^d, \mathbb{R}^m):$$

$$\|F - N\|_{\Omega, \infty} = \sup_{x \in \Omega} \|F(x) - N(x)\|_{\mathbb{R}^m} < \varepsilon.$$

*Beweis.* ① Wir zeigen zuerst, dass wir ohne Einschränkung  $\ell = 0$  annehmen dürfen. Dazu behaupten wir, dass für eine beliebige diskriminatorische und nicht-konstante Funktion  $\sigma: \mathbb{R} \rightarrow \mathbb{R}$  stets  $\sigma + \ell$  ebenfalls diskriminatorisch ist und zwar für beliebiges  $\ell \in \mathbb{R}$ . Sei dazu  $\Omega \subseteq \mathbb{R}$  kompakt und  $\mu \in M(\Omega)$  derart dass

$$\forall w \in \mathbb{R}^n, b \in \mathbb{R}: \int_{\Omega} \sigma(\langle w, x \rangle + b) + \ell \, d\mu(x) = 0$$

gilt. Wir spezialisieren in der obigen Bedingung  $w = 0$  und nutzen aus, dass dann der Integrand konstant wird. Es folgt

$$\forall b \in \mathbb{R}: (\sigma(b) + \ell) \cdot \mu(\Omega) = 0.$$

Weil  $\sigma$  per Voraussetzung nicht konstant ist, können wir ein  $b \in \mathbb{R}$  finden, sodass  $\sigma(b) + \ell \neq 0$  gilt. Aus obigem folgt also, dass  $\mu(\Omega) = 0$  sein muss. (Beachte, dass hieraus im Allgemeinen noch nicht  $\mu = 0$  folgt, da  $\mu \in M(\Omega)$  möglicherweise nicht positiv ist!) Benutzen wir aber nochmal die allererste Bedingung, so erhalten wir

$$\begin{aligned} \forall w, b \in \mathbb{R}: 0 &= \int_{\Omega} \sigma(wx + b) + \ell \, d\mu(x) \\ &= \int_{\Omega} \sigma(wx + b) \, d\mu(x) + \ell \cdot \mu(\Omega) \\ &= \int_{\Omega} \sigma(wx + b) \, d\mu(x). \end{aligned}$$

Da  $\sigma$  diskriminatorisch ist, folgt hieraus jetzt wie gewünscht  $\mu = 0$ .

② Sei jetzt  $\sigma$  wie im Satz, wobei wir nach dem ersten Teil ohne Einschränkung voraussetzen, dass  $\ell = 0$  ist. Für  $\Omega \subseteq \mathbb{R}$  kompakt und  $\mu \in M(\Omega)$  gelte

$$\forall w \in \mathbb{R}^d, b \in \mathbb{R}: \int_{\Omega} \sigma(wx + b) \, d\mu(x) = 0.$$

Wir behaupten, dass  $\mu(\Omega \cap I_{w,b}) = 0$  ist für alle  $w, b \in \mathbb{R}$ , wobei  $I_{w,b} := \{x \in \mathbb{R} \mid wx + b > 0\}$  für  $w > 0$  gleich dem Intervall  $(-b/w, \infty)$  ist, während sich für  $w < 0$  das Intervall  $(-\infty, -b/w)$  ergibt. Für  $w = 0$  ist  $I_{w,b}$  entweder leer und es ist nichts zu zeigen, oder  $I_{w,b} = \mathbb{R}$  und  $\mu(\Omega) = 0$  folgt sofort mit Argumenten wie in ①. Für

$w, b, a \in \mathbb{R}$  berechnen wir

$$\lim_{j \rightarrow \infty} \sigma(j(wx + b) + a) = \begin{cases} r & \text{falls } wx + b > 0 \\ \sigma(a) & \text{falls } wx + b = 0 \\ 0 & \text{sonst} \end{cases} = r \cdot \mathbb{1}_{I_{w,b}} + \sigma(a) \cdot \mathbb{1}_{\{-b/w\}}$$

und sehen, dass per Voraussetzung für jedes  $j \in \mathbb{N}$

$$\int_{\Omega} \sigma(j(wx + b) + a) d\mu(x) = \int_{\Omega} \sigma(jwx + jb + a) d\mu(x) = 0$$

gilt. Nach dem Lebesgueschen Konvergenzsatz erhalten wir aus beidem zusammen, dass bei fixierten  $w, b \in \mathbb{R}$  für jedes  $a \in \mathbb{R}$

$$r\mu(I_{w,b} \cap \Omega) + \sigma(a)\mu(\{-b/w\} \cap \Omega) = 0$$

gilt. Da der erste Summand unabhängig von  $a$  ist und  $\sigma$  nicht konstant ist, muss  $\mu(P\{-b/w\} \cap \Omega) = 0$  sein und wir erhalten

$$\forall w, b \in \mathbb{R}: \mu(I_{w,b} \cap \Omega) = 0.$$

③ Jetzt behaupten wir schließlich, dass  $\mu = 0$  gilt. Für fixiertes  $w \in \mathbb{R}$  wählen wir ein kompaktes Intervall  $J \supseteq \{wx \mid x \in \Omega\}$  und betrachten

$$\varphi_w: L^\infty(J) \rightarrow \mathbb{C}, \quad \varphi_w(f) := \int_{\Omega} f(wx) d\mu(x).$$

Für reelle  $a < b$  mit  $(a, b] \subseteq \Omega$  berechnen wir

$$\begin{aligned} \varphi_w(\mathbb{1}_{(a,b]}) &= \int_{\Omega} \mathbb{1}_{(a,b]}(wx) d\mu(x) \\ &= \int_{\Omega} \mathbb{1}_{(a,\infty)}(wx) d\mu(x) - \int_{\Omega} \mathbb{1}_{(b,\infty)}(wx) d\mu(x) \\ &= \mu(\Omega \cap I_{w,-a}) - \mu(\Omega \cap I_{w,-b}) = 0 \end{aligned}$$

Da die einfachen Funktionen in  $L^\infty(J)$  dicht liegen, folgt  $\varphi_w \equiv 0$ . Wir betrachten nun einerseits  $f \in L^\infty(J)$  definiert durch  $f(t) = e^{-it}$ , und andererseits  $\mu \in M(\mathbb{R})$ , indem wir  $\mu \in M(\Omega)$  per Null fortsetzen. Jetzt berechnen wir die Fouriertransformierte  $\hat{\mu} \in C_b(\mathbb{R})$  und erhalten

$$\hat{\mu}(w) = \frac{1}{\sqrt{2\pi}} \int_{\Omega} e^{-iwx} d\mu(x) = \frac{1}{\sqrt{2\pi}} \varphi_w(f) = 0.$$

Mit Lemma 16.14 folgt  $\mu = 0$  und wir haben bewiesen, dass  $\sigma$  diskriminatorisch ist.

④ Lemma 16.17 impliziert nun, dass  $\mathcal{S}^\sigma(\mathbb{R}) \subseteq C(\mathbb{R})$  gleichmäßig dicht auf Kompakta ist und die Propositionen 16.8 und 16.12 zeigen, dass daraus die Aussage des Satzes folgt.  $\square$

Als Korollar des obigen Satzes sowie unserer Vorarbeiten erhalten wir ein analoges Approximationsresultat für die nicht-abgeschnittene ReLU-Funktion.

**Korollar 16.21.** *Der Unterraum  $\mathcal{S}^{\text{ReLU}}(\mathbb{R}^d, \mathbb{R}^m) \subseteq C(\mathbb{R}^d, \mathbb{R}^m)$  ist gleichmäßig dicht auf Kompakta.*

*Beweis.* Es genügt zu zeigen, dass ReLU diskriminatorisch ist. Sei  $\Omega \subseteq \mathbb{R}$  kompakt und  $\mu \in \mathcal{M}(\Omega)$  sei derart dass  $\int_{\Omega} \text{ReLU}(wx + b) d\mu(x) = 0$  für alle  $w, b \in \mathbb{R}$  gilt. Für beliebige  $w, b \in \mathbb{R}$  gilt per Definition

$$\overline{\text{ReLU}}(wx + b) = \begin{cases} 0, & \text{falls } wx + b < 0, \\ wx + b, & \text{falls } 0 \leq wx + b < 1, \\ 1, & \text{falls } 1 \leq wx + b, \end{cases}$$

womit man per Fallunterscheidung sieht, dass

$$\begin{aligned} \overline{\text{ReLU}}(wx + b) &= \max(0, wx + b) - \max(0, wx + b - 1) \\ &= \text{ReLU}(wx + b) - \text{ReLU}(wx + b - 1) \end{aligned}$$

für alle  $x \in \mathbb{R}$  gilt. Durch zweimalige Anwendung der Voraussetzung erhalten wir  $\int_{\Omega} \overline{\text{ReLU}}(wx + b) d\mu(x) = 0$  und dies für beliebige  $w, b \in \mathbb{R}$ . Da der Beweis von 16.20 insbesondere gezeigt hat, dass  $\overline{\text{ReLU}}$  diskriminatorisch ist, folgt hieraus  $\mu = 0$ .  $\square$

Als nächstes widmen wir uns unstetigen Aktivierungsfunktionen. Um zu adressieren, dass für solche nicht mehr  $\mathcal{S}^{\sigma}(\mathbb{R}^d, \mathbb{R}^m) \subseteq C(\mathbb{R}^d, \mathbb{R}^m)$  gilt, ersetzen wir in Definition 16.6 das Supremum über  $\Omega$  durch das essentielle Supremum, betrachten also

$$\|f\|_{\Omega, \text{ess sup}} = \text{ess sup}_{x \in \Omega} \|f(x)\|_{\mathbb{R}^m} := \inf_{\substack{N \in \Sigma \\ \lambda(N)=0}} \|f|_{\Omega \setminus N}\|_{\infty}$$

wobei  $\lambda$  das Lebesguemaß auf  $\Omega$  ist und  $\Sigma$  die  $\sigma$ -Algebra der Borelmengen von  $\Omega$ . Mit dieser Anpassung erhalten wir den folgenden Satz.

**Satz 16.22.** (von Leshno, Lin, Pinkus, Schocken) *Sei  $t_0 \in \mathbb{R}$  und  $\sigma: \mathbb{R} \rightarrow \mathbb{R}$  sei stetig auf  $[-s, t_0) \cup (t_0, s]$  für ein  $s > 0$  mit reellen Grenzwerten  $\lim_{t \nearrow t_0} \sigma(t) \neq \lim_{t \searrow t_0} \sigma(t)$ . Dann gilt:*

$$\begin{aligned} \forall F \in C(\mathbb{R}^d, \mathbb{R}^m), \Omega \subseteq \mathbb{R}^d \text{ kompakt}, \varepsilon > 0 \exists N \in \mathcal{S}^{\sigma}(\mathbb{R}^d, \mathbb{R}^m): \\ \|F - N\|_{\Omega, \text{ess sup}} < \varepsilon. \end{aligned}$$

*Beweis.* Nach den Propositionen 16.8 und 16.12 genügt es den Fall  $d = m = 1$  zu betrachten. Hier haben wir

$$\mathcal{S}^{\sigma}(\mathbb{R}) = \text{span} \{ \sigma(w \cdot + b) \mid w, b \in \mathbb{R} \}.$$

Für  $w = \pm 1$  wählen wir  $b$  derart, dass  $\sigma(b) \neq 0$  ist. Nehmen wir dann  $1/\sigma(b)$  als

Koeffizient, so sehen wir, dass  $\mathbf{1}_{\mathbb{R}} \in \mathcal{S}^\sigma(\mathbb{R})$  gilt. Für  $w = 1$  und  $a \in \mathbb{R}$  geeignet gilt

$$\ell := \lim_{t \nearrow 0} a\sigma(t - t_0) < \lim_{t \searrow 0} a\sigma(t - t_0) =: r$$

und folglich erfüllt  $\tau \in \mathcal{S}^\sigma(\mathbb{R})$ ,  $\tau(t) := \frac{1}{r-\ell}(a\sigma(t - t_0) - \ell \mathbf{1}_{\mathbb{R}})$

$$\lim_{t \nearrow 0} \tau(t) = 0 \quad \text{und} \quad \lim_{t \searrow 0} \tau(t) = 1.$$

Wir betrachten nun  $\mathcal{S}^\tau(\mathbb{R}) \subseteq \mathcal{S}^\sigma(\mathbb{R})$  und sehen, dass für  $c \in \mathbb{R}$  und  $a < b$  in  $\mathbb{R}$

$$\|\tau(w \cdot + c) - \mathbf{1}_{[c, \infty)}\|_{[a, b], \infty} = \sup_{t \in [a, b]} |\tau(wt + c) - \mathbf{1}_{[c, \infty)}| \xrightarrow{w \rightarrow 0} 0$$

gilt, also  $\lim_{w \rightarrow 0} \tau(w \cdot + c) = \mathbf{1}_{[c, \infty)}$  gleichmäßig auf kompakten Teilmengen von  $\mathbb{R}$  gilt. Damit ist

$$\begin{aligned} \mathbf{1}_{[c, \infty)} \in \overline{\mathcal{S}^\sigma(\mathbb{R})} &:= \{F: \mathbb{R} \rightarrow \mathbb{R} \mid \forall \Omega \subseteq \mathbb{R} \text{ kompakt, } \varepsilon > 0 \\ &\exists N \in \mathcal{S}^\sigma(\mathbb{R}): \|F - N\|_{\Omega, \text{ess sup}} < \varepsilon\} \end{aligned}$$

für beliebiges  $c \in \mathbb{R}$ . Da  $\overline{\mathcal{S}^\sigma(\mathbb{R})}$  ein Vektorraum ist, enthält letzterer dann alle Treppenfunktionen (mit endlich vielen Stufen!). Aus der Analysis ist wohlbekannt, dass auf einem festen kompakten Intervall jede stetige Funktion gleichmäßig durch Treppenfunktionen approximiert werden kann. Damit folgt die Aussage des Satzes.  $\square$

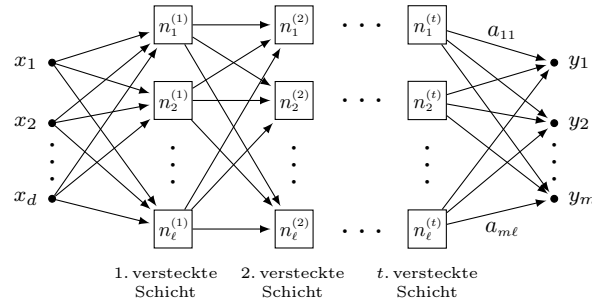
Satz 16.22 deckt insbesondere flache neuronale Netze mit Heaviside- oder Signum-Aktivierung ab. In der Tat kann man in diesen Spezialfällen sogar den ersten Teil des Beweises weglassen bzw. deutlich kürzen.

In der Praxis stellt sich heraus, dass die Approximation mit einem flachen neuronalen Netz sehr viele Neuronen erfordert, während man im Vergleich mit deutlich weniger Neuronen auskommt, wenn man letztere nur anders anordnet. Wir verweisen auf Aufgabe 16.8 für ein Beispiel, welches das obige unterlegt und fahren hier fort mit der formalen Definition dessen was sich hinter dem Begriff ‘Deep Learning’ verbirgt.

**Definition 16.23.** Sei  $\sigma: \mathbb{R} \rightarrow \mathbb{R}$ . Eine Funktion  $N: \mathbb{R}^d \rightarrow \mathbb{R}^m$  der Form

$$N = A \circ n^{(t)} \circ \dots \circ n^{(1)} \quad \text{mit} \quad n^{(k)} = \begin{bmatrix} n_1^{(k)} \\ \vdots \\ n_\ell^{(k)} \end{bmatrix}$$

mit  $\sigma$ -aktivierten Neuronen  $n_1^{(1)}, \dots, n_\ell^{(1)}: \mathbb{R}^d \rightarrow \mathbb{R}$  und  $n_1^{(j)}, \dots, n_\ell^{(j)}: \mathbb{R}^\ell \rightarrow \mathbb{R}$  für  $k = 2, \dots, t$ , und  $A \in \mathbb{R}^{\ell \times m}$  heißt *tiefes, vollständig verbundenes, vorwärtspropagierendes neuronales Netz der Breite  $\ell$  und Tiefe  $t \geq 2$  mit linearem Ausgang*. Ein solches Netz kann durch das folgende Bild beschrieben werden.



Wir bezeichnen mit  $\mathcal{D}^{\sigma, \ell, t}(\mathbb{R}^d, \mathbb{R}^m)$  die Menge der neuronalen Netze obigen Typs, welche wir kurz als *tiefe Netze* bezeichnen. Wir definieren überdies

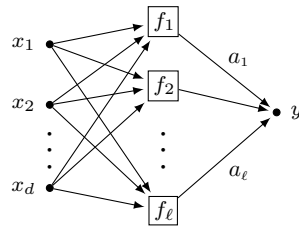
$$\mathcal{D}^{\sigma, \ell}(\mathbb{R}^d, \mathbb{R}^m) = \bigcup_{t=1}^{\infty} \mathcal{D}^{\sigma, \ell, t}(\mathbb{R}^d, \mathbb{R}^m),$$

also die Menge der beliebig tiefen neuronalen Netze mit fester Breite  $\ell$ . Wir setzen  $\mathcal{D}^{\sigma, \ell, t}(\mathbb{R}^d) := \mathcal{D}^{\sigma, \ell, t}(\mathbb{R}^d, \mathbb{R})$  und  $\mathcal{D}^{\sigma, \ell}(\mathbb{R}^d) := \mathcal{D}^{\sigma, \ell}(\mathbb{R}^d, \mathbb{R})$ .

Der nächste Satz zeigt, dass, unter bestimmten Voraussetzungen, ein Trade-off zwischen Breite und Tiefe möglich ist. Wir behandeln zunächst den skalarwertigen Fall.

**Satz 16.24.** (von Hanin) *Sei  $F \in \mathcal{S}^{\text{ReLU}}(\mathbb{R}^d)$  ein flaches Netz mit Breite  $\ell$  und  $\Omega \subseteq \mathbb{R}_{\geq 0}^d$  kompakt. Dann gibt es ein tiefes Netz  $N \in \mathcal{D}^{\text{ReLU}, d+3, \ell}(\mathbb{R}^d)$  mit  $F(x) = N(x)$  für alle  $x \in \Omega$ .*

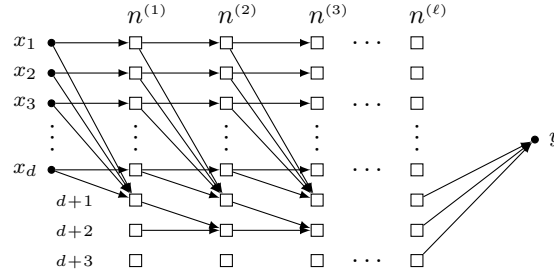
*Beweis.* Das flache Netz  $F: \mathbb{R}^d \rightarrow \mathbb{R}$  kann durch das folgende Bild illustriert werden



und wird formal durch die Vorschrift

$$F = [a_1 \ \cdots \ a_\ell] \circ \begin{bmatrix} f_1 \\ \vdots \\ f_\ell \end{bmatrix} = \sum_{i=1}^{\ell} a_i f_i$$

mit  $a_i \in \mathbb{R}$  und  $\sigma$ -aktivierten Neuronen  $f_i: \mathbb{R}^d \rightarrow \mathbb{R}$  dargestellt. Da  $\Omega$  kompakt und  $F$  stetig ist, existiert  $c_0 \geq 0$  mit  $F + c_0 \geq 0$  auf  $\Omega$ . Wir konstruieren das tiefe Netz  $N$  nun entsprechend dem folgenden Bild, bei welchem nicht eingezeichnete Pfeile einer Null in Gewichtsvektor bzw. Ausgangsmatrix entsprechen.



Die Neuronen in den Zeilen 1 bis  $d$  duplizieren die Eingänge. Genauer setzen wir  $n_j^{(k)} := \text{ReLU}(\langle e_j, \cdot \rangle) + 0$  für  $k = 1, \dots, \ell$  und  $j = 1, \dots, d$ , wobei  $e_j$  der  $j$ -te Einheitsvektor in  $\mathbb{R}^d$  bzw. in  $\mathbb{R}^{d+3}$  ist. Wir erhalten folglich  $n_j^{(k)}(x) = x_j$  für jedes  $x = (x_1, \dots, x_d) \in \Omega$ .

Die  $(d+1)$ -te Zeile soll nun die Ausgabe des gegebenen flachen Netzes  $F$  replizieren. Dazu setzen wir  $n_{d+1}^{(1)} := f_1$ . Für  $k = 2, \dots, \ell$  müssen ist  $n_{d+1}^{(k)}$  eine Funktion von  $d+3$  Variablen. Ist also  $w_k$  der Gewichtsvektor und  $b_k$  das Bias von  $f_k$ , so setzen wir

$$n_{d+1}^{(k)} := \text{ReLU}\left(\left\langle \begin{bmatrix} w_k \\ 0 \\ 0 \\ 0 \end{bmatrix}, \cdot \right\rangle + b_k\right)$$

für  $k \geq 2$  und erhalten das gewünschte. In der  $(d+2)$ -ten Zeile setzen wir

$$n_{d+2}^{(1)} = \text{ReLU}(\langle 0, \cdot \rangle + c_0) \quad \text{und} \quad n_{d+2}^{(k)} = \text{ReLU}(\langle [a_{k-1}^1], \cdot \rangle + 0) \quad \text{für } k \geq 2.$$

Sukzessives Einsetzen zeigt

$$\begin{aligned} n_{d+2}^{(1)} &= \text{ReLU}(\langle 0, \cdot \rangle + c_0) = c_0 \\ n_{d+2}^{(2)} \circ n^{(1)} &= \text{ReLU}(\langle [a_1^1], n^{(1)} \rangle) = \text{ReLU}(a_1 f_1 + c_0) = a_1 f_1 + c_0 \\ &\vdots \\ n_{d+2}^{(\ell)} \circ n^{(\ell-1)} \circ \dots \circ n^{(1)} &= \text{ReLU}(\langle [a_{\ell-1}^1], n^{(\ell-1)} \circ \dots \circ n^{(1)} \rangle) = \sum_{i=1}^{\ell-1} a_i f_i + c_0. \end{aligned}$$

auf  $\Omega$ . Es fehlt nun einerseits noch der  $\ell$ -te Summand und andererseits haben wir mit  $c_0$  einen Term zuviel. Um Letzteres zu korrigieren, definieren wir die Neuronen der  $(d+3)$ -ten Zeile per

$$n_{d+3}^{(k)} = \text{ReLU}(\langle 0, \cdot \rangle + 0) \quad \text{für } k = 1, \dots, \ell-1 \quad \text{und} \quad n_{d+3}^{(\ell)} := \text{ReLU}(\langle 0, \cdot \rangle + c_0) = c_0.$$

Für das letzte Neuron in der  $(d+1)$ -ten Zeile haben wir per vorangegangener Konstruktion gerade  $n_{d+1}^{(\ell)} \circ n^{(\ell-1)} \circ \dots \circ n^{(1)}(x) = f_\ell(x)$  für  $x \in \Omega$ . Wir wählen nun  $B = [0 \dots 0 \ a_\ell \ 1 \ -1]$  als Ausgangsmatrix und erhalten

$$N = B \circ n^{(\ell)} \circ \dots \circ n^{(1)} = a_\ell f_\ell + \sum_{i=1}^{\ell-1} a_i f_i + c_0 - c_0 = F$$

auf  $\Omega$ . □

Im obigen Beweis sind wir bei der Definition des tiefen Netzes  $N$  sehr verschwenderisch vorgegangen; in der Tat kopieren dort  $(n \cdot \ell)$ -viele Neuronen nur die Eingänge und  $(\ell - 1)$ -weitere sind gleich der Nullfunktion. Dies suggeriert einerseits, dass wir bei gegebener Breite mit weitaus geringerer Tiefe auskommen sollten. Andererseits benötigen wir für ein Dichtheitsresultat beliebig viele Neuronen, müssen also bei fester Breite dann beliebige Tiefe zulassen, siehe das folgende Korollar.

**Korollar 16.25.** *Sei  $\Omega \subseteq \mathbb{R}^n$  kompakt und  $\sigma = \cdot$ . Dann ist  $\mathcal{D}^{\sigma, n+3m}(\Omega, \mathbb{R}^m) \subseteq C(\Omega, \mathbb{R}^m)$  dicht.*

*Beweis.* Sei  $F \in C(\mathbb{R}^d, \mathbb{R}^m)$ ,  $\Omega \subseteq \mathbb{R}^d$  sei kompakt und  $\varepsilon > 0$ .

① Wir überlegen uns zuerst, dass wir ohne Einschränkung  $\Omega \subseteq \mathbb{R}_{\geq 0}^d$  annehmen dürfen. Dazu wählen wir  $x_0 \in \mathbb{R}^d$  derart, dass  $\Omega_0 := x_0 + \Omega \subseteq \mathbb{R}_{\geq 0}^d$  gilt und setzen  $F_0 := F(\cdot - x_0)$ . Sei jetzt  $N_0 \in \mathcal{D}^{\sigma, n+3m}(\mathbb{R}^d, \mathbb{R}^m)$  derart dass  $\|F_0 - N_0\|_{\Omega_0, \infty} < \varepsilon$  gilt. Mit  $N := N_0(\cdot + x_0)$  haben wir also  $\|F - N\|_{\Omega, \infty} < \varepsilon$  und wir müssen noch zeigen, dass  $N$  ein tiefes neuronales Netz ist, was wir durch Abänderung der Biase in der ersten versteckten Schicht von  $N_0$  sehen können. Für  $x \in \mathbb{R}^d$  haben wir nämlich

$$N(x) = N_0(x + x_0) = [A \circ n^{(t)} \circ \cdots \circ n^{(1)}](x - x_0) = [A \circ n^{(t)} \circ \cdots \circ \tilde{n}^{(1)}](x)$$

wobei die  $n^{(i)}$  die Schichten von  $N$  sind und  $\tilde{n}_i^{(1)}(x) = \text{ReLU}(\langle w_i^{(1)}, x \rangle + \tilde{b}_i)$  mit  $\tilde{b}_i^{(1)} = \langle w_i^{(1)}, x_0 \rangle + b_i^{(1)}$ . Ab jetzt sei also ohne Einschränkung  $\Omega \subseteq \mathbb{R}_{\geq 0}^d$ .

② Als nächstes bemerken wir, dass wir Satz 16.24  $\mathbb{R}^m$ -wertig machen können, indem wir wie im Beweis von Satz 16.24 beginnen und dann für jede weitere Koordinatenfunktion dem Netz weitere drei Zeilen hinzufügt. In diesem Sinne erhalten wir

$$\mathcal{S}^{\text{ReLU}, \ell}(\mathbb{R}^d, \mathbb{R}^m) \subseteq \mathcal{D}^{\text{ReLU}, d+3m, \ell}(\mathbb{R}^d, \mathbb{R}^m)$$

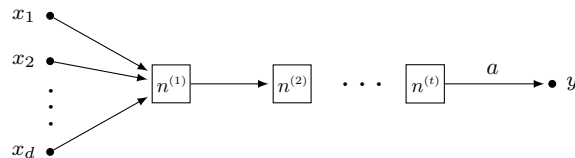
für jedes  $\ell \in \mathbb{N}$  und Vereinigen über  $\ell$  liefert

$$\mathcal{S}^{\text{ReLU}}(\mathbb{R}^d, \mathbb{R}^m) \subseteq \mathcal{D}^{\text{ReLU}, d+3m}(\mathbb{R}^d, \mathbb{R}^m) \subseteq C(\mathbb{R}^d, \mathbb{R}^m).$$

Nach Korollar 16.21 finden wir ein  $N \in \mathcal{S}^{\text{ReLU}}(\mathbb{R}^d, \mathbb{R}^m)$  mit  $\|F - N\|_{\Omega, \infty} < \varepsilon$  und sind wegen der obigen Inklusion fertig. □

Das folgende einfache Beispiel zeigt, dass wir nicht beliebig kleine Breiten durch größere Tiefe ausgleichen können.

**Beispiel 16.26.** Wir betrachten ein neuronales Netz mit Breite  $\ell = 1$ , Tiefe  $t \geq 2$  und ReLU-Aktivierung.



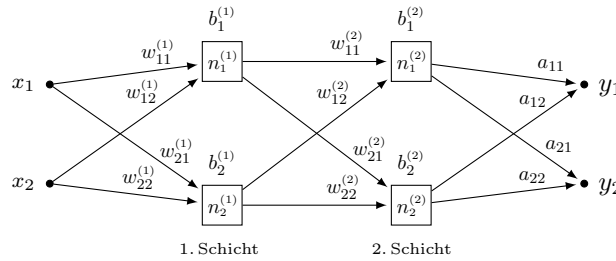


Schreibt man die entsprechende Funktion auf, so sieht man, dass diese von der Form  $N = a\text{ReLU}(\dots)$  ist und folglich je nach Vorzeichen von  $a$  entweder auf ganz  $\mathbb{R}^d$  größer gleich Null oder auf ganz  $\mathbb{R}^d$  kleiner gleich Null ist.

Es gibt weitere Resultate zur Expressivität, die, teils unter zusätzlichen Annahmen an die zu approximierenden Funktionen, asymptotische Schranken für die nötige Breite und Tiefe angeben. Wir verweisen auf das Kapitelende für Referenzen.

## 16.2 Rückwärtspropagation

Kapitel 16.1 liefert reine Existenzaussagen und verrät nichts darüber wie Gewichte und Biase zu wählen sind, wenn zu einer vorgegebene Datenmenge ein neuronales Netz als Klassifizierer oder Regressor bestimmt werden soll. Im folgenden befassen wir uns mit dieser Frage und bemerken zunächst, dass wir ähnliches bereits im allerersten Kapitel über affin-lineare Regression behandelt haben. Danach, im Kontext der logistischen Regression, haben wir sogar bereits (einzelne) Neuronen als Prediktoren/Approximanden verwendet. Die Vorgehensweise bei neuronalen Netzen ist prinzipiell die gleiche: Wir definieren eine Kostenfunktion, welche für die Featurevektoren einer gegebene Datenmenge die Abweichungen der Label von den Werten der Features unter dem Netz berechnet und minimieren dann über die Parameter des Netzes. Wir betrachten das folgende Beispiel



eines Netzes wie in Definition 16.23 mit einer differenzierbaren Aktivierungsfunktion  $\sigma: \mathbb{R} \rightarrow \mathbb{R}$ , Gewichten und Biases wie angegeben und linearem Ausgang  $A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$ . Wir führen im Folgenden die für neuronale Netze übliche Terminologie ein, weisen aber darauf hin, dass diese zum Teil schlicht aus neuen Namen für wohlbekannte mathematische Konzepte besteht.

Es sei eine Datenmenge

$$D = \{(x^{(p)}, y^{(p)}) \mid p = 1, \dots, n\} \subseteq \mathbb{R}^2 \times \mathbb{R}^2$$

mit  $x^{(p)} = (x_1^{(p)}, x_2^{(p)})$  und  $y^{(i)} = (y_1^{(p)}, y_2^{(p)})$  gegeben, deren Elemente wir im folgenden als Trainingsdaten verwenden. Das Netz im obigen Bild bezeichnen wir mit  $N \in \mathcal{D}^{\sigma, 2, 2}(\mathbb{R}^2, \mathbb{R}^2)$ . Wir betrachten die *Kostenfunktion*

$$C: \mathbb{R}^{16} \rightarrow \mathbb{R}, \quad C(w_{ij}^{(k)}, b_i^{(k)}, a_{ij}) := \frac{1}{2n} \sum_{p=1}^n \|N(x^{(p)}) - y^{(p)}\|^2$$

und betonen, dass die  $x^{(p)}$ ,  $y^{(p)}$  konstant sind, während  $C$  eine Funktion der Gewichte, Biase und Matrixeinträge ist, die das neuronale Netz  $N$  parametrisieren. Der Grund, warum wir diese rechts nicht notieren, ist offensichtlich: Bereits unser sehr kleines neuronales Netz führt auf einen länglichen Ausdruck, wenn wir dieses explizit aufschreiben wollten. In der Tat werden wir unten oft weiter abkürzen und z.B. ‘ $C(w, b, a) = \dots$ ’ schreiben oder auch die Argumente, der Übersichtlichkeit wegen, ganz weglassen, also einfach ‘ $C = \dots$ ’ schreiben, statt ‘ $C(w_{ij}^{(k)}, b_i^{(k)}, a_{ij}) = \dots$ ’.

Wir legen jetzt die Heuristik zugrunde, dass unser Netz  $N$  die Datenmenge gut approximiert, falls  $C$  klein ist, suchen also einen Minimierer

$$(w^*, b^*, a^*) \in \underset{(w, b, a) \in \mathbb{R}^{16}}{\operatorname{argmin}} C(w, b, a)$$

und verwenden dazu die Gradientenmethode, die in Kapitel 17 genauer diskutiert wird. Dabei initialisieren wir die Gewichte, Biase und Matrixeinträge zunächst beliebig und updaten diese dann, indem wir einen Schritt in Richtung  $-\nabla C$  machen.

**Algorithmus 16.27.** *Der folgende Pseudocode gibt das Gradientenverfahren für die in diesem Kapitel diskutierte Kostenfunktion wieder:*

```

1: function GRADIENTENVERFAHREN ( $D, \gamma > 0, \varepsilon > 0, T$ )
2:    $(w, b, a) \leftarrow$  beliebiger Punkt in  $\mathbb{R}^{16}$ 
3:   for  $\tau \leftarrow 1$  to  $T$  do
4:      $w \leftarrow w - \gamma \cdot \frac{\partial C}{\partial w}(w, b, a)$ 
5:      $b \leftarrow b - \gamma \cdot \frac{\partial C}{\partial b}(w, b, a)$ 
6:      $a \leftarrow a - \gamma \cdot \frac{\partial C}{\partial a}(w, b, a)$ 
7:     if  $C(w, b, a) \leq \varepsilon$  then break
8:   return  $(w, b, a)$ 

```

Den Parameter  $\gamma > 0$ , welcher in Kapitel 17 als *Schrittweite* bezeichnet wird, nennt man im Kontext neuronaler Netze oft *Lernrate*. Letztere steuert, wie weit man in die Richtung  $-\nabla C$  geht, und kann wie oben konstant sein, aber auch in jeder Iteration verändert werden, vergleiche Beispiel 17.6 und die nachfolgende Diskussion. Im Pseudocode sind  $w$ ,  $b$  und  $a$  natürlich von der Iteration  $\tau$  abhängig. Um Verwechslungen mit den oben auch nicht angegebenen Indizes  $i$ ,  $j$  und  $k$  vorzubeugen, verzichten wir darauf, die Abhängigkeit von  $\tau$  explizit anzugeben. Neben einem Abbruch beim Erreichen einer vorgegebenen Genauigkeit  $\varepsilon > 0$ , ist die maximale Anzahl durchzuführender Iterationen durch  $T \in \mathbb{N}$  beschränkt.

Im Kapitel 17 zur Gradientenmethode wenden wir diese auf konvexe Funktionen an und beweisen Sätze zur Konvergenz. Unter den Voraussetzungen dieses Kapitels ist  $C$  evtl. nicht konvex. Die Resultate aus Kapitel 17 können also nur so gelesen werden, als dass sie nahelegen, dass die beschriebene Vorgehensweise zur Bestimmung eines Netzes  $N$ , welches im durch  $C$  formalisierten Sinne gut zu den Trainingsdaten passt, erfolgsversprechend ist — einen Satz über Konvergenz werden wir im aktuellen Kapitel nicht beweisen.

Was wir im folgenden beweisen werden, sind Rekursionsformeln zur Berechnung der Ableitungen von  $C$ , die unter dem Namen *Rückwärtspropagation* bekannt sind. Die folgenden Punkte werden sich hierfür als hilfreich erweisen.

**Bemerkung 16.28.** (i) Zuerst stellen wir fest, dass  $C$  differenzierbar ist; dafür sorgen einerseits die Quadrate an den euklidischen Normen und andererseits die spezielle Form tiefer neuronaler Netze und unsere Annahme, dass die Aktivierungsfunktion  $\sigma$  differenzierbar ist. Darüber hinaus gilt

$$\frac{\partial C}{\partial w_{ij}^{(k)}} = \frac{\partial}{\partial w_{ij}^{(k)}} \left( \frac{1}{2n} \sum_{p=1}^n \|N(x^{(p)}) - y^{(p)}\|^2 \right) = \frac{1}{n} \sum_{p=1}^n \frac{\partial}{\partial w_{ij}^{(k)}} \underbrace{\left( \frac{1}{2} \|N(x^{(p)}) - y^{(p)}\|^2 \right)}_{=: C_p}$$

und entsprechende Formeln gelten für die partiellen Ableitungen nach  $b^{(k)}$  und  $a_{ij}$ . Es genügt also, wenn wir die Ableitungen aller  $C_p$  bestimmen. Um die Lesbarkeit zu erhöhen werden wir dabei den Index  $p$  weglassen und also von nun annehmen, dass die Datenmenge  $D = \{(x, y) = ((x_1, x_2), (y_1, y_2))\} \in \mathbb{R}^2 \times \mathbb{R}^2$  nur einen Punkt enthält und die Kostenfunktion sich zu

$$C = \frac{1}{2} \|N(x) - y\|^2$$

vereinfacht.

(ii) Wenn Gewichte und Biase gegeben sind, dann kann  $N(x)$  berechnet werden indem man Schicht-für-Schicht von links nach rechts geht. Dies nennt man *Vorwärtspropagation*. Wir werden zeigen, dass man

$$\frac{\partial C}{\partial w_{ij}^{(k)}}, \frac{\partial C}{\partial b_i^{(k)}} \text{ und } \frac{\partial C}{\partial a_{ij}}$$

rekursiv und zwar Schicht-für-Schicht von rechts nach links berechnen kann. Dies nennt man *Rückwärtspropagation*.

(iii) Um letzteres effizient aufzuschreiben, definieren wir

$$z^{(k)} := \begin{bmatrix} z_1^{(k)} \\ z_2^{(k)} \end{bmatrix} \text{ mit } z_i^{(k)} = \left\langle \begin{bmatrix} w_{i1}^{(k)} \\ w_{i2}^{(k)} \end{bmatrix}, \begin{bmatrix} n_1^{(k-1)} \\ n_2^{(k-1)} \end{bmatrix} \right\rangle + b_i^{(k)}$$

für  $k = 1, 2, i = 1, 2$  wobei  $n_i^{(0)} = x_i$ . Die  $z_i^{(k)}$  sind also jeweils das, was wir ‘innerhalb’ eines Neurons ausrechnen *bevor* wir die Aktivierungsfunktion anwenden.

Wir behandeln zunächst die Ausgangsmatrix und die 2. Schicht.

**Proposition 16.29.** Sei  $\sigma: \mathbb{R} \rightarrow \mathbb{R}$  differenzierbar,  $N \in \mathcal{D}^{\sigma,2,2}(\mathbb{R}^2, \mathbb{R}^2)$  ein neuronales Netz,  $D = \{(x, y)\} \subseteq \mathbb{R}^2 \times \mathbb{R}^2$  eine einpunktige Datenmenge und  $C: \mathbb{R}^{16} \rightarrow \mathbb{R}$ ,  $C(w, b, a) = \|N(x) - y\|^2$ . Dann gelten

$$\frac{\partial C}{\partial a_{ij}} = (N_i - y_i) \cdot n_j^{(2)}, \text{ sowie } \frac{\partial C}{\partial w_{ij}^{(2)}} = \delta_i^{(2)} \cdot n_j^{(1)} \text{ und } \frac{\partial C}{\partial b_i^{(2)}} = \delta_i^{(2)}$$

mit  $\delta_i^{(2)} = \sum_{\mu=1}^2 a_{\mu i} (N_\mu - y_\mu) \cdot \sigma'(z_i^{(2)})$  und  $N = \begin{bmatrix} N_1 \\ N_2 \end{bmatrix}$ .

*Beweis.* Es gilt

$$N = \begin{bmatrix} N_1 \\ N_2 \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} n_1^{(2)} \\ n_2^{(2)} \end{bmatrix} = \begin{bmatrix} a_{11}n_1^{(2)} + a_{12}n_2^{(2)} \\ a_{21}n_1^{(2)} + a_{22}n_2^{(2)} \end{bmatrix}$$

und mit der Kettenregel

$$\begin{aligned} \frac{\partial C}{\partial a_{ij}} &= \frac{\partial}{\partial a_{ij}} \left( \frac{1}{2} \|N(x) - y\|^2 \right) = \frac{1}{2} \frac{\partial}{\partial a_{ij}} \left\| \begin{bmatrix} a_{11}n_1^{(2)} + a_{12}n_2^{(2)} - y_1 \\ a_{21}n_1^{(2)} + a_{22}n_2^{(2)} - y_2 \end{bmatrix} \right\|^2 \\ &= \frac{1}{2} \frac{\partial}{\partial a_{ij}} \left[ (a_{11}n_1^{(2)} + a_{12}n_2^{(2)} - y_1)^2 + (a_{21}n_1^{(2)} + a_{22}n_2^{(2)} - y_2)^2 \right] \\ &= (a_{i1}n_1^{(2)} + a_{i2}n_2^{(2)} - y_i) \cdot n_j^{(2)} = (N_i - y_i) \cdot n_j^{(2)} \end{aligned}$$

folgt zunächst die letzte Ableitung. Um die ersten zwei Ableitungen einzusehen, benutzen wir die Hilfsfunktionen  $z_i^{(2)}$  aus Bemerkung 16.28(iii), mit denen wir die Neuronen der 2. Schicht als

$$\begin{bmatrix} n_1^{(2)} \\ n_2^{(2)} \end{bmatrix} = \begin{bmatrix} \sigma \left( \left\langle \begin{bmatrix} w_{11}^{(2)} \\ w_{12}^{(2)} \end{bmatrix}, \begin{bmatrix} n_1^{(1)} \\ n_2^{(1)} \end{bmatrix} \right\rangle + b_1^{(2)} \right) \\ \sigma \left( \left\langle \begin{bmatrix} w_{21}^{(2)} \\ w_{22}^{(2)} \end{bmatrix}, \begin{bmatrix} n_1^{(1)} \\ n_2^{(1)} \end{bmatrix} \right\rangle + b_2^{(2)} \right) \end{bmatrix} = \begin{bmatrix} \sigma(z_1^{(2)}) \\ \sigma(z_2^{(2)}) \end{bmatrix}$$

schreiben können. Dann beginnen wir wie oben, müssen jetzt aber berücksichtigen, dass beide Summanden in der eckigen Klammer von  $z_i^{(2)}$  abhängen. Wir erhalten daher

$$\begin{aligned} \frac{\partial C}{\partial z_i^{(2)}} &= \frac{1}{2} \frac{\partial}{\partial z_i^{(2)}} \left[ (a_{11}n_1^{(2)} + a_{12}n_2^{(2)} - y_1)^2 + (a_{21}n_1^{(2)} + a_{22}n_2^{(2)} - y_2)^2 \right] \\ &= \sum_{\mu=1}^2 (a_{\mu 1}n_1^{(2)} + a_{\mu 2}n_2^{(2)} - y_\mu) \cdot \frac{\partial}{\partial z_i^{(2)}} (a_{\mu 1}n_1^{(2)} + a_{\mu 2}n_2^{(2)} - y_\mu) \\ &= \sum_{\mu=1}^2 (N_\mu - y_\mu) \cdot a_{\mu i} \sigma'(z_i^{(2)}) = \sum_{\mu=1}^2 a_{\mu i} (N_\mu - y_\mu) \cdot \sigma'(z_i^{(2)}) = \delta_i^{(2)} \end{aligned}$$

aus. Per Kettenregel folgt dann

$$\frac{\partial C}{\partial w_{ij}^{(2)}} = \frac{\partial C}{\partial z_i^{(2)}} \cdot \frac{\partial z_i^{(2)}}{\partial w_{ij}^{(2)}} \underset{\text{s.o.}}{\uparrow} \delta_i^{(2)} \cdot n_j^{(1)} \quad \text{und} \quad \frac{\partial C}{\partial b_i^{(2)}} = \frac{\partial C}{\partial z_i^{(2)}} \cdot \frac{\partial z_i^{(2)}}{\partial b_i^{(2)}} \underset{\text{s.o.}}{\uparrow} \delta_i^{(2)} \cdot 1$$

was den Beweis beendet.  $\square$

Führen wir eine Iteration der Gradientenmethode durch, so sind uns zunächst die Werte der Parameter  $w_{ij}^{(k)}$ ,  $b_i^{(k)}$  und  $a_{ij}$  aus der vorangegangenen Iteration bekannt. Durch Vorwärtspropagation erhalten wir die zugehörigen Werte der  $z_j^{(k)}$  und der  $n_i^{(k)}$  und der  $N_i$ . Daraus können wir  $\delta_i^{(2)}$  berechnen und mit Proposition 16.29 dann alle

drei Ableitungen

$$\frac{\partial C}{\partial w_{ij}^{(2)}}(w_{ij}^{(k)}, b_i^{(k)}), \quad \frac{\partial C}{\partial b_i^{(2)}}(w_{ij}^{(k)}, b_i^{(k)}) \quad \text{und} \quad \frac{\partial C}{\partial a_{ij}} = (N_i - y_i) \cdot n_j^{(2)}$$

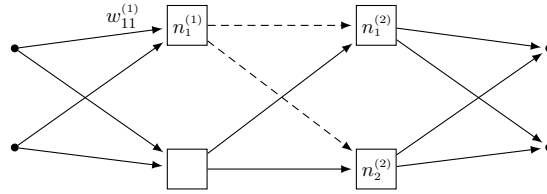
bestimmen. Damit haben wir die 2. Schicht abgearbeitet und kommen zur 1. Schicht.

**Proposition 16.30.** *Sei  $\sigma: \mathbb{R} \rightarrow \mathbb{R}$  differenzierbar,  $N \in \mathcal{D}^{\sigma,2,2}(\mathbb{R}^2, \mathbb{R}^2)$  ein neuronales Netz,  $D = \{(x, y)\} \subseteq \mathbb{R}^2 \times \mathbb{R}^2$  eine einpunktige Datenmenge und  $C: \mathbb{R}^{16} \rightarrow \mathbb{R}$ ,  $C(w, b, a) = \|N(x) - y\|^2$ . Dann gelten*

$$\frac{\partial C}{\partial w_{ij}^{(1)}} = \delta_i^{(1)} \cdot x_j \quad \text{und} \quad \frac{\partial C}{\partial b_i^{(1)}} = \delta_i^{(1)}$$

mit  $\delta_i^{(1)} := (\delta_1^{(2)} w_{1i}^{(2)} + \delta_2^{(2)} w_{2i}^{(2)}) \cdot \sigma'(z_i^{(1)})$ .

*Beweis.* Wir beginnen wie im Beweis von Proposition 16.29, jedoch werden die Abhängigkeiten nochmal etwas komplizierter. Beispielsweise geht das Gewicht  $w_{11}^{(1)}$  in beide Neuronen der 2. Schicht ein



und damit in  $z_1^{(2)}$  und in  $z_2^{(2)}$ . Lesen wir  $C = C(w_{11}^{(1)})$  als Verkettung  $\mathbb{R} \rightarrow \mathbb{R}^2 \rightarrow \mathbb{R}$ ,  $C = C(z_1^{(2)}(w_{11}^{(1)}), z_2^{(2)}(w_{11}^{(1)}))$ , so folgt mit der entsprechenden mehrdimensionalen Kettenregel und den Abkürzungen und Rechnungen aus Proposition 16.29

$$\begin{aligned} \frac{\partial C}{\partial w_{1j}^{(1)}} &= \frac{\partial C}{\partial z_1^{(2)}} \frac{\partial z_1^{(2)}}{\partial w_{1j}^{(1)}} + \frac{\partial C}{\partial z_2^{(2)}} \frac{\partial z_2^{(2)}}{\partial w_{1j}^{(1)}} \\ &= \delta_1^{(2)} \frac{\partial}{\partial w_{1j}^{(1)}} \left( \left\langle \begin{bmatrix} w_{11}^{(2)} \\ w_{12}^{(2)} \end{bmatrix}, \begin{bmatrix} n_1^{(1)} \\ n_2^{(1)} \end{bmatrix} \right\rangle + b_1^{(2)} \right) + \delta_2^{(2)} \frac{\partial}{\partial w_{1j}^{(1)}} \left( \left\langle \begin{bmatrix} w_{21}^{(2)} \\ w_{22}^{(2)} \end{bmatrix}, \begin{bmatrix} n_1^{(1)} \\ n_2^{(1)} \end{bmatrix} \right\rangle + b_2^{(2)} \right) \\ &= \delta_1^{(2)} w_{11}^{(2)} \frac{\partial n_1^{(1)}}{\partial w_{1j}^{(1)}} + \delta_2^{(2)} w_{21}^{(2)} \frac{\partial n_1^{(1)}}{\partial w_{1j}^{(1)}} \\ &= (\delta_1^{(2)} w_{11}^{(2)} + \delta_2^{(2)} w_{21}^{(2)}) \cdot \frac{\partial}{\partial w_{1j}^{(1)}} \underbrace{\left( \left\langle \begin{bmatrix} w_{11}^{(1)} \\ w_{12}^{(1)} \end{bmatrix}, \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \right\rangle + b_1^{(1)} \right)}_{=z_1^{(1)}} \\ &= (\delta_1^{(2)} w_{11}^{(2)} + \delta_2^{(2)} w_{21}^{(2)}) \cdot \sigma'(z_1^{(1)}) \cdot x_j. \end{aligned}$$

Die Gleichung für die Ableitung nach  $w_{2j}^{(1)}$  sieht man völlig analog ein. Für die Ableitung nach  $b_i^{(1)}$  geht man wiederum analog zu oben vor, sieht dann aber in der vorletzten Zeile, dass sich

$$\frac{\partial C}{\partial b_i^{(1)}} = (\delta_1^{(2)} w_{1i}^{(2)} + \delta_2^{(2)} w_{2i}^{(2)}) \cdot \sigma'(z_1^{(1)}) \cdot \frac{\partial z_1^{(1)}}{\partial b_i^{(1)}} = (\delta_1^{(2)} w_{1i}^{(2)} + \delta_2^{(2)} w_{2i}^{(2)}) \cdot \sigma'(z_1^{(1)}) \cdot 1$$

ergibt. □

Da wir die  $\delta_i^{(2)}$  bereits bei der Bearbeitung der 2. Schicht berechnet haben, bekommen wir mit Proposition 16.30 die Ableitungen

$$\frac{\partial C}{\partial w_{ij}^{(1)}}(w_{ij}^{(k)}, b_i^{(k)}) \quad \text{und} \quad \frac{\partial C}{\partial b_i^{(1)}}(w_{ij}^{(k)}, b_i^{(k)})$$

wobei die  $w_{ij}^{(k)}$ ,  $b_i^{(k)}$  und  $a_{ij}$  immer noch die Werte der Parameter der vorangegangenen Iteration sind. Da wir nun alle 16 Ableitungen an der entsprechenden Stelle  $(w_{ij}^{(k)}, b_i^{(k)}, a_{ij})$  kennen, können wir das Update entsprechend des Gradientenverfahrens durchführen.

In Satz 16.31 verallgemeinern wir das obige für Netze mit mehr als zwei Schichten und mehr als zwei Neuronen pro Schicht. Um die entsprechende Rekursionsformel für die  $\delta_i^{(2)}$  elegant formulieren zu können, brauchen wir die folgende Notation.

**Definition 16.31.** Für  $x, y \in \mathbb{R}^d$  heißt

$$x \odot y = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix} \odot \begin{bmatrix} y_1 \\ \vdots \\ y_d \end{bmatrix} := \begin{bmatrix} x_1 y_1 \\ \vdots \\ x_d y_d \end{bmatrix}$$

das *Hadamard-Produkt* von  $x$  und  $y$ .

Wir betrachten im Folgenden ein tiefes vollständig verbundenes vorwärtspropagierendes neuronales Netz  $N \in \mathcal{D}^{\sigma, \ell, t}(\mathbb{R}^d, \mathbb{R}^m)$  mit linearem Ausgang wie durch das Bild in Definition 16.23 dargestellt. Wir fassen die Gewichte jeder Schicht zu einer Matrix zusammen und die Biase jeder Schicht zu einem Vektor. Da im Allgemeinen  $d \neq \ell$  ist, hat die Gewichtsmatrix der ersten Schicht dabei eventuell ein anderes Format als die restlichen. Das Netz ist also durch die Aktivierungsfunktion und die Matrizen

$$W^{(1)} = \begin{bmatrix} w_{11}^{(1)} & \cdots & w_{1d}^{(1)} \\ \vdots & & \vdots \\ w_{\ell 1}^{(1)} & \cdots & w_{\ell d}^{(1)} \end{bmatrix}, \quad W^{(k)} = \begin{bmatrix} w_{11}^{(k)} & \cdots & w_{1\ell}^{(k)} \\ \vdots & & \vdots \\ w_{\ell 1}^{(k)} & \cdots & w_{\ell \ell}^{(k)} \end{bmatrix}, \quad b^{(k)} = \begin{bmatrix} b_1^{(k)} \\ \vdots \\ b_\ell^{(k)} \end{bmatrix}, \quad A = \begin{bmatrix} a_{11} & \cdots & a_{1\ell} \\ \vdots & & \vdots \\ a_{m1} & \cdots & a_{m\ell} \end{bmatrix}$$

für  $k = 2, \dots, t$  bzw.  $k = 1, \dots, t$  gegeben. Insgesamt enthalten die obigen Matrizen  $\nu := (\ell \cdot d + (t-1) \cdot \ell^2 + m \cdot \ell)$ -viele skalare Parameter. Die für die Gradientenmethode nötigen Ableitungen können per Rückwärtspropagation mithilfe des folgenden Satzes bestimmt werden, in welchem wir die  $\delta_i^{(k)}$  zu Vektoren  $\delta^{(k)} = (\delta_1^{(k)}, \dots, \delta_\ell^{(k)})$  zusammenfassen,  $\sigma'(\cdot)$  koordinatenweise lesen und die Hilfsfunktionen aus Bemerkung 16.28(iii)

$$z^{(k)} := \begin{bmatrix} z_1^{(k)} \\ \vdots \\ z_\ell^{(k)} \end{bmatrix} \quad \text{mit} \quad z_i^{(k)} = \left\langle \begin{bmatrix} w_{i1}^{(k)} \\ \vdots \\ w_{i\ell}^{(k)} \end{bmatrix}, \begin{bmatrix} n_1^{(k-1)} \\ \vdots \\ n_\ell^{(k-1)} \end{bmatrix} \right\rangle + b_i^{(k)}$$

entsprechend in höherer Dimension verwenden.

**Satz 16.32.** Sei  $\sigma: \mathbb{R} \rightarrow \mathbb{R}$  differenzierbar,  $N \in \mathcal{D}^{\sigma, \ell, t}(\mathbb{R}^d, \mathbb{R}^m)$  ein neuronales Netz,  $D = \{(x, y)\} \subseteq \mathbb{R}^d \times \mathbb{R}^m$  eine einpunktige Datenmenge und  $C: \mathbb{R}^\nu \rightarrow \mathbb{R}$ ,  $C(w, b, a) =$

$\|N(x) - y\|^2$ . Mit  $n_j^{(0)} := x_j$  gilt

$$\frac{\partial C}{\partial a_{ij}} = (N_i - y_i) \cdot n_j^{(t)}, \text{ sowie } \frac{\partial C}{\partial w_{ij}^{(k)}} = \delta_i^{(k)} n_j^{(k-1)} \text{ und } \frac{\partial C}{\partial b_i^{(k)}} = \delta_i^{(k)}$$

für alle Schichten  $k = 1, \dots, t$ , und die  $\delta_i^{(k)}$  können rekursiv per

$$\delta^{(t)} = A^\top (N - y) \odot \sigma'(z^{(t)}) \text{ und } \delta^{(k)} = (W^{(k+1)})^\top \delta^{(k+1)} \odot \sigma'(z^{(k)}),$$

berechnet werden, wobei  $k = t, \dots, 1$ , d.h. wir beginnen mit der letzten Schicht und arbeiten uns rückwärts vor bis zur ersten Schicht. Die Auswertungen  $z^{(k)} = W^{(k)} n^{(k-1)} + b^{(k)}$  sind hierbei jeweils aus der Vorwärtspropagation bereits bekannt.

*Beweis.* Anhand der Propositionen 16.29 und 16.30 und an deren Beweisen liest man ab, dass sich in der allgemeinen Situation die oben angegebenen drei Formeln für die Ableitungen ergeben, wenn man

$$\delta_i^{(t)} := \sum_{\mu=1}^m a_{\mu i} (N_\mu - y_\mu) \cdot \sigma'(z_i^{(t)})$$

und für  $k = 2, \dots, t$

$$\delta_i^{(k)} := (\delta_1^{(k)} w_{1i}^{(k)} + \dots + \delta_\ell^{(k)} w_{\ell i}^{(k)}) \cdot \sigma'(z_i^{(k-1)})$$

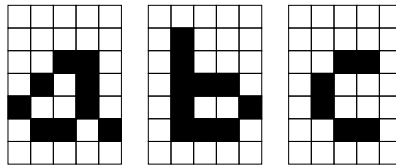
definiert. Obiges sind aber genau die jeweils  $i$ -ten Einträge der im Satz angegebenen Hadamard-Produkte.  $\square$

Wir erinnern an dieser Stelle nochmal daran, dass bei einer Datenmenge  $D = \{(x^{(p)}, y^{(p)}) \mid p = 1, \dots, n\}$  mit  $p \geq 2$  die Kostenfunktion

$$C: \mathbb{R}^\nu \rightarrow \mathbb{R}, \quad C(w_{ij}^{(k)}, b_i^{(k)}, a_{ij}) := \frac{1}{2n} \sum_{p=1}^n \|N(x^{(p)}) - y^{(p)}\|^2$$

verwendet werden muss, dass sich aber die Ableitungen per Linearität sofort aus Satz 16.32 ergeben, vergleiche auch Bemerkung 16.28(ii).

Wie am Anfang dieses Kapitels angekündigt, kommen wir zum Schluss nochmal auf Klassifizierer zurück und zwar insbesondere auf solche mit möglicherweise mehr als zwei Klassen. Ein anschauliches Beispiel hierfür ist die Handschrifterkennung. Wie im folgenden Bild dargestellt



arbeiten wir hierbei mit Daten, bei denen der Featurevektor ein handschriftlich geschriebener Buchstabe ist, oben z.B. gegeben als Element von  $\mathbb{R}^{7 \times 5} \cong \mathbb{R}^{35}$ . Als Label

verwenden wir Einheitsvektoren, d.h. wir weisen jedem Buchstaben einen Einheitsvektor zu, z.B. dem Buchstaben ‘a’ den Vektor  $e_1$ , ‘b’ den Vektor  $e_2$ , ‘c’ den Vektor  $e_3$  usw. Letzteres wird auch *One-hot-Kodierung* genannt. Bei den Trainingsdaten ist dann bekannt, welcher Buchstabe handschriftlich notiert wurde, was zu einer Datenmenge

$$D = \{(x^{(p)}, y^{(p)}) \mid p = 1, \dots, n\} \subseteq \mathbb{R}^d \times \{e_1, \dots, e_m\} \quad (16.1)$$

führt. Das naheliegende Ziel ist die Bestimmung eines Klassifizierers  $K: \mathbb{R}^d \rightarrow \{0, 1\}^m$ . Wie bereits im Kapitel 2 über die logistische Regression diskutiert, ist es allerdings sinnvoller einen Regressor  $R: \mathbb{R}^d \rightarrow (0, 1)^m$  zu bestimmen: Erstens kann man sich  $K$  jederzeit durch Rundung verschaffen, wenn man  $R$  einmal bestimmt hat. Zweitens kann man die Werte von  $R$  als Wahrscheinlichkeiten interpretieren; hat man im obigen Beispiel etwa  $R(x) = (0.88, 0.01, 0.10, \dots)$  für  $x \in \mathbb{R}^{35} \setminus D$ , so kann man dies so lesen, als dass es sich bei  $x$  mit Wahrscheinlichkeit 0.88 um ein ‘a’ handelt, mit Wahrscheinlichkeit 0.01 um ein ‘b’, mit Wahrscheinlichkeit 0.10 um ein ‘c’ usw. Drittens wird sich zeigen, dass der Ansatz mit  $R$  auf ein Optimierungsproblem über eine Klasse  $\mathcal{R}$  von differenzierbaren Funktionen führt und daher eine Anwendung der Rückwärtspropagation wie im ersten Teil dieses Unterkapitels möglich ist. Um letztere Klasse formal zu definieren, benötigen wir die sogenannte *Softmaxfunktion*

$$\text{softmax}: \mathbb{R}^m \rightarrow \mathbb{R}^m, \text{softmax}(z_1, \dots, z_m) = \frac{(e^{z_1}, \dots, e^{z_m})}{e^{z_1} + \dots + e^{z_m}}$$

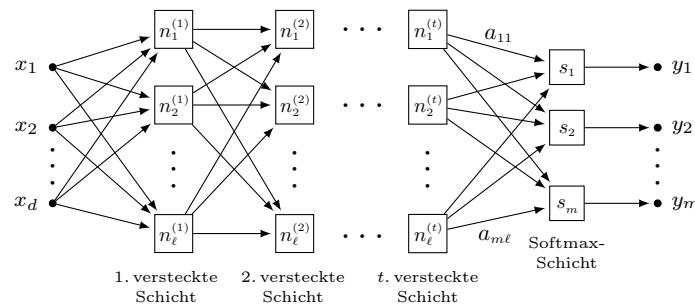
und definieren die Elemente von  $\mathcal{R}$  durch Nachschaltung derselben hinter ein tiefes neuronales Netz.

**Definition 16.33.** Sei  $\sigma: \mathbb{R} \rightarrow \mathbb{R}$  eine Aktivierungsfunktion. Wir setzen

$$\mathcal{R}^{\sigma, \ell, t}(\mathbb{R}^d, \mathbb{R}^m) := \{\text{softmax} \circ N \mid N \in \mathcal{D}^{\sigma, \ell, t}(\mathbb{R}^d, \mathbb{R}^m)\}.$$

und bezeichnen die Elemente von  $\mathcal{R}^{\sigma, \ell, t}(\mathbb{R}^d, \mathbb{R}^m)$  als *tiefe neuronale Netze mit Softmaxausgang*.

Bezeichnen wir die Koordinatenfunktionen der Softmaxfunktion mit  $s_i$ , so kann man  $R \in \mathcal{R}^{\sigma, \ell, t}(\mathbb{R}^d, \mathbb{R}^m)$  durch das folgende Bild veranschaulichen. Wir weisen allerdings darauf hin, dass es sich wegen  $s_i: \mathbb{R}^m \rightarrow \mathbb{R}$  bei den Boxen in der Softmax-Schicht nicht um Neuronen im Sinne der Definition 16.1 handelt.



Bevor wir zur Rückwärtspropagation kommen, wollen wir noch folgendes bemer-



ken: Ist eine Datenmenge mit kategoriellen Labels wie in (16.1) gegeben, so könnte man die Approximation auch mit einem Netz  $N \in \mathcal{D}^{\sigma,\ell,t}(\mathbb{R}^d, \mathbb{R}^m)$  versuchen. Man muss dann darauf bauen, dass der Optimierungsprozess eine Funktion  $N$  liefert, die ‘sinnvolle’ Featurevektoren, im Beispiel oben etwa solche  $x \in \mathbb{R}^{35}$ , die tatsächlich einem handgeschriebenen Buchstaben entsprechen, auf Label schickt die ‘in der Nähe’ eines Einheitsvektors liegen. Obwohl letzteres natürlich auch für  $R \in \mathcal{R}^{\sigma,\ell,t}(\mathbb{R}^d, \mathbb{R}^m)$  fehlschlagen kann, erzwingt doch der oben vorgestellte Ansatz immerhin einige wünschenswerte Eigenschaften. Bezeichnen wir nämlich mit  $R_i$  die  $i$ -te Koordinatenfunktion von  $R$ , so erhalten wir per Definition, dass stets

$$\forall x \in \mathbb{R}^d: R_i \in (0, 1) \text{ und } \sum_{i=1}^m R_i(x) = 1$$

gilt. Insbesondere können wir für jedes  $R \in \mathcal{R}^{\sigma,\ell,t}(\mathbb{R}^d, \mathbb{R}^m)$  und jedes  $x \in \mathbb{R}^d$  die  $R_i(x)$  als Wahrscheinlichkeiten auffassen und analog zur logistischen Regression die Maximum Likelihood Methode anwenden, um eine zum Problem passende Kostenfunktion zu definieren, nämlich die *Likelihood-Funktion*

$$L: \mathcal{R}^{\sigma,\ell,t}(\mathbb{R}^d, \mathbb{R}^m) \rightarrow \mathbb{R}, \quad L(R) := \prod_{p=1}^n \prod_{i=1}^m R_i(x^{(p)})^{y_i^{(p)}}.$$

Da jedes  $y^{(p)}$  gleich einem Einheitsvektor ist, ist für fixiertes  $p$  im inneren Produkt genau derjenige Faktor mit  $y^{(p)} = e_i$  gleich  $R_i(x^{(p)})$ . Alle anderen Faktoren sind gleich Eins. Es folgt, dass  $L$  groß ist, wenn  $R_i(x^{(p)}) \approx 1$  für genau diejenigen Paare  $(i, p)$  mit  $y^{(p)} = e_i$  gilt — und dann zwangsläufig  $R_i(x^{(p)}) \approx 0$  für Paare  $(i, p)$  mit  $y^{(p)} \neq e_i$ . Zusammengefasst ist  $L$  also groß, wenn  $R(x^{(p)}) \approx y^{(p)}$  komponentenweise gilt.

Um  $L$  zu maximieren, parametrisieren wir wie gehabt  $R \in \mathcal{R}^{\sigma,\ell,t}(\mathbb{R}^d, \mathbb{R}^m)$  durch  $\nu = (\ell \cdot d + (t-1) \cdot \ell^2 + m \cdot \ell)$ -viele reelle Parameter und gehen zur Minimierung der *negativen Log-Likelihood-Funktion* über, die wir hier wieder mit  $C$  bezeichnen

$$C: \mathbb{R}^\nu \rightarrow \mathbb{R}, \quad \ell(w_{ij}^{(k)}, b_i^{(k)}, a_{ij}) := -\log \prod_{p=1}^n \prod_{i=1}^m R_i(x^{(p)})^{y_i^{(p)}}.$$

Drch Ausrechnen des Logarithmus’ erhalten wir

$$C(w_{ij}^{(k)}, b_i^{(k)}, a_{ij}) = -\sum_{p=1}^n \sum_{i=1}^m y_i^{(p)} \log R_i(x^{(p)}) \quad (16.2)$$

und wollen diese Formel verwenden, um per Gradientenverfahren, siehe Algorithmus 16.27, einen Minimierer zu bestimmen. Im nächsten Satz geben wir, analog zu Satz 16.32, an wie die hierfür nötigen Ableitungen per Rückwärtspropagation berechnet werden können. Entsprechend Bemerkung 16.28(i) formulieren wir den Satz nur für den Fall einer einpunktigen Datenmenge.

**Satz 16.34.** Sei  $\sigma: \mathbb{R} \rightarrow \mathbb{R}$  differenzierbar,  $R \in \mathcal{R}^{\sigma,\ell,t}(\mathbb{R}^d, \mathbb{R}^m)$  ein neuronales Netz mit Softmaxausgang,  $D = \{(x, y)\} \subseteq \mathbb{R}^d \times \{0, 1\}^m$  eine one-hot-kodierte einpunktige

Datenmenge und  $C: \mathbb{R}^p \rightarrow \mathbb{R}$ ,  $C(w, b, a) = -(y_1 \log R_1(x) + \dots + y_m \log R_m(x))$ . Mit  $n_j^{(0)} := x_j$  gilt

$$\frac{\partial C}{\partial a_{ij}} = (R_i - y_i) \cdot n_j^{(t)}, \text{ sowie } \frac{\partial C}{\partial w_{ij}^{(k)}} = \delta_i^{(k)} n_j^{(k-1)} \text{ und } \frac{\partial C}{\partial b_i^{(k)}} = \delta_i^{(k)}$$

für alle Schichten  $k = 1, \dots, t$ , und die  $\delta_i^{(k)}$  können rekursiv per

$$\delta^{(t)} = A^\top(R - y) \odot \sigma'(z^{(t)}) \text{ und } \delta^{(k)} = (W^{(k+1)})^\top \delta^{(k+1)} \odot \sigma'(z^{(k)}),$$

berechnet werden.

*Beweis.* Wir berechnen zunächst die partiellen Ableitungen der Softmaxfunktion, welche wir in diesem Beweis zur Abkürzung mit  $s$ , bzw.  $s_i$  für die  $i$ -te Koordinatenfunktion, abkürzen. Da  $\text{ran } s_i \subseteq (0, 1)$  gilt, können wir

$$\log s_i = \log \frac{e^{z_i}}{e^{z_1} + \dots + e^{z_m}} = z_i - \log(e^{z_1} + \dots + e^{z_m})$$

betrachten. Differenzieren liefert einerseits

$$\frac{\partial}{\partial z_j} \log s_i = \frac{1}{s_i} \cdot \frac{\partial s_i}{\partial z_j}$$

und, mit Kroneckersymbol  $\delta_{ij}$ , andererseits

$$\begin{aligned} \frac{\partial}{\partial z_j} \log s_i &= \delta_{ij} - \frac{1}{e^{z_1} + \dots + e^{z_m}} \cdot \frac{\partial}{\partial z_j} (e^{z_1} + \dots + e^{z_m}) \\ &= \delta_{ij} - \frac{e^{z_j}}{e^{z_1} + \dots + e^{z_m}} = \delta_{ij} - s_j. \end{aligned}$$

Durch Umstellen und Einsetzen ergibt sich

$$\frac{\partial s_i}{\partial z_j} = s_i \cdot \frac{\partial}{\partial z_j} \log s_i = s_i (\delta_{ij} - s_j).$$

Per Definition gilt  $R = s \circ N$  mit  $N \in \mathcal{D}^{\sigma, \ell, t}(\mathbb{R}^m, \mathbb{R}^n)$ . Wir lesen  $C = C(N_1, \dots, N_m)$  und berechnen die entsprechenden Ableitungen

$$\begin{aligned} \frac{\partial C}{\partial N_j} &= \frac{\partial}{\partial N_j} \left( - \sum_{i=1}^m y_i \log s_i \right) = - \sum_{i=1}^m y_i \frac{\partial}{\partial N_j} \log s_i - \sum_{i=1}^m \frac{y_i}{s_i} \cdot \frac{\partial s_i}{\partial N_j} \\ &= - \sum_{i=1}^m \frac{y_i}{s_i} \cdot s_i (\delta_{ij} - s_j) = - \sum_{i=1}^m y_i \delta_{ij} - y_i s_j s_j \left( \sum_{i=1}^m y_i \right) - y_j \\ &\stackrel{\substack{\uparrow \\ y \in \{0,1\}^m}}{=} s_j - y_j \end{aligned}$$

und erhalten für deren Auswertungen

$$\frac{\partial C}{\partial N_j}(N_1, \dots, N_m) = s_j(N_1, \dots, N_m) - y_j = R_j - y_j.$$

Weiter gilt

$$\frac{\partial N_\mu}{\partial a_{ij}} = \frac{\partial}{\partial a_{ij}} [a_{\mu 1} n_1^{(t)} + \cdots + a_{\mu \ell} n_\ell^{(t)}] = \begin{cases} n_j^{(t)} & \text{falls } \mu = i, \\ 0 & \text{sonst,} \end{cases}$$

und damit per mehrdimensionaler Kettenregel

$$\frac{\partial C}{\partial a_{ij}} = \sum_{\mu=1}^m \frac{\partial C}{\partial N_\mu} \cdot \frac{\partial N_\mu}{\partial a_{ij}} = (R_i - y_i) \cdot n_j^{(t)}.$$

Wir gehen jetzt ähnlich zu den Beweisen der Propositionen 16.29 und 16.30 vor, benutzen die Hilfsfunktionen

$$z^{(k)} := \begin{bmatrix} z_1^{(k)} \\ \vdots \\ z_\ell^{(k)} \end{bmatrix} \quad \text{mit} \quad z_i^{(k)} = w_{i1}^{(k)} n_1^{(k-1)} + \cdots + w_{i\ell}^{(k)} n_\ell^{(k-1)} + b_i^{(k)}.$$

und beginnen mit den Ableitungen nach den Gewichten und Biasen der letzten Schicht. Zunächst stellen wir fest, dass

$$\frac{\partial N_\mu}{\partial z_i^{(t)}} = \frac{\partial}{\partial z_i^{(t)}} [a_{\mu 1} n_1^{(t)} + \cdots + a_{\mu \ell} n_\ell^{(t)}] = a_{\mu i} \cdot \frac{\partial}{\partial z_i^{(t)}} n_i^{(t)} = a_{\mu i} \cdot \sigma'(z_i^{(t)})$$

gilt, weil  $z_i^{(t)}$  nur in  $n_i^{(t)}$  eingeht. Da  $w_{ij}^{(t)}$  und  $b_i^{(t)}$  nur in  $z_i^{(t)}$  vorkommen, erhalten wir weiter

$$\frac{\partial C}{\partial w_{ij}^{(t)}} = \sum_{\mu=1}^m \frac{\partial C}{\partial N_\mu} \cdot \frac{\partial N_\mu}{\partial z_i^{(t)}} \cdot \frac{\partial z_i^{(t)}}{\partial w_{ij}^{(t)}} = \sum_{\mu=1}^m (R_\mu - y_\mu) \cdot a_{\mu i} \cdot \sigma'(z_i^{(t)}) \cdot n_j^{(t-1)},$$

und

$$\frac{\partial C}{\partial b_i^{(t)}} = \sum_{\mu=1}^m (R_\mu - y_\mu) \cdot a_{\mu i} \cdot \sigma'(z_i^{(t)}) \cdot 1,$$

woraus die Gleichungen für  $\delta^{(t)}$  und für die Ableitungen im Fall  $k = t$  folgen. Wir gehen nun Rückwärts per Induktion vor und müssen für  $1 \leq k < t$  zusätzlich berücksichtigen, dass  $w_{ij}^{(k)}$  und  $b_i^{(k)}$  in alle Neuronen der  $(k+1)$ -ten Schicht eingehen und damit auch in alle  $z_v^{(k+1)}$  für  $v = 1, \dots, \ell$ . Wir erhalten nach der Kettenregel für Funktionen  $\mathbb{R} \rightarrow \mathbb{R}^\ell \rightarrow \mathbb{R}$  also

$$\begin{aligned} \frac{\partial C}{\partial w_{ij}^{(k)}} &= \sum_{v=1}^{\ell} \frac{\partial C}{\partial z_v^{(k+1)}} \cdot \frac{\partial z_v^{(k+1)}}{\partial w_{ij}^{(k)}} \\ &= \sum_{v=1}^{\ell} \frac{\partial C}{\partial z_v^{(k+1)}} \cdot \frac{\partial}{\partial w_{ij}^{(k)}} (w_{v1}^{(k+1)} n_1^{(k)} + \cdots + w_{v\ell}^{(k+1)} n_\ell^{(k)} + b_v^{(k+1)}) \\ &= \sum_{v=1}^{\ell} \frac{\partial C}{\partial z_v^{(k+1)}} \cdot w_{vi}^{(k+1)} \cdot \frac{\partial n_i^{(k)}}{\partial w_{ij}^{(k)}} \end{aligned}$$

$$\begin{aligned}
&= \sum_{v=1}^{\ell} \frac{\partial C}{\partial z_v^{(k+1)}} \cdot w_{vi}^{(k+1)} \cdot \frac{\partial}{\partial w_{ij}^{(k)}} \sigma(\underbrace{w_{i1}^{(k)} n_1^{(k-1)} + \dots + w_{i\ell}^{(k)} n_{\ell}^{(k-1)} + b_i^{(k)}}_{=z_i^{(k)}}) \\
&= \sum_{v=1}^{\ell} \frac{\partial C}{\partial z_v^{(k+1)}} \cdot w_{vi}^{(k+1)} \cdot \sigma'(z_i^{(k)}) \cdot n_j^{(k-1)} \\
&= \left( \sum_{v=1}^{\ell} w_{vi}^{(k+1)} \frac{\partial C}{\partial z_v^{(k+1)}} \cdot \sigma'(z_i^{(k)}) \right) \cdot n_j^{(k-1)}.
\end{aligned}$$

Für  $k = t - 1$  haben wir

$$\frac{\partial C}{\partial z_v^{(k+1)}} = \frac{\partial C}{\partial z_v^{(t)}} = \sum_{\mu=1}^m \frac{\partial C}{\partial N_{\mu}} \cdot \frac{\partial N_{\mu}}{\partial z_v^{(t)}} = \sum_{\mu=1}^m (R_{\mu} - y_{\mu}) \cdot a_{\mu v} \cdot \sigma'(z_v^{(t)}) = \delta_v^{(t)}$$

und wir erhalten für  $k = t - 1, t - 2, \dots, 1$  induktiv die Formel

$$\frac{\partial C}{\partial w_{ij}^{(k)}} = \left( \sum_{v=1}^{\ell} w_{vi}^{(k+1)} \cdot \delta_v^{(k+1)} \cdot \sigma'(z_i^{(k)}) \right) \cdot n_j^{(k-1)}$$

in welcher die große Klammer gerade der  $i$ -te Eintrag von  $(W^{(k+1)})^T \delta^{(k+1)} \odot \sigma'(z^{(k)})$  ist. Für die Ableitung nach  $b_i^{(k)}$  geht man analog vor, erhält aber nach der fünften Zeile in der obigen langen Rechnung

$$\frac{\partial C}{\partial b_i^{(k)}} = \sum_{v=1}^{\ell} \frac{\partial C}{\partial z_v^{(k+1)}} \cdot w_{vi}^{(k+1)} \cdot \sigma'(z_i^{(k)}) \cdot 1$$

und damit alle noch ausstehenden Gleichungen.  $\square$

Als wirklich letzten Punkt dieses Kapitels betrachten wir noch den Fall binärer Label. Hier kann man natürlich obiges für  $m = 2$  anwenden, es ist aber in dieser Situation naheliegend, statt der Einheitsvektoren  $e_1$  und  $e_2$  die Zahlen 0 und 1 zu verwenden, vergleiche die logistische Regression in Kapitel 2, und entsprechend statt des 2-dimensionalen Softmaxausganges einen 1-dimensionalen Sigmoidausgang. Die Klasse der Approximanden besteht dann aus Funktionen der Form

$$S := \text{sig} \circ N : \mathbb{R}^{\nu} \rightarrow (0, 1)$$

mit  $N \in \mathcal{D}^{\sigma, \ell, t}(\mathbb{R}^d, \mathbb{R}^1)$ . Ersetzt man in der Kostenfunktion (16.2) bei  $m = 2$ ,  $R_1$  durch  $S$  und  $R_2$  durch  $1 - S$ , so erhält man

$$C : \mathbb{R}^{\nu} \rightarrow \mathbb{R}, \quad C(a, w, b) = - \sum_{p=1}^n y^{(p)} \log S(x^{(p)}) + (1 - y^{(p)}) \log(1 - S(x^{(p)}))$$

und damit genau das Analogon der negativen Log-Likelihood-Funktion, die wir bei der logistischen Regression minimiert haben. In Aufgabe 16.5 besprechen wir eine praktische Anwendung und in Aufgabe 16.9 zeigen wir, dass man beim Übergang von  $\mathcal{R}^{\sigma, \ell, t}(\mathbb{R}^d, \mathbb{R}^2)$  zu den Funktionen  $S$  wie oben weder Expressivität gewinnt noch einbüßt.

## Referenzen

Wir notieren, dass die erste Version des Satzes von Stone-Weierstraß 16.10 z.B. in [Wer18, Satz VIII.4.7] nachgelesen werden kann. Die zweite Version 16.15 folgt aus der ersten per 1-Punkt-Kompaktifizierung von  $\mathbb{R}$ . Den Rieszsche Darstellungssatz in der Version 16.16 kann man in [Wer18, Theorem II.2.5] nachlesen, die zweite Version, oft auch Riesz-Markov Theorem genannt, findet man in [Rud87, Theorem 6.19]. Für Satz 16.18 verweisen wir auf [Wer18, Korollar III.1.8]. Ferner verweisen wir auf [Wer18, Anhang A] für Definitionen und Resultate welche über den klassischen Inhalt einer Maßtheorievorlesung hinausgehen, wie z.B. die Räume  $M(\Omega)$  und  $M(\mathbb{R})$  von komplexen Maßen, das Integral bzgl. eines komplexen Maßes und den Satz von Fubini für komplexe Maße.

Die Resultate am Kapitelanfang zur Darstellung logischer Funktionen sind allesamt Standard. Unsere Version von Satz 16.4 ist eine Umformulierung von [SSBD14, Claim 20.1]. Die Geschichte der kontinuierlichen Approximationsresultate beginnt mit der Arbeit [Cyb89] von 1989, in welcher Lemma 16.17 und Satz 16.20 für  $\Omega = [0, 1]^d$  direkt mit einer multivariablen Variante der von uns gegebenen Beweise, gezeigt wurden. Beachte das Korrigendum [Cyb92]. Es folgten mehrere Arbeiten [HSW89, CL92, Hor91, HSW90, LLPS93, LP93, SC92] zwischen 1989 und 1993, die Cybenko's Ergebnis verallgemeinert und ausgebaut haben im Sinne, dass (a) die Voraussetzungen an die Aktivierungsfunktion abgeschwächt wurden, und (b) die Approximation von anderen Klassen von Funktionen (z.B.  $L^p$ ,  $C^m$ ) betrachtet wurde. Bei der Behandlung der ReLU-Funktion folgen wir dem Zugang in [Gui18, Lemma 3.15]. Das von uns als Proposition 16.12 gebrachte Reduktionsargument basiert auf [CL92] und [LP93]. Der Beweis von 16.22 folgt [LLPS93, Remark 3]. Was die Approximation stetiger Funktionen im Sinne gleichmäßiger Konvergenz auf Kompakta angeht, erwähnen wir, dass in [LLPS93] für  $\sigma \in L_{\text{loc}}^\infty(\mathbb{R})$  derart, dass der Abschluss der Menge der Unstetigkeitsstellen von  $\sigma$  eine Lebesgue-Nullmenge ist, die Äquivalenz

$$C(\mathbb{R}^n, \mathbb{R}^m) \subseteq \overline{S_\sigma(\mathbb{R}^n, \mathbb{R}^m)}^{L_{\text{loc}}^\infty(\mathbb{R}^n, \mathbb{R}^m)} \iff \sigma \notin \mathbb{R}[X]$$

gezeigt wird, vergleiche mit Beispiele 16.7(i). Der das Trade-off zwischen flachen und tiefen Netzen erlaubende Satz 16.24 basiert [Han19], wurde hier leicht abgewandelt um Netze mit linearem Ausgang zu behandeln. Im Text haben wir nur sehr einfache Architekturen neuronaler Netze behandelt und verweisen für eine umfassende Behandlung auf [Cal20].

Die Technik der Rückwärtspropagation hat ihren Ursprung in den 1960er Jahren; unsere Darstellung basiert auf dem Buch [Nie15] wurde aber leicht abgewandelt und um die Behandlung eines Netzes mit Softmax-Ausgang erweitert.

## Aufgaben

**Aufgabe 16.1.** In dieser Aufgabe betrachten wir Neuronen mit Heaviside-Aktivierung.

- (i) Zeigen Sie, dass die Entweder-Oder-Funktion XOR:  $\{0, 1\}^2 \rightarrow \{0, 1\}$  nicht durch ein einzelnes Neuron mit zwei Eingängen dargestellt werden kann.
- (ii) Geben Sie eine Darstellung von XOR durch ein neuronales Netz mit (höchstens) drei Neuronen und Heaviside-Aktivierung an.

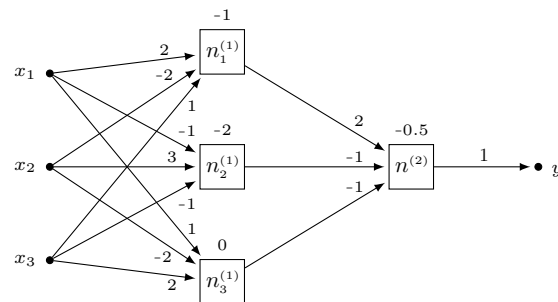
**Aufgabe 16.2.** Nutzen Sie die konjunktive Normalform um die durch die folgende Tabelle definierte Funktion  $f: \{0, 1\}^3 \rightarrow \{0, 1\}$

$x_1$	$x_2$	$x_3$	$f(x_1, x_2, x_3)$
0	0	0	1
0	0	1	1
0	1	0	1
0	1	1	1
1	0	0	0
1	0	1	1
1	1	0	0
1	1	1	0

durch ein Heaviside-aktiviertes neuronales Netz darzustellen.

**Aufgabe 16.3.** Zeigen Sie durch Anwendung des Satzes von Stone-Weierstraß, dass der Raum der flachen neuronalen Netze mit Exponentialaktivierung  $\mathcal{S}^{\text{exp}}(\mathbb{R}) \subseteq C(\mathbb{R})$  in den stetigen Funktionen gleichmäßig dicht auf Kompakta ist.

**Aufgabe 16.4.** Wir betrachten das folgende neuronale Netz



bei welchem die Zahlen über den Pfeilen die Gewichte angeben und die Zahlen über den Boxen jeweils das Bias des Neurons. Die Neuronen  $n_1^{(1)}, n_2^{(1)}, n_3^{(1)}$  haben die ReLU-Funktion und  $n^{(2)}$  die Sigmoidfunktion als Aktivierung. Berechnen Sie für jeden in der folgenden Tabelle gegebenen Eingang sukzessive die Ausgaben der Neuronen  $n_1^{(1)}, n_2^{(1)}, n_3^{(1)}$ , und dann die des Neurons  $n^{(2)}$ , welche gleich der Ausgabe des gesamten Netzes ist.

$x_1$	$x_2$	$x_3$	$n_1^{(1)}$	$n_2^{(1)}$	$n_3^{(1)}$	$n^{(2)}$	$y$
0	0	0					
0	0	1					
0	1	0					
0	1	1					
1	0	0					
1	0	1					
1	1	0					
1	1	1					

**Aufgabe 16.5.** Trainieren Sie mit der Funktion `MLPClassifier` aus dem Pythonpaket `sklearn` ein neuronales Netz mit einer versteckten Schicht, sodass die Funktion

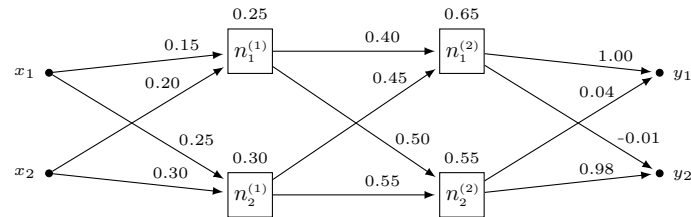
$$f: \{0, 1\}^3 \rightarrow \{0, 1\}, f(x_1, x_2, x_3) := \text{XOR}(x_1, x_2) \vee (x_1 \wedge x_2 \wedge x_3)$$

möglichst gut approximiert wird. Wie breit muss die versteckte Schicht sein?

*Hinweis:* In `MLPClassifier` müssen Sie nur die Breite der versteckten Schichten angeben und deren Aktivierung. Mit `.predict()` können Sie dann die  $\{0, 1\}$ -wertige Ausgabe für alle Trainingsdaten auf

einmal ausdrucken passen. Im Hintergrund benutzt `MLPClassifier` allerdings eine Ausgangsschicht wie in Aufgabe 16.4 und berechnet ein  $p \in [0, 1]$ , sodass sich die Einträge des Tupels  $(p, 1 - p)$  als Wahrscheinlichkeiten für die Klassifizierung mit 0 bzw. 1 interpretieren lassen. Diese Tupel können per `.predict_proba()` ausgedruckt werden.

**Aufgabe 16.6.** Wir betrachten das folgende Netz



bei welchem die Zahlen über den Pfeilen die Gewichte angeben und die Zahlen über den Boxen jeweils das Bias des Neurons. Alle Neuronen seien Sigmoid-aktiviert. Es sei ein einzelner Trainingsdatenpunkt  $(x, y)$  gegeben mit  $x = (0.05, 0.10)$  und  $y = (0.05, 1.01)$ . Führen Sie *per Hand* eine Vorwärts- und eine Rückwärtspropagation durch und bestimmen Sie das erste Update der Gewichte, Biase, Matrixelemente bei einer Lernrate von 0.5.

**Aufgabe 16.7.** Verketteten Sie das Netz aus Aufgabe 16.6 mit der Softmaxfunktion und führen Sie wieder eine Vorwärts- und eine Rückwärtspropagation durch.

**Aufgabe 16.8.** Schreiben Sie ein Python-Programm, welches handschriftlich notierte Zahlen mittels eines tiefen neuronalen Netzes klassifiziert. Nutzen Sie dafür die MNIST-Datenmenge und Satz 16.34. Verwenden Sie einen Teil der Datenmenge als Trainingsdaten und dann testen Sie dann den Klassifizierer auf zufällig ausgewählten Datenpunkten, die nicht zu den Trainingsdaten gehört haben.

*Hinweis:* Dies ist eine eher aufwändige Aufgabe, wenn Sie alles ‘from scratch’ programmieren, aber auch eine Gelegenheit, zumindest einmal ein neuronales Netz von Grund auf selber zu erstellen und zu trainieren — bevor Sie dann in Zukunft auf fertige Pakete zurückgreifen, die das für Sie erledigen.

**Aufgabe 16.9.** Zeigen Sie, dass gilt

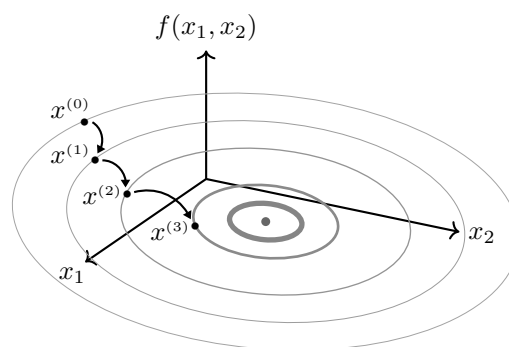
$$\mathcal{R}^{\sigma, \ell, t}(\mathbb{R}^d, \mathbb{R}^2) = \left\{ \begin{bmatrix} S \\ 1-S \end{bmatrix} \mid S = \text{sig} \circ N \mid N \in \mathcal{D}^{\sigma, \ell, t}(\mathbb{R}^d, \mathbb{R}^1) \right\}.$$

## Kapitel 17

# Gradientenverfahren für konvexe Funktionen

In diesem Kapitel geben wir eine Einführung in das Gradientenverfahren. Letzteres ist eine klassische Methode der Optimierung mit welcher Extrema reellwertiger Funktionen numerisch bestimmt werden können. Damit ist das Gradientenverfahren vielfältig einsetzbar und wird in der Numerik- und Optimierungsliteratur ausgiebig behandelt. Da das Verfahren in vielen Data Science und Machine Learning Problemen Anwendung findet und wie wir es in vorhergehenden Kapiteln auch schon mehrfach notiert haben, soll hier das Verfahren unter gutartigen Bedingungen untersucht werden.

Die der Gradientenmethode zugrundeliegende Idee kann anschaulich wie folgt beschrieben werden. Angenommen wir suchen das Minimum einer Funktion  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  wie im folgenden Bild.



Dann kann man versuchen, dieses zu approximieren, indem man an einem, z.B. zufällig gewählten, Punkt  $x^{(0)}$  startet und dann schrittweise jeweils in die Richtung geht, in der es am steilsten bergab geht. Formal definiert man per

$$x^{(k+1)} := x^{(k)} - \gamma_k \nabla f(x^{(k)}) \quad (17.1)$$

die Folge  $(x^{(k)})_{k \in \mathbb{N}}$  wobei  $\gamma_k > 0$  *Schrittweite* oder, im Kontext von maschinellem



Lernen, auch *Lernrate*, genannt wird. Letztere kann konstant sein, d.h.,  $\gamma_k \equiv \gamma$ , sie kann durch eine vorgegebene Folge variiert werden, z.B.,  $\gamma_k = 1/k$ , oder in jedem Schritt neu gewählt werden, z.B. so, dass  $f(x^{(k)} - \gamma_k \nabla f(x^{(k)}))$  möglichst klein wird. Erstmal besteht natürlich keine Garantie, dass man auf diese Weise mit den  $x^{(k)}$  in die Nähe einer Minimalstelle, bzw. mit  $f(x^{(k)})$  in die Nähe des minimalen Wertes von  $f$  gelangt. Wir werden in diesem Kapitel zeigen, dass dies unter einer geeigneten Kombination von Konvexitäts- und Differenzierbarkeitseigenschaften an  $f$ , und geeignet kleiner und konstanter Schrittweite, der Fall ist, und dass unter diesen Voraussetzungen quantifiziert werden kann, wieviele Iterationen nötig sind um eine gewisse Genauigkeit der Approximation zu erreichen. Damit sichergestellt ist, dass (17.1) wohldefiniert ist, schränken wir uns für das ganze Kapitel auf Funktionen ein, die differenzierbar, und der Einfachheit halber, auf ganz  $\mathbb{R}^n$  definiert sind.

Als erstes präzisieren wir unsere Notation.

**Definition 17.1.** Sei  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  eine Funktion. Ein Punkt  $x^* \in \mathbb{R}^n$  heißt *Minimalstelle* oder *Minimierer* von  $f$ , falls  $f(x^*) \leq f(x)$  für alle  $x \in \mathbb{R}^d$  gilt. Wir bezeichnen mit

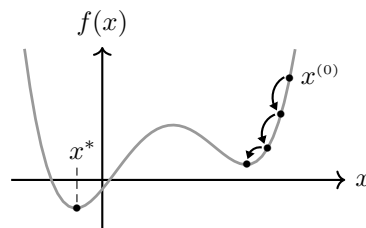
$$X^* = \operatorname{argmin}_{x \in \mathbb{R}^n} f(x)$$

die Menge aller Minimierer von  $f$ . Falls  $X^* = \{x^*\}$  einelementig ist, schreiben wir auch  $x^* = \operatorname{argmin}_{x \in \mathbb{R}^n} f(x)$ . Insbesondere sind im folgenden Minimierer immer globale Minimalstellen. Den *minimalen Wert* oder das *Minimum* von  $f$  bezeichnen wir mit

$$f^* := \min_{x \in \mathbb{R}^n} f(x)$$

falls dieses existiert.

Im Allgemeinen kann es natürlich sein, dass  $f$  gar keine Minimierer hat. Obwohl man in solchen Fällen eventuell trotzdem das Gradientenverfahren sinnvoll verwenden kann, siehe Aufgabe 17.1, werden wir im folgenden entweder explizit voraussetzen, dass Minimierer existieren, oder durch andere Voraussetzungen an  $f$  garantieren, dass dem so ist. Ist  $X^* \neq \emptyset$  gesichert, so wird ein weiteres zentrales Problem des Gradientenverfahrens durch das folgende Bild verdeutlicht, für ein konkretes Beispiel siehe Aufgabe 17.3.



Die dargestellte Funktion hat einen sogar eindeutig bestimmten Minimierer und die Folge  $(f(x^{(k)}))_{k \in \mathbb{N}}$  konvergiert, aber nicht gegen das Minimum  $f^*$ . Auch für derartige Funktionen kann die Gradientenmethode erfolgreich angewandt werden, z.B. könnte

man es bei der Funktion im Bild mit einem negativen Startwert versuchen. Für unsere theoretischen Resultate in diesem Kapitel wollen wir aber solche Effekte ausschließen und tun dies durch die folgende Definition.

**Definition 17.2.** Eine Funktion  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  heißt *konvex*, falls

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y)$$

für alle  $\lambda \in [0, 1]$  und  $x, y \in \mathbb{R}^n$  gilt. Beachte, dass die Ungleichung für  $\lambda \in \{0, 1\}$  oder  $x = y$  trivialerweise gilt.

Für konvexe Funktionen gilt das folgende *Lokal-zu-Global-Prinzip*.

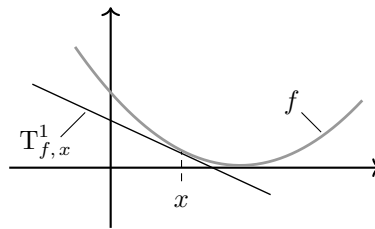
**Proposition 17.3.** Sei  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  konvex. Wenn  $x$  eine lokale Minimalstelle von  $f$  ist, so ist  $x \in X^*$ , d.h.,  $x$  ist ein (globaler) Minimierer.

*Beweis.* Sei  $x$  eine lokale Minimalstelle. Für beliebiges  $y \in \mathbb{R}^n$  existiert dann  $\lambda_0 \in (0, 1)$  sodass für alle  $\lambda \in (0, \lambda_0)$  wegen der Konvexität die Abschätzung

$$f(x) \leq f((1 - \lambda)x + \lambda y) \leq (1 - \lambda)f(x) + \lambda f(y) = f(x) - \lambda f(x) + \lambda f(y)$$

gilt. Hieraus folgt  $f(x) \leq f(y)$  wegen  $\lambda > 0$ . Da  $y \in \mathbb{R}^n$  beliebig war, zeigt dies, dass  $x \in X^*$  gilt.  $\square$

Als nächstes zeigen wir, dass für konvexe differenzierbare Funktionen die notwendige Bedingung  $\nabla f(x) = 0$  für einen Minimierer bereits hinreichend ist. Zuerst charakterisieren wir die Konvexität mithilfe des Gradienten. Beachte, dass im folgenden Lemma die kleinere Seite der Abschätzung gerade das erste Taylorpolynom  $T_{f,x}^1(y)$  von  $f$  in  $x$ , ausgewertet in  $y$ , ist.



Die Abschätzung im Lemma besagt also anschaulich, dass die Funktionswerte von  $f$  oberhalb der Tangente an den Graph von  $f$  im Punkt  $x$  liegen.

**Proposition 17.4.** Sei  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  differenzierbar. Dann ist  $f$  konvex genau dann wenn  $f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle$  für alle  $x, y \in \mathbb{R}^d$  gilt.

*Beweis.* “ $\implies$ ” Aus der Definition der Konvexität erhalten wir durch Abziehen von  $f(x)$ , Division durch  $\lambda$ , und Grenzübergang

$$f(y) - f(x) \geq \frac{f(x + \lambda(y - x)) - f(x)}{\lambda} \xrightarrow{\lambda \searrow 0} f'(x)(y - x) = \langle \nabla f(x), y - x \rangle.$$

“ $\Leftarrow$ ” Seien  $x, y \in \mathbb{R}^n$  und  $\lambda \in (0, 1)$ . Wir setzen  $z := \lambda x + (1 - \lambda)y$ . Dann gelten per Voraussetzung  $f(x) \geq f(z) + \langle \nabla f(z), x - z \rangle$  und  $f(y) \geq f(z) + \langle \nabla f(z), y - z \rangle$ . Multiplikation mit  $\lambda$  bzw. mit  $(1 - \lambda)$  und Addition beider Gleichungen liefert

$$\begin{aligned} \lambda f(x) + (1 - \lambda)f(y) &\geq \lambda f(z) + \lambda \langle \nabla f(z), x - z \rangle + (1 - \lambda)f(z) + (1 - \lambda) \langle \nabla f(z), y - z \rangle \\ &= f(z) + \langle \nabla f(z), \underbrace{\lambda x - \lambda z + (1 - \lambda)y - (1 - \lambda)z}_{=z - \lambda z - (1 - \lambda)z = 0} \rangle \\ &= f(\lambda x + (1 - \lambda)y) \end{aligned}$$

und damit die Konvexität von  $f$ .  $\square$

Wir notieren als Folgerung.

**Folgerung 17.5.** *Sei  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  konvex und differenzierbar. Dann gilt  $x \in X^*$  genau dann wenn  $\nabla f(x) = 0$ .*

*Beweis.* Es muss nur gezeigt werden, dass  $\nabla f(x) = 0$  hinreichend für ein lokales Minimum ist, denn dann folgt das Gewünschte aus Proposition 17.3. Wegen der Konvexität von  $f$  impliziert Lemma 17.4 mit  $\nabla f(x) = 0$  aber, dass  $f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle = f(x)$  für alle  $y \in \mathbb{R}^n$  gilt und wir sind fertig.  $\square$

Sei nun eine Funktion  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  gegeben, die konvex und differenzierbar ist, und sei zusätzlich für diese Funktion  $X^* \neq \emptyset$ . Für einen gegebenen Startwert  $x^{(0)}$  betrachten wir die Folge  $(x^{(k)})_{k \in \mathbb{N}}$  definiert per (17.1), d.h.

$$x^{(k+1)} = x^{(k)} - \gamma_k \nabla f(x^{(k)}).$$

In erster Linie interessiert uns, ob, und wenn ja wie schnell, die Bildfolge  $(f(x^{(k)}))_{k \in \mathbb{N}}$  gegen  $f^*$  konvergiert, vergleiche die in Kapitel 2 und 16 erläuterten Anwendungen. Das folgende Beispiel illustriert zunächst, dass Konvexität und Differenzierbarkeit alleine nicht ausreichen, um die Konvergenz der Folge  $(f(x^{(k)}))_{k \in \mathbb{N}}$  zu garantieren, sondern dass die Wahl der Schrittweiten essentiell ist.

**Beispiel 17.6.** Sei  $f: \mathbb{R} \rightarrow \mathbb{R}$ ,  $f(x) = x^2$ , also  $x^* := \operatorname{argmin}_{x \in \mathbb{R}} f(x) = 0$ , und sei  $x^{(0)} \neq 0$ . Dann gilt  $x^{(k+1)} = (1 - 2\gamma_k)x^{(k)}$  und wir haben die folgenden Effekte je nach Wahl der  $\gamma_k$ .

- (i) Für  $\gamma_k \equiv 1$  ist  $x^{(k)} = \pm x^{(0)}$ , und die Bildfolge  $f(x^{(k)}) \equiv f(x^{(0)})$  ist konstant und ungleich dem Minimum.
- (ii) Für  $\gamma_k \equiv \gamma > 1$  ist  $(x^{(k)})_{k \in \mathbb{N}}$  divergent und  $f(x^{(k)}) \rightarrow \infty$ .
- (iii) Für  $\gamma_k \equiv \gamma < 1$  haben wir  $x^{(k)} \rightarrow 0 = x^*$  sowie  $f(x^{(k)}) \rightarrow 0 = f^*$  und damit das gewünschte Ergebnis.
- (iv) Für  $\gamma_k = 1/k$  haben wir  $x^{(k)} = (1 - \frac{2}{k})^{k-1} x^{(0)} \rightarrow x^{(0)} e^{-2}$  sowie  $f(x^{(k)}) \rightarrow f(x^{(0)} e^{-2}) = e^{-4} (x^{(0)})^2$ , d.h. die Bildfolge konvergiert gegen einen Wert, welcher, abhängig vom Startpunkt, beliebig viel größer als  $f^*$  sein kann.

In Beispiel 17.6(i)–(ii) kommt  $f(x^{(k)}) \not\rightarrow f^*$  natürlich nur dadurch zustande, dass mutwillig eine zu große Schrittweite gewählt wird. In der Tat kann man bei jeder konvexen und differenzierbaren Funktion  $f$  im  $k$ -ten Schritt des Gradientenverfahrens ein  $\gamma_k > 0$  wählen, sodass

$$f(x^{(k+1)}) = f(x^{(k)} - \gamma_k \nabla f(x^{(k)})) < f(x^{(k)})$$

gilt, solange nur  $\nabla f(x^{(k)}) \neq 0$  ist: Wäre nämlich  $f(x^{(k+1)}) = f(x^{(k)} - \gamma \nabla f(x^{(k)})) \geq f(x^{(k)})$  für alle  $\gamma > 0$ , so würde für die Richtungsableitung

$$-\langle \nabla f(x^{(k)}), \nabla f(x^{(k)}) \rangle = \partial_{-\nabla f(x^{(k)})} f(x^{(k)}) = \lim_{h \searrow 0} \frac{f(x^{(k)} - h \nabla f(x^{(k)})) - f(x^{(k)})}{h} \geq 0$$

folgen, was nur geht, wenn der Gradient in  $x^{(k)}$  verschwindet. Es gibt also stets eine Folge von Schrittweiten  $(\gamma_k)_{k \in \mathbb{N}}$  sodass die Bildfolge  $(f(x^{(k+1)}))_{k \in \mathbb{N}}$  monoton fällt und daher konvergiert. Solange wir nicht mit den  $x^{(k)}$  einen Minimierer genau treffen, fällt die Folge sogar strikt, was bedeutet, dass unsere Approximation von  $f^*$  mit jeder Iteration besser wird. Andererseits zeigt Beispiel 17.6(iv), dass bei zu schnell fallenden  $\gamma_k$  wieder  $f(x^{(k)}) \not\rightarrow f^*$  eintreten kann.

Die Schrittweiten dürfen also einerseits nicht konstant zu groß sein, aber andererseits auch nicht zu schnell gegen Null gehen. Unsere vorhergehenden Argumente und das Beispiel legen nahe, dass man mit einem konstanten  $\gamma > 0$  auskommen sollte, wenn die Ableitung von  $f$  beschränkt ist. In der Tat fordern wir nun etwas stärkeres, vergleiche aber Aufgabe 17.10.

**Definition 17.7.** Sei  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  differenzierbar und sei  $L > 0$ . Dann heißt  $f$  *L-glatt*, falls  $\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|$  für alle  $x, y \in \mathbb{R}^n$  gilt.

Äquivalent ausgedrückt bedeutet dies, dass  $\nabla f$  Lipschitz-stetig ist. Wir benötigen die zwei folgenden Konsequenzen.

**Lemma 17.8.** Sei  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  eine *L-glatte Funktion* mit  $L > 0$ . Dann gilt  $f(x) - f(y) - \langle \nabla f(y), x - y \rangle \leq \frac{L}{2}\|x - y\|^2$  für  $x, y \in \mathbb{R}^n$ .

*Beweis.* Seien  $x, y \in \mathbb{R}^n$  fest. Wir definieren die Hilfsfunktion  $h: \mathbb{R} \rightarrow \mathbb{R}$ ,  $h(t) = f(y + t(x - y))$ . Dann gilt  $h'(t) = \langle \nabla f(y + t(x - y)), x - y \rangle$  nach der Kettenregel. Nach dem Hauptsatz gilt weiter

$$f(x) - f(y) = h(1) - h(0) = \int_0^1 h'(t) dt = \int_0^1 \langle \nabla f(y + t(x - y)), x - y \rangle dt.$$

Damit erhalten wir, unter Beachtung dass das Integral von 0 bis 1 läuft,

$$\begin{aligned} f(x) - f(y) - \langle \nabla f(y), x - y \rangle &= \int_0^1 \langle \nabla f(y + t(x - y)), x - y \rangle - \langle \nabla f(y), x - y \rangle dt \\ &\leq \int_0^1 |\langle \nabla f(y + t(x - y)) - \nabla f(y), x - y \rangle| dt \\ &\leq \int_0^1 \|\nabla f(y + t(x - y)) - \nabla f(y)\| \|x - y\| dt \end{aligned}$$

$$\begin{aligned}
&\leq \int_0^1 L \|y + t(x - y) - y\| \|x - y\| dt \\
&= \int_0^1 Lt \|x - y\|^2 dt = \frac{L}{2} \|x - y\|^2
\end{aligned}$$

wobei wir zuerst die Cauchy-Schwarz-Bunjakowski-Ungleichung und dann die  $L$ -Glattheit benutzt haben.  $\square$

Wir fügen jetzt die Voraussetzung hinzu, dass  $f$  konvex ist.

**Lemma 17.9.** *Sei  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  konvex und  $L$ -glatt mit  $L > 0$ . Dann gilt  $f(x) - f(y) \leq \langle \nabla f(x), x - y \rangle - \frac{1}{2L} \|\nabla f(x) - \nabla f(y)\|^2$  für  $x, y \in \mathbb{R}^n$ .*

*Beweis.* Seien  $x, y \in \mathbb{R}^n$  gegeben. Wir setzen  $z := y - \frac{1}{L}(\nabla f(y) - \nabla f(x))$  und schätzen in der folgenden Rechnung  $f(x) - f(z)$  mithilfe von Lemma 17.4 nach oben ab und  $f(z) - f(y)$  mithilfe von Lemma 17.8. Dann setzen wir  $z$  ein und vereinfachen. Das ergibt

$$\begin{aligned}
f(x) - f(y) &= f(x) - f(z) + f(z) - f(y) \\
&\leq \langle \nabla f(x), x - z \rangle + \langle \nabla f(y), z - y \rangle + \frac{L}{2} \|z - y\|^2 \\
&= \langle \nabla f(x), x - y + \frac{1}{L}(\nabla f(y) - \nabla f(x)) \rangle - \langle \nabla f(y), y - y \\
&\quad + \frac{1}{L}(\nabla f(y) - \nabla f(x)) \rangle + \frac{L}{2} \|y - \frac{1}{L}(\nabla f(y) - \nabla f(x)) - y\|^2 \\
&= \langle \nabla f(x), x - y \rangle + \langle \nabla f(x) - \nabla f(y), \frac{1}{L}(\nabla f(y) - \nabla f(x)) \rangle \\
&\quad + \frac{1}{2L} \|\nabla f(y) - \nabla f(x)\|^2 \\
&= \langle \nabla f(x), x - y \rangle - \frac{1}{2L} \|\nabla f(y) - \nabla f(x)\|^2
\end{aligned}$$

wie behauptet.  $\square$

Bevor wir das erste Hauptresultat zum Gradientenverfahren beweisen, notieren wir noch die folgende Äquivalenz, die in der Literatur manchmal auch als Definition von konvexen  $L$ -glatten Funktionen verwendet wird.

**Folgerung 17.10.** *Sie  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  differenzierbar und konvex. Sei  $L > 0$ . Dann sind folgende Aussagen äquivalent.*

(i)  $f$  ist  $L$ -glatt und konvex.

(ii)  $\forall x, y \in \mathbb{R}^n: 0 \leq f(x) - f(y) - \langle \nabla f(y), x - y \rangle \leq \frac{L}{2} \|x - y\|^2$ .

*Beweis.* (i)  $\implies$  (ii): Aus der Konvexität folgt die erste Ungleichung mit Lemma 17.4 und die zweite Ungleichung folgt aus der  $L$ -Glattheit mit Lemma 17.8.

(ii)  $\implies$  (i): Die erste Ungleichung impliziert Konvexität mit Lemma 17.4. Der Beweis von Lemma 17.9 zeigt, dass Konvexität und die zweite Ungleichung, d.h. gerade die Ungleichung in Lemma 17.8, implizieren, dass

$$\frac{1}{2L} \|\nabla f(x) - \nabla f(y)\|^2 \leq f(y) - f(x) + \langle \nabla f(x), x - y \rangle$$

$$\begin{aligned}
&= f(y) - f(x) - \langle \nabla f(x), y - x \rangle \\
&\leq \frac{L}{2} \|x - y\|^2
\end{aligned}$$

gilt, wobei wir im letzten Schritt nochmal die zweite Ungleichung aus (ii) benutzt haben, aber mit vertauschten Rollen von  $x$  und  $y$ . Multiplikation mit  $2L$  und Wurzelziehen zeigt nun  $\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|$  und damit ist  $f$   $L$ -glatt.  $\square$

Jetzt sind wir bereit für unser erstes Resultat im Zusammenhang mit dem Gradientenverfahren. Dieses wird zeigen, dass für eine  $L$ -glatte und konvexe Funktion  $f$  die Folge  $(f(x^{(k)}))_{k \in \mathbb{N}}$  konvergiert und zwar für geeignet kleine, aber dann konstante, Schrittweite  $\gamma$  in (17.1).

**Satz 17.11.** *Sei  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  konvex und  $L$ -glatt mit  $L > 0$ . Seien  $x^* \in X^*$  und  $x^{(0)} \in \mathbb{R}^n$  beliebig,  $(x^{(k)})_{k \in \mathbb{N}}$  bezeichne die Folge aus (17.1) wobei die Schrittweite  $\gamma := 1/L$  konstant gewählt wird. Dann gilt*

$$0 \leq f(x^{(k)}) - f^* \leq \frac{2L}{k} \|x^{(0)} - x^*\|^2$$

für  $k \geq 1$  und insbesondere  $f(x^{(k)}) \rightarrow f^*$  für  $k \rightarrow \infty$ .

*Beweis.* ① Für  $j \geq 0$  setzen wir  $x = x^{(j+1)} = x^{(j)} - \frac{1}{L} \nabla f(x^{(j)})$  und  $y = x^{(j)}$  in die Ungleichung aus Lemma 17.8 ein und erhalten

$$f(x^{(j+1)}) - f(x^{(j)}) - \langle \nabla f(x^{(j)}), -\frac{1}{L} \nabla f(x^{(j)}) \rangle \leq \frac{L}{2} \left\| 0 - \frac{1}{L} \nabla f(x^{(j)}) \right\|^2.$$

Ausrechnen und umstellen liefert

$$f(x^{(j+1)}) - f(x^{(j)}) \leq \frac{1}{2L} \|\nabla f(x^{(j)})\|^2 - \frac{1}{L} \langle \nabla f(x^{(j)}), \nabla f(x^{(j)}) \rangle = -\frac{1}{2L} \|\nabla f(x^{(j)})\|^2$$

woran wir sehen, dass für diese Wahl von  $\gamma$  die Folge  $(f(x^{(j)}))_{j \in \mathbb{N}}$  monoton fällt und damit konvergiert.

② Aus Lemma 17.9 und Lemma 17.4 folgern wir

$$\begin{aligned}
\langle \nabla f(x) - \nabla f(y), x - y \rangle &= \langle f(x), x - y \rangle - \langle \nabla f(y), x - y \rangle \\
&\stackrel{\text{Lem. 17.9}}{\geq} f(x) - f(y) + \frac{1}{2L} \|\nabla f(x) - \nabla f(y)\|^2 - \langle \nabla f(y), x - y \rangle \\
&\stackrel{\text{Lem. 17.4}}{\geq} \frac{1}{2L} \|\nabla f(x) - \nabla f(y)\|^2
\end{aligned}$$

für beliebige  $x, y \in \mathbb{R}^n$ . Wir setzen nun  $x = x^{(j)}$  und  $y = x^*$ , also  $\nabla f(x^*) = 0$  nach Folgerung 17.5, und erhalten

$$\langle \nabla f(x^{(j)}), x^{(j)} - x^* \rangle \geq \frac{1}{2L} \|\nabla f(x^{(j)})\|^2.$$

Nach diesen Vorbereitungen folgt

$$\|x^{(j+1)} - x^*\|^2 = \|x^{(j)} - \frac{1}{L} \nabla f(x^{(j)}) - x^*\|^2$$

$$\begin{aligned}
&= \|x^{(j)} - x^*\|^2 + 2\left\langle -\frac{1}{L}\nabla f(x^{(j)}), x^{(j)} - x^* \right\rangle + \left\| -\frac{1}{L}\nabla f(x^{(j)}) \right\|^2 \\
&\leq \|x^{(j)} - x^*\|^2 - \frac{2}{L^2}\|\nabla f(x^{(j)})\|^2 + \frac{1}{L^2}\|\nabla f(x^{(j)})\|^2 \\
&\stackrel{\text{s.o.}}{\uparrow} \\
&= \|x^{(j)} - x^*\|^2.
\end{aligned}$$

③ Jetzt bemerken wir, dass wir die Abschätzung im Satz nur für  $k < k_0 := \inf\{j \geq 0 \mid x^{(j)} \in X^*\}$  zeigen müssen. Ist nämlich  $k_0 < \infty$  und  $k \geq k_0$ , so ist die linke Seite der Ungleichung Null. Wir nehmen an, dass  $k_0 \geq 1$  ist, insbesondere gilt also  $x^{(0)} \neq x^*$ , wir fixieren  $0 \leq k < k_0$  und setzen  $\delta_k := f(x^{(k)}) - f(x^*)$ . Nach dem gerade notierten sind dann  $\delta_0, \dots, \delta_k > 0$ . Für  $0 \leq j < k$  gilt nach ①

$$\delta_{j+1} = f(x^{(j+1)}) - f(x^{(j)}) + f(x^{(j)}) - f(x^*) \leq \delta_j - \frac{1}{2L}\|\nabla f(x^{(j)})\|^2$$

und weiter

$$\begin{aligned}
\delta_j &= f(x^{(j)}) - f(x^*) \leq -\langle \nabla f(x^{(j)}), x^* - x^{(j)} \rangle \\
&\stackrel{\text{Lem. 17.4}}{\uparrow} \\
&\leq \|\nabla f(x^{(j)})\| \|x^{(j)} - x^*\| \leq \|\nabla f(x^{(j)})\| \|x^{(0)} - x^*\|. \\
&\stackrel{\text{CSB-Ungl.}}{\uparrow} \qquad \qquad \qquad \stackrel{\text{Teil ②}}{\uparrow}
\end{aligned}$$

Es folgt  $\|\nabla f(x^{(j)})\| \geq \delta_j / \|x^{(0)} - x^*\|$  und wir erhalten

$$\delta_{j+1} \leq \delta_j - \frac{1}{2L}\|\nabla f(x^{(j)})\|^2 \leq \delta_j - \frac{1}{2L\|x^{(0)} - x^*\|^2}\delta_j^2 = \delta_j - \omega\delta_j^2,$$

mit  $\omega := (2L\|x^{(0)} - x^*\|^2)^{-1}$ . Dies bedeutet  $\omega\delta_j^2 + \delta_{j+1} \leq \delta_j$  woraus per Division durch  $\delta_{j+1} \cdot \delta_j$  zunächst die Ungleichung  $\omega\frac{\delta_j}{\delta_{j+1}} + \frac{1}{\delta_j} \leq \frac{1}{\delta_{j+1}}$  folgt und wir dann durch Umstellen  $\frac{1}{\delta_{j+1}} - \frac{1}{\delta_j} \geq \omega$  erhalten. Für  $k \geq 1$  summieren wir nun über die Ungleichung ab und erhalten per Teleskopsumme

$$\frac{1}{\delta_k} \geq \frac{1}{\delta_k} - \frac{1}{\delta_0} = \sum_{j=0}^{k-1} \frac{1}{\delta_{j+1}} - \frac{1}{\delta_j} \geq k\omega$$

also  $\delta_k \leq \frac{1}{k\omega}$ . Daraus folgt, zusammen mit  $f^* - f(x^*) \leq f(x^{(k)})$ , dass

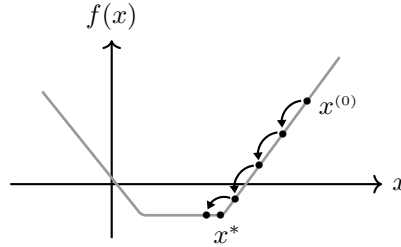
$$0 \leq f(x^{(k)}) - f^* = \delta_k \leq \frac{1}{k\omega} = \frac{2L\|x^{(0)} - x^*\|^2}{k} \leq \frac{2L}{k} \text{dist}(x^{(0)}, x^*)^2$$

gilt wie gewünscht.  $\square$

In der obigen Situation ist  $2L\|x^{(k)} - x^*\|^2$  eventuell zwar nicht explizit bekannt, aber doch konstant. Die Anzahl der Iterationen von (17.1), die durchgeführt werden müssen, um eine Approximationsqualität von  $|f(x^{(k)}) - f^*| < \varepsilon$  zu erreichen, ist dann durch  $\mathcal{O}(\frac{1}{\varepsilon})$  beschränkt. Wählt man oben  $\gamma \in (0, 1/L]$  fest, so erhält man eine analoge Abschätzung wie im Satz und weiß somit, dass bei der Anwendung der Gradientenmethode auf eine Funktion, bei der  $L$  nicht explizit bekannt ist, geeignet

kleine und konstante  $\gamma$  zum Erfolg führen werden, vgl. Aufgabe 17.6.

Satz 17.11 garantiert nicht die Konvergenz  $x^{(k)} \rightarrow x^*$  und dies kann man auch nicht erwarten, betrachte das Trivialbeispiel  $f \equiv 0$  und  $x^* \neq x^{(0)}$  oder eine Situation wie im folgenden Bild, welches zeigt, dass man eventuell mit den  $x^{(k)}$  nicht den am nächsten am Startpunkt liegenden Minimierer approximiert.



Andererseits sichern die Voraussetzungen von Satz 17.11 durchaus die Konvergenz der Folge  $(x^{(k)})_{k \in \mathbb{N}}$  gegen einen Minimierer – dieser kann allerdings von  $x^*$  verschieden sein.

**Korollar 17.12.** *Unter den Voraussetzungen von Satz 17.11 konvergiert  $(x^{(k)})_{k \in \mathbb{N}}$  gegen einen Minimierer  $y^*$  von  $f$ .*

*Beweis.* Aus Teil ② des vorhergehenden Beweises folgt

$$\forall x^* \in X^*, k \geq 0: \|x^{(k+1)} - x^*\| \leq \|x^{(k)} - x^*\|$$

Die Folge  $(x^{(k)})_{k \in \mathbb{N}}$  ist also beschränkt und enthält daher eine konvergente Teilfolge  $(x^{(k_j)})_{j \in \mathbb{N}}$  deren Grenzwert wir  $y^*$  nennen. Da  $f$  stetig ist und  $f(x^{(k)}) \rightarrow f^*$  nach Satz 17.11 gilt, haben wir

$$f(y^*) = \lim_{j \rightarrow \infty} f(x^{(k_j)}) = f^*$$

und damit  $y^* \in X^*$ . Nach dem obigem ist also  $(\|x^{(k+1)} - y^*\|)_{k \in \mathbb{N}}$  monoton fallend mit einer Teilfolge die gegen Null konvergiert. Das heisst aber, dass bereits die ganze Folge gegen Null konvergiert, also  $x^{(k)} \rightarrow y^*$ .  $\square$

Die folgenden Definition schließt die Existenz von mehreren Minimierern aus.

**Definition 17.13.** Eine Funktion  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  heißt *strikt konvex*, falls

$$f(\lambda x + (1 - \lambda)y) < \lambda f(x) + (1 - \lambda)f(y)$$

für alle  $\lambda \in (0, 1)$  und  $x \neq y \in \mathbb{R}^n$ .

Offenbar ist jede strikt konvexe Funktion auch konvex. Beispiele für konvexe und nicht strikt konvexe Funktionen sind z.B. der Betrag  $|\cdot|: \mathbb{R} \rightarrow \mathbb{R}$  oder jede beliebige affin-lineare Funktion. Wie angekündigt, genügt strikte Konvexität, um die Eindeutigkeit von Minimierern zu garantieren, wenn solche existieren.



**Proposition 17.14.** *Sei  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  strikt konvex und sei  $x^* \in X^*$ . Dann ist  $x^*$  der einzige Minimierer von  $f$ .*

*Beweis.* Sei  $x \in X^*$  und  $x \neq x^*$ . Dann gilt für  $\lambda \in (0, 1)$

$$f(\lambda x + (1 - \lambda)x^*) < \lambda f(x) + (1 - \lambda)f(x^*) = \lambda f(x^*) + (1 - \lambda)f(x^*) = f(x^*).$$

Für  $\lambda \rightarrow 0$  erhalten wir Punkte  $y = \lambda x + (1 - \lambda)x^* \rightarrow x^*$ , also beliebig nah an  $x^*$ , mit  $f(y) < f(x^*)$  im Widerspruch dazu, dass  $x^* \in X^*$  ein Minimierer ist.  $\square$

Wir betrachten das folgende Beispiel, welches sowohl für die bereits skizzierten Anwendungen, als auch für die folgenden theoretischen Resultate zum Gradientenverfahren, zentral ist.

**Beispiel 17.15.** Das Quadrat der euklidischen Norm  $\|\cdot\|^2: \mathbb{R}^n \rightarrow \mathbb{R}$  ist strikt konvex. In der Tat gilt für  $\lambda \in (0, 1)$  und  $x \neq y \in \mathbb{R}^n$

$$\begin{aligned} \lambda\|x\|^2 + (1 - \lambda)\|y\|^2 - \|\lambda x + (1 - \lambda)y\|^2 &= \lambda\|x\|^2 + (1 - \lambda)\|y\|^2 - \langle \lambda x + (1 - \lambda)y, \lambda x + (1 - \lambda)y \rangle \\ &= \lambda\|x\|^2 + (1 - \lambda)\|y\|^2 - \|\lambda x\|^2 - 2\langle \lambda x, (1 - \lambda)y \rangle - \|(1 - \lambda)y\|^2 \\ &= (\lambda - \lambda^2)\|x\|^2 + [(1 - \lambda) - (1 - \lambda)^2]\|y\|^2 - \lambda(1 - \lambda)2\langle x, y \rangle \\ &= \lambda(1 - \lambda)[\|x\|^2 + \|y\|^2 - 2\langle x, y \rangle] \\ &= \lambda(1 - \lambda)\|x - y\|^2 > 0 \end{aligned}$$

und Umstellen liefert genau die nötige Abschätzung.

In Kombination mit dem folgende Lemma ergibt sich aus Beispiel 17.15 eine Klasse von (strikt) konvexen Funktionen.

**Lemma 17.16.** *Sei  $g: \mathbb{R}^n \rightarrow \mathbb{R}$  konvex und  $h: \mathbb{R}^m \rightarrow \mathbb{R}^n$  affin-linear. Dann ist  $g \circ h$  konvex. Ist  $g$  strikt konvex und  $h$  affin-linear und injektiv, dann ist  $g \circ h$  sogar strikt konvex.*

*Beweis.* ① Sei  $h$  gegeben durch  $h(x) = Ax + b$  mit  $A \in \mathbb{R}^{n \times m}$  und  $b \in \mathbb{R}^n$  und sei  $\lambda \in (0, 1)$  und  $x \neq y \in \mathbb{R}^m$ . Dann gilt wegen

$$\begin{aligned} (g \circ h)(\lambda x + (1 - \lambda)y) &= g(\lambda Ax + (1 - \lambda)Ay + b) \\ &= g(\lambda Ax + \lambda b + (1 - \lambda)Ay + (1 - \lambda)b) \\ &= g(\lambda(Ax + b) + (1 - \lambda)(Ay + b)) \\ &\leq \lambda g(h(x)) + (1 - \lambda)g(h(y)) \end{aligned}$$

wobei wir in der zweiten Gleichung  $b = \lambda b + (1 - \lambda)b$  und in der Ungleichung am Ende die Konvexität von  $g$  benutzt haben.

② Wenn  $h$  injektiv ist, dann ist oben  $Ax + b \neq Ay + b$  und daher die Abschätzung strikt wegen der strikten Konvexität von  $g$ .  $\square$

Wir werden nun die strikte Konvexität weiter verschärfen und zeigen, dass dann die Existenz und Eindeutigkeit eines Minimierers  $x^*$  automatisch folgt, dass  $x^{(k)} \rightarrow x^*$  gilt und für die Genauigkeit  $|f(x^{(k)}) - f^*| < \varepsilon$  nur  $\mathcal{O}(\log \frac{1}{\varepsilon})$ -viele Iterationen nötig sind.

**Definition 17.17.** Sei  $\mu > 0$ . Die Funktion  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  heißt  $\mu$ -konvex, falls

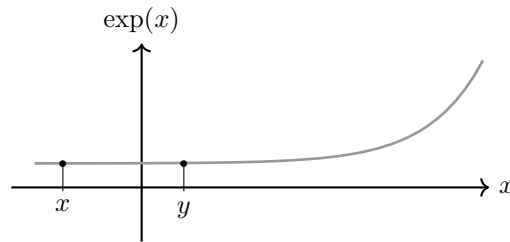
$$f(\lambda x + (1 - \lambda)y) + \lambda(1 - \lambda)\frac{\mu}{2}\|x - y\|^2 \leq \lambda f(x) + (1 - \lambda)f(y)$$

gilt für alle  $x \neq y \in \mathbb{R}^n$  und  $\lambda \in (0, 1)$ .

Beispiele für  $\mu$ -konvexe Funktionen werden wir weiter unten geben. Wir bemerken jetzt erstmal, dass gilt

$$\mu\text{-konvex} \begin{array}{c} \xRightarrow{\quad} \\ \nleftarrow{\quad} \end{array} f \text{ strikt konvex} \begin{array}{c} \xRightarrow{\quad} \\ \nleftarrow{\quad} \end{array} f \text{ konvex}.$$

Ein Beispiel für die erste durchgestrichene Implikation ist die Exponentialfunktion: Für  $x \rightarrow -\infty$ ,  $y := x + 1$  und  $\mu > 0$  beliebig gehen in Definition 17.17 die rechte Seite und der erste Summand links gegen Null, während der zweite Summand auf der linken Seite echt positiv und unabhängig von  $x$  ist.



Als nächstes charakterisieren wir  $\mu$ -konvexe Funktionen wie folgt.

**Lemma 17.18.** Eine Funktion  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  ist  $\mu$ -konvex mit  $\mu > 0$  genau dann wenn die Hilfsfunktion  $\varphi: \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $\varphi(x) := f(x) - \frac{\mu}{2}\|x\|^2$  konvex ist. Hierbei bezeichnet  $\|\cdot\|$  die euklidische Norm auf  $\mathbb{R}^n$ .

*Beweis.* Wir bemerken zunächst, dass Multiplikation der Gleichung aus Beispiel 17.15 mit  $\mu/2$  auf die Gleichung

$$\lambda\frac{\mu}{2}\|x\|^2 + (1 - \lambda)\frac{\mu}{2}\|y\|^2 - \frac{\mu}{2}\|\lambda x + (1 - \lambda)y\|^2 = \frac{\mu}{2}\lambda(1 - \lambda)\|x - y\|^2$$

führt, die für alle  $x, y \in \mathbb{R}^n$  und  $\lambda \in [0, 1]$  gilt. Einsetzen der Formel für  $\varphi$  in die Definition von Konvexität zeigt, dass  $\varphi$  konvex ist genau dann wenn für alle  $x \neq y$  und  $\lambda \in (0, 1)$  die Abschätzung

$$f(\lambda x + (1 - \lambda)y) - \frac{\mu}{2}\|\lambda x + (1 - \lambda)y\|^2 \leq \lambda(f(x) + \frac{\mu}{2}\|x\|^2) + (1 - \lambda)(f(y) - \frac{\mu}{2}\|y\|^2)$$

gilt. Sammelt man hier nun alle Terme die eine Norm enthalten auf der linken Seite, so lassen sich diese entsprechend der ersten Gleichung des aktuellen Beweises ersetzen. Dies führt genau auf die Bedingung für  $\mu$ -Konvexität von  $f$ .  $\square$

Als Konsequenz erhalten wir, dass das Quadrat der euklidischen Norm  $\|\cdot\|^2: \mathbb{R}^n \rightarrow \mathbb{R}$  nicht nur strikt konvex ist, wie in Beispiel 17.15 angegeben, sondern sogar  $\mu$ -konvex für  $\mu \in (0, 2]$ . Weiter können wir nun zeigen, dass  $\mu$ -Konvexität die Existenz eines (dann wegen der strikten Konvexität eindeutigen) Minimierers impliziert.

**Proposition 17.19.** *Sei  $f: \mathbb{R}^n \rightarrow \mathbb{R}$   $\mu$ -konvex mit  $\mu > 0$  und differenzierbar. Dann existiert ein eindeutig bestimmter Minimierer  $x^* = \operatorname{argmin}_{x \in \mathbb{R}^n} f(x)$ .*

*Beweis.* Wir zeigen, dass  $f(y) \rightarrow \infty$  gilt für  $\|y\| \rightarrow \infty$ . In der Tat ist nach Lemma 17.18 die Hilfsfunktion  $\varphi = f - \frac{\mu}{2} \|\cdot\|^2$  konvex. Für  $x, y \in \mathbb{R}^n$  gilt daher mit Lemma 17.4

$$\begin{aligned} f(y) - \frac{\mu}{2} \|y\|^2 = \varphi(y) &\geq \varphi(x) + \langle \nabla \varphi(x), y - x \rangle \\ &= f(x) - \frac{\mu}{2} \|x\|^2 + \langle \nabla f(x) - \frac{\mu}{2} \cdot 2x, y - x \rangle. \end{aligned}$$

Wir fixieren nun  $x \in \mathbb{R}^d$ , setzen  $\alpha := \nabla f(x) - \mu x \in \mathbb{R}^n$ , addieren auf beiden Seiten  $\frac{\mu}{2} \|y\|^2$  und erhalten

$$\begin{aligned} f(y) &\geq \underbrace{f(x) - \frac{\mu}{2} \|x\|^2 - \langle \alpha, x \rangle + \langle \alpha, y \rangle}_{=: K} + \frac{\mu}{2} \|y\|^2 \\ &\geq \underbrace{K}_{\substack{\uparrow \\ \text{CSB-Ungl.}}} + \underbrace{\left(\frac{\mu}{2} \|y\| - \|\alpha\|\right)}_{\substack{> 0 \text{ für } \|y\| \\ \text{groß genug}}} \|y\| \xrightarrow{\|y\| \rightarrow \infty} \infty. \end{aligned}$$

Da  $f$  insbesondere stetig ist, folgt die Existenz eines Minimierers aus dem Extremalsatz. Die Eindeutigkeit haben wir bereits in Proposition 17.14 unter schwächeren Bedingungen gezeigt.  $\square$

Das nächste Lemma ist eine Charakterisierung von  $\mu$ -Konvexität für differenzierbare Funktionen. Die entsprechende Bedingung wird in der Literatur auch manchmal als Definition für  $\mu$ -Konvexität verwendet.

**Lemma 17.20.** *Sei  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  differenzierbar. Dann ist  $f$   $\mu$ -konvex mit  $\mu > 0$  genau dann wenn  $f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\mu}{2} \|y - x\|^2$  für alle  $x, y \in \mathbb{R}^d$  gilt.*

*Beweis.* Nach Lemma 17.18 ist  $f$   $\mu$ -konvex genau dann wenn  $\varphi = f - \frac{\mu}{2} \|\cdot\|^2$  konvex ist. Dies ist nach Lemma 17.4 äquivalent zu

$$\forall x, y \in \mathbb{R}^n: \varphi(y) \geq \varphi(x) + \langle \nabla \varphi(x), y - x \rangle.$$

Einsetzen von  $\varphi$  zeigt, dass die Ungleichung zu

$$f(y) - \frac{\mu}{2} \|y\|^2 \geq f(x) - \frac{\mu}{2} \|x\|^2 + \langle \nabla f(x), y - x \rangle - \mu \langle x, y - x \rangle$$

äquivalent ist und damit zu

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\mu}{2} \|y\|^2 - \frac{\mu}{2} \|x\|^2 + \mu \|x\|^2 - \mu \langle x, y \rangle$$

$$\begin{aligned}
&= f(x) + \langle \nabla f(x), y - x \rangle + \frac{\mu}{2} [\|y\|^2 - 2\langle x, y \rangle + \|x\|^2] \\
&= f(x) + \langle \nabla f(x), y - x \rangle + \frac{\mu}{2} \|x - y\|^2
\end{aligned}$$

wie behauptet.  $\square$

Ein Vergleich mit der Abschätzung in Lemma 17.4 zeigt nochmal, dass  $\mu$ -Konvexität (strikte) Konvexität verschärft, indem auf der kleineren Seite der Konvexitätsabschätzung ein für  $x \neq y$  stets echt positiver Term addiert wird.

**Bemerkung 17.21.** Vertauschen von  $x$  und  $y$  in Lemma 17.20 zeigt, dass  $f: \mathbb{R}^n \rightarrow \mathbb{R}$   $\mu$ -konvex mit  $\mu > 0$  ist genau dann wenn

$$\forall x, y \in \mathbb{R}^n: f(x) - f(y) - \langle \nabla f(y), x - y \rangle \geq \frac{\mu}{2} \|x - y\|^2$$

gilt. Diese Abschätzung sieht derjenigen in Lemma 17.8 und Folgerung 17.10 sehr ähnlich. Es ist aber zu beachten, dass wir hier *nach unten* abschätzen — und zwar mit dem gleichen Term  $\|x - y\|^2$  wie bei der  $L$ -Glattheit, jedoch mit einer anderen multiplikativen Konstante. Kombiniert man Folgerung 17.10 mit Lemma 17.20, so sieht man leicht, dass eine differenzierbare Funktion  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  genau dann gleichzeitig  $\mu$ -konvex und  $L$ -glatt ist, wenn

$$\forall x, y \in \mathbb{R}^n: \frac{\mu}{2} \|x - y\|^2 \leq f(x) - f(y) - \langle \nabla f(y), x - y \rangle \leq \frac{L}{2} \|x - y\|^2$$

gilt. Hieraus folgt auch sofort, dass in diesem Fall notwendigerweise  $L \geq \mu$  gelten muss.

Das folgende Lemma ist unsere letzte Vorbereitung, bevor wir unseren finalen Satz zum Gradientenverfahren beweisen können.

**Lemma 17.22.** Sei  $f: \mathbb{R}^n \rightarrow \mathbb{R}$   $\mu$ -konvex und  $L$ -glatt mit  $\mu, L > 0$ . Dann gilt  $\langle \nabla f(y) - \nabla f(x), y - x \rangle \geq \frac{\mu L}{\mu + L} \|x - y\|^2 + \frac{1}{\mu + L} \|\nabla f(x) - \nabla f(y)\|^2$  für alle  $x, y \in \mathbb{R}^n$ .

*Beweis.* Seien  $x, y \in \mathbb{R}^n$  beliebig und sei  $\varphi = f - \frac{\mu}{2} \|\cdot\|^2$  die bereits in vorhergehenden Beweisen verwendete Hilfsfunktion. Wir benutzen zunächst dieselben Argumente und Rechnungen wie im Beweis des vorhergehenden Lemmas, aber nutzen dann die  $L$ -Glattheit aus und schätzen weiter nach oben ab. Es ergibt sich

$$\begin{aligned}
0 &\leq \varphi(x) - \varphi(y) - \langle \nabla \varphi(y), x - y \rangle \\
&\quad \uparrow \text{Lemmas 17.18 u. 17.4} \\
&= f(x) - f(y) - \langle \nabla f(y), x - y \rangle - \frac{\mu}{2} \|x - y\|^2 \\
&\quad \uparrow \text{wie in Lem. 17.20} \\
&\leq \frac{L}{2} \|x - y\|^2 - \frac{\mu}{2} \|x - y\|^2 \\
&\quad \uparrow \text{Lem. 17.8} \\
&\leq (L - \mu) \|x - y\|^2.
\end{aligned}$$

① Wir behandeln nun zuerst den Fall, dass  $L = \mu$  ist. Oben werden dann alle Ungleichungen zu Gleichungen und es folgt  $f(x) = f(y) + \langle \nabla f(y), x - y \rangle + \frac{\mu}{2} \|x - y\|^2$  für alle  $x, y \in \mathbb{R}^n$ . Fixiert man  $y$  und differenziert nach  $x$ , so erhält man  $\nabla f(x) = \nabla f(y) + \mu(x - y)$  und durch Umstellen

$$\forall x, y \in \mathbb{R}^n: \nabla f(x) - \nabla f(y) = \mu(x - y).$$

Setzt man  $L = \mu$  in die Aussage des Lemmas ein, so sieht man, dass sich die Behauptung zu

$$\forall x, y \in \mathbb{R}^n: \langle \nabla f(y) - \nabla f(x), y - x \rangle \geq \frac{\mu}{2} \|x - y\|^2 + \frac{1}{2\mu} \|\nabla f(x) - \nabla f(y)\|^2$$

reduziert. Einsetzen von  $\nabla f(x) - \nabla f(y)$  zeigt, dass auf beiden Seiten  $\mu \|x - y\|^2$  herauskommt.

② Sei jetzt  $L > \mu$ . Dann besagt unsere anfängliche Ungleichungskette zusammen mit Folgerung 17.10, angewandt auf die per Konstruktion differenzierbare Funktion  $\varphi$ , dass diese  $(L - \mu)$ -glatt ist. Zweimalige Anwendung (einmal mit  $x$  und  $y$  vertauscht) von Lemma 17.9 liefert die Ungleichungen

$$\begin{aligned} \varphi(y) - \varphi(x) &\leq \langle \nabla \varphi(x), y - x \rangle - \frac{1}{2(L - \mu)} \|\nabla \varphi(y) - \nabla \varphi(x)\|^2, \\ \varphi(x) - \varphi(y) &\leq \langle \nabla \varphi(y), x - y \rangle - \frac{1}{2(L - \mu)} \|\nabla \varphi(x) - \nabla \varphi(y)\|^2. \end{aligned}$$

Jetzt addieren wir beide Ungleichungen und setzen  $\varphi = f - \frac{\mu}{2} \|\cdot\|^2$  ein. Es folgt

$$\begin{aligned} 0 &\leq \langle \nabla \varphi(x) - \nabla \varphi(y), x - y \rangle - \frac{1}{L - \mu} \|\nabla \varphi(x) - \nabla \varphi(y)\|^2 \\ &= \langle \nabla f(x) - \nabla f(y) - \mu(x - y), x - y \rangle - \frac{1}{L - \mu} \|\nabla f(x) - \nabla f(y) - \mu(x - y)\|^2 \\ &= \langle \nabla f(x) - \nabla f(y), x - y \rangle - \mu \|x - y\|^2 - \frac{1}{L - \mu} \left( \|\nabla f(x) - \nabla f(y)\|^2 \right. \\ &\quad \left. - 2\mu \langle \nabla f(x) - \nabla f(y), x - y \rangle + \mu^2 \|x - y\|^2 \right) \\ &= \left( 1 + \frac{2\mu}{L - \mu} \right) \langle \nabla f(x) - \nabla f(y), x - y \rangle - \left( \mu + \frac{\mu^2}{L - \mu} \right) \|x - y\|^2 \\ &\quad - \frac{1}{L - \mu} \|\nabla f(x) - \nabla f(y)\|^2 \\ &= \frac{L + \mu}{L - \mu} \langle \nabla f(x) - \nabla f(y), x - y \rangle - \frac{L}{L - \mu} \|x - y\|^2 - \frac{1}{L - \mu} \|\nabla f(x) - \nabla f(y)\|^2. \end{aligned}$$

Umstellen und Multiplikation mit  $\frac{L - \mu}{L + \mu}$  ergibt die behauptete Ungleichung.  $\square$

Der folgende Satz zeigt, wie sich die Aussagen, die über die rekursiv definierte Folge (17.1) gemacht werden können, signifikant verbessern lassen wenn die zu minimierende Funktion die oben definierten Eigenschaften hat.

**Satz 17.23.** *Sei  $f: \mathbb{R}^n \rightarrow \mathbb{R}$   $\mu$ -konvex und  $L$ -glatt mit  $\mu, L > 0$ . Sei  $x^{(0)} \in \mathbb{R}^n$  beliebig und  $(x^{(k)})_{k \geq 0}$  die Folge aus (17.1) mit Schrittweite  $\gamma := \frac{2}{\mu + L}$ . Dann existiert gemäß Proposition 17.19 ein eindeutig bestimmter Minimierer  $x^* = \operatorname{argmin}_{x \in \mathbb{R}^n} f(x)$ , es gilt*

$L \geq \mu$  und für  $k \geq 1$  haben wir

$$(i) \quad \|x^{(k)} - x^*\| \leq \left| \frac{L-\mu}{L+\mu} \right|^k \|x^{(0)} - x^*\| \quad \text{und}$$

$$(ii) \quad 0 \leq f(x^{(k)}) - f(x^*) \leq \frac{L}{2} \left| \frac{L-\mu}{L+\mu} \right|^{2k} \|x^{(0)} - x^*\|^2.$$

Insbesondere gilt also  $x^{(k)} \rightarrow x^*$  sowie  $f(x^{(k)}) \rightarrow f(x^*)$  für  $k \rightarrow \infty$ .

*Beweis.* Zunächst bemerken wir, dass  $L \geq \mu$  aus Bemerkung 17.21 folgt.

(i) Wir setzen die Rekursionsvorschrift (17.1) ein und rechnen dann die linke Seite aus, wobei wir den Term  $-\nabla f(x^*)$  künstlich hinzufügen um Lemma 17.22 anwenden zu können:

$$\begin{aligned} \|x^{(k+1)} - x^*\|^2 &= \|x^{(k)} - \gamma \nabla f(x^{(k)}) - x^*\|^2 \\ &= \|x^{(k)} - x^*\|^2 - 2\gamma \underbrace{\langle \nabla f(x^{(k)}) - \nabla f(x^*), x^{(k)} - x^* \rangle}_{=0 \text{ wegen Folg. 17.5}} + \gamma^2 \|\nabla f(x^{(k)})\|^2 \\ &\stackrel{\uparrow \text{Lem. 17.22}}{\leq} \|x^{(k)} - x^*\|^2 + \frac{2}{\mu+L} \cdot \left[ \frac{\mu L}{\mu+L} \|x^{(k)} - x^*\|^2 + \right. \\ &\quad \left. \frac{1}{\mu+L} \|\nabla f(x^{(k)}) - \underbrace{\nabla f(x^*)}_{=0}\|^2 \right] + \frac{1}{(\mu+L)^2} \|\nabla f(x^{(k)})\|^2 \\ &= \left(1 - \frac{4\mu L}{(\mu+L)^2}\right) \|x^{(k)} - x^*\|^2 + \underbrace{\left(-\frac{4}{(\mu+L)^2} + \frac{1}{(\mu+L)^2}\right) \|\nabla f(x^{(k)})\|^2}_{\leq 0} \\ &\leq \frac{\mu^2 + 2\mu L + L^2 - 4\mu L}{(\mu+L)^2} \|x^{(k)} - x^*\|^2 = \left(\frac{L-\mu}{L+\mu}\right)^2 \|x^{(k)} - x^*\|^2. \end{aligned}$$

Wurzelziehen und Iteration liefert die behauptete Ungleichung.

(ii) Nach der Vorarbeit von oben gilt nun

$$\begin{aligned} f(x^{(k)}) - f(x^*) &\stackrel{\uparrow \text{Lem. 17.8}}{\leq} \langle \nabla f(x^*), x^{(k)} - x^* \rangle + \frac{L}{2} \|x^{(k)} - x^*\|^2 \\ &= \frac{L}{2} \|x^{(k)} - x^*\|^2 \stackrel{\uparrow \nabla f(x^*)=0}{\leq} \left(\frac{L-\mu}{L+\mu}\right)^{2k} \|x^{(k)} - x^*\|^2. \end{aligned}$$

Die Konvergenz von Folge und Bildfolge erhält man, da  $\frac{L-\mu}{L+\mu} \in [0, 1)$  liegt.  $\square$

Setzen wir  $\kappa := L/\mu$ , dann ist  $\kappa \geq 1$  und es gilt

$$\left(\frac{L-\mu}{L+\mu}\right)^2 = \left(\frac{\frac{\kappa}{\mu} - 1}{\frac{\kappa}{\mu} + 1}\right)^2 = \left(\frac{\kappa - 1}{\kappa + 1}\right)^2 = \left(1 - \frac{2}{\kappa + 1}\right)^2 \leq e^{-\frac{4}{\kappa + 1}}.$$

Unter den Voraussetzungen von Satz 17.23 erhalten wir also

$$0 \leq f(x^{(k)}) - f(x^*) \leq \frac{L}{2} \|x^{(0)} - x^*\|^2 e^{-\frac{4}{\kappa+1}k} =: C e^{-ck}$$

und weil  $Ce^{-ck} < \varepsilon$  genau dann gilt wenn  $k > \frac{1}{C}(\log c + \log \frac{1}{\varepsilon})$  ist, sehen wir, dass die Anzahl Iterationen von (17.1), die durchgeführt werden müssen um die Genauigkeit  $|f(x^{(k)}) - f(x^*)| < \varepsilon$  zu erreichen, durch  $K \cdot \log \frac{1}{\varepsilon}$  beschränkt sind, wobei  $K$  eine von  $\varepsilon$  unabhängige Konstante ist, vergleiche mit der deutlich langsameren Konvergenz in Satz 17.11.

**Bemerkung 17.24.** In vielen Anwendungen minimiert man eine Funktion der Form

$$f: \mathbb{R}^n \rightarrow \mathbb{R}, f(x) = \sum_{j=1}^N f_j(x)$$

mit sehr vielen Funktionen  $f_j$ , vergleiche Kapitel 16.2. Aufgrund der speziellen Form ergibt sich  $\nabla f$  als die Summe der  $\nabla f_j$ . Zur Beschleunigung von (17.1) kann man zum *Stochastischen Gradientenverfahren* übergehen und

$$x^{(k+1)} := x^{(k)} - \gamma \sum_{i=1}^m \nabla f_{j_i}(x^{(k)})$$

setzen, wobei die  $j_1, \dots, j_m \in \{1, \dots, N\}$  in jedem Schritt zufällig gewählt werden.

## Referenzen

Unsere Darstellung basiert hauptsächlich auf [Bub13, Bub15, Lot22, Bec17, Nes18, BV04, NW06]. Aufgabe 17.10 ist ein Variation der sogenannten *Projizierten Gradientenmethode*. Den Inhalt der Bemerkung findet man z.B. in [Roc70, Corollary 25.5.1].

Alternativ zum Zugang in diesem Kapitel setzen viele Autoren die zu minimierende Funktion  $f$  von vorn herein als zweimal differenzierbar voraus.  $L$ -Glattheit und  $\mu$ -Konvexität können dann durch Abschätzungen von  $\langle \nabla^2 f(\cdot)h, h \rangle$  nach oben bzw. unten beschrieben werden. Für den genauen Zusammenhang zu den oben benutzten Begriffen verweisen wir auf [Nes18].

Dieses Kapitel stellt nur einen sehr kleinen Ausschnitt zum Thema Gradientenverfahren dar. Für weitere Algorithmen, z.B. zur Beschleunigung von (17.1), Strategien zur Schrittwahl, klassische (Gegen-)Beispiele, Verallgemeinerungen auf nicht differenzierbare Funktionen mittels Subgradienten, die Behandlung von Funktion  $f: A \subset \mathbb{R}^n \rightarrow \mathbb{R}$ , und auch die Frage wie Gradienten numerisch berechnet werden können, verweisen wir auf die oben genannte Literatur.

Auch zum Stochastischen Gradientenverfahren gibt es theoretische Resultate, die dann den Erwartungswert von  $\|x^{(k)} - x^*\|$  abschätzen. Wir verweisen auf [Wen17] und die darin zitierte Originalliteratur.

## Aufgaben

**Aufgabe 17.1.** Sei  $f: \mathbb{R}^2 \rightarrow \mathbb{R}$ ,  $f(x) = e^x$ . Durch welche Wahl der  $\gamma_k$  können Sie erreichen, dass  $(f(x^{(k)}))_{k \in \mathbb{N}}$  mit  $x^{(k+1)} = x^{(k)} - \gamma_k f'(x^{(k)})$  gegen  $\inf_{x \in \mathbb{R}} f(x)$  konvergiert?

**Aufgabe 17.2.** Sei  $f: \mathbb{R}^2 \rightarrow \mathbb{R}$ ,  $f(x) = x^2$ . Untersuchen Sie die Folge  $(x^{(k)})_{k \in \mathbb{N}}$  mit  $x^{(k+1)} = x^{(k)} - 2^{-k} f'(x^{(k)})$  auf Konvergenz. Überrascht Sie das Ergebnis?

**Aufgabe 17.3.** Zeigen Sie für die Funktion  $f: \mathbb{R} \rightarrow \mathbb{R}$ ,  $f(x) = 3x^4 + 4x^3 - 36x^2 + 42$ , durch Ausführung des Gradientenverfahrens mit konstanter Schrittweite, dass selbiges, bei ungünstig gewähltem Startwert gegen eine lokale Minimalstelle konvergiert die kein (globaler) Minimierer ist.

**Aufgabe 17.4.** Zeigen Sie, dass die Ableitung der folgenden Funktion für  $x \searrow 0$  exponentiell schnell fällt. Die Gradientenmethode kommt also, wenn Sie erstmal in der Nähe der Null ist, nur noch sehr langsam voran.

$$f: \mathbb{R} \rightarrow \mathbb{R}, f(x) = \begin{cases} e^{-1/x} & \text{falls } x \in (0, \frac{1}{4}], \\ (x + 8e^{-4} - \frac{1}{4})^2 - 64e^{-8} + e^{-4} & \text{falls } x \in (\frac{1}{4}, \infty), \\ 0 & \text{sonst.} \end{cases}$$

Glauben Sie, dass man dennoch durch (17.1) mit konstanter Schrittweite eine Folge bekommen kann, die unabhängig vom Startwert gegen einen Minimierer konvergiert?

**Aufgabe 17.5.** Zeigen Sie, dass  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $f(x) = \frac{1}{2}\langle x, Ax \rangle + \langle b, x \rangle + c$  mit  $A \in \mathbb{R}^{n \times n}$  positiv definit, sowie  $b \in \mathbb{R}^n$  und  $c \in \mathbb{R}$  beliebig,  $\mu$ -konvex und  $L$ -glatt ist und stellen Sie fest für welche  $\mu$  und  $L$  dies gilt.

**Aufgabe 17.6.** Sei  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  konvex und  $L$ -glatt mit  $L > 0$ . Sei  $x^* \in X^*$  und  $x^{(0)}$  beliebig,  $(x^{(k)})_{k \in \mathbb{N}}$  sei die Folge aus (17.1) mit konstanter Schrittweite  $\gamma \in (0, 1/L)$ . Zeigen Sie, dass dann  $f(x^{(k)}) - f(x^*) \leq \frac{2}{k\gamma} \|x^{(0)} - x^*\|^2$  gilt für  $k \geq 0$ .

**Aufgabe 17.7.** In Lemma 17.22 ist die Hilfsfunktion  $\varphi$  in Wirklichkeit sogar  $\frac{L-\mu}{2}$ -glatt, da wir in der allerersten Abschätzung einen Faktor  $1/2$  ‘verschenkt’ haben. Gehen Sie den Beweis noch einmal durch und finden Sie heraus, auf welche Weise sich die Abschätzung in Satz 17.23 verbessert, wenn man das  $1/2$  mitnimmt.

**Aufgabe 17.8.** Bei der Implementierung des Gradientenverfahrens kann man einerseits ein Abbruchkriterium einbauen, sodass der Algorithmus solange läuft, bis sich der Abstand zwischen  $x^{(k+1)}$  und  $x^{(k)}$  nur noch wenig ändert. Andererseits kann es sinnvoll sein, die Anzahl der Iterationen von vornherein zu beschränken. Dies gilt insbesondere, wenn man Funktionen behandelt, bei denen keine Schranke für die Laufzeit bekannt ist.

```

1: function GD( $f, (\gamma_k)_{k \in \mathbb{N}}, K, \varepsilon$ )
2:    $x^{(0)} \leftarrow$  zufälliger Punkt in  $\mathbb{R}^d$ 
3:   for  $k \leftarrow 0$  to  $K$  do
4:      $x^{(k+1)} \leftarrow x^{(k)} - \gamma \nabla f(x^{(k)})$ 
5:     if  $\|x^{(k+1)} - x^{(k)}\| < \varepsilon$  then
6:       break
7:   return  $x^{(k+1)}$ 
```

Implementieren Sie das Gradientenverfahren entsprechend dem angegebenen Pseudocode für die Funktion  $f: \mathbb{R}^2 \rightarrow \mathbb{R}$ ,  $f(x, y) = x^2 + xy + y^2$ , wobei Sie den Gradienten von Hand analytisch ausrechnen. Testen Sie es für verschiedene konstante Schrittweiten  $\gamma_k \equiv \gamma > 0$  und Nullfolgen  $\gamma_k \rightarrow 0$ . Berechnen Sie manuell das eindeutig bestimmte Minimum und vergleichen Sie.

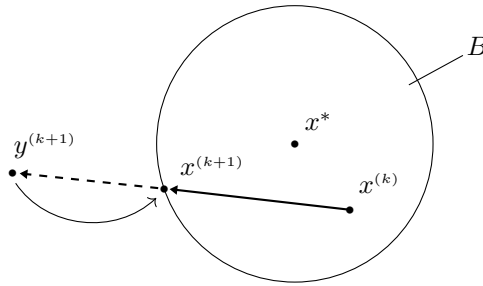
**Aufgabe 17.9.** Zeigen Sie, dass die Funktion  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $f(x) = \|x\|^{1+\alpha}$ ,  $\alpha \in (0, 1)$ , konvex und stetig differenzierbar ist, aber für jedes  $L > 0$  nicht  $L$ -glatt ist. Führen Sie dann für  $n = 1$  und  $\alpha = 1/2$  das Gradientenverfahren mit  $\gamma_k \equiv 1$  aus, und zwar mit verschiedenen Startwerten und verschieden vielen Iterationen. Stellen Sie eine Vermutung zum Konvergenzverhalten der Folge  $(x^{(k)})_{k \in \mathbb{N}}$  auf und beweisen Sie diese.



**Aufgabe 17.10.** Sei  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  stetig differenzierbar und konvex, sei  $x^*$  ein Minimierer und sei  $x^{(0)} \in \mathbb{R}^n$  beliebig. In dieser Aufgabe zeigen wir, dass stets eine Folge von Schrittweiten  $(\gamma_k)_{k \in \mathbb{N}} \subseteq [0, \infty)$  existiert, sodass die durch (17.1) definierte Folge  $(f(x^{(k)}))_{k \in \mathbb{N}}$  gegen  $f^*$  konvergiert. Dazu sei  $\eta_0 := 1$  und  $\eta_k := 1/\sqrt{k}$  für  $k \geq 1$ , sowie  $B := \overline{B}_R(x^*)$  mit  $R > \|x^* - x^{(0)}\|$ . Jetzt gehen wir wie folgt vor: Ist  $x^{(k)}$  gegeben, so setzen wir

$$y^{(k+1)} := x^{(k)} - \eta_k \nabla f(x^{(k)}).$$

Falls dann  $y^{(k+1)} \in B$  gilt, so setzen wir  $\gamma_k := \eta_k$ , also  $x^{(k+1)} = y^{(k+1)}$ . Falls nicht, so wählen wir  $\gamma_k < \eta_k$  derart, dass  $x^{(k+1)} = x^{(k)} - \gamma_k \nabla f(x^{(k)}) \in \partial B$  gilt.



- (i) Zeigen Sie mittels Proposition 17.4 und via der Identität  $2\langle a, b \rangle = \|a\|^2 + \|b\|^2 - \|a - b\|^2$  und für  $k \geq 0$  die Abschätzung

$$f(x^{(k)}) - f(x^*) \leq \frac{1}{2\eta_k} (\|x^{(k)} - x^*\|^2 + \|y^{(k)} - y^{(k+1)}\|^2 - \|x^{(k+1)} - x^*\|^2).$$

- (ii) Zeigen Sie, dass die Reihe  $\sum_{k=1}^{\infty} \frac{1}{k} (f(x^{(k)}) - f(x^*))$  konvergiert.

- (iii) Folgern Sie mit einem Widerspruchsargument, dass  $f(x^{(k)}) \rightarrow f^*$  gilt.

*Bemerkung:* Man kann zeigen, dass eine differenzierbare und konvexe Funktion  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  automatisch stetig differenzierbar ist. Damit zeigt die Aufgabe, dass die Bedingungen mit denen wir das Kapitel begonnen haben, in der Tat ausreichen, um die Konvergenz  $f(x^{(k)}) \rightarrow f^*$  zu garantieren.

## Anhang A

# Ausgewählte Resultate der Wahrscheinlichkeitstheorie

Wir präsentieren in diesem Anhang zunächst eine Zusammenfassung derjenigen Resultate aus der Wahrscheinlichkeitstheorie, die wir im Haupttext benutzen. Insbesondere fixieren wir unsere Notation zu Zufallsvariablen, Erwartungswerten, Varianzen usw. Im zweiten Teil beweisen wir die in den Kapiteln 8–12 benutzte Charakterisierung normalverteilter Zufallsvektoren, sowie die Rechenregeln für Linearkombinationen, die wir in Fakt 10.1 ohne Beweis notiert hatten.

### A.1 Wahrscheinlichkeit

Für Notation aus der Maßtheorie verweisen wir auf die am Kapitelende empfohlene Literatur. Die folgenden Definitionen sind allesamt Standard, auch wenn der Erwartungswert in Grundvorlesungen oftmals anders definiert wird.

**Definition A.1.** Sei  $\Omega \neq \emptyset$  eine nichtleere Menge. Ein Tripel  $(\Omega, \Sigma, P)$  heißt *Wahrscheinlichkeitsraum*, falls die folgenden Bedingungen erfüllt sind.

(W1)  $\Sigma$  ist eine  $\sigma$ -Algebra über  $\Omega$ , d.h.  $\Sigma \subseteq \mathcal{P}(\Omega)$ , sodass

- (i)  $\Omega \in \Sigma$ ,
- (ii)  $A \in \Sigma \implies \Omega \setminus A \in \Sigma$ ,
- (iii)  $(A_n)_{n \in \mathbb{N}} \subseteq \Sigma \implies \bigcup_{n=1}^{\infty} A_n \in \Sigma$ .

(W2)  $P: \Sigma \rightarrow [0, 1]$  ist ein *Wahrscheinlichkeitsmaß*, d.h.

- (iv)  $P(\Omega) = 1$ ,
- (v)  $P$  ist  $\sigma$ -additiv, d.h. für  $(A_n)_{n \in \mathbb{N}} \subseteq \Sigma$  mit paarweise disjunkten  $A_n$  gilt

$$P\left(\bigcup_{n=1}^{\infty} A_n\right) = \sum_{n=1}^{\infty} P(A_n).$$

Ist  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum, so heißt  $\Omega$  der *Ergebnisraum* und die Elemente von  $\Sigma$  heißen *Ereignisse*.

**Definition A.2.** Sei  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum.

- (i) Zwei Ereignisse  $A, B \in \Sigma$  heißen *unabhängig*, falls  $P(A \cap B) = P(A) \cdot P(B)$  gilt.
- (ii) Für  $A, B \in \Sigma$  heißt  $P(A|B) := \frac{P(A \cap B)}{P(B)}$  die *bedingte Wahrscheinlichkeit von A gegeben B*.

**Proposition A.3.** (Satz von der totalen Wahrscheinlichkeit) Sei  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum und seien  $A, B \in \Sigma$ . Dann gilt

$$P(A) = P(A|B) \cdot P(B) + P(A|B^c) \cdot P(B^c).$$

**Definition A.4.** Sei  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum. Eine messbare Abbildung  $X: \Omega \rightarrow \mathbb{R}^d$  heißt *Zufallsvektor* bzw. *Zufallsvariable* falls  $d = 1$ . Hierbei ist  $\mathbb{R}^d$  mit der  $\sigma$ -Algebra  $\mathcal{B}^d$  der Borelmengen ausgestattet. Wir verwenden die folgende Notation

$$P[X \in B] := P(X^{-1}(B)) = P(\{\omega \in \Omega \mid X(\omega) \in B\}) \text{ für } B \in \mathcal{B}^d,$$

$$P[X \geq a] := P(\{\omega \in \Omega \mid X(\omega) \geq a\}) \text{ für } a \in \mathbb{R},$$

$$P[a \leq \|X\| \leq b] := P(\{\omega \in \Omega \mid a \leq \|X(\omega)\| \leq b\}) \text{ für } a, b \in \mathbb{R},$$

$$P[X_i \geq c_i \text{ für alle } i] := P(\{\omega \in \Omega \mid X_i(\omega) \geq c_i \text{ für alle } i\}) \text{ für } c_i \in \mathbb{R}$$

für Zufallsvariablen und -vektoren. Ist  $X$  ein Zufallsvektor und  $\omega \in \Omega$ , so nennen wir  $X(\omega)$  eine *Realisierung* desselben.

**Definition A.5.** Sei  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum,  $X, Y: \Omega \rightarrow \mathbb{R}$  seien Zufallsvariablen und  $k \geq 2$ .

- (i) Falls das folgende Integral existiert, so heißt

$$E(X) := \int_{\Omega} X \, dP$$

der *Erwartungswert* von  $X$ . Analog heißt  $E(X^k)$  das *k-te Moment* von  $X$ , falls das entsprechende Integral existiert.

- (ii) Ist  $X$  integrierbar, und damit  $E(X) \in \mathbb{R}$ , so heißen

$$V(X) := E((X - E(X))^2) \quad \text{und} \quad \sigma(X) := \sqrt{V(X)}$$

die *Varianz* von  $X$ , und falls  $V(X) < \infty$ , die *Standardabweichung* von  $X$ .

- (iii) Sind  $X$  und  $Y$  integrierbar und die dann wohldefinierten Zufallsvariablen  $X - E(X)$  und  $Y - E(Y)$  seien auch wieder integrierbar. Dann heißt

$$\text{Cov}(X, Y) := E(X - E(X)) E(Y - E(Y))$$

die *Kovarianz* von  $X$  und  $Y$ .

**Proposition A.6.** Sei  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum,  $X, Y: \Omega \rightarrow \mathbb{R}$  seien integrierbare Zufallsvariablen mit  $V(X), V(Y), \text{Cov}(X, Y) < \infty$  und  $a, b \in \mathbb{R}$ . Dann gelten

- (i)  $E(aX + bY) = aE(X) + bE(Y)$ ,
- (ii)  $X \leq Y$  punktweise  $\implies E(X) \leq E(Y)$ ,
- (iii)  $V(X) = E(X^2) - E(X)^2$ ,
- (iv)  $V(aX + b) = a^2 V(X)$ .
- (v)  $V(aX + bY) = a^2 V(X) + b^2 V(Y) + 2ab \text{Cov}(X, Y)$ .
- (vi)  $|\text{Cov}(X, Y)| \leq \sqrt{V(X)} \sqrt{V(Y)}$ .

**Definition A.7.** Sei  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum,  $X: \Omega \rightarrow \mathbb{R}^d$  ein Zufallsvektor und  $\rho: \mathbb{R}^d \rightarrow [0, \infty)$  eine meßbare Funktion. Falls gilt

$$\forall A \in \mathcal{B}^d: P[X \in A] = \int_A \rho(x) d\lambda^d(x)$$

so sagen wir, die Zufallsvariable ist  $\rho$ -verteilt. Wichtige Beispiele sind gaußverteilte Zufallsvektoren und gleichmäßig verteilte Zufallsvektoren mit

$$\rho(x) = \frac{1}{(2\pi\sigma^2)^{d/2}} e^{-\frac{\|x-\mu\|^2}{2\sigma^2}} \quad \text{bzw.} \quad \rho(x) = \frac{1}{\lambda(B)} \cdot \mathbb{1}_B(x)$$

wobei  $\mu \in \mathbb{R}^d$ ,  $\sigma > 0$  und  $B \in \mathcal{B}^d$  mit  $\lambda^d(B) \in (0, \infty)$  fest sind.

**Proposition A.8.** Sei  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum und  $X: \Omega \rightarrow \mathbb{R}^d$ , eine  $\rho$ -verteilte Zufallsvariable und  $k \geq 1$ . Dann gilt für Momente und Erwartungswert

$$E(X^k) = \int_{\mathbb{R}} x^k \rho(x) d\lambda(x).$$

In der Tat gilt die Verallgemeinerung  $E(f(X)) = \int_{\mathbb{R}} f(x) \rho(x) d\lambda(x)$  für eine große Klasse von Funktionen  $f$  und wird manchmal als das ‘Gesetz des unbewussten Statistikers’ bezeichnet.

**Satz A.9.** (Markov-Ungleichung) Es sei  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum und  $X: \Omega \rightarrow [0, \infty)$  eine Zufallsvariable. Dann gilt für  $a > 0$

$$P[X \geq a] \leq \frac{E(X)}{a}.$$

**Satz A.10.** (Tschebyscheff-Ungleichung) Sei  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum und  $X: \Omega \rightarrow \mathbb{R}$  eine Zufallsvariable, sodass  $E(X)$  und  $V(X)$  endlich sind und  $V(X)$  von Null verschieden ist. Dann gilt für  $a > 0$

$$P[|X - E(X)| \geq a] \leq \frac{V(X)}{a^2}.$$

**Definition A.11.** Sei  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum und  $X_1, \dots, X_n: \Omega \rightarrow \mathbb{R}$  seien Zufallsvariablen. Die  $X_1, \dots, X_n$  heißen *unabhängig*, falls die Gleichung

$$P\left[(X_1, \dots, X_n) \in \prod_{i=1}^n A_i\right] = \prod_{i=1}^n P[X_i \in A_i].$$

für beliebige Borelmengen  $A_i \in \mathcal{B}$  gilt. Eine Folge  $(X_i)_{i \in \mathbb{N}}$  von Zufallsvariablen heißt *unabhängig*, falls jede endliche Teilfolge unabhängig ist.

**Proposition A.12.** Sei  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum und  $X, Y: \Omega \rightarrow \mathbb{R}$  unabhängige Zufallsvariablen mit endlichem Erwartungswert und endlicher Varianz. Dann gelten

- (i)  $E(X \cdot Y) = E(X) \cdot E(Y)$ ,
- (ii)  $V(X + Y) = V(X) + V(Y)$ .

**Satz A.13.** (Klonsatz) Für  $i \in \mathbb{N}$  seien Wahrscheinlichkeitsräume  $(\Omega_i, \Sigma_i, P_i)$  und Zufallsvariablen  $Y_i: \Omega_i \rightarrow \mathbb{R}$  gegeben. Dann existiert ein Wahrscheinlichkeitsraum  $(\Omega, \Sigma, P)$  und unabhängige Zufallsvariablen  $X_1, X_2, \dots: \Omega \rightarrow \mathbb{R}$  derart dass jedes  $X_i$  die gleiche Verteilung hat wie  $Y_i$ , also  $P[X_i \in A] = P_i[Y_i \in A]$  für jede Borelmenge  $A \in \mathcal{B}$  und jedes  $i \in \mathbb{N}$  gilt.

**Satz A.14.** (Schwaches Gesetz der großen Zahl) Sei  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum und  $X: \Omega \rightarrow \mathbb{R}$  eine Zufallsvariable, sodass  $E(X)$  und  $V(X)$  endlich sind. Sei  $(X_i)_{i \in \mathbb{N}}$  eine Folge unabhängiger Kopien von  $X$ . Dann gilt für  $\varepsilon > 0$  und  $n \in \mathbb{N}$

$$P\left[\left|\frac{X_1 + \dots + X_n}{n} - E(X)\right| \geq \varepsilon\right] \leq \frac{V(X)}{n\varepsilon^2}.$$

Für  $n \rightarrow \infty$  geht also  $P[\dots]$  gegen Null. Man sagt, dass

$$\frac{X_1 + \dots + X_n}{n} \xrightarrow{n \rightarrow \infty} E(X)$$

in Wahrscheinlichkeit konvergiert, wobei  $E(X)$  auf der rechten Seite als konstante Zufallsvariable zu verstehen ist.

**Bemerkung A.15.** Die Elemente  $X_i$  eine Folge wie in Satz A.14 nennt man auch *Stichproben* der Zufallsvariable  $X$ . Beachte, dass man für *theoretische* Aussagen über die mehrfache Durchführung eines durch  $X$  beschriebenen Zufallsexperimentes Stichproben in diesem Sinne betrachtet und dann z.B. wie in Satz A.14 deren Mittelwert (als neue Zufallsvariable!) untersucht. Für *Simulationen* benutzt man Realisierungen  $X(\omega_1), \dots, X(\omega_n)$ , vgl. Definition A.4, von  $X$  und berechnet dann z.B. deren Mittelwert (als reelle Zahl!).

**Definition A.16.** Seien  $\rho, \tau \in L^1(\mathbb{R})$ . Dann existiert

$$\rho * \tau: \mathbb{R} \rightarrow \mathbb{R}, (\rho * \tau)(s) = \int_{\mathbb{R}} \rho(s - t) \tau(t) \, d\lambda(t)$$

für fast alle  $s \in \mathbb{R}$  und  $\rho * \tau \in L^1(\mathbb{R})$  heißt die *Faltung* von  $\rho$  und  $\tau$ .

**Satz A.17.** Sei  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum und  $X, Y: \Omega \rightarrow \mathbb{R}$  unabhängige  $\rho$ - bzw.  $\tau$ -verteilte Zufallsvariablen mit  $\rho, \tau \in L^1(\mathbb{R})$ . Dann ist  $X + Y: \Omega \rightarrow \mathbb{R}$  eine  $\rho * \tau$ -verteilte Zufallsvariable.

## A.2 Gaußverteilung

Im folgenden stellen wir für den Haupttext unverzichtbare Fakten über gaußverteilte Zufallsvariablen und -vektoren zusammen. Als erstes zeigen wir, dass die Bezeichnungen *Mittelwert* und *Varianz* für die Parameter einer gaußverteilten Zufallsvariable überhaupt gerechtfertigt sind.

**Satz A.18.** Sei  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum und sei  $X: \Omega \rightarrow \mathbb{R}$  eine gaußverteilte Zufallsvariable, d.h.  $X \sim \mathcal{N}(\mu, \sigma^2)$  mit  $\mu \in \mathbb{R}$  und  $\sigma > 0$ . Dann gelten

$$E(X) = \mu \quad \text{und} \quad V(X) = \sigma^2.$$

*Beweis.* In der Tat haben wir für den Erwartungswert

$$\begin{aligned} E(X) &= \int_{\mathbb{R}} x \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} d\lambda(x) \stackrel{u:=\frac{x-\mu}{\sqrt{2}\sigma}}{=} \frac{1}{\sqrt{2\pi}\sigma} \int_{\mathbb{R}} (\sqrt{2}\sigma u + \mu) e^{-u^2} \sqrt{2}\sigma d\lambda(u) \\ &= \frac{1}{\sqrt{\pi}} \left[ \underbrace{\sqrt{2}\sigma \int_{\mathbb{R}} u e^{-u^2} d\lambda(u)}_{=0} + \mu \underbrace{\int_{\mathbb{R}} e^{-u^2} d\lambda(u)}_{=\sqrt{\pi}} \right] = \mu \end{aligned}$$

wobei das erste Integral aus Symmetriegründen Null ist und man das zweite z.B. durch Quadrieren, dann Anwendung von Fubini und schließlich als uneigentliches Riemannintegral in Polarkoordinaten ausrechnen kann. Für die Varianz ergibt sich

$$\begin{aligned} V(X) &= E(X^2) - E(X)^2 = \int_{\mathbb{R}} x^2 \rho(x) d\lambda(x) - \mu^2 \\ &\stackrel{u:=\frac{x-\mu}{\sqrt{2}\sigma}}{=} \frac{1}{\sqrt{2\pi}\sigma} \int_{\mathbb{R}} (\sqrt{2}\sigma u + \mu)^2 e^{-u^2} \sqrt{2}\sigma d\lambda(u) - \mu^2 \\ &= \frac{1}{\sqrt{\pi}} \left[ \underbrace{2\sigma^2 \int_{\mathbb{R}} u^2 e^{-u^2} d\lambda(u)}_{=\sqrt{\pi}/2} + \underbrace{2\sqrt{2}\sigma\mu \int_{\mathbb{R}} u e^{-u^2} d\lambda(u)}_{=0} + \underbrace{\mu^2 \int_{\mathbb{R}} e^{-u^2} d\lambda(u)}_{=\sqrt{\pi}} \right] - \mu^2 \\ &= \sigma^2 \end{aligned}$$

wobei man das erste Integral durch partielle Integration ausrechnet und die beiden anderen wie im ersten Teil.  $\square$

Wir notieren das folgende sehr einfache Resultat.

**Proposition A.19.** Sei  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum und sei  $X: \Omega \rightarrow \mathbb{R}$  eine gaußverteilte Zufallsvariable, d.h.  $X \sim \mathcal{N}(\mu, \sigma^2)$ . Sei  $a \in \mathbb{R}$  gegeben. Dann gilt  $a + X \sim \mathcal{N}(a + \mu, \sigma^2)$ .

*Beweis.* Für  $A \in \mathcal{B}$  gilt

$$P[a + X \in A] = \frac{1}{\sqrt{2\pi}\sigma} \int_{A-a} e^{\frac{(x-\mu)^2}{2\sigma^2}} d\lambda(x) = \frac{1}{\sqrt{2\pi}\sigma} \int_A e^{\frac{(y-(a+\mu))^2}{2\sigma^2}} d\lambda(y)$$

per Substitution  $y := x + a$ .  $\square$

Es folgt die im Haupttext vielfach benutzte Charakterisierung (sphärisch) gaußverteilter Zufallsvektoren via deren Koordinatenfunktionen.

**Satz A.20.** Sei  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum und  $X: \Omega \rightarrow \mathbb{R}^d$  ein Zufallsvektor mit Koordinatenfunktionen  $X_1, \dots, X_d$ . Dann gilt  $X \sim \mathcal{N}(\mu, \sigma^2, \mathbb{R}^d)$  genau dann, wenn  $X_i \sim \mathcal{N}(\mu, \sigma^2)$  für  $i = 1, \dots, d$  gilt und die  $X_1, \dots, X_d$  unabhängig sind.

*Beweis.* Für  $A = A_1 \times \dots \times A_d \subseteq \mathbb{R}^d$  mit  $A_i \in \mathcal{B}^d$  berechnen wir

$$\begin{aligned} P[X \in A] &\stackrel{X \sim \mathcal{N}(\mu, \sigma^2, \mathbb{R}^d)}{=} \frac{1}{(2\pi\sigma^2)^{d/2}} \int_A e^{-\frac{\|x-\mu\|^2}{2\sigma^2}} d\lambda^d(x) \\ &= \frac{1}{(2\pi\sigma^2)^{d/2}} \int_A e^{\frac{(x_1-\mu_1)^2}{2\sigma^2}} \dots e^{-\frac{(x_d-\mu_d)^2}{2\sigma^2}} d\lambda(x_1, \dots, x_d) \\ &= \frac{1}{\sqrt{2\pi}\sigma} \int_{A_1} e^{\frac{(x_1-\mu_1)^2}{2\sigma^2}} d\lambda(x_1) \dots \frac{1}{\sqrt{2\pi}\sigma} \int_{A_d} e^{\frac{(x_d-\mu_d)^2}{2\sigma^2}} d\lambda(x_d) \\ &\stackrel{X_i \sim \mathcal{N}(\mu, \sigma^2)}{=} P[X_1 \in A_1] \dots P[X_d \in A_d] \end{aligned}$$

wobei die erste bzw. letzte Gleichung unter der dort notierten Voraussetzung gilt. Die behauptete Äquivalenz folgt jetzt aus der obigen Rechnung.

“ $\implies$ ” Für  $B \in \mathcal{B}$  und fixiertes  $1 \leq i \leq d$  setze oben  $A = \mathbb{R} \times \dots \times B \times \dots \times \mathbb{R}$ , wobei das  $B$  in der  $i$ -ten Koordinate steht. Dann folgt  $P[X_i \in B] = P[X \in \mathbb{R}^{i-1} \times B \times \mathbb{R}^{d-i}]$  und die für die Unabhängigkeit benötigte Gleichung.

“ $\impliedby$ ” Sei  $A = A_1 \times \dots \times A_d \subseteq \mathbb{R}^d$  ein Quader. Da die  $X_i$ ’s unabhängig sind, folgt  $P[X \in A] = P[X_1 \in A_1] \dots P[X_d \in A_d]$ . Lesen wir nun die obige Rechnung rückwärts, so erhalten wir

$$P[X \in A] = \frac{1}{(2\pi\sigma^2)^{d/2}} \int_A e^{-\frac{\|x-\mu\|^2}{2\sigma^2}} d\lambda(x)$$

für Quader  $A$  und damit dann auch für beliebige Borelmengen in  $\mathbb{R}^d$ .  $\square$

Als letztes widmen wir uns Linearkombinationen unabhängiger normalverteilter Zufallsvariablen.

**Satz A.21.** Sei  $(\Omega, \Sigma, P)$  ein Wahrscheinlichkeitsraum, seien  $X_1, \dots, X_n \rightarrow \mathbb{R}$  unabhängige Zufallsvariablen mit  $X_i \sim \mathcal{N}(0, 1)$  und seien  $\lambda_1, \dots, \lambda_n \in \mathbb{R} \setminus \{0\}$ . Dann ist

die Zufallsvariable

$$X := \sum_{i=1}^n \lambda_i X_i \sim \mathcal{N}(0, \sigma^2)$$

gaußverteilt mit Mittelwert Null und Varianz  $\sigma^2 := \lambda_1^2 + \dots + \lambda_d^2 > 0$ .

*Beweis.* Wir notieren zuerst, dass für  $\lambda \neq 0$ ,  $X \sim \mathcal{N}(0, 1)$  und  $A \in \mathcal{B}$

$$\mathbb{P}[\lambda X \in A] = \mathbb{P}[X \in \frac{1}{\lambda} A] = \int_{\frac{1}{\lambda} A} \rho(x) dx = \int_A \rho(\frac{x}{\lambda}) \frac{1}{\lambda} dx$$

gilt, wobei  $\rho(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$  ist. Damit  $\lambda X$  also  $\rho(\frac{\cdot}{\lambda}) \frac{1}{\lambda}$ -verteilt. Einsetzen von  $\rho$  zeigt

$$\rho(\frac{t}{\lambda}) \frac{1}{\lambda} = \frac{1}{\sqrt{2\pi}} e^{-(\frac{t}{\lambda})^2/2} \frac{1}{\lambda} = \frac{1}{\sqrt{2\pi}\lambda} e^{-\frac{t^2}{2\lambda^2}}$$

für  $t \in \mathbb{R}$  und damit  $\lambda X \sim \mathcal{N}(0, \lambda^2)$ . Wenden wir dies auf unsere Ausgangssituation an, so haben wir  $\lambda_i X_i \sim \mathcal{N}(0, \lambda_i^2)$  für  $i = 1, \dots, d$ . Man überlegt sich leicht, dass die Faltung von  $L^1$ -Funktionen assoziativ ist und damit dann, dass Satz A.17 auch für mehr als zwei Zufallsvariablen gilt. Um unseren Beweis zu vervollständigen, müssen wir also

$$(\rho_1 * \dots * \rho_d)(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{x^2}{2\sigma^2}} \quad \text{mit} \quad \rho_i(x) = \frac{1}{\sqrt{2\pi}\lambda_i} e^{-\frac{x^2}{2\lambda_i^2}}$$

zeigen. Wir machen dies für  $d = 2$ , kürzen  $a := \lambda_1^2$ ,  $b := \lambda_2^2$  und  $c := a + b$  ab. Damit folgt

$$\begin{aligned} (\rho_1 * \rho_2)(s) &= \frac{1}{2\pi\sqrt{ab}} \int_{\mathbb{R}} \exp\left(-\frac{(s-t)^2}{2a}\right) \exp\left(-\frac{t^2}{2b}\right) d\lambda(t) \\ &= \frac{1}{2\pi\sqrt{ab}} \int_{\mathbb{R}} \exp\left(-\frac{b(s^2 - 2st + t^2) + at^2}{2ab}\right) d\lambda(t) \\ &= \frac{1}{2\pi\sqrt{ab}} \int_{\mathbb{R}} \exp\left(-\frac{t^2(b+a) - 2stb + bs^2}{2ab}\right) d\lambda(t) \\ &= \frac{1}{2\pi c^2 \sqrt{ab/c}} \int_{\mathbb{R}} \exp\left(-\frac{t^2(b+a)/c - 2stb/c + bs^2/c}{2ab/c}\right) d\lambda(t) \\ &= \frac{1}{2\pi c^2 \sqrt{ab/c}} \int_{\mathbb{R}} \exp\left(-\frac{(t - (bs)/c)^2 - (sb/c)^2 + s^2(b/c)}{2ab/c}\right) d\lambda(t) \\ &= \frac{1}{\sqrt{2\pi}c} \exp\left(-\frac{(sb/c)^2 - s^2(b/c)}{2ab/c}\right) \frac{1}{\sqrt{2\pi(ab/c)}} \int_{\mathbb{R}} \exp\left(-\frac{(t - (bs)/c)^2}{2ab/c}\right) d\lambda(t) \\ &= \frac{1}{\sqrt{2\pi}c} \exp\left(-\frac{(sb/c)^2 c^2 - s^2(b/c)c^2}{2abc}\right) \\ &= \frac{1}{\sqrt{2\pi}c} \exp\left(-\frac{s^2(b^2 - bc)}{2abc}\right) \\ &= \frac{1}{\sqrt{2\pi}c} \exp\left(-\frac{s^2}{2c}\right) \end{aligned}$$

wie gewünscht. □



## Referenzen

Wir verweisen auf [Beh13] für eine Einführung und auf das Buch [Wen08], in welchem die Wahrscheinlichkeitslehre von Anfang an auf Maßtheorie aufgebaut wird.

# Literaturverzeichnis

- [AA84] A. Albert and J. A. Anderson. On the existence of maximum likelihood estimates in logistic regression models. *Biometrika*, 71:1–10, 1984.
- [Agg16] C. C. Aggarwal. *Recommender Systems*. Springer Cham, 2016.
- [AM12] Y. Abu-Mostafa. Lecture 14—Support Vector Machines. Youtube-Video, Caltech, [Link](#), 2012.
- [BC99] C. Burges and D. Crisp. Uniqueness of the svm solution. In *NIPS Vol. 99*, 1999.
- [Bec17] A. Beck. *First-order Methods in Optimization*, volume 25 of *MOS/SIAM Ser. Optim.* 2017.
- [Beh13] E. Behrends. *Elementare Stochastik. Ein Lernbuch – von Studierenden mitentwickelt*. Heidelberg: Springer Spektrum, 2013.
- [Ber] B. Bernstein. On the uniqueness of the SVM solution. [Link](#), abgerufen am 30. September 2023.
- [BHK20] A. Blum, J. Hopcroft, and R. Kannan. *Foundations of Data Science*. Cambridge University Press, Cambridge, 2020.
- [Bis06] C. M. Bishop. *Pattern Recognition and Machine Learning*. Inf. Sci. Stat. Springer New York, 2006.
- [Bub13] S. Bubeck. ORF523: The complexities of optimization. Blog, [Link](#), 2013.
- [Bub15] S. Bubeck. Convex optimization: algorithms and complexity. *Found. Trends Mach. Learn.*, 8(3-4):231–357, 2015.
- [BV04] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
- [Cal20] O. Calin. *Deep Learning Architectures. A Mathematical Approach*. Springer Ser. Data Sci. Springer, Cham, 2020.
- [Chu07] Fan Chung. Random walks and local cuts in graphs. *Linear Algebra and its Applications*, 423(1):22–32, 2007.
- [CL92] C. K. Chui and X. Li. Approximation by ridge functions and neural networks with one hidden layer. *J. Approx. Theory*, 70(2):131–141, 1992.
- [CLRS22] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein. *Introduction to Algorithms*. Cambridge, MIT Press, 4th edition edition, 2022.
- [Col12] M. Collins. Convergence proof for the perceptron algorithm. Lecture Notes, Columbia University, [Link](#), 2012.

- [Cyb89] G. Cybenko. Approximation by superpositions of a sigmoidal function. *Math. Control Signals Syst.*, 2(4):303–314, 1989.
- [Cyb92] G. Cybenko. Correction to: Approximation by superpositions of a sigmoidal function. *Math. Control Signals Syst.*, 5(4):455, 1992.
- [DH08] P. Deuffhard and A. Hohmann. *Numerische Mathematik 1. Eine algorithmisch orientierte Einführung*. Berlin: de Gruyter, 4th revised and extended ed. edition, 2008.
- [Don04] D. L. Donoho. High-dimensional data analysis: The curses and blessings of dimensionality. [Link](#), 2004.
- [DS07] S. Dasgupta and L. Schulman. A probabilistic analysis of EM for mixtures of separated, spherical gaussians. *The Journal of Machine Learning Research*, 8:203–226, 2007.
- [For] M. Fornasier. Foundations of Data Analysis. Aufgaben zur Vorlesung, TU München, [Link](#), abgerufen am 6. September 2020.
- [Fun06] S. Funk. Netflix update: Try this at home. Blog Post, [Link](#), 2006.
- [Giu03] E. Giusti. *Direct methods in the calculus of variations*. Singapore: World Scientific, 2003.
- [GK02] C. Geiger and C. Kanzow. *Theorie und Numerik restringierter Optimierungsaufgaben*. Springer Berlin, 2002.
- [Gow14] S. Gower. Netflix Prize and SVD. [Link](#), 2014.
- [Gui18] L. F. Guilhoto. An overview of artificial neural networks for mathematicians. [Link](#), 2018.
- [Ham] M. Hampton. SVD computation example. Lecture Notes, University of Minnesota Duluth, [Link](#), abgerufen am 30. September 2023.
- [Han19] B. Hanin. Universal function approximation by deep neural nets with bounded width and relu activations. *Mathematics*, 7(10), 2019.
- [Hor91] K. Hornik. Approximation capabilities of multilayer feedforward networks. *Neural Networks*, 4(2):251–257, 1991.
- [HSW89] K. Hornik, M. Stinchcombe, and H. White. Multilayer feedforward networks are universal approximators. *Neural Netw.*, 2(5):359–366, 1989.
- [HSW90] K. Hornik, M. Stinchcombe, and H. White. Universal approximation of an unknown mapping and its derivatives using multilayer feedforward networks. *Neural Networks*, 3(5):551–560, 1990.
- [JL84] W. B. Johnson and J. Lindenstrauss. Extensions of Lipschitz mappings into a Hilbert space. *Contemp. Math.* 26, 189–206 (1984)., 1984.
- [Kal96] D. Kalman. A singularly valuable decomposition: The SVD of a matrix. *The College Mathematics Journal*, 27(1):2–23, 1996.
- [Köp13] M. Köppen. The curse of dimensionality. In *5th Online World Conference on Soft Computing in Industrial Applications*, 2013.
- [Lef19] M. Lefkowitz. Professor’s perceptron paved the way for AI — 60 years too soon. *Cornell Chronicle*, 2019. [Link](#).

- [LL16] I. E. Leonard and J. E. Lewis. *Geometry of Convex Sets*. Hoboken, NJ: John Wiley & Sons, 2016.
- [LLPS93] M. Leshno, V. Ya. Lin, A. Pinkus, and S. Schocken. Multilayer feedforward networks with a nonpolynomial activation function can approximate any function. *Neural Networks*, 6(6):861–867, 1993.
- [Lot22] M. Lotz. Mathematics of Machine Learning. Lecture Notes, The University of Warwick, [Link](#), 2022.
- [LP93] V. Ya. Lin and A. Pinkus. Fundamentality of ridge functions. *J. Approx. Theory*, 75(3):295–311, 1993.
- [LRU12] J. Leskovec, A. Rajaraman, and J. D. Ullman. *Mining of Massive Datasets*. Cambridge University Press, Cambridge, 2012.
- [MC] MIT-CSAIL. Quadratic programming with Python and CVXOPT. Lecture Notes, MIT-CSAIL, [Link](#), abgerufen am 30. September 2023.
- [Nes18] Yu. Nesterov. *Lectures on Convex Optimization*, volume 137 of *Springer Optim. Appl.* Springer, Cham, 2nd edition edition, 2018.
- [Nie15] M. A. Nielsen. *Neural Networks and Deep Learning*. Determination Press, 2015.
- [Nos16] J. Noss. 6.034 Recitation 7: Support Vector Machines (SVMs). Youtube-Video, MIT, [Link](#), 2016.
- [NW06] J. Nocedal and S. J. Wright. *Numerical optimization*. Springer Ser. Oper. Res. Financ. Eng. Springer, New York, 2nd ed. edition, 2006.
- [Ola96] M. Olazara. A sociological study of the official history of the perceptrons controversy. *Social Studies of Science*, 26(3):611–659, 1996.
- [Pin20] I. Pinelis. Asymptotically tight concentration of norms of subgaussian random vectors with independent coordinates, as the dimension  $n \rightarrow \infty$ ? Mathoverflow Post, [Link](#), 2020.
- [Rin08] C. M. Ringel. Funktionen (GHR). Vorlesungsmanuskript, Uni Bielefeld, [Link](#), 2008.
- [Roc70] R. T. Rockafellar. *Convex Analysis*, volume 28 of *Princeton Math. Ser.* Princeton University Press, Princeton, NJ, 1970.
- [Ros59] F. Rosenblatt. The design of an intelligent automaton. *Research Trends, Cornell Aeronauticak Laboratory, inc.*, VI(2), 1959.
- [Rud87] W. Rudin. *Real and Complex Analysis*. McGraw-Hill, New York, 3rd ed. edition, 1987.
- [Rud91] W. Rudin. *Functional Analysis*. New York, NY: McGraw-Hill, 2nd ed. edition, 1991.
- [SC92] X. Sun and E. W. Cheney. The fundamentality of sets of ridge functions. *Aequationes Math.*, 44(2-3):226–235, 1992.
- [SC08] I. Steinwart and A. Christmann. *Support Vector Machines*. Inf. Sci. Stat. Springer, New York, 2008.
- [Sha08] A. Shashua. Introduction to Machine Learning. arXiv: 0904.3664v1, 2008.

- [Sha15] C. Shalizi. Lecture Notes on Modern Regression. Lecture Notes, Carnegie Mellon University, [Link](#), 2015.
- [SKKR01] B. Sarwar, G. Karypis, J. Konstan, and J. Riedl. Item-based collaborative filtering recommendation algorithms. In *WWW10, May 1–5*, Hong Kong, 2001.
- [SSBD14] S. Shalev-Shwartz and S. Ben-David. *Understanding Machine Learning. From Theory to Algorithms*. Cambridge University Press, Cambridge, 2014.
- [Ver18] R. Vershynin. *High-Dimensional Probability*, volume 47 of *Cambridge Series in Statistical and Probabilistic Mathematics*. Cambridge University Press, Cambridge, 2018.
- [Wen08] J. Wengenroth. *Wahrscheinlichkeitstheorie*. de Gruyter, Berlin, 2008.
- [Wen17] C. Wendler. Das Stochastische Gradientenverfahren. Bachelorarbeit, Universität Innsbruck, [Link](#), 2017.
- [Wer18] D. Werner. *Funktionalanalysis*. Springer-Lehrb. Springer Spektrum, Berlin, 8th revised edition edition, 2018.