

Machine Learning für das KMU

Martin Sterchi

2023-12-15

Contents

Über das Buch	5
Zielgruppe	6
Aufbau des Buchs	6
Weiterführende Literatur	8
Lizenz	10
Kontakt	10
1 Einführung	11
1.1 Was ist Machine Learning?	11
1.2 Wann macht es Sinn ML einzusetzen?	13
1.3 Anwendungsfälle von ML	15
1.4 Supervised vs. Unsupervised Learning	16
1.5 Regression vs. Klassifikation	18
1.6 Parametrische vs. nicht-parametrische Modelle	18
1.7 Machine Learning Pipeline	20
2 Mathematik- und Statistik-Grundlagen	21
2.1 Funktionen	21
2.2 Integral- und Differentialrechnung	29
2.3 Lineare Algebra	29
2.4 Wahrscheinlichkeitsrechnung	29
2.5 Verteilungen	30
3 Einführung in das Programmieren mit R	31

4	Lineare Regression	33
5	Lineare Klassifikation	35
6	Machine Learning Pipeline	37
7	Decision Trees	39
8	Ensembles	41
9	Support Vector Machines	43
10	Artificial Neural Networks	45
11	Convolutional Neural Networks	47
12	Recurrent Neural Networks	49
13	Generative AI	51

Über das Buch

Die Motivation für dieses Buch kam aus der Erkenntnis, dass viele kleine und mittelgrosse Unternehmen (KMU) in der Schweiz zwar über grosse Datenmengen verfügen, aber nicht das nötige Knowhow haben, um die Daten zu analysieren und für die Optimierung von Entscheidungsprozessen zu nutzen. Mit diesem Buch möchte ich einen kleinen Beitrag leisten, den Knowhow Transfer von Fachhochschulen in die Unternehmen zu katalysieren.

Das Buch versucht, sowohl die klassischen Machine Learning Methoden als auch neueste Entwicklungen im Deep Learning mit einem Fokus auf die Anwendung zu vermitteln. Deep Learning kann als eine Teilmenge des Machine Learnings gesehen werden. Das heisst, jede Deep Learning Methode ist automatisch auch eine Machine Learning Methode. Machine Learning enthält jedoch weitere Methoden, welche nicht dem Deep Learning zugeordnet werden können. Das Gebiet Machine Learning ist wiederum eine Teilmenge der Methoden der Künstlichen Intelligenz. Letztere enthält weitere Methoden, welche nicht dem Machine Learning zuzuordnen sind. Abbildung 1 versucht diesen Sachverhalt schematisch darzustellen.

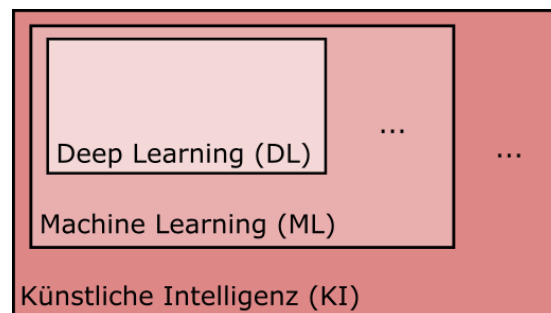


Figure 1: Unterscheidung zwischen KI, ML und DL.

Wir werden im ganzen Buch die folgenden (üblichen) Abkürzungen verwenden:

- Künstliche Intelligenz = KI (oft spricht man auch von AI, was die Abkürzung für den englischen Begriff *Artificial Intelligence* ist).

- Machine Learning = ML
- Deep Learning = DL

Obwohl das Buch einen anwendungsorientierten Ansatz verfolgt, soll die mathematisch-statistische Intuition hinter den beschriebenen Modellen und Methoden nicht zu kurz kommen. Diese Intuition ist aus meiner Sicht zwingend, um beurteilen zu können, ob sich ein Modell überhaupt für ein gegebenes Problem eignet. Am Schluss geht es nämlich darum, dass wir mit dem Einsatz von Machine Learning einen Mehrwert für ein Unternehmen oder für die Gesellschaft schaffen können. Das erfordert, dass wir uns eingehend und kritisch mit den Modellen und deren Eignung für ein gegebenes Problem auseinander setzen.

Zielgruppe

Das Buch richtet sich insbesondere an Fachhochschulstudierende in der deutschsprachigen Schweiz mit einem intrinsischen Interesse an quantitativen Methoden im Allgemeinen und Machine Learning im Besonderen. Vorausgesetzt werden Mathematikkenntnisse auf Stufe Mittelschule (Berufs- oder gymnasiale Matur), d.h. Sie sollten vertraut sein mit den Grundlagen bezüglich mathematischer Funktionen, der Integral- und Differentialrechnung sowie den wichtigsten Resultaten aus der Algebra. Ausserdem gehe ich davon aus, dass Sie bereits eine Einführung in das Thema Statistik besucht haben und Konzepte aus der deskriptiven Statistik (Mittelwert, Median, Varianz, Quantile, etc.) sowie aus der Inferenzstatistik (Verteilungen, statistisches Testen, etc.) bekannt sind.

Bevor Sie sich aber nun Sorgen machen: Kapitel 2 enthält eine Einführung in die wichtigsten Mathematik- und Statistikgrundlagen, die nötig sind für das Verständnis von Machine Learning Modellen.

Da ich mit diesem Buch einen anwendungsorientierten Ansatz verfolge, werden wir auch in das Programmieren einsteigen. Dazu verwenden wir in diesem Buch die Programmiersprache R. Es werden keine Vorkenntnisse vorausgesetzt. Kapitel 3 enthält eine kurze Einführung in die Programmiersprache R und verweist Sie auf weiterführende Ressourcen zum Thema Programmieren. Jedes Modell, das wir uns anschauen werden, ist mit R-Code dokumentiert, so dass Sie lernen, wie die Modelle in der Praxis angewendet werden können.

Aufbau des Buchs

Das Buch enthält folgende Kapitel:

- Kapitel 1: Einführung in das Thema Machine Learning mit **Definitionen** sowie Anwendungsbeispielen.
- Kapitel 2: Wichtigste **Mathematik- und Statistikgrundlagen**, die für das Verständnis der Modelle in den späteren Kapitel elementar sind.
- Kapitel 3: Einführung in das **Programmieren** mit **R** sowie Überblick über die wichtigsten **R-Packages**, die wir verwenden werden.
- Kapitel 4: Hier erlernen wir die Grundmodelle, um **Regressionsprobleme** zu lösen. Es sind lineare Modelle, was bedeutet, dass die funktionale Form der Modelle linear von den Parametern des Modells abhängen. Grafisch bedeutet dies, dass ein solches Modell im einfachsten Fall durch eine Gerade beschrieben werden kann.
- Kapitel 5: In diesem Kapitel lernen wir die Grundmodelle für das **Klassifikationsproblem** kennen. Diese Modelle führen typischerweise zu einer linearen Entscheidungsgrenze (engl. *Decision Boundary*) zwischen den verschiedenen Klassen, die wir unterscheiden oder klassifizieren wollen.
- Kapitel 6: Damit wir ML in der Praxis anwenden können, lernen wir hier die typische **ML-Pipeline** kennen. Sie werden die Techniken und Methoden kennen lernen, die es braucht, um überhaupt erst an den Punkt zu kommen, um ein ML-Modell rechnen zu können. Oft werden diese Techniken und Methoden unter dem Begriff Preprocessing der Daten zusammengefasst. Doch die Pipeline endet nicht mit dem Rechnen eines ML-Modells. Danach muss ein Modell evaluiert werden und wenn Sie als Analyst*in zufrieden sind, müssen Sie sich Gedanken machen, wie das Deployment des Modells aussehen soll. Das heisst, wie kann Ihr Modell Dritten zur Verfügung gestellt werden? Wir werden uns hier auch kurz mit den wichtigsten Techniken aus dem Unsupervised Learning befassen.
- Kapitel 7: Nach den ersten linearen Modellen für das Regressions- und Klassifikationsproblem lernen wir hier ein flexibleres Modell kennen, nämlich den **Entscheidungsbaum** (engl. *Decision Tree*). Entscheidungsbäume eignen sich sowohl für das Regressions- als auch für das Klassifikationsproblem. Obwohl sie in realen Projekten typischerweise anderen Modellen unterlegen sind, wenn es um die Vorhersagequalität geht, sind sie trotzdem attraktive Modelle, da sie gut visualisierbar sind.
- Kapitel 8: Aufbauend auf den Entscheidungsbäumen aus dem vorherigen Kapitel können sehr mächtige Modelle erstellt werden, die in der Praxis oft die besten Vorhersagen liefern. Weil es sich dabei üblicherweise um eine clevere Aggregation der Resultate einer grossen Anzahl individueller Entscheidungsbäume handelt, werden diese Modelle **Ensembles** genannt. Wie die individuellen Entscheidungsbäume eignen sich Ensembles sowohl für das Regressions- als auch für das Klassifikationsproblem.
- Kapitel 9: Ein weiteres mächtiges Modell, das sich sowohl für das Regressions- als auch für das Klassifikationsproblem eignet, sind die **Support Vector Machines**. Ihre Popularität ist mit dem Aufstieg von Deep Learning etwas verblasst. Es lohnt sich aber immer noch allemal, diese Familie von Modellen kennen zu lernen, insbesondere auch weil sie nicht als Blackbox-Modelle gelten und theoretisch gut fundiert sind.

- Kapitel 10: Ab diesem Kapitel steigen wir in das Thema Deep Learning ein. Sie werden die Architektur von einfachen **Artificial Neural Networks** (ANNs) kennen lernen. Ausserdem schauen wir uns in diesem Kapitel den genialen Backpropagation Algorithmus anhand eines einfachen linearen Regressionsproblems an. Dieser Algorithmus ist der Schlüssel für die viel diskutierten Fortschritte im Bereich der künstlichen Intelligenz, weil er das Trainieren von riesigen Modellen überhaupt erst möglich macht.
- Kapitel 11: Hier lernen wir sogenannte **Convolutional Neural Networks** (CNNs) kennen. Sie sind die Basis für die Fortschritte auf dem Gebiet Computer Vision und erlauben beispielsweise Anwendungen im Bereich automatische Gesichtserkennung in Bildern oder Videos.
- Kapitel 12: Nach ANNs und CNNs lernen wir hier **Recurrent Neural Networks** (RNNs) kennen. Diese Modelle bilden die Basis für Probleme, in denen die Daten als Sequenzen vorliegen. Das können einfache Zeitreihen (z.B. Börsenkurse) sein, aber auch komplexere Sequenzdaten wie beispielsweise geschriebene oder gesprochene Sprache oder Tonaufnahmen.
- Kapitel 13: In diesem letzten Kapitel geht es schliesslich um **Generative KI**. Wir beschäftigen uns hier also mit Modellen, die nicht nur einfach ein Vorhersageprobleme lösen können, sondern auch neue Inhalte (z.B. Texte, Musik, Bilder) generieren können. Abbildung 2 enthält als Beispiel den Output einer generativen Software, die basierend auf einem Prompt ein Bild erstellt. Nach dem Lesen dieses Kapitels sollten Sie ein grundlegendes Verständnis für die Funktionsweise von Modellen wie Chat-GPT haben.

Weiterführende Literatur

Ein grosser Teil des vorliegenden Buchs baut auf bestehenden Büchern zum Thema Machine Learning auf. Ich werde im Buch immer wieder auf die Quellen verweisen. Die wichtigsten Referenzen für dieses Buch sind folgende:

- Gareth James, Daniela Witten, Trevor Hastie, Robert Tibshirani. (2021). An Introduction to Statistical Learning: with Applications in R. New York: Springer. 2nd Edition.
- Aurélien Géron. (2019). Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems. Sebastopol: O'Reilly Media Inc. 3rd Edition.
- Christopher M. Bishop. (2006). Pattern Recognition and Machine Learning. Berlin, Heidelberg: Springer.
- Kevin P. Murphy. (2012). Machine Learning A Probabilistic Perspective. The MIT Press.

Die ersten beiden Referenzen sind einführende Texte und können parallel zum vorliegenden Buch gelesen werden. Die letzten zwei Referenzen sind fortgeschrit-

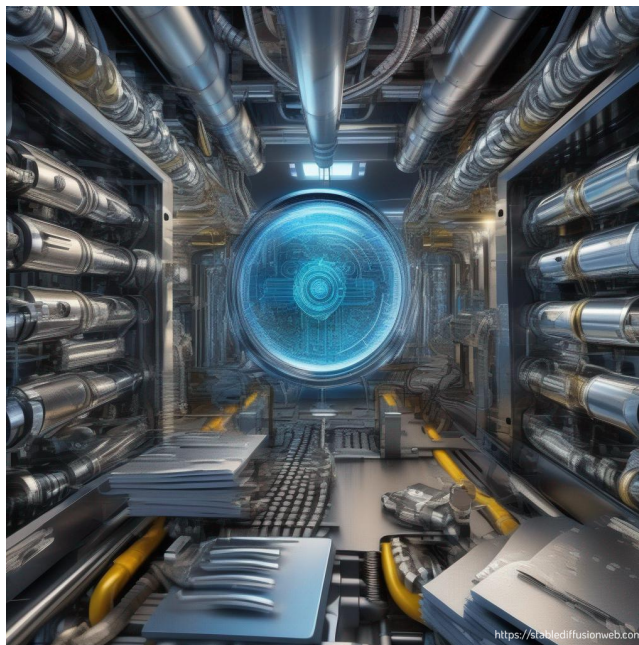


Figure 2: Beispielsoutput einer generativen Bildgenerierungssoftware (<https://stablediffusionweb.com/>) basierend auf dem Prompt "A title image for a textbook about Machine Learning targeting small and medium companies."

tene Texte und ich empfehle, sie erst nach dem vollständigen Verständnis des vorliegenden Buchs oder der ersten beiden Referenzen zu lesen.

Lizenz

Das vorliegende Buch ist unter Lizenz CC BY-NC-SA 4.0 DEED (Namensnennung, nicht-kommerziell, Weitergabe unter gleichen Bedingungen 4.0 International) lizenziert. Bitte halten Sie sich an die Lizenzbedingungen.

Kontakt

Für Fragen und Anregungen zum Buch stehe ich gerne zur Verfügung:

Martin Sterchi
Riggenbachstrasse 16
4600 Olten
martin.sterchi@fhnw.ch

Chapter 1

Einführung

In diesem Kapitel geht es darum zu verstehen, was ML überhaupt ist, warum es nützlich sein kann und was typische Anwendungsfälle von ML sind. Wir werden ausserdem verschiedene Unterkategorien von ML kennen lernen.

1.1 Was ist Machine Learning?

Im Prinzip geht die Geschichte des MLs weit zurück, nämlich zu den Anfängen der Statistik. Viele Modelle, die heutzutage im ML angewendet werden sind nämlich eigentlich von Statistiker*innen erfundene Modelle. Die Geschichte des MLs und der Statistik sind darum eng verknüpft. Einen eigentlichen Startpunkt des MLs könnte man vielleicht in den 1960er Jahren ausmachen, mit den Arbeiten von Frank Rosenblatt¹, welcher das sogenannte **Perceptron** und einen dazugehörigen Lernalgorithmus prägte (dazu später mehr). Danach blieb es aber rund 20 Jahre relativ ruhig bis die Forschung im Bereich Machine Learning so richtig Fahrt aufnahm. Ein grosser Schub für die Entwicklung von ML ging vom Aufkommen von extrem grossen Datenmengen (**Big Data**) und dem Internet aus. Das führte nämlich dazu, dass sich immer mehr Leute aus den Fachbereichen Informatik und Computer Science mit dem Thema ML befassten und effiziente Hard- und Software sowie algorithmische Kniffs und Tricks beisteuerten. Ausserdem ermöglichte das Internet den Zugang zu gewaltigen Datenmengen an Bildern, Videos, Klicks, etc. - denken Sie beispielsweise nur schon an die Informationen, die jede*r von uns tagtäglich im Internet hinterlässt. Ein weiterer Schub für das Machine Learning war (und ist) zudem die immer besser werdende Rechenleistung von Computern. Diese Entwicklungen haben sich im November 2022 kulminiert in der erstmaligen breiten öffentlichen Wahrnehmung von sogenannten **Large Language Models** wie ChatGPT.

¹https://en.wikipedia.org/wiki/Frank_Rosenblatt

Wie der Name sagt, geht es im ML darum, dass eine Maschine (oder präziser, ein Computer) aus einem gegebenen Datensatz automatisch Muster lernt, ohne dass ein Mensch dem Computer (explizit) sagen muss, was er lernen soll. Der Mensch gibt jedoch dem Computer die Rahmenbedingungen für das selbständige Lernen vor. Die erlernten Muster sind selbstverständlich nur nützlich, wenn sie genereller Natur sind und auch für neue bzw. zukünftige Beobachtungen gelten. Beispiel: ein Spital hat während der Corona Pandemie ein Modell trainiert, um den täglichen Pflegebedarf je nach Wochentag, Saison, und weiteren Indikatoren vorherzusagen. Das Modell funktioniert nun nach der Pandemie aber nicht wunschgemäß und prognostiziert in den Tendenz einen zu hohen Pflegebedarf. Das Problem ist, dass die erlernten Muster nicht auf eine Zeit nach der Pandemie übertragbar sind. Die Trainingsdaten waren nicht repräsentativ genug.

Bevor wir etwas konkreter anschauen, wie genau ein Computer selbständig aus Daten lernen kann, schauen wir uns die Definitionen von zwei Experten im Gebiet ML an:

“[Machine Learning is the] field of study that gives computers the ability to learn without being explicitly programmed.” Arthur Samuel, 1959

“Machine Learning is the science (and art) of programming computers so they can learn from data.” Aurélien Géron²

Zusammenfassend lässt sich sagen, dass wir mit ML dem Computer die Möglichkeit geben, automatisch und selbständig aus Daten zu lernen. Nichtsdestotrotz braucht es Sie als ML-Expert*in, und zwar wie folgt:

1. Sie entscheiden sich für ein spezifisches ML Modell. Typischerweise kann ein ML Modell durch eine mathematische Funktion (siehe Kapitel 2) charakterisiert werden. ML Modelle können unterschiedlich flexibel sein und es liegt im Ermessen von Ihnen, wie flexibel das Modell sein soll. Sie müssen bei der Wahl des Modells die Komplexität des Problems berücksichtigen. Grundsätzlich gilt bei der Wahl des Modells, dass flexiblere Modelle komplexere Sachverhalte abbilden können. Ein zu flexibles Modell kann aber zu Overfitting führen, aber dazu später mehr. Dieser Schritt wird im Fachjargon typischerweise **Model Selection** (Modelauswahl) genannt.
2. Sobald Sie das Modell ausgewählt haben, übergeben Sie dem Computer (etwas vereinfacht gesagt) das Modell, einen Datensatz sowie einen Lernalgorithmus. Nun hat der Computer alle Zutaten, um automatisch zu lernen. Doch was lernt er eigentlich? Der Computer lernt die Parameter Ihres gewählten Modells, so dass das Modell sich optimal an die Daten anpasst. Dieser Schritt wird im Fachjargon **Model Training** (Trainieren des Modells) genannt.

²Aurélien Géron. (2019). Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems. Sebastopol: O'Reilly Media Inc. 3rd Edition.

3. Falls Sie mit dem erlernten Modell zufrieden sind, können Sie es nun entweder dazu verwenden Vorhersagen zu machen oder um Zusammenhänge in den Daten zu interpretieren und daraus wertvolle Einsichten gewinnen. Dieser Schritt wird im Fachjargon als **Model Inference** (Modellinferenz) zusammengefasst. Typischerweise sind Sie in der Realität mit dem ersten erlernten Modell allerdings noch nicht zufrieden und gehen zurück zu Schritt 1 und wählen ein anderes Modell.

Es handelt sich bei dieser Vorgehensweise um eine sehr allgemeine Beschreibung des Machine Learning Prozesses. Wie diese drei Schritte konkret funktionieren, werden Sie in den nachfolgenden Kapiteln dieses Buchs erfahren.

1.2 Wann macht es Sinn ML einzusetzen?

Ein ML Modell zu trainieren kann viel Zeit und Geld kosten. Zum Beispiel müssen Sie unter Umständen überhaupt erst die Daten sammeln (oder von einem Datendienstleister kaufen), um ein Modell zu trainieren. Oder das Projekt ist so komplex, dass Sie als Analyst*in unzählige Stunden benötigen, um die Daten überhaupt erst in eine Form zu bringen, die es erlaubt ein Modell zu trainieren. Für neuartige DL Modelle oder Generative KI kann das Trainieren bzw. Lernen eines Modells durch den reinen Stromverbrauch bzw. die vom Cloud-Betreiber in Rechnung gestellten Kosten so hoch sein, dass sich Ihr ursprüngliches Vorhaben nicht mehr lohnt. Es ist also ungemein wichtig, dass Sie sich vor Projektbeginn gut überlegen, ob ML für Ihr vorliegendes Problem überhaupt Sinn macht und einen Mehrwert generieren kann.

Folgende Daumenregeln³ können Ihnen dabei helfen, zu entscheiden, ob ML für Ihr Projekt Sinn macht:

- Ihr Problem entspricht einem Standard ML-Problem, das bereits mehrfach gelöst wurde und für das es sogenannte “off-the-shelf” Lösungen gibt. Beispiel: Sie wollen das Sentiment (positive vs. negative Grundhaltung) von Social Media Posts über Ihr Unternehmen automatisch klassifizieren. Dazu gibt es viele vortrainierte Modelle, die teilweise gratis verwendet werden können.
- Der manuelle Arbeitsaufwand ist sehr gross, wenn das Problem durch Menschen gelöst werden soll. Das Problem ist aber ansonsten klar strukturiert und benötigt keinen grossen kognitiven Einsatz eines Menschen. Beispiel: In den Post-Verteilzentren werden die von Hand geschriebenen PLZ problemlos mittels Computer bzw. ML Modellen “gelesen” und die Briefe und Pakete entsprechend sortiert.

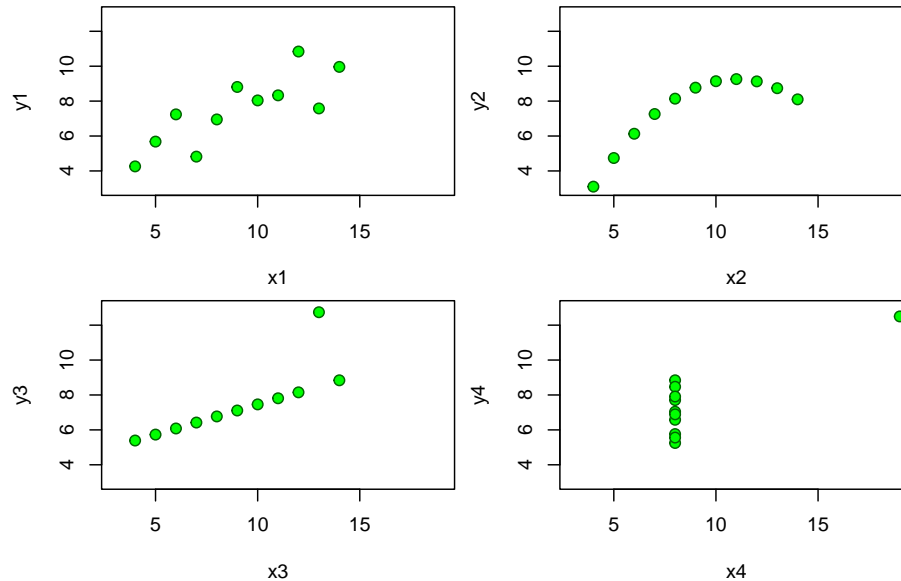
³siehe auch Seiten 6 - 7 in Aurélien Géron. (2019). Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems. Sebastopol: O'Reilly Media Inc. 3rd Edition.

- Komplexe Probleme, in denen ein Mensch keinen Überblick hat, weil so grosse und komplexe Datenmengen vorhanden sind. Wir Menschen haben grosse Mühe damit, in Rohdaten (reinen Datentabellen) irgendwelche Muster zu erkennen. In diesem Fall können wir entweder versuchen, die Daten zu visualisieren oder mithilfe von ML Zusammenhänge zu lernen, die wir sonst nicht erkennen könnten. Ein illustratives Beispiel ist das Anscombe Quartett⁴, das vier kleine Stichproben mit jeweils elf Datenpunkten enthält. Jeder Datenpunkt wird durch eine x und eine y Variable beschrieben. Die vier x - sowie die vier y -Variablen haben identische Mittelwerte. Erst eine einfache Visualisierung der vier Stichproben mithilfe eines Streudiagramms zeigt die Muster sowie die Unterschiede zwischen den vier Stichproben deutlich auf.

TODO: DL nur mit grossen Datenmengen

```
#>      x1 x2 x3 x4      y1  y2    y3    y4
#> 1   10 10 10  8   8.04 9.14  7.46  6.58
#> 2     8  8  8  8   6.95 8.14  6.77  5.76
#> 3   13 13 13  8   7.58 8.74 12.74  7.71
#> 4     9  9  9  8   8.81 8.77  7.11  8.84
#> 5   11 11 11  8   8.33 9.26  7.81  8.47
#> 6   14 14 14  8   9.96 8.10  8.84  7.04
#> 7     6  6  6  8   7.24 6.13  6.08  5.25
#> 8     4  4  4 19   4.26 3.10  5.39 12.50
#> 9   12 12 12  8  10.84 9.13  8.15  5.56
#> 10    7  7  7  8   4.82 7.26  6.42  7.91
#> 11    5  5  5  8   5.68 4.74  5.73  6.89
```

⁴<https://de.wikipedia.org/wiki/Anscombe-Quartett>



1.3 Anwendungsfälle von ML

Stellen Sie sich vor, wir haben einen Datensatz mit 300 Spam Emails und 700 “Ham” Emails (kein Spam). Ohne Machine Learning müssten wir nun von Hand die 300 Spam Emails mit den 700 Ham Emails vergleichen und versuchen, Muster zu finden, die es uns erlauben Regeln aufzustellen, um die Spam Emails korrekt zu klassifizieren (z.B. Spam enthält tendenziell eher Geldbeträge oder Preise als Ham). Danach könnten wir die Regeln mit R implementieren. Dann stellt sich aber auch noch die Frage, wie die verschiedenen Regeln miteinander kombiniert werden, um eine Klassifikation zu machen. Dieses Vorgehen würde sehr viel zu tun geben und es würde gezwungenermaßen zu willkürlichen Entscheidungen führen.

Machine Learning führt zu i) weniger Aufwand und ii) besseren Lösungen, indem wir in einem R-Skript ein Modell (z.B. logistische Regression) aufsetzen und dann dem Modell die Daten in geeigneter Form füttern. Danach lernt der Computer selbständig, wie er die Emails bestmöglich in Spam und Ham klassifiziert.

ML Beispiele - Spam Filter - ChatGPT - Face Recognition in Fotos - Predictive Maintenance von Maschinen - Fraisa Fotos - Summarizing documents - Customer Service Chatbots - Credit Card Fraud Detection - Recommender Systems - Optimierungsprobleme (RL) - Anomaly/Novelty Detection

1.4 Supervised vs. Unsupervised Learning

Den Unterschied zwischen dem Supervised Learning und dem Unsupervised Learning können wir am besten erklären, indem wir uns mit ein paar mathematischen Grundlagen des Machine Learnings befassen. Keine Sorge, diese Grundlagen sind sehr einfach, aber versuchen Sie, diese bereits gut zu verstehen, denn wir bauen später darauf auf.

Im **Supervised Learning** haben wir einerseits sogenannte Input-Daten und andererseits einen Output, den wir vorhersagen wollen. Für die Input-Daten gibt es ganz viele verschiedene Begriffe, die synonym verwendet werden: z.B. Features, unabhängige Variablen, Attribute, Prädiktoren. Dasselbe gilt für den Output, hier gibt es folgende Synonyme: Zielvariable, abhängige Variable, Label, oder auch einfach y . Unsere Konvention hier ist aber folgende: es gibt Input-Daten (oder Input-Variablen) und einen Output (oder Output-Variable).

Die Input-Daten für eine Beobachtung i schreiben wir mathematisch wie folgt:

$$\mathbf{x}_i = \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{ip} \end{pmatrix},$$

Diese Notation bedarf ein paar Erklärungen:

- Den Index i brauchen wir, um die verschiedenen Beobachtungen zu kennzeichnen. i kann eine Ganzzahl zwischen 1 und n annehmen, wobei n die Anzahl Beobachtungen im Datensatz bezeichnet. Wenn wir zum Beispiel etwas über die Input-Daten der dritten Beobachtung sagen wollen, dann können wir die Notation \mathbf{x}_3 verwenden.
- Für jede Beobachtung i haben wir insgesamt p Variablen, welche die verschiedenen Attribute einer Beobachtung enthalten. x_{i1} bezeichnet also die erste Variable der i -ten Beobachtung, x_{i2} die zweite Variable der i -ten Beobachtung und x_{ip} die p -te (letzte) Variable der i -ten Beobachtung.
- Was Sie oben sehen, ist aus mathematischer Sicht ein Spaltenvektor. Im Moment reicht es, wenn Sie wissen, dass wir mit diesem Spaltenvektor die Input-Daten einer Beobachtung *kompakt* darstellen können.

Neben den Input-Daten haben wir im Supervised Learning aber wie erwähnt auch einen Output und den bezeichnen wir üblicherweise mit y_i . Auch hier hilft uns der Index i dabei, die Beobachtungen eindeutig zu kennzeichnen. Schauen wir uns am besten kurz ein konkretes Beispiel an:

Aufgabe

Wichtig: Beim Supervised Learning geht es um ML Probleme, in denen sowohl Input-Daten als auch ein Output vorhanden ist. Ziel beim Supervised Learning

ist es, ein Modell zu trainieren, das basierend auf den Input-Daten möglichst gute Vorhersagen für den Output macht. Es geht also hier um Vorhersageprobleme. In einem gewissen Sinn ist der Output die überwachend Instanz (engl. Supervisor), welche den Lernprozess des Modells kontrolliert.

Im Gegensatz zum Supervised Learning haben wir im **Unsupervised Learning** nur Input-Daten und *keinen Output*. Im Unsupervised Learning geht es darum, aus den Input-Daten interessante Muster zu lernen, welche für bessere unternehmerische Entscheidungen verwendet werden können. Ein einfaches Beispiel ist das Clustering von Kundinnen und Kunden eines Unternehmens in ähnliche Kundengruppen, so dass die verschiedenen Kundengruppen gezielter mit Marketingaktionen angesprochen werden können. Techniken, um komplexe Datensätze zu visualisieren, werden typischerweise auch zum Unsupervised Learning gezählt.

Neben dem Supervised und dem Unsupervised Learning gibt es noch eine dritte Kategorie von Machine Learning, nämlich das **Reinforcement Learning** (RL). Dieser Kategorie gehören Modelle an, die (virtuelle) Agenten so trainieren, dass sie langfristig möglichst optimal handeln. Das bekannteste Beispiel aus dem RL ist Googles AlphaGo Agent, welcher den menschlichen Go Weltmeister im Jahr 2017 schlug.⁵ Reinforcement Learning ist aber auch eine wichtige Komponente in der Optimierung von grossen Sprachmodellen wie ChatGPT. In einer ersten Fassung dieses Buchs werden wir uns nicht (oder nur am Rande) mit RL befassen.

Die Unterscheidung zwischen den drei Arten von Machine Learning ist im oberen Teil der Abbildung 1.1 visualisiert:

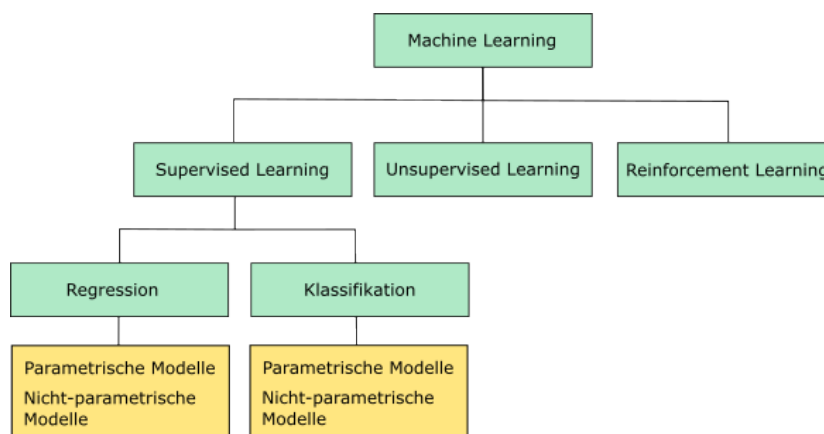


Figure 1.1: Die verschiedenen Kategorien des Machine Learnings und deren Hierarchie.

⁵<https://deepmind.google/technologies/alphago/>

1.5 Regression vs. Klassifikation

In der Kategorie des Supervised Learnings unterscheiden wir weiter zwischen Regressions- und Klassifikationsproblemen (siehe auch Abbildung 1.1).

Im Regressionsproblem ist der Output eine **stetige** Variable (Intervall- oder Verhältnisskalierung), d.h. die Variable enthält numerische Werte.

Im Klassifikationsproblem ist der Output bzw. die Zielvariable eine **kategorische** Variable (Nominal- oder Ordinalskalierung).

Aufgabe

1.6 Parametrische vs. nicht-parametrische Modelle

Ein ML Modell gehört entweder der Familie **parametrischer** Modelle oder der Familie **nicht-parametrischer** Modelle an. Dabei spielt es keine Rolle, ob wir mit dem Modell ein Regressions- oder ein Klassifikationsproblem lösen wollen.

Womöglich sind Sie in Ihrer Ausbildung bereits **parametrischen Modellen** begegnet, denn das einfache lineare Regressionsmodell ist ein typisches Beispiel für ein parametrisches ML Modell. Das Modell ist vollkommen charakterisiert durch die beiden lernbaren (optimierbaren) Parameter w_0 und w_1 ⁶ und kann wie folgt (mathematisch) aufgeschrieben werden:

$$\hat{y}_i = f(x_i) = w_0 + w_1 \cdot x_i$$

Wenn Ihnen der obige Ausdruck noch fremd vorkommt, dann ist das nicht schlimm. Wir werden im Kapitel 4 ausführlich auf lineare Regressionsmodelle eingehen. Im Moment müssen Sie nur wissen, dass ein parametrisches Modell wie oben mit einer mathematischen Funktion beschrieben werden kann und dass diese Funktion durch lernbare **Parameter** (hier w_0 und w_1) charakterisiert wird.

Nicht-parametrische Modelle wiederum sind Modelle, welche nicht (oder zumindest nicht explizit) durch Parameter charakterisiert sind. Am besten schauen wir uns gleich ein einfaches nicht-parametrisches Modell an, nämlich das **K-Nearest-Neighbors** (KNN) Modell. Stellen Sie sich vor, Sie haben einen Datensatz mit 55 Produkten aus Ihrem Sortiment. Sie haben jedes dieser 55 Produkte auf Instagram und auf Tiktok durch Influencer*innen bewerben lassen. Für jedes der 55 Produkte hatten Sie ein Werbebudget für Instagram (x_{i1}) und ein Werbebudget für Tiktok (x_{i2}). Am Ende des Geschäftsjahrs haben

⁶In Statistikvorlesungen werden die beiden Parameter oft eher mit b_0 und b_1 oder mit β_0 und β_1 bezeichnet. Im Machine Learning nennt man Parameter oft Gewichte (engl. Weights), weshalb die Parameter typischerweise mit w bezeichnet werden.

Sie für jedes der 55 Produkte bestimmt, ob die Absatzziele erreicht wurden oder nicht (Output y_i). Die erfolgreichen Produkte (= Absatzziel erreicht) sind in untenstehender App als blaue Punkte eingezeichnet. Die roten Dreiecke repräsentieren die nicht-erfolgreichen Produkte. Sie sehen, dass erfolgreiche Produkte tendenziell höhere Instagram und Tiktok Werbebudgets aufwiesen als nicht-erfolgreiche Produkte. Sie möchten nun ein Modell schätzen, dass die Produkte automatisch klassifizieren kann. Dazu verwenden Sie das KNN Modell, das die K nächsten Nachbarn unter den 55 gegebenen Produkten sucht und dann die häufigste Beobachtung unter den K nächsten Nachbarn vorhersagt. In anderen Worten: wir suchen die **K ähnlichsten** Beobachtungen und nutzen diese, um eine Vorhersage zu machen.

Selbstverständlich spielt der konkrete Wert von K hier eine grosse Rolle - sollen wir nur $K = 1$ Nachbarn berücksichtigen? Oder $K = 10$ Nachbarn? Die erste Abbildung in der App zeigt nicht nur die 55 Datenpunkte, sondern auch die **Entscheidungsgrenze** (in schwarz). Untersuchen Sie kurz, wie sich diese Entscheidungsgrenze verändert, wenn Sie K erhöhen oder reduzieren.

Ausserdem können Sie in der ersten Abbildung auch den schwarzen Punkt mit der Maus setzen, wodurch Ihnen die K nächsten Punkte des schwarzen Punkts angezeigt werden.

Die zweite Abbildung zeigt die Entscheidungsregionen mit unterschiedlicher Intensität je nachdem wie sicher sich das Modell ist. In einer Region, in der alle K Nachbarn nicht-erfolgreiche Produkte sind, sind wir uns eher sicher bezüglich der Vorhersage als in einer Region, in der die Anteile zwischen erfolgreichen und nicht-erfolgreichen Produkten ausgeglichen sind.

Um die K nächsten Nachbarn zu finden, müssen wir die Distanzen zwischen Punkten rechnen können. Dazu verwenden wir die Euklidische Distanz, welche wir in Kapitel 2 kennen lernen werden.

Das KNN Modell ist ein sehr einfaches ML Modell, welches in der Praxis allerdings nicht allzu häufig angewendet wird. Warum nicht? Weil es am sogenannten **Fluch der Dimensionalität** (engl. Curse of Dimensionality) leidet. Doch was bedeutet das? Je mehr Input-Variablen wir haben, desto weiter entfernt sind Datenpunkte voneinander (das ist etwas, das man sich nur schwer vorstellen kann, aber Sie können es mir für den Moment einfach mal glauben). Das KNN beruht auf der Grundidee, dass wir K nahe, ähnliche Beobachtungen für die Vorhersage verwenden. Wenn diese K nahen Beobachtungen im hochdimensionalen Raum (= viele Input-Variablen) nicht mehr nahe sind, dann funktioniert auch das Modell nicht mehr gut.

Ev. Market-Basket Analyse

Zwei Fragen mit Bildern (siehe Rapp)

1.7 Machine Learning Pipeline

Pipeline zeigen

Chapter 2

Mathematik- und Statistik-Grundlagen

In diesem Kapitel repetieren wir die wichtigsten Grundlagen aus der Mathematik und Statistik, die es braucht, um Machine Learning Modelle zu verstehen. Das Thema *Lineare Algebra* wird für die meisten von Ihnen wahrscheinlich Neuland sein.

2.1 Funktionen

Eine Funktion, die wir in der Mathematik typischerweise mit f bezeichnen, ordnet jedem **Argument** x aus dem Definitionsbereich D (engl. *Domain*) **genau einen Wert** y aus dem Wertebereich W (engl. *Codomain*) zu. Oft sind D und W die Menge der reellen Zahlen, also \mathbb{R} . Die Menge der reellen Zahlen enthält alle möglichen Zahlen, die Sie sich vorstellen können.¹ Zum Beispiel die Zahlen 3, -4.247 , $\sqrt{14}$, $5/8$, etc.

Wie eine Funktion grafisch aussieht, ist aus Panel (a) der Abbildung 2.1 ersichtlich. Hier zeigen wir die Form einer Funktion in einem kartesischen Koordinatensystem. Die Funktionskurve weist jedem Wert x auf der x-Achse genau einen Wert y auf der y-Achse zu. Der wichtigste Teil der oben aufgeführten Definition ist der Teil “genau einen Wert”, denn eine Funktion kann einem Element x nicht zwei oder mehr Werte zuweisen, sondern nur genau einen. Genau aus diesem Grund handelt es sich bei Panel (b) in Abbildung 2.1 *nicht* um eine Funktion, da gewissen x -Werten mehrere Werte y zugeordnet werden. *Wichtig*: das heisst aber nicht, dass zwei verschiedenen x -Werten, nennen wir sie x' und x'' , derselbe y -Wert zugeordnet werden kann (vgl. Panel (a)).

¹Einzige Ausnahme sind die komplexen Zahlen.

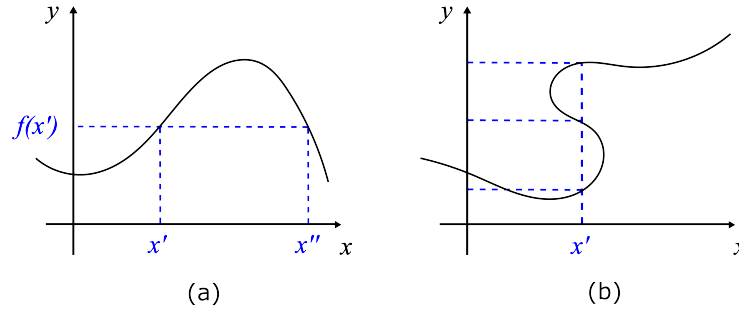


Figure 2.1: (a) Eine Funktion, die jedem x -Wert genau einen y -Wert zuweist. (b) Keine Funktion.

Mathematisch wird diese allgemeine Definition einer Funktion häufig wie folgt beschrieben:

$$f : x \mapsto y$$

Wir haben also eine Funktion f , die jedem Element x genau einen Wert y zuweist. Der Pfeil in obiger mathematischer Schreibweise beschreibt genau dieses Mapping. Wie genau dieses Mapping einem Argument x den entsprechenden y -Wert zuordnet, wird durch die Funktion $f(x)$ beschrieben. In den folgenden Abschnitten schauen wir uns typische Beispiele von Funktionen an, angefangen mit linearen Funktionen. Doch vorher wollen wir uns kurz überlegen, warum Funktionen für das Machine Learning überhaupt wichtig sind. Ein grosser Teil des Machine Learnings, der **Supervised Learning** genannt wird, befasst sich mit dem Problem, wie eine Zielvariable y mithilfe von einem oder mehreren Prädiktoren x vorhergesagt werden kann. Ein Machine Learning Modell ist darum nichts anderes als eine Funktion $y = f(x)$, die basierend auf den Prädiktoren x die Zielvariable y möglichst gut beschreiben kann.²

2.1.1 Lineare Funktionen

Nun schauen wir uns an, wie eine **lineare** Funktion aussieht. Eine lineare Funktion kann allgemein wie folgt geschrieben werden:

$$y = f(x) = a \cdot x + b$$

Obige Funktionsgleichung besagt, dass wir den entsprechenden y -Wert kriegen, indem wir den Wert des Arguments x mit a multiplizieren und danach eine Konstante b addieren. a und b sind die **Parameter** dieser Funktion. Die konkreten Zahlenwerte dieser beiden Parameter definieren, wie die Funktion am Schluss genau aussieht.

²Zumindest aus einer nicht-probabilistischen Perspektive.

Eine lineare Funktion hat auch eine geometrische Interpretation und zwar entspricht eine lineare Funktion einer Gerade. Das ist auch der Grund, warum wir diese Funktionen **linear** nennen, sie können graphisch durch eine “Linie” dargestellt werden. Der Parameter a ist die Steigung dieser Geraden und der Parameter b entspricht dem Ort, wo die Gerade die y -Achse schneidet (sogenannter y -Achsenabschnitt).

Am besten schauen wir uns ein paar konkrete Beispiele an (Abb. 2.2).

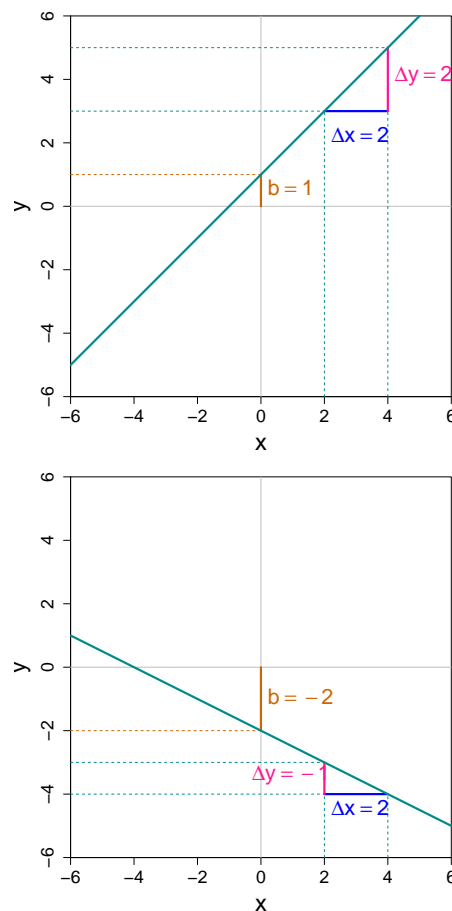


Figure 2.2: Beispiele linearer Funktionen.

Aus der linken Abbildung können wir ablesen, dass die Steigung dieser Geraden $\frac{\Delta y}{\Delta x} = \frac{2}{2} = 1$ ist und dass die Gerade die y -Achse am Ort 1 schneidet. Die entsprechende lineare Funktion kann dementsprechend als $y = x + 1$ geschrieben werden.³

³Wir müssen hier die Steigung 1 nicht explizit schreiben, aber selbstverständlich ist es nicht

Aus der rechten Abbildung können wir ablesen, dass die Steigung $\frac{\Delta y}{\Delta x} = \frac{-1}{2} = -0.5$ ist und dass die Gerade die y-Achse am Ort -2 schneidet. Die entsprechende lineare Funktion kann dementsprechend als $y = -0.5 \cdot x - 2$ geschrieben werden.

Es ist wichtig zu sehen, dass der Effekt einer Veränderung von x (also Δx) auf y überall derselbe ist. Es spielt also keine Rolle, ob wir von $x = -2$ zu $x = -1$ gehen oder von $x = 100$ zu $x = 101$, die entsprechende Veränderung in y (also Δy) wird dieselbe sein. Das muss so sein, denn die Gerade steigt (oder sinkt) mit konstanter Steigung.

Aufgaben

1. Zeichnen Sie die Funktion $y = 2 \cdot x$ in ein Koordinatensystem ein. Warum fehlt der Parameter b ?
2. Zeichnen Sie die Funktion $y = -3$ in ein Koordinatensystem ein. Ist das überhaupt eine Funktion nach obiger Definition?

2.1.2 Quadratische Funktionen

Nun wollen wir uns eine etwas interessantere (und flexiblere) Familie von Funktionen anschauen, nämlich **quadratische** Funktionen. Auch hier wollen wir die Funktion erstmal allgemein aufschreiben:

$$y = f(x) = a \cdot x^2 + b \cdot x + c$$

Eine quadratische Funktion hat drei **Parameter**, nämlich a , b und c . Grafisch entspricht die quadratische Funktion einer **Parabel** (vgl. Abb. 2.3). Die Parameter sind hier nicht mehr so einfach grafisch zu interpretieren, aber die vier Beispiele in unten stehender Abbildung geben Anhaltspunkte, was passiert, wenn die Parameterwerte sich ändern.

Aufgaben

1. Sie haben folgende quadratische Gleichung: $y = 2 \cdot x^2 + x - 2$. Berechnen Sie mit der bekannten Lösungsformel $x_{1,2} = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$ die Orte auf der x-Achse, wo die Parabel die Achse schneidet (oder einfacher gesagt die Nullstellen).
2. Verwenden Sie folgenden R-Code, um beliebige quadratische Funktionen grafisch darzustellen, indem Sie die Parameterwerte auf der ersten Code-Zeile verändern.

falsch die lineare Funktion als $y = 1 \cdot x + 1$ zu schreiben.

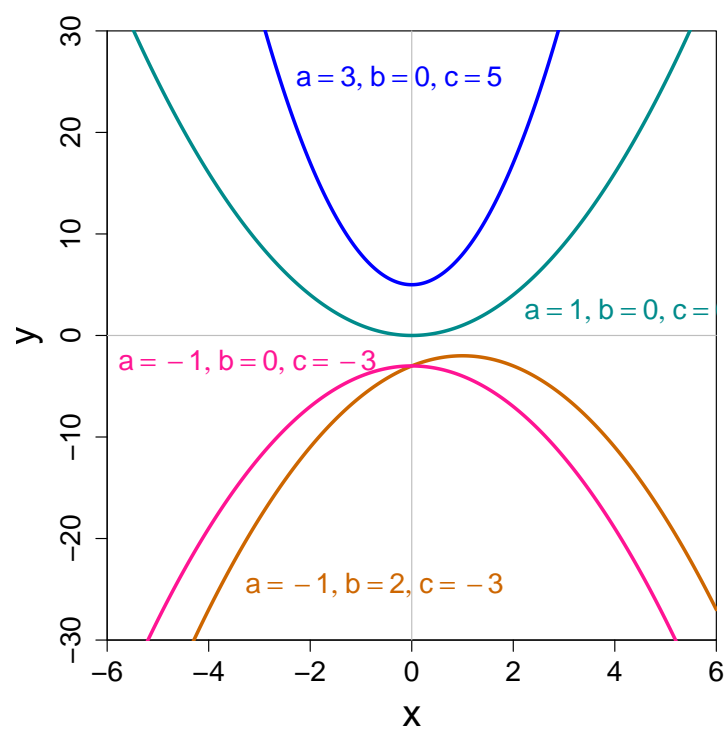


Figure 2.3: Beispiele quadratischer Funktionen.

```

# Parameter setzen
a <- 2; b <- 0; c <- 1
# Quadratische Funktion
quad <- function(x, a, b, c) {a * x^2 + b * x + c}
# x-Werte
x <- seq(-6, 6, 0.01)
# y-Werte
y <- quad(x, a, b, c)
# Plot
plot(x, y, type = "l", lwd = 2, col = "darkcyan")

```

Sie wundern sich nun vielleicht, könnte man nicht auch eine Funktion antreffen, in der x^3 , x^4 , etc. vorkommen? Das ist selbstverständlich möglich. In diesem Fall spricht man dann von einem sogenannten **Polynom**. Die höchste Potenz des Arguments x definiert den Grad des Polynoms.

Schauen wir uns doch am besten gleich wieder ein Beispiel an:

$$y = f(x) = 1 \cdot x^4 - 2 \cdot x^3 - 5 \cdot x^2 + 8 \cdot x - 2$$

Die Visualisierung dieser Funktion ist in Abb. 2.4 gegeben. Diese Funktion ist nun bereits enorm flexibel und kann je nach Parameterwerten ganz unterschiedliche Zusammenhänge abbilden.

Aufgaben

1. Eine quadratische Funktion ist ein Polynom welchen Grades?
2. Handelt es sich bei der Funktion $y = 2x^5 + x + 1$ immer noch um ein Polynom? Falls ja, ein Polynom welchen Grades?
3. Handelt es sich bei der Funktion $y = x^{0.5} + 2$ um ein Polynom?

2.1.3 Funktionen mehrerer Argumente

Bisher haben wir nur Funktionen mit **einem Argument** x angeschaut, doch die meisten für das Machine Learning interessanten Funktionen sind Funktionen **mehrerer Argumente**.

Der Einfachheit halber schauen wir uns hier nur mal eine **lineare** Funktion zweier Argumente, nennen wir sie x_1 und x_2 , an, denn diese können wir in 3D immer noch visualisieren. Wir betrachten folgende Funktion: $y = f(x_1, x_2) = 1 \cdot x_1 + 0.5 \cdot x_2 + 5$.

Aha! Während eine lineare Funktion eines Arguments grafisch einer Gerade entspricht, sehen wir nun, dass eine lineare Funktion zweier Argumente nichts anderes als eine Ebene darstellt. Wir sehen, dass die Ebene die y-Achse am Punkt 5 schneidet. Etwas schwieriger zu sehen ist die Steigung der Ebene in

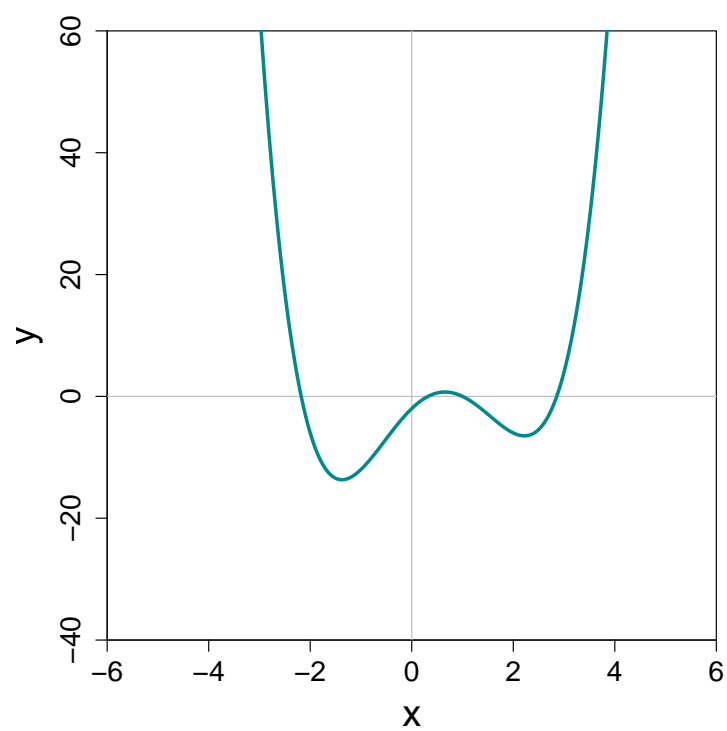


Figure 2.4: Beispiel einer polynomischen Funktion vierten Grades.

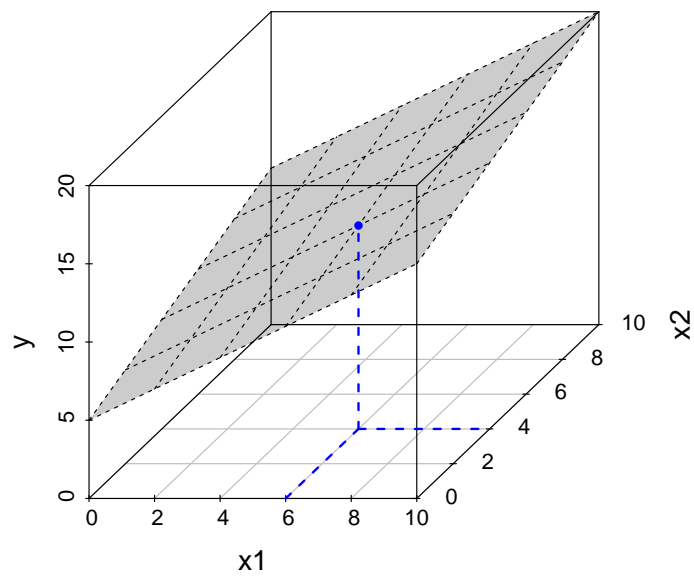


Figure 2.5: Lineare Funktion zweier Argumente (Ebene).

die Richtung der x_1 -Achse und in die Richtung der x_2 -Achse. Sie können aber vielleicht bereits erraten, dass die (partiellen) Steigungen 1 und 0.5 betragen.

Die Funktion ordnet jeden möglichen Punkt (x_1, x_2) einem Punkt auf der Ebene zu. Wir können zum Beispiel für den in Abb. 2.5 eingezeichneten Punkt $(6, 4)$ den entsprechenden Punkt auf der Ebene ausrechnen:

$$\begin{aligned} y &= 1 \cdot x_1 + 0.5 \cdot x_2 + 5 \\ &= 1 \cdot 6 + 0.5 \cdot 4 + 5 \\ &= 13 \end{aligned}$$

Selbstverständlich könnten wir uns nun auch quadratische Funktionen oder Polynome mehrerer Argumente anschauen, aber darauf verzichten wir vorerst.

2.1.4 Potenzen und Logarithmen

Blabla...

2.2 Integral- und Differentialrechnung

Olteanu materials: Local vs. global minima From a maximization to a minimization problem Basic definition of derivative Differentiation rules local min., max. and saddle point Second derivative test Partial derivatives What is a gradient? What is Hessian? What is Jacobian? Chain rules Lagrange optimization

2.3 Lineare Algebra

Olteanu materials: What is a scalar? What is a vector? What is a matrix? Vector norms Inner products Symmetric, diagonal, square and identity matrix Associative, commutative laws for matrices Matrix addition and multiplication Matrix inversion Eigenvectors and eigenvalues Quadratic form and positive (semi-) definiteness Differentiation rules for matrices

2.4 Wahrscheinlichkeitsrechnung

Olteanu materials: Sample space and axioms of probability Conditional probability definition Discrete vs. continuous random variables Joint probability distributions Expectation and variance, covariance (always for discrete and continuous) Bernoulli, Binomial, Normal, Multivariate Normal, Laplace

2.4.1 Diskrete Zufallsvariablen

Wir werden später sehen, dass im Machine Learning oftmals Dinge als **Zufallsvariablen** modelliert werden. Eine Zufallsvariable X ist eine Variable, für die der konkrete Wert nicht von vornherein klar ist. Wir können mit X zum Beispiel das Resultat eines Münzwurfs modellieren. Die zwei möglichen Resultate sind Kopf und Zahl. Vor dem Münzwurf ist nicht klar, ob Kopf oder Zahl erscheinen wird. Genau darum modellieren wir das Resultat des Münzwurfs als Zufallsvariable.

Es gibt in diesem einfachen Beispiel nur zwei mögliche Resultate (Kopf und Zahl), d.h. die Anzahl möglicher Resultate ist endlich (= nicht unendlich). Darum handelt es sich in diesem Fall um eine **diskrete** Zufallsvariable.

2.5 Verteilungen

Chapter 3

Einführung in das Programmieren mit R

leaRn Materialien

tidymodels

Referenzen auf andere Ressourcen (Hadley et al.)

Chapter 4

Lineare Regression

Chapter 5

Lineare Klassifikation

Chapter 6

Machine Learning Pipeline

Chapter 7

Decision Trees

Chapter 8

Ensembles

Chapter 9

Support Vector Machines

Chapter 10

Artificial Neural Networks

Chapter 11

Convolutional Neural Networks

Chapter 12

Recurrent Neural Networks

Chapter 13

Generative AI