

Maschinelles Lernen verstehen. Ein Einstieg.

Maschinelles Lernen verstehen. Ein Einstieg.



Dr. Christian Richter
Referent für Medienbildung (und digitale Kompetenzen)
am Landesinstitut für Schule und Medien Berlin-Brandenburg (LISUM)
Medienwissenschaftler

In Zusammenarbeit mit Cornelia Brückner (LISUM)

Beispiele für KI im Alltag

- Empfehlungen bei Amazon und Spotify und Netflix, Chatbots,
- Suchanfragen bei Google, Navigation
- Gesichts- und Bilderkennung
- Empfehlungen medizinischer Behandlungen
- Kreditvergaben, Auswahl von Bewerber:innen
- Einschätzung Gefahr einer wiederholten Straffälligkeit von Kriminellen

Maschinelles Lernen

- Maschinen finden ohne explizite Programmierung von Regeln selbstständig Muster in Datensätzen
- treffen auf Basis dieser Analyse dann Vorhersagen und Entscheidungen

Voraussetzung ist die Verarbeitung großer Datenmengen

„Lernstile“

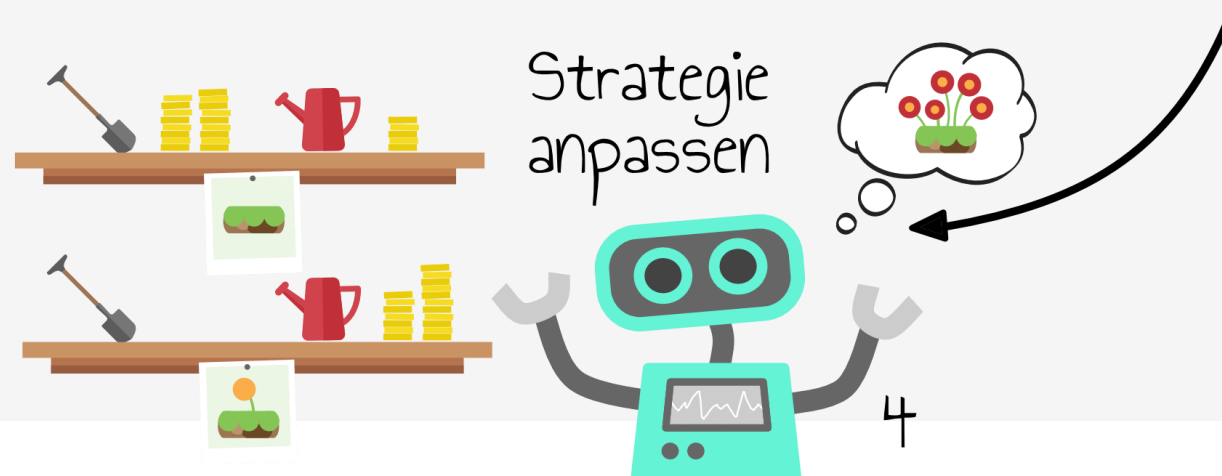
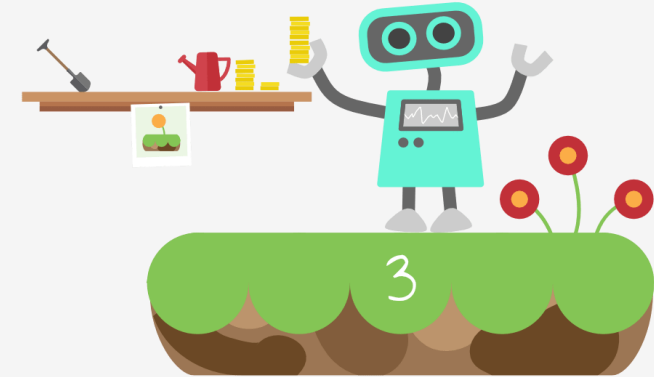
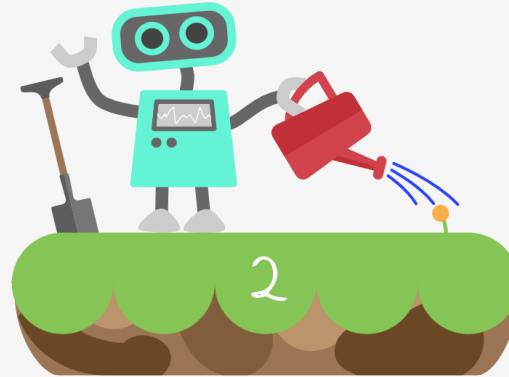
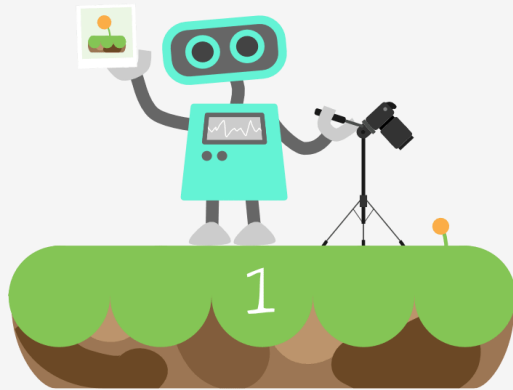
- Reinforced Learning (verstärkendes L)
- Supervised Learning (überwachtes L.)

Reinforced Learning

Zustand erfassen

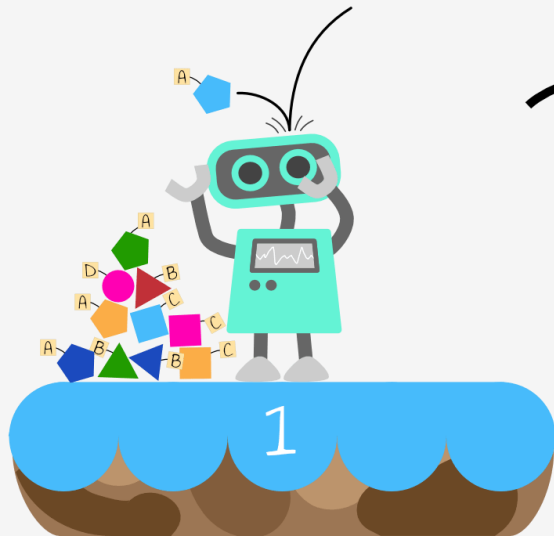
Aktion wählen und durchführen

Belohnung oder Bestrafung erhalten

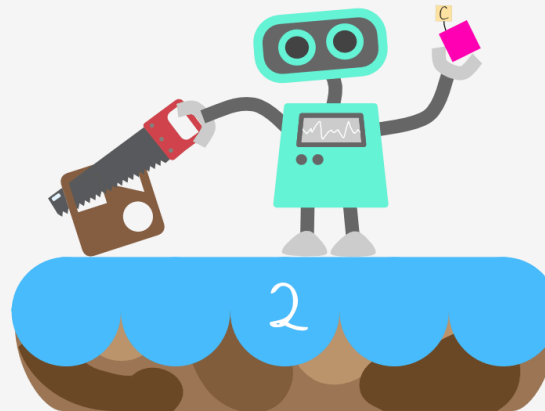


Supervised Learning

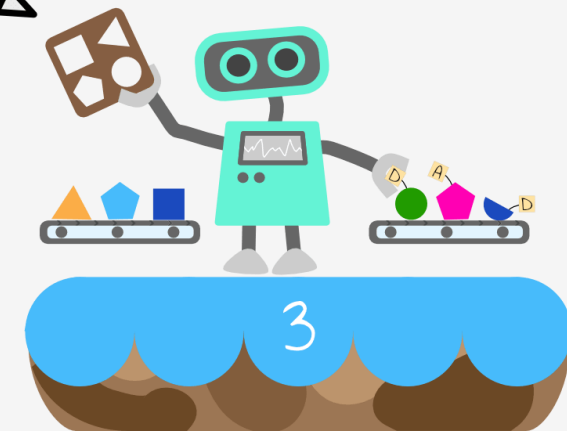
Beschriftete Eingaben
erhalten



Regeln finden, die
bekannte Eingaben
richtig beschriften



Neue Eingaben
entsprechend der gefundenen
Regeln beschriften



Supervised Learning

- „beschriftete Daten“
- Klassifikation nach Prinzip: „a“ oder „nicht a“
- in der Trainingsphase wird selbständig eine Verbindung zwischen Trainingsdaten und Beschriftung hergestellt
- Ergebnisse werden als Wahrscheinlichkeiten ausgegeben

Datensatz der Passagiere auf der Titanic

- Mit folgenden Angaben

SURVIVED	Überlebensvariable mit 0=Nein und 1=Ja
PCLASS	Klasse auf dem Schiff 1= Erste, 2=Zweite, 3= Dritte
NAME	Nachname, Title, Vorname(n)
SEX	Geschlecht „male“=männlich, „female“=weiblich
AGE	Alter in Jahren
SIBSP	Anzahl der Geschwister bzw. Partner, die mit an Bord waren
PARCH	Anzahl der Eltern bzw. Kinder, die mit an Bord waren
TICKET	ID Nummer des Tickets
FARE	Preis des Tickets
CABIN	Kabinennummer
EMBARKED	Einstiegshafen C= Cherbourg, Q= Queenstown, S=Southampton

Datensatz der Passagiere auf der Titanic



KI-generiertes Bild

Wie waren die Überlebenschancen von Jack & Rose?

Datensatz der Passagiere auf der Titanic



Wie waren die Überlebenschancen von Jack & Rose?

Modell prognostizierte:

Jack stirbt = 88,9% Wahrscheinlichkeit

Rose überlebt = 96,5% Wahrscheinlichkeit

Teachable Machine

The screenshot displays the Teachable Machine interface. On the left, two class cards are visible: 'Kuh' (Cow) and 'Schaf' (Sheep). Each card shows '5 Bild-Beispiele' (5 image examples) and options for 'Webcam' and 'Hochladen' (Upload). A 'Klasse hinzufügen' (Add class) button is at the bottom. In the center, a 'Training' panel indicates 'Modell ist trainiert' (Model is trained) and 'Erweitert' (Advanced) settings. On the right, a 'Vorschau' (Preview) panel shows 'Modell exportieren' (Export model), an 'Eingabe' (Input) section with a toggle for 'AN' (ON) and a 'Datei' (File) dropdown, and an 'Ausgabe' (Output) section. The output shows a 99% confidence for 'Kuh' and a lower confidence for 'Schaf'. A large image of a sheep is shown in the preview area.

<https://teachablemachine.withgoogle.com/>

Verzerrung im System

Bias und Fairness in algorithmischen Entscheidungssystemen

Unfaire Systeme

- treffen Entscheidungen, bei denen Personen aufgrund ihrer angeborenen oder erworbenen Eigenschaften bevorzugt oder benachteiligt werden
- eine Form von Diskriminierung

Begriffe

Faire Systeme

- stellen eine Gleichbehandlung durch KI-Systeme sicher
- Funktionieren insbesondere ohne Diskriminierung auf Basis von ethnischer Herkunft, Geschlecht, Alter, Religion /Weltanschauung oder sonstiger Indikatoren
- relevant vor allem für KI-Systeme, die Entscheidungen über Personen treffen

Begriffe

Bias

- Verzerrungseffekt in der Datenerhebung oder – verarbeitung
- kann Entscheidungen oder Wahrnehmungen des Systems beeinflussen
- bewusst und unbewusst
- z.B. durch Stereotypen

Erste Annäherung

- Facebook und LinkedIn setzen KI-Systeme ein, um fragwürdige Inhalte zu filtern. (z.B. gewaltverherrlichende, pornografische oder politisch extreme Bilder, Texte und Videos)
- Welche Schwierigkeiten könnten sich hierbei ergeben?

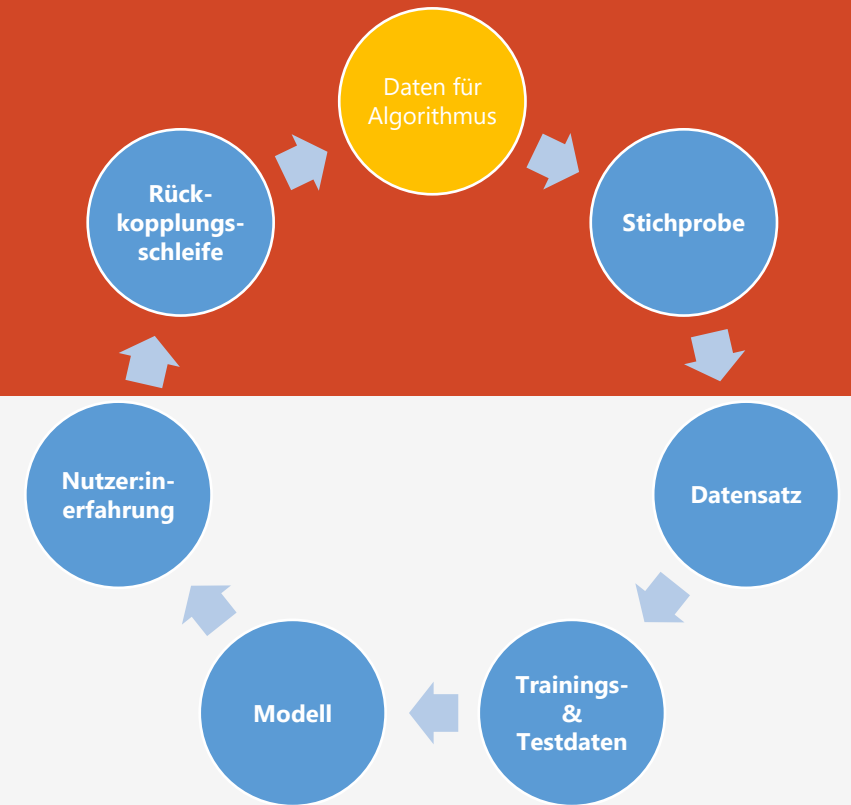
ML-Prozesses



Arten von Bias

Historical Bias

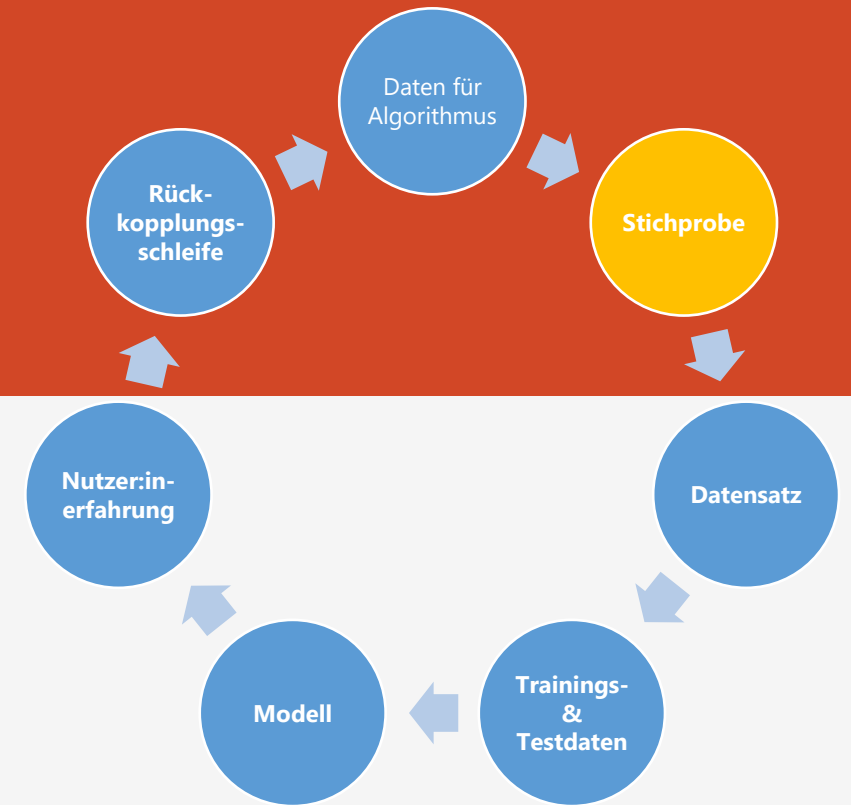
- historische Diskriminierungsmuster
- sind diese Muster im Datensatz enthalten, reproduzieren Systeme diese Muster
- Personen, die in der Vergangenheit Schwierigkeiten hatten, werden durch System weiter benachteiligt
- Bsp.: People of Color in US-Kriminalstatistiken



Arten von Bias

Representation Bias

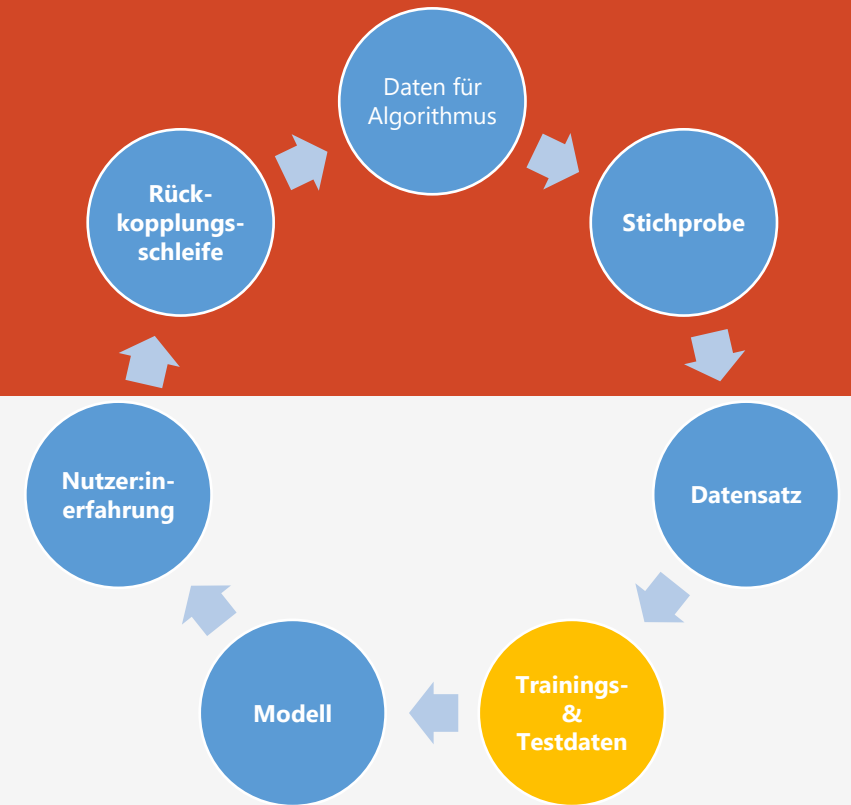
- Nutzung von nicht repräsentativen Stichproben
- Unterrepräsentierte Bereiche verschwinden im System
- Bsp.: beim Training von StableDiffusion & GPT3 wurden nur englischsprachige Quellen genutzt



Arten von Bias

Evaluation Bias

- Testdaten sind nicht repräsentativ
- KI wird unfair trainiert
- Bsp.: Vergleichsdatensatz von Gesichtserkennungssystemen besteht hauptsächlich aus hellhäutigen Menschen



Diskussionsbeispiele

- Ein Unternehmen entwickelte ein System, um aus den vielen eingehenden Bewerbungen die geeignetsten automatisiert herauszufiltern. Dafür wurde das System mit den Daten der bereits eingestellten Bewerber:innen trainiert. So sollte es lernen, welche Eigenschaften von der Leitung bei Einstellungen bevorzugt werden. Bei einer Überprüfung kam heraus, dass Frauen durch das System diskriminiert wurden. Wie kam das zustande?

Diskussionsbeispiele

- Ein anderes Unternehmen entwickelte ebenfalls ein System, um aus den vielen eingehenden Bewerbungen die geeignetsten automatisiert herauszufiltern. Das System stellt fest, dass Arbeitnehmende, die weiter entfernt vom Firmensitz arbeiten, häufig kündigen. Daher werden Bewerbungen von Personen bevorzugt, die kurze Anfahrtswege haben. Welche Folge hatte diese Entscheidung?

Content-Filter

- Facebook und LinkedIn setzen KI-Systeme ein, um fragwürdige Inhalte zu filtern. (z.B. gewaltverherrlichende, pornografische oder politisch extreme Bilder, Texte und Videos)
- Welche Schwierigkeiten könnten sich hierbei ergeben?

Ermittlung der Fairness

- Zuordnung der Ergebnisse nach Korrektheit
- also die Zuordnung der Anzahl der richtig und falsch klassifizierten Testdaten
- Testdaten müssen vorher verifiziert sein, Zuordnung muss vorab zweifelsfrei feststehen
- Darstellung in Confusion Matrix

Ermittlung der

		Vom System zugeordnet	
		gutartig	bösartig
Tatsächliche Zuordnung	gutartig	True Positives	False Negatives
	bösartig	False Positives	True Negatives

Wie kann man Fairness gewährleisten?

- Diversität in Programmier-Teams
- Transparenz des KI-Systems, sodass Menschen es nachvollziehen können
- Einsatz von Algorithmen zur Korrektur

Der Abschnitt „Verzerrung im System“ basiert zu großen Teilen auf dem Online-Angebot:

Bias & Fairness algorithmischer Entscheidungssysteme
des **Instituts für Business Analytics** der Universität Ulm

Prof. Dr. Mathias Klier

<https://bias-and-fairness-in-ai-systems.de/>

Weitere Quellen

- Beispiele zu Bias in Bewerbungsprozessen:
<https://www.zeit.de/arbeit/2018-10/bewerbungsroboter-kuenstliche-intelligenz-amazon-frauen-diskriminierung/komplettansicht>
- Schlechtes Abschneiden von People of Color bei Gesichtserkennungssoftware
<http://gendershades.org/>
- COMPAS - algorithmische Systeme, um Gefahr der wiederholten Straffälligkeit von Kriminellen einzuschätzen
<https://netzpolitik.org/2017/radiofeature-algorithmen-als-schicksalsmaschinen/>
- Verzerrungen bei Kreditvergabe (Beispiel für Historical Bias)
<https://safe-frankfurt.de/de/aktuelles/safe-finance-blog/details/ki-basiert-ermittelte-kreditausfallrisiken-sind-mit-vorsicht-zu-geniessen.html>
- Titanic-Datensatz - Überleben Rose und Jack?
<https://www.eoda.de/wissen/blog/oeffentlich-verfuegbare-datensaetze-titanic/>

Materialien

verwendetes Material

<https://computingeducation.de/proj-it2school/>

weitere Materialien:

<https://www.medien-in-die-schule.de/unterrichtseinheiten/machine-learning-intelligente-maschinen/modul-2-wie-funktioniert-machine-learning/>

Deep Learning

https://www.science-on-stage.de/sites/default/files/material/machine-learning-in-der-schule_3.-auflage.pdf

nach Datensätzen suchen <https://datasetsearch.research.google.com/>

Lizenzhinweis:

„Maschinelles Lernen verstehen. Ein Einstieg.“

von Christian Richter und Cornelia Brückner

ist, wenn nicht anders gekennzeichnet, lizenziert unter [CC-BY 4.0](https://creativecommons.org/licenses/by/4.0/)