

Maschinelles Lernen 01

Prof. Dr. Christoph Böhm

Hochschule München

3. Januar 2024

Einführung

Einführung

Einführung

Was ist maschinelles Lernen?

Beispiele:

- ▶ Spracherkennung und Gesichtserkennung
- ▶ Gesichtserkennung
- ▶ Künstliche Intelligenz / Bots in Computerspielen
- ▶ Autonomes Fahren
- ▶ Kreditausfälle oder Betrugsversuche vorhersagen
- ▶ Medizinische Diagnosen

Frage:

Was haben alle diese Beispiele **gemeinsam**?

Einführung

Was ist maschinelles Lernen?

Paradigmenwechsel

Für alle diese Beispiele ist es relativ schwierig, entsprechenden Programmcode manuell zu schreiben. Beim **maschinellen Lernen** (ML) wird daher ein anderes Paradigma verwendet.

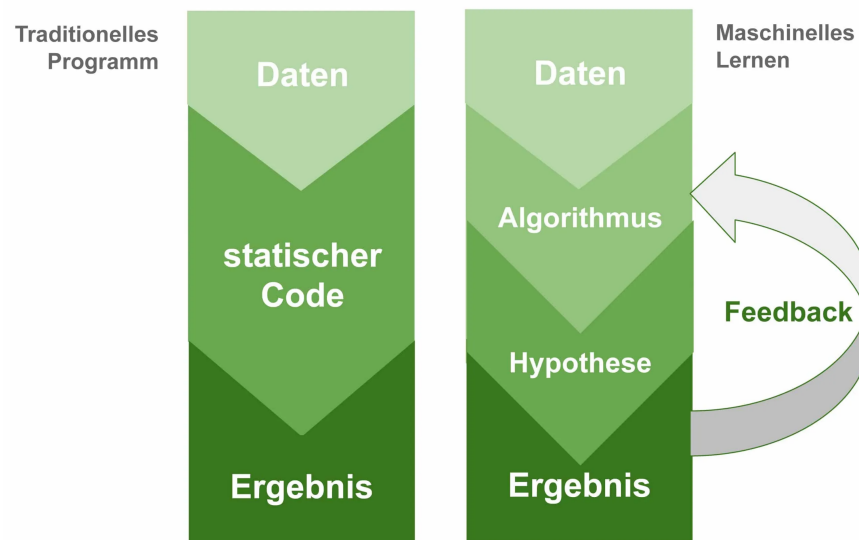
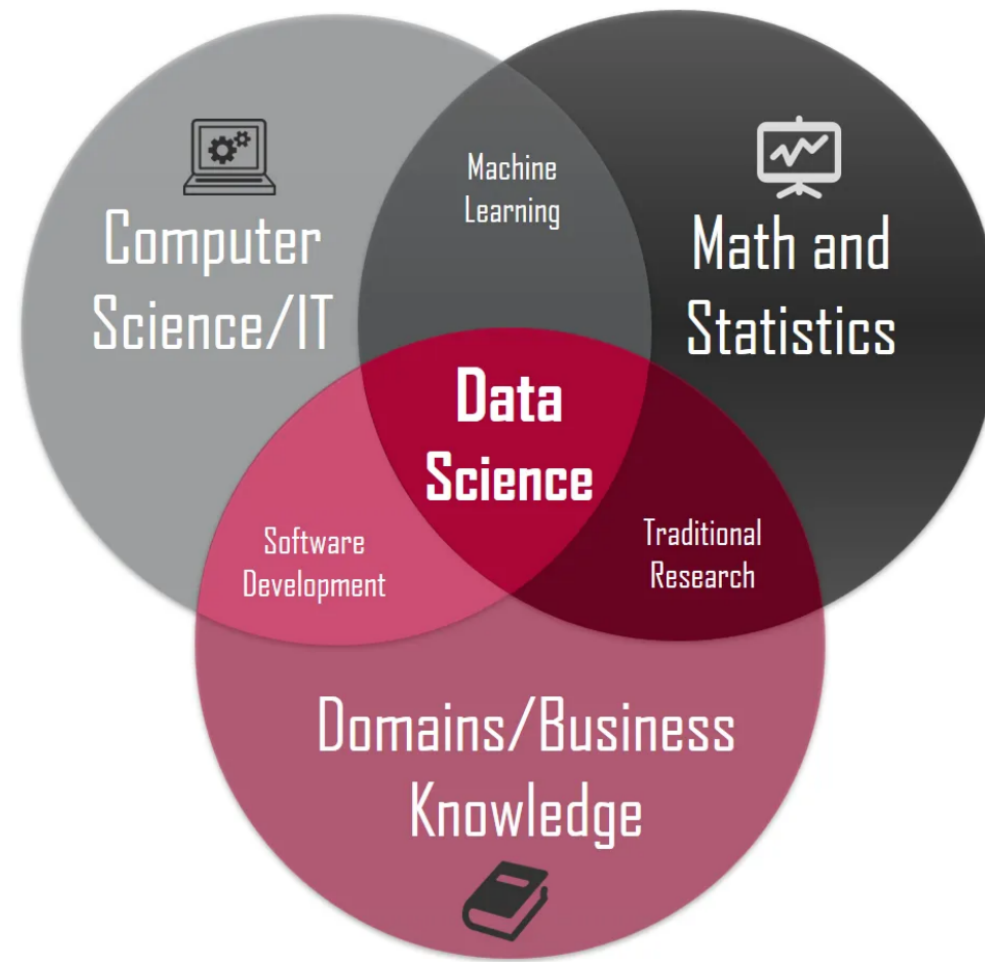


Abbildung 1: Paradigmenwechsel von manuell geschriebenem Code zu trainierten Modellen.

Einführung

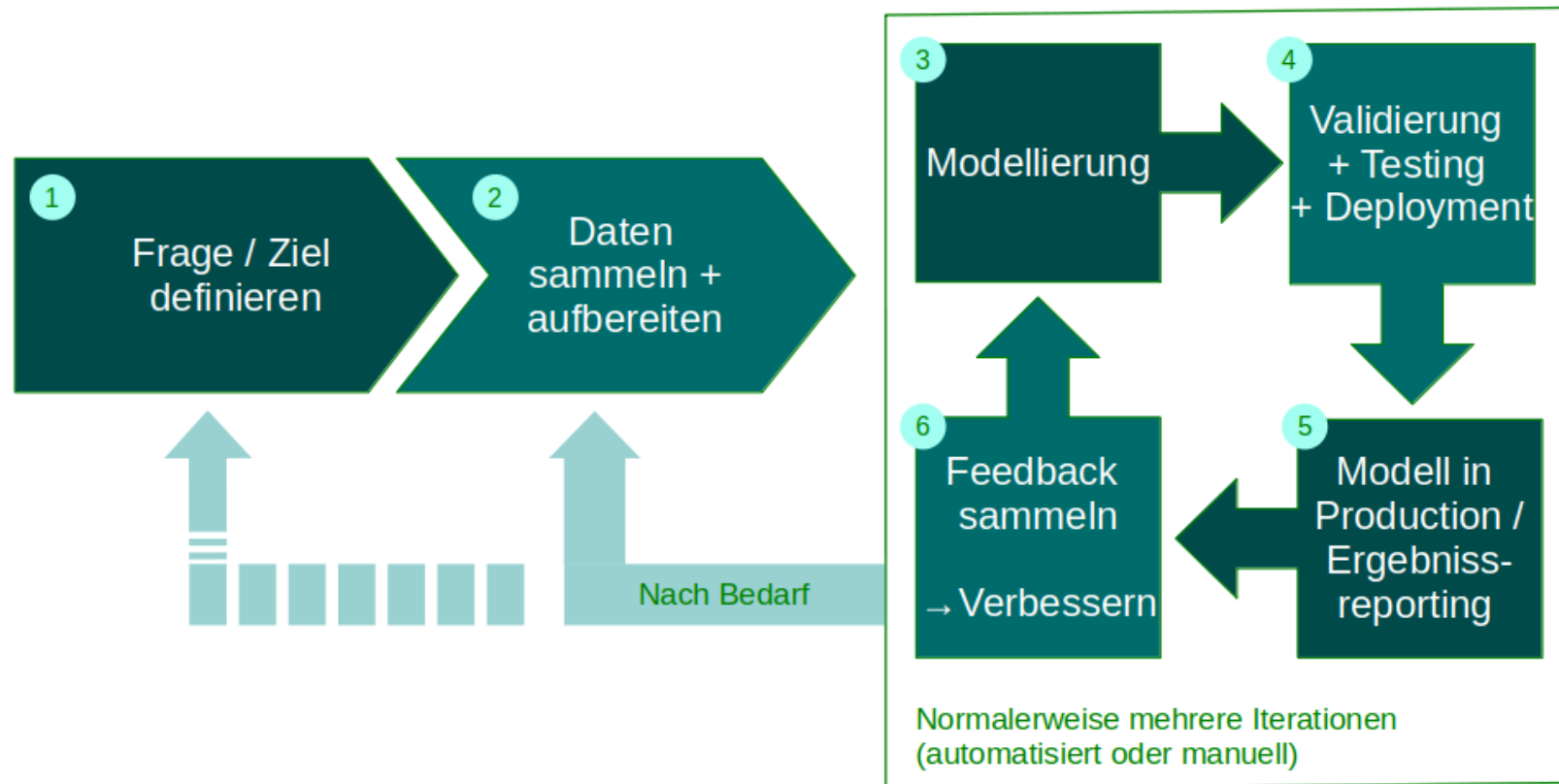
Was ist maschinelles Lernen?



Einführung

Was ist maschinelles Lernen?

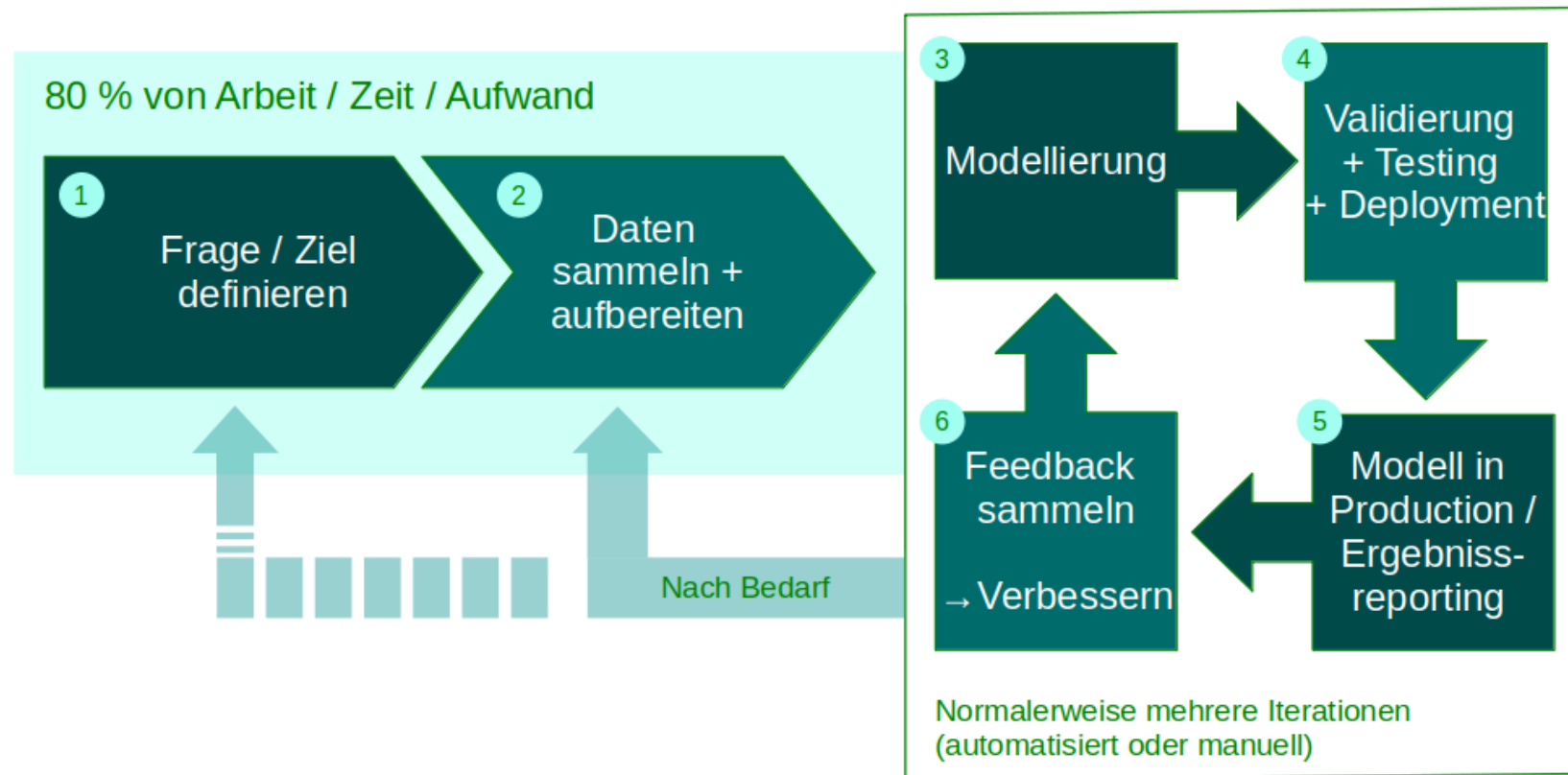
Überblick Data Science Pipeline



Einführung

Was ist maschinelles Lernen?

Überblick Data Science Pipeline



Einführung

Was ist maschinelles Lernen?

ChatGPT 4 (7.11.23) über den Data Science Prozess:

To tackle the problem or the question,
We first set the right direction.

Data big and small we garner,
The quest for answers, we charter.

Prepping data, the next phase,
This step can truly amaze.

Exploration is profound,
Patterns, trends, relations found.

Models built and evaluated,
With predictions that are calculated.

Last is sharing this creation,
Our model's story, its narration.

In the real world to be deployed,
It's the knowledge we enjoyed.

Einführung

Was ist maschinelles Lernen?

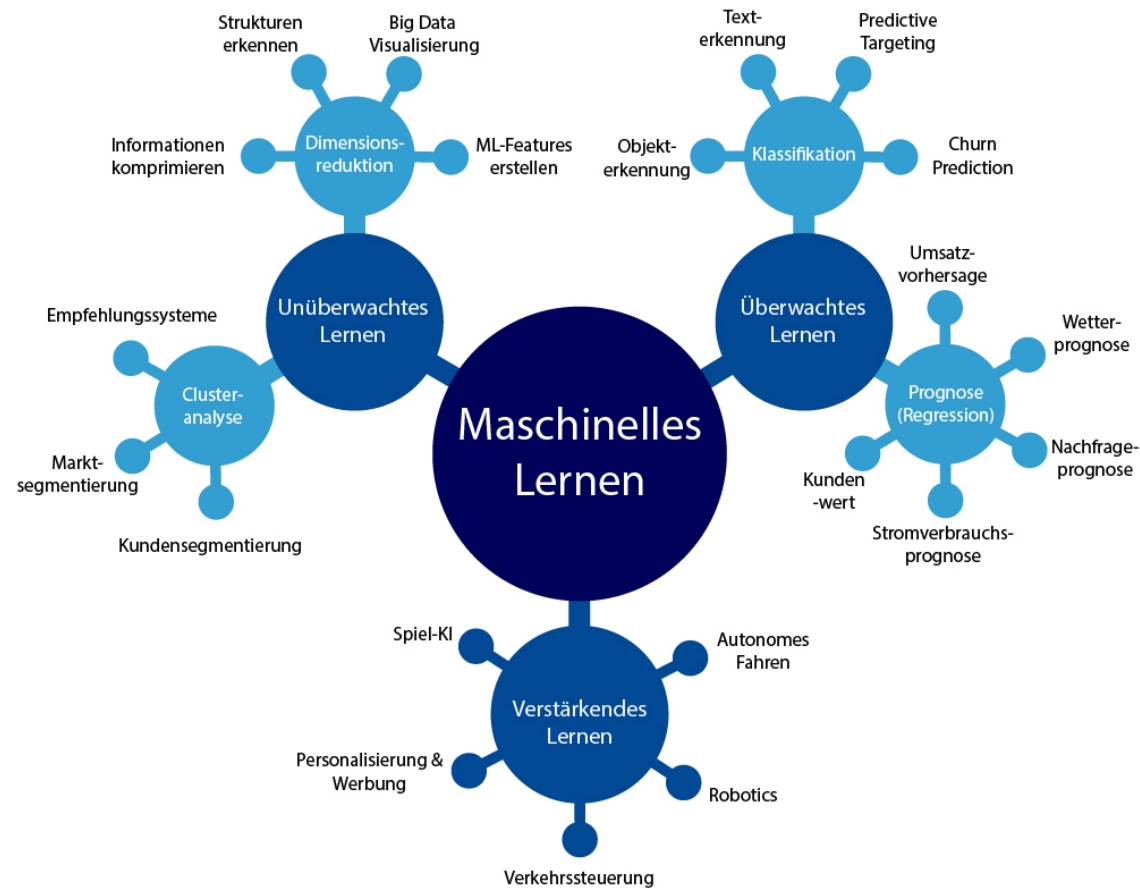


Abbildung 2: Überblick ML mit Anwendungsbeispielen.

Quelle: <https://datasolut.com/was-ist-machine-learning/> (20.02.2023)

Einführung

Was ist maschinelles Lernen?

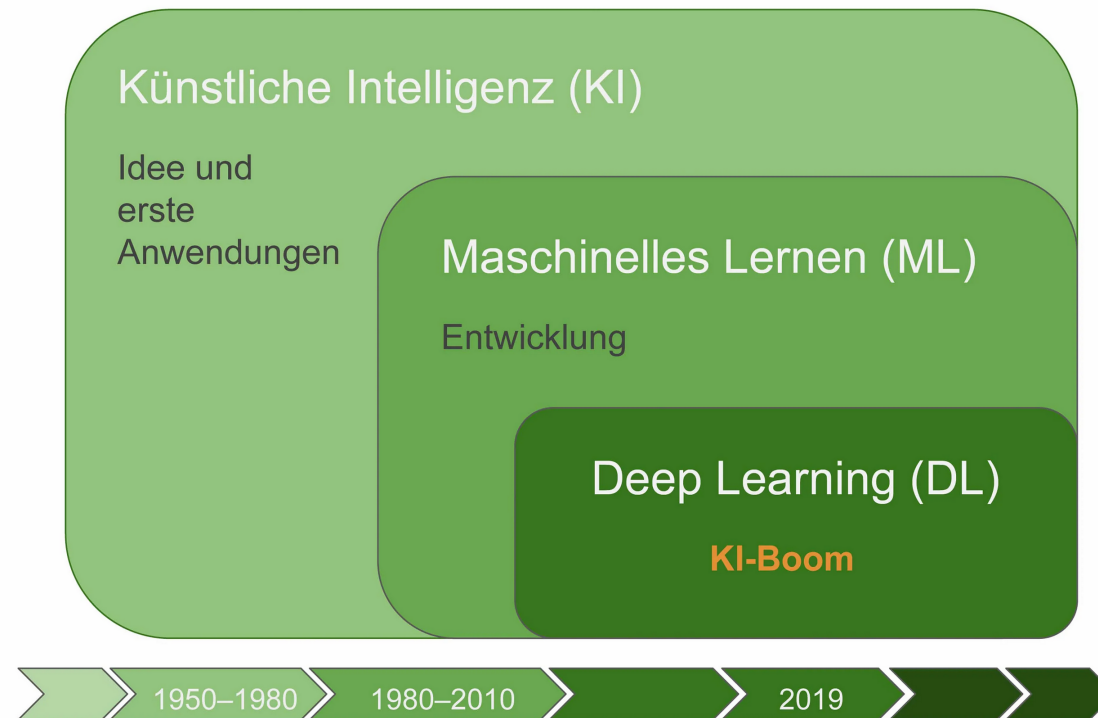


Abbildung 3: Historische Entwicklung von Künstlicher Intelligenz und Machine Learning.

Quelle: [https://medium.com/@friedrich.seck/forschungsfeld-ki-k%C3%](https://medium.com/@friedrich.seck/forschungsfeld-ki-k%C3%9C)

BCnstliche-intelligenz-maschinelles-lernen-deep-learning-und-knn-959c21715b20 (20.02.2023)

Einführung

Was ist maschinelles Lernen?

Ziel des maschinellen Lernens ist es, Verständnis über Daten zu gewinnen und Vorhersagen bzgl. potentiell neuartiger Daten treffen zu können. Grundsätzlich gibt es drei verschiedene Lernmethoden

- ▶ Überwachtes Lernen (Supervised Learning)
- ▶ Unüberwachtes Lernen (Unsupervised Learning)
- ▶ Verstärkendes Lernen (Reinforcement Learning)

In diesem Kurs werden wir uns mit den ersten beiden Methoden beschäftigen.

Einführung

Was ist maschinelles Lernen?

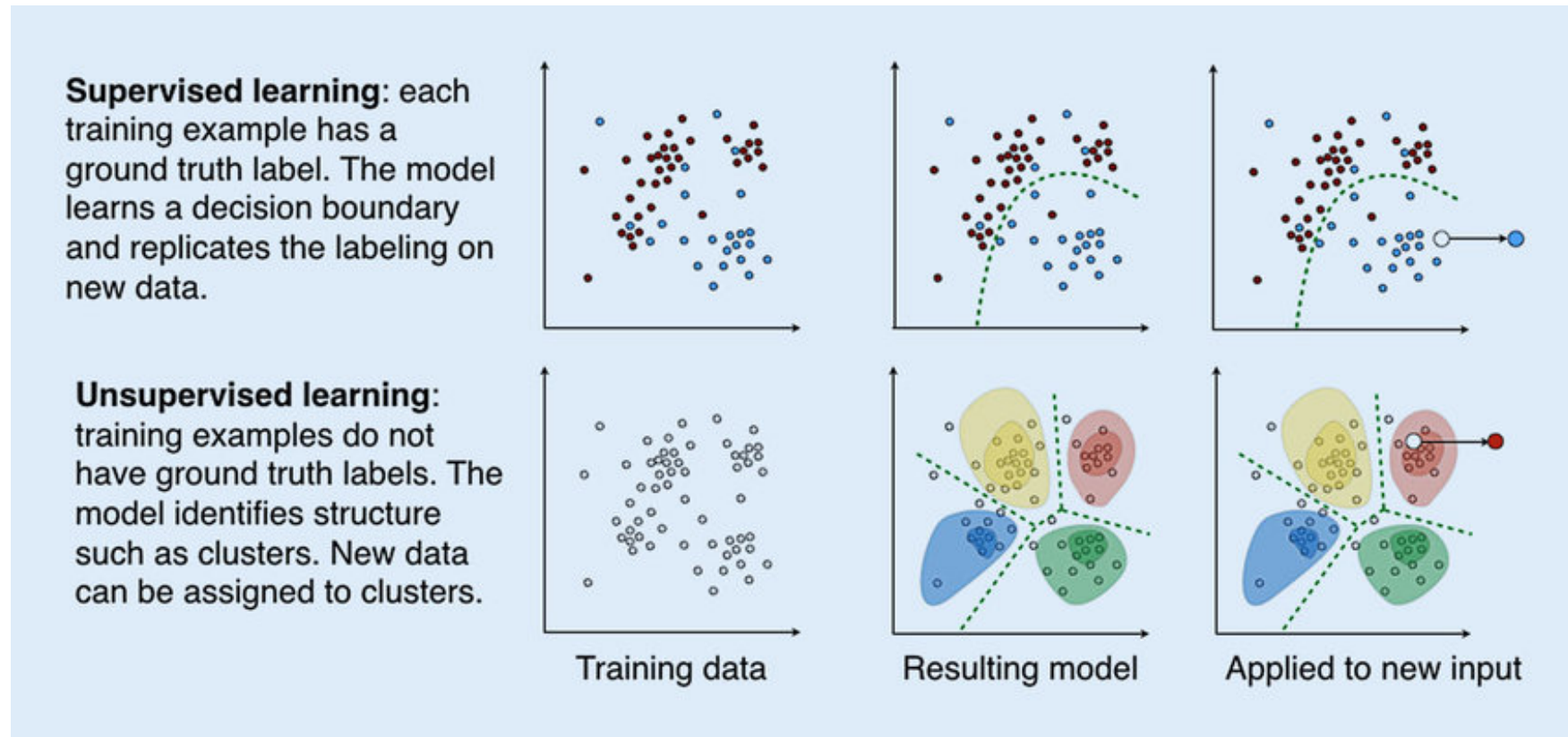


Abbildung 4: Überwachtes und Unüberwachtes Lernen.

Quelle:

https://www.researchgate.net/figure/Supervised-and-unsupervised-machine-learning_fig2_325867536

(20.02.2023)

Einführung

Überwachtes Lernen

Beim **überwachten Lernen** versuchen wir eine Funktion

$$f : \mathcal{X} \rightarrow \mathcal{Y}$$

zu finden, welche den Zusammenhang zwischen den potentiell mehrdimensionalen Mengen \mathcal{X} und \mathcal{Y} möglichst gut repräsentiert, denn meistens werden wir eine perfekte Abbildung aufgrund von statistischen Effekten nicht erreichen. Dabei gibt es zwei Arten von **Fehlern**:

- ▶ **reduzierbar** z.B. durch eine bessere Funktion f
- ▶ **nicht reduzierbar** z.B. aufgrund von Messfehlern in den Daten

Einführung

Überwachtes Lernen

Modell

Wir nennen eine Repräsentation von f mathematisch aber auch als Datenstruktur im Computer **Modell**.

Die Dimensionen von

- ▶ \mathcal{X} werden **Eingabevariablen, Prädiktoren, unabhängige Variablen** oder **Features**
- ▶ \mathcal{Y} werden **Ausgabevariablen, Responses** oder **abhängige Variablen**

genannt.

Einführung

Überwachtes Lernen

Grundsätzlich gibt es beim überwachten Lernen zwei grobe Zielsetzungen zwischen denen meist abgewogen werden muss:

- ▶ **Vorhersage**: Gewünscht ist eine möglichst gute Vorhersage $y = f(\mathbf{x})$ wobei die Funktionsweise von f im Extremfall eine Blackbox sein kann.
- ▶ **Inferenz**: Hier steht die **Interpretierbarkeit** von f im Vordergrund, z.B. Aussagen welche Prädiktoren für welchen Response relevant sind oder auch welcher Zusammenhang (linear, quadratisch, etc.) genau besteht.

Einführung

Überwachtes Lernen

Auch für die Herangehensweise gibt es im Großen und Ganzen zwei Möglichkeiten:

- ▶ **Parametrische** Methoden: Hier wird zunächst eine Annahme bzgl. einer parametrisierten Struktur von f gemacht und diese Parameter werden schließlich mit Hilfe von Daten bestimmt.
- ▶ **Nicht-parametrische** Methoden: Es wird keine Annahme bzgl. der Struktur von f gemacht und es wird versucht f möglichst direkt mit Hilfe von Daten zu definieren.

Einführung

Überwachtes Lernen

Üblicherweise kennen wir die Mengen \mathcal{X} und \mathcal{Y} , aber die genaue Abbildung f können wir trotzdem nur anhand von vielen Beispielen

$$\mathcal{D} = \{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \mid \mathbf{x}^{(i)} \in \mathcal{X}, \mathbf{y}^{(i)} \in \mathcal{Y}, 1 \leq i \leq n\}$$

erahnen.

Trainingsdatensatz

Wir nennen eine solche Menge an Beispielen, die wir für den Lernprozess verwenden **Trainingsdatensatz**.

Wir sprechen bei der Menge \mathcal{D} auch von **gelabelten Daten**. Oft muss ein großer (manueller) Aufwand investiert werden, um an solche Daten zu gelangen.

Einführung

Überwachtes Lernen

Üblicherweise ist \mathcal{X} ein d -dimensionaler reellwertiger Vektorraum, im allgemeinen ist also $\mathcal{X} = \mathbb{R}^d$ für ein $d \in \mathbb{N}$.

Beispiele

- ▶ $\mathcal{X} = \mathbb{R}$: Temperatur in $^{\circ}\text{C}$
- ▶ $\mathcal{X} = \mathbb{R}^2$: Temperatur in $^{\circ}\text{C}$ und Windgeschwindigkeit in $\frac{\text{m}}{\text{s}}$
- ▶ $\mathcal{X} = \mathbb{R}^{16384}$: Graustufenbild 128×128 Pixel (Grauwerte von 0.0 bis 1.0)

Hier wird auch klar, warum wir meist (außer für Beispiele zu Illustrationszwecken) keine einfachen Wertetabellen für f verwenden können.

Einführung

Überwachtes Lernen

Ist \mathcal{Y} eine diskrete Menge, das heißt $\mathcal{Y} = \{C_1, \dots, C_k\}$ für ein $k \in \mathbb{N}$, dann handelt es sich um ein **Klassifikationsproblem**. Bei der Klassifikation sind wir an **qualitativen** Aussagen interessiert. Die einzelnen Objekte C_1, \dots, C_k werden **Klassen** oder **Kategorien** genannt.

Einführung

Überwachtes Lernen

Beispiel: Binäre Klassifikation mit $|\mathcal{Y}| = 2$

Temperaturklassifikation nach menschlichem Empfinden:

$$f : \mathbb{R} \rightarrow \{\text{angenehm}, \text{unangenehm}\}$$

$$f(x) = \begin{cases} \text{angenehm} & \text{falls } x \in [18.0, 25.0] \\ \text{unangenehm} & \text{andernfalls.} \end{cases}$$

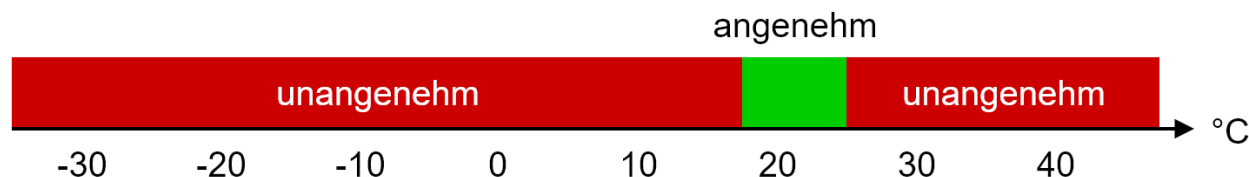


Abbildung 5: Temperaturklassifikation.

Einführung

Überwachtes Lernen

Natürlich kann es wie in der Definition beschrieben auch mehrere Klassen geben.

Beispiel: Mehrklassen-Klassifikation mit $|\mathcal{Y}| = 5$

Temperaturklassifikation nach menschlichem Empfinden:

$$f : \mathbb{R} \rightarrow \{\text{frostig, kalt, angenehm, warm, heiß}\}$$

$$f(x) = \begin{cases} \text{frostig} & \text{falls } x \in (-\infty, 4.0) \\ \text{kalt} & \text{falls } x \in [4.0, 18.0) \\ \text{angenehm} & \text{falls } x \in [18.0, 25.0) \\ \text{warm} & \text{falls } x \in [25.0, 35.0) \\ \text{heiß} & \text{falls } x \in [35.0, \infty) \end{cases}$$

Einführung

Überwachtes Lernen

Ist \mathcal{Y} eine kontinuierliche Menge, das heißt $\mathcal{Y} \subseteq \mathbb{R}$, dann handelt es sich um ein **Regressionsproblem**. Bei der Regression sind wir an **quantitativen** Aussagen interessiert.

Einführung

Überwachtes Lernen

Ein Beispiel einer Regression ist ein *linearer Zusammenhang* zwischen der Temperatur und der Anzahl der verkauften Eiskugeln in einer Eisdiele.

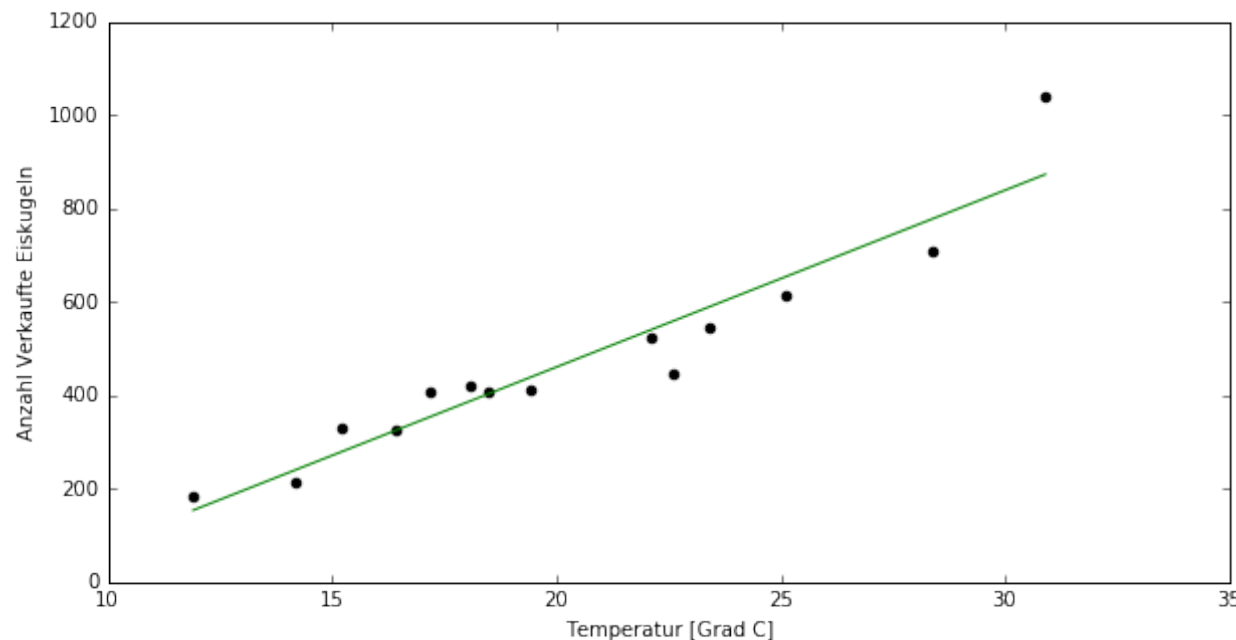


Abbildung 6: Linearer Zusammenhang zwischen der Temperatur und der Anzahl der verkauften Eiskugeln, $f : \mathbb{R} \rightarrow \mathbb{R}$ mit $f(x) = -300 + 40x$.

Einführung

Überwachtes Lernen

Die Ausgabemenge \mathcal{Y} kann prinzipiell auch mehrdimensional sein.

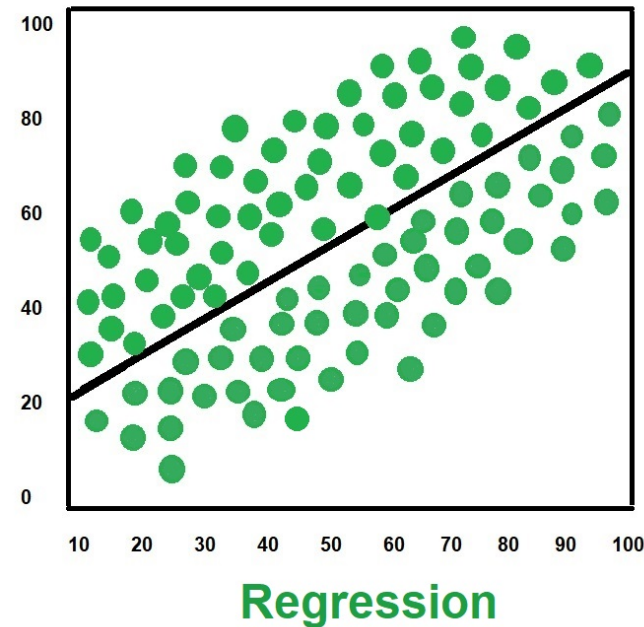
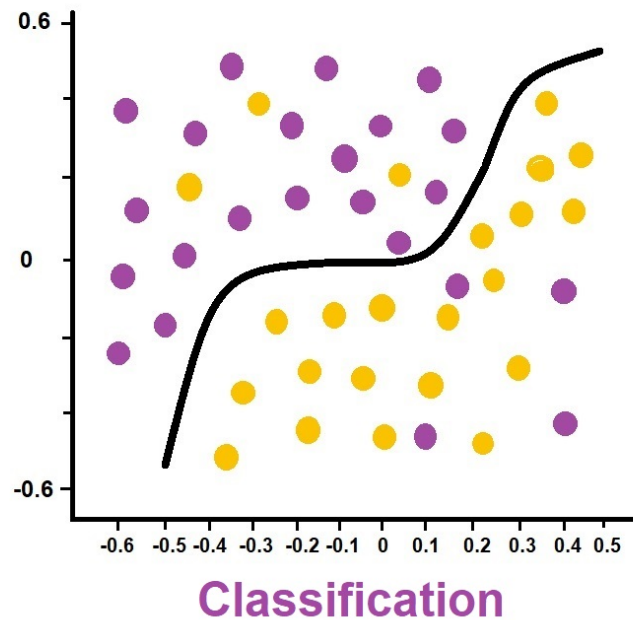
Beispiele

- ▶ $\mathcal{Y} = \{\text{gut, schlecht}\} \times \{\text{günstig, normal, teuer}\}$
- ▶ $\mathcal{Y} = \mathbb{R}^2$: Anzahl verkaufte Eiskugeln, Anzahl verkaufte Pizzen

Einführung

Überwachtes Lernen

Übersicht: Klassifikation (kategoriale Zielgröße) vs.
Regression (metrische Zielgröße)



Quelle: <https://www.ejable.com/tech-corner/ai-machine-learning-and-deep-learning/a-guide-to-linear-regression-and-logistic-regression-in-machine-learning/> (20.02.2023)

Einführung

Unüberwachtes Lernen

Beim **unüberwachten Lernen** versucht man ohne Zuhilfenahme von gelabelten Daten einen Mehrwert zu erhalten. Das Ziel ist daher ausgehend von einer Menge von Daten

$$\mathcal{D} = \{\mathbf{x}^{(i)} \mid \mathbf{x}^{(i)} \in \mathcal{X}, 1 \leq i \leq n\}$$

mehr über die Beschaffenheit von \mathcal{X} herauszubekommen, um dieses Wissen dann direkt oder indirekt anwenden zu können.

Einführung

Unüberwachtes Lernen

Beispiele

- ▶ Lernen der **Verteilung** von \mathcal{X} z.B. bei Sprachmodellen (Welche Wörter folgen auf ein bestimmtes Wort oder einen Satz).
- ▶ **Dimensionsreduktion** zur Verbesserung von überwachten Lernverfahren, z.B. $\mathbf{X} = \mathbb{R}^{10}$ statt $\mathbf{X} = \mathbb{R}^{100}$ für $f : \mathcal{X} \rightarrow \mathcal{Y}$
- ▶ Finden von Ähnlichkeitsstrukturen durch **Clustering**

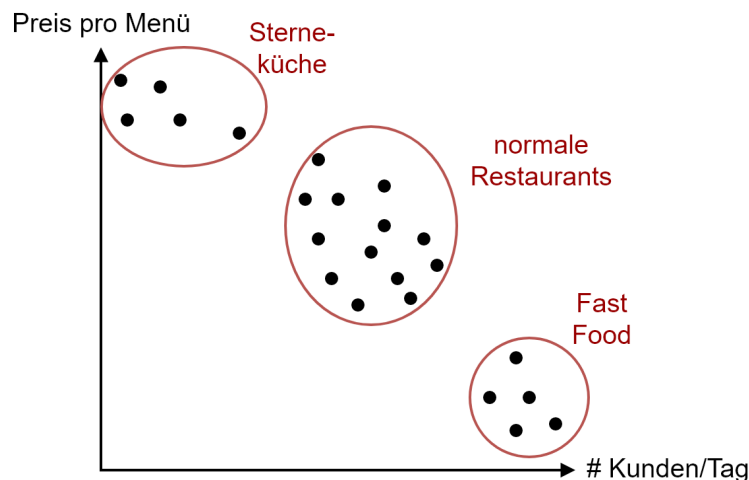


Abbildung 7: Clustering von Restaurants.

Einführung

Datenvisualisierung

Wenn man ein Projekt mit maschinellen Lernmethoden beginnt, ist es ratsam, sich zunächst einen **Überblick** über die Daten zu verschaffen. Meist gelingt dies am besten, wenn man die Daten geeignet **visualisiert**. Im Folgenden finden Sie einige Beispiele verschiedener Diagrammtypen.

Einführung

Datenvisualisierung

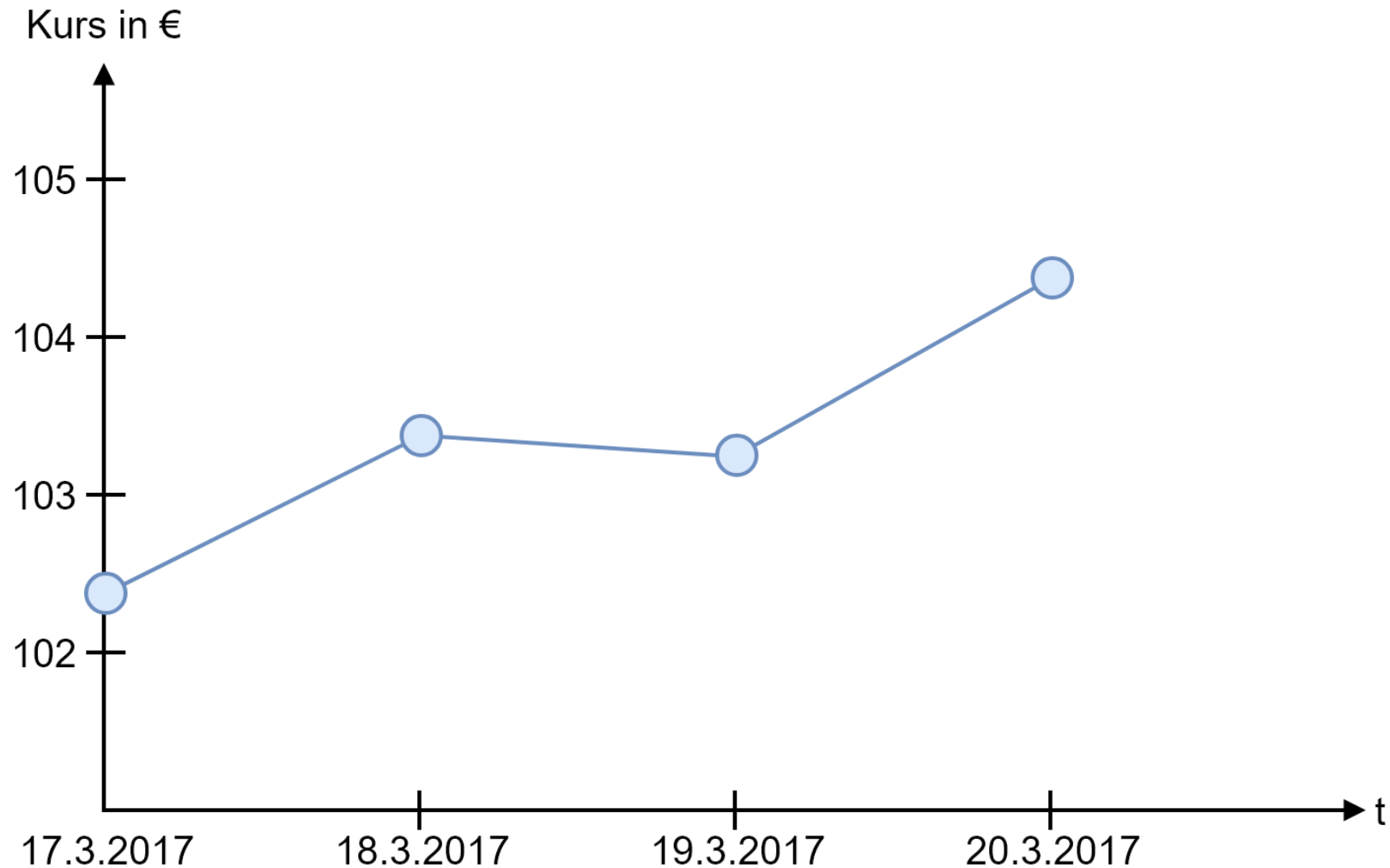


Abbildung 8: Beispiel eines Liniendiagramms.

Einführung

Datenvisualisierung

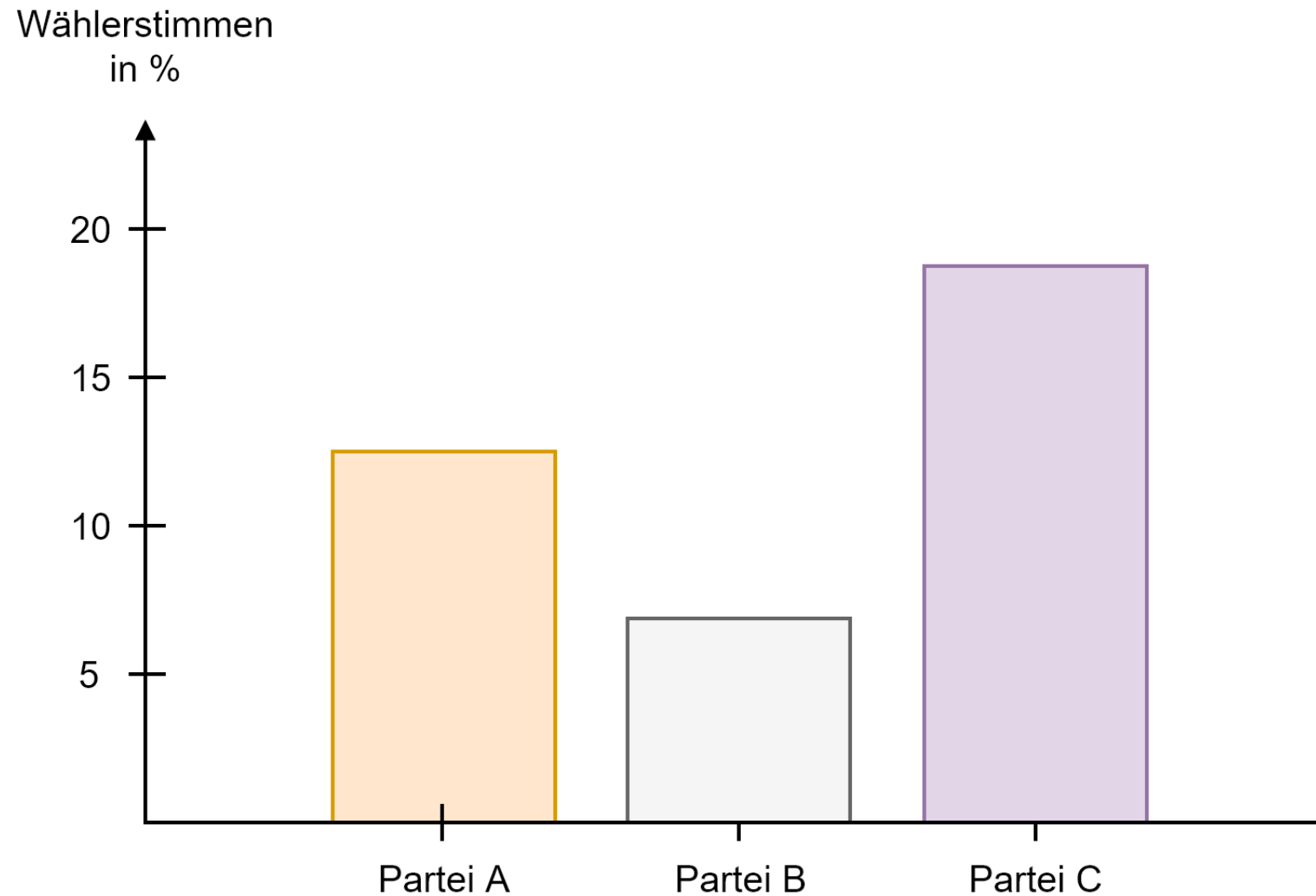


Abbildung 9: Beispiel eines Balkendiagramms.

Einführung

Datenvisualisierung

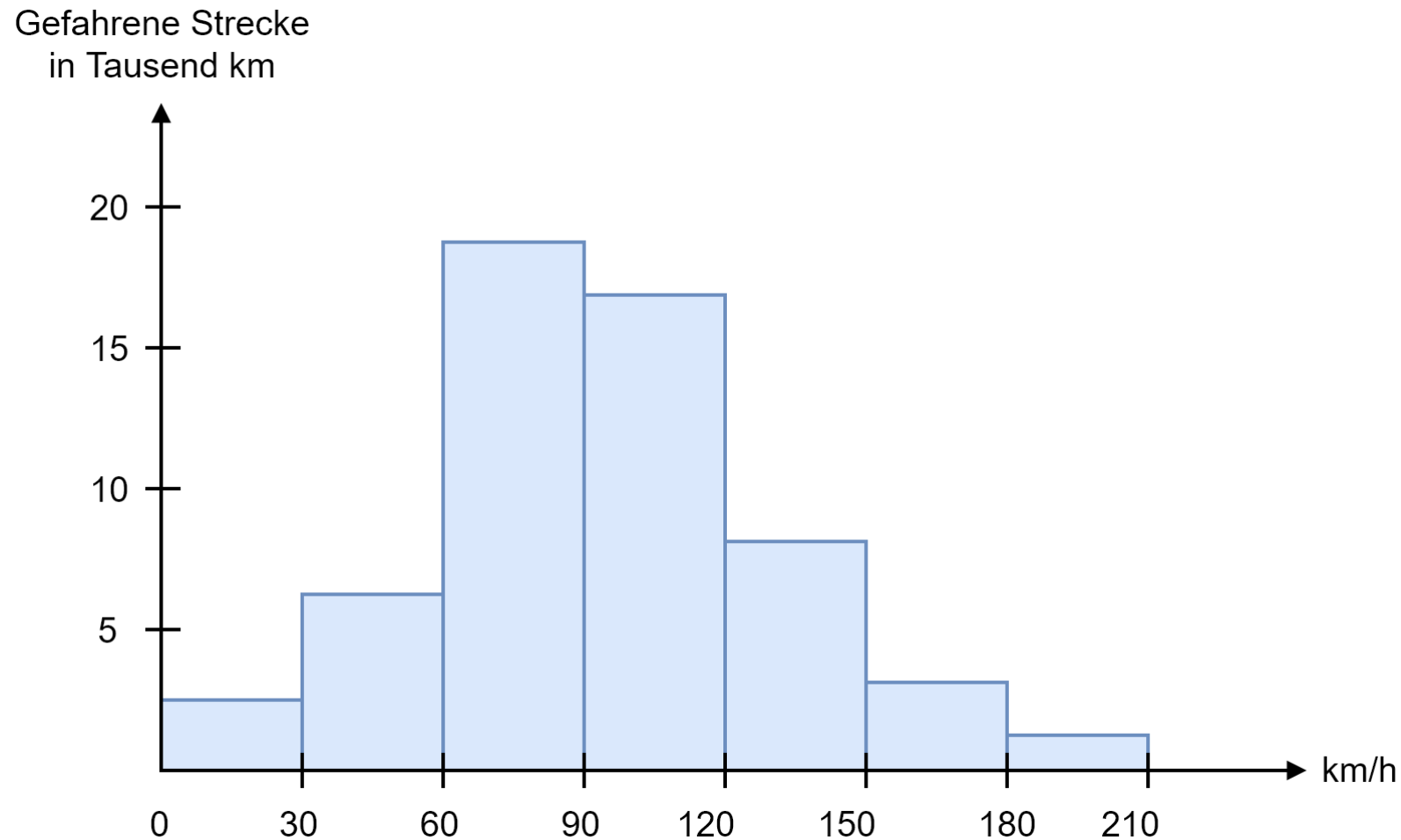


Abbildung 10: Beispiel eines Histogramms – eines speziellen Balkendiagramms.

Einführung

Datenvisualisierung

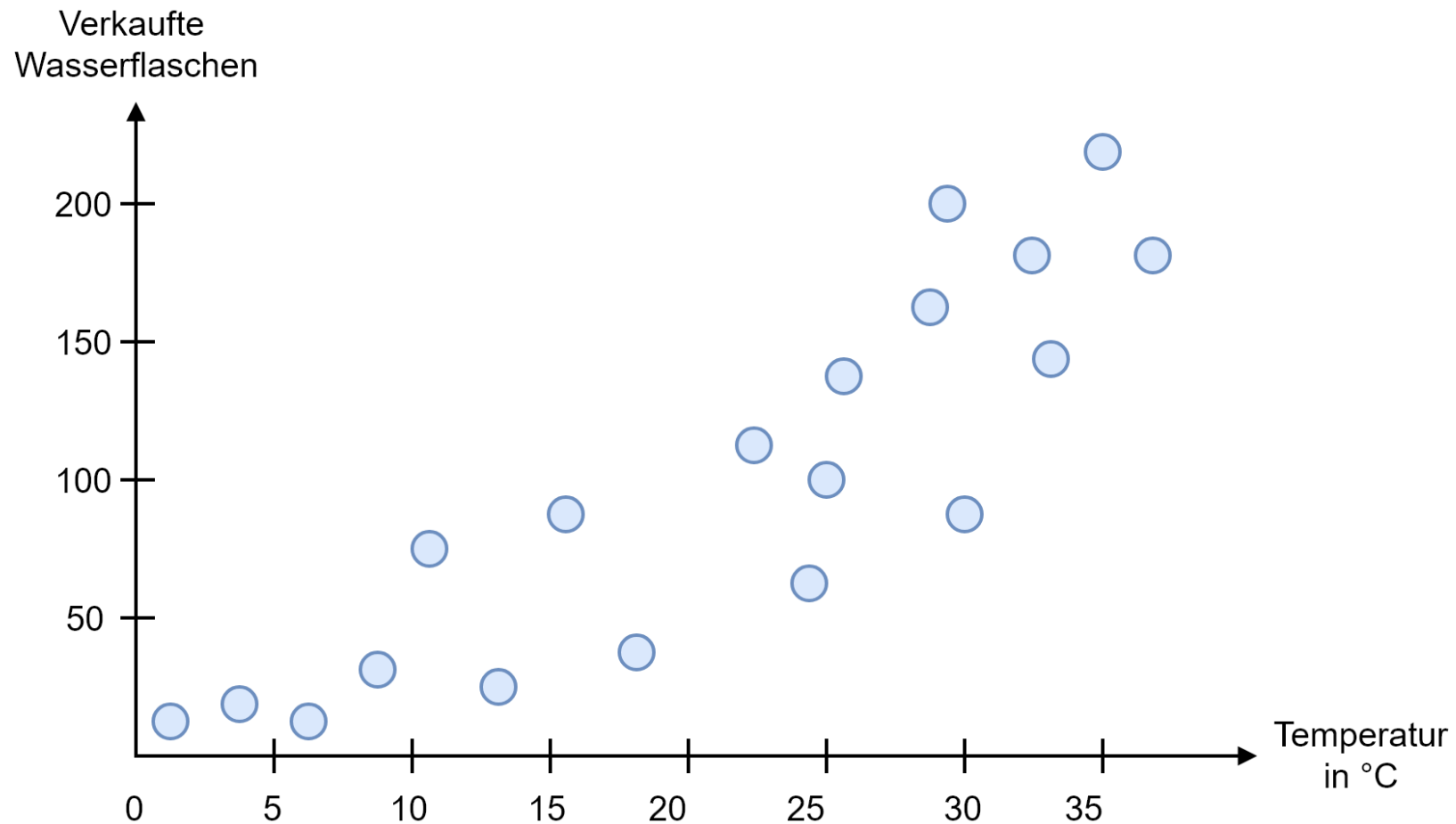


Abbildung 11: Beispiel eines Streudiagramms.

Einführung

Datenvorverarbeitung

Bevor tatsächlich ein ML Modell erstellt und trainiert wird, müssen die entsprechenden Daten **vorverarbeitet** werden. Dazu gehören grundsätzlich drei Schritte

1. Auswahl
2. Aufbereitung
3. Transformation

der Daten. Oftmals muss auch aufgrund neuer Erkenntnisse zwischen den Schritten hin und her gewechselt werden.

Einführung

Datenvorverarbeitung

Auswahl: Nicht immer sind mehr Daten auch wirklich besser, d.h. es sollte darauf geachtet werden, dass nur für den Anwendungszweck **relevante** Daten verwendet werden, um die Rechen- und Speicheranforderungen im Rahmen zu halten. Auch die Leistung des Systems könnte u.U. unter zu vielen bzw. den falschen Daten leiden – natürlich auch unter zu wenig.

Einführung

Datenvorverarbeitung

Fragestellungen, die bzgl. der Auswahl helfen:

- ▶ Auf welche Daten hat man Zugriff?
- ▶ Welche Daten kann man mit welchem Aufwand erstellen bzw. simulieren?
- ▶ Auf welchen Teil der Daten kann/sollte man verzichten?

Starthilfe

Im Rahmen von Wettbewerben und Benchmarks werden immer wieder Datensätze veröffentlicht, die zum Lernen von ML Techniken verwendet werden können. Ein Beispiel ist <https://www.kaggle.com/datasets>.

Einführung

Datenvorverarbeitung

Aufbereitung:

- ▶ **Definition** eines geeigneten Formats (Tabellen, Big Data Formate wie Parquet, CSV, Bilder, etc.) und **Umwandlung** der Daten
- ▶ **Bereinigung**, d.h. Entfernung von *unvollständigen* oder *ungültigen* Daten oder aufgrund von rechtlichen Bestimmungen (Datenschutz)
- ▶ **Unterauswahl** der Daten (lange Laufzeit, großer Speicheraufwand). Hier muss auf eine *repräsentative* Auswahl (Zeit, Ort, Gruppen, etc.) geachtet werden, um keinen systematischen Fehler einzuführen.

Einführung

Datenvorverarbeitung

Transformation:

- ▶ **Skalierung**: Features in den geeigneten Wertebereich für ML Methode bringen, z.B. auf Wertebereiche $[0, 1]$ oder $[-1, 1]$. Auch eine Normierung auf Mittelwert 0 und Standardabweichung 1 kann notwendig sein.
- ▶ **Zerlegung** in sinnvolle Features, z.B. Extraktion der Zeit und des Fehlercodes aus Logfile-Einträgen
- ▶ **Aggregation** mehrerer Features, z.B. Gesamtzahl der Aktienverkäufe an einem Tag statt jede Einzeltransaktion