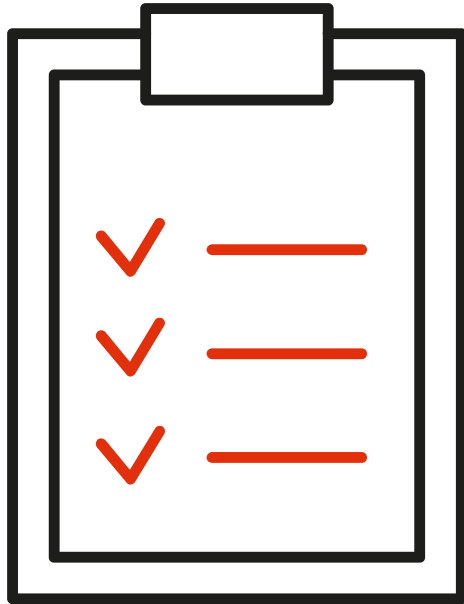


Vertrauenswürdige KI-Anwendungen

Thementag am 03.05.2023

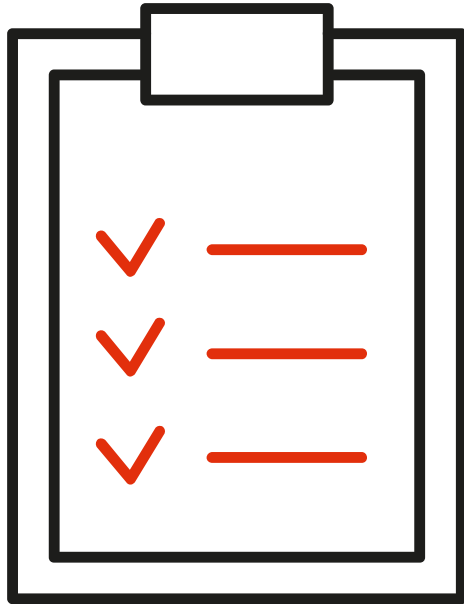
Martin Folz M.Sc., Dr.-Ing. Alexander Dementyev, Ass. iur. Michael Rätze

Agenda



- Vorstellung des Mittelstand-Digitalzentrum Chemnitz
- Einführung in erklärbare KI
- Technische Möglichkeiten der XAI mit Fokus auf virtuelle Sensorik
- Beleuchtung und Diskussion über rechtliche Auswirkungen beim Einsatz von KI

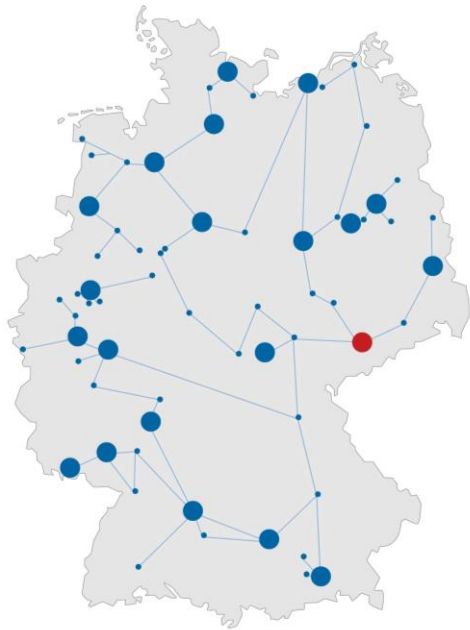
Agenda



- Vorstellung des Mittelstand-Digitalzentrum Chemnitz
- Einführung in erklärbare KI
- Technische Möglichkeiten der XAI mit Fokus auf virtuelle Sensorik
- Beleuchtung und Diskussion über rechtliche Auswirkungen beim Einsatz von KI

Das Mittelstand-Digital Netzwerk

Regionale Zentren und Zentren mit Themenschwerpunkten



- Mit dem Mittelstand-Digital Netzwerk unterstützt das Bundesministerium für Wirtschaft und Klimaschutz die Digitalisierung in kleinen und mittleren Unternehmen.
- Das Zentrum in Chemnitz fokussiert die Digitalisierung in sächsischen Betrieben.
- bundesweite Vernetzung

Unsere Experten

Starkes Partnernetzwerk vor Ort

Mittelstand-Digital Zentrum Chemnitz
Geschäftsstelle c/o TU Chemnitz

Unterstützer
Verbände, Kammern,
öffentliche Hand



Konsortium


TECHNISCHE UNIVERSITÄT
IN DER KULTURHAUPTSTADT EUROPAS
CHEMNITZ

- Prof. Fabrikplanung und Intralogistik
- Prof. Arbeitswissenschaft und Innovationsmanagement
- Prof. für Privatrecht und Recht des geistigen Eigentums

 **Fraunhofer**
IWU

an den Standorten
Chemnitz und Dresden



 **IHK** Industrie- und Handelskammer
Chemnitz

 **WfE** WIRTSCHAFTS
FÖRDERUNG
ERZGEBIRGE

 **tti**
TECHNOLOGIETRANSFER UND
INNOVATIONSFÖRDERUNG
MAGDEBURG GMBH

Unser Ziel: Digitalisierung unterstützen

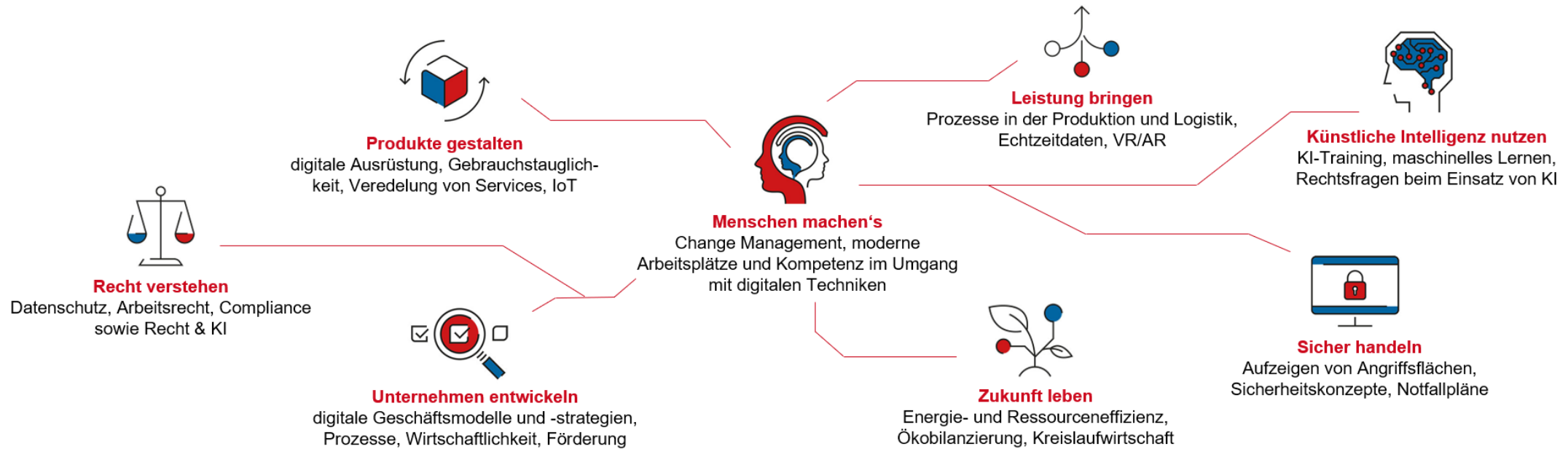
Sächsischer Mittelstand in Industrie, Handwerk und Handel als Zielgruppe

→ Kostenfreier und anbieterneutraler Wissens- und Technologietransfer

- Veranstaltungen wie Workshops, Seminare und Expertenrunden anbieten
- Fachwissen zugänglich machen
- Digitalisierungsprojekte begleiten
- Trainings- und Testumgebungen zur Verfügung stellen
- Lösungen mit Hilfe von Demonstratoren veranschaulichen
- Erfahrungsaustausch zwischen Unternehmen fördern

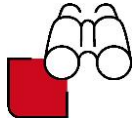
Unsere Themen

Im Fokus steht der Mensch



Praxisnahe Unterstützung

Angebote für Einsteiger in die Digitalisierung und Erfahrene



Potentiale entdecken

- Potentialanalysen
- Selbstchecks und Reifegradmodelle



Wissen vermitteln

- Technologie- und Trendthemen
- Unternehmerisches Fachwissen
- Sprechstunden
- Erfolgsgeschichten



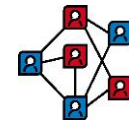
Mitarbeitende qualifizieren

- Interaktive Basis- und Fachworkshops
- Onlineseminare
- Selbstlernangebote
- Thementage



Projekte begleiten

- Potentialanalysen
- Einführung von digitalen Prozessen und Technologien
- Entwicklung digitaler Geschäftsmodelle unterstützen



Netzwerk ausbauen

- Partnernetzwerk
- Unternehmerforen
- Erfahrungskreise



Projekte begleiten

Impulsprojekte



- wenige Wochen Projektdauer
- Potenziale finden
- Impulse setzen
- Strategieentwicklung unterstützen

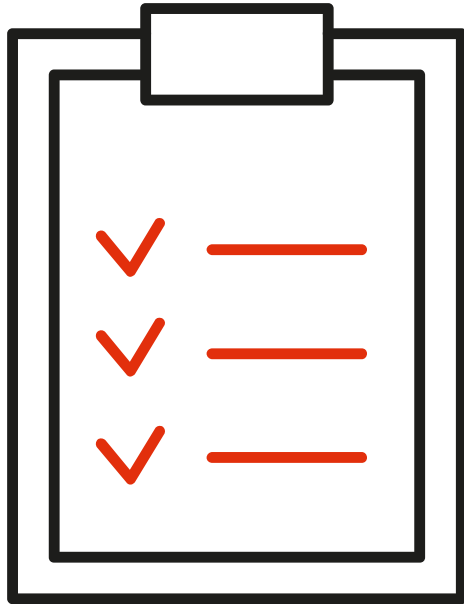
Projekte begleiten

Digitalisierungsprojekte



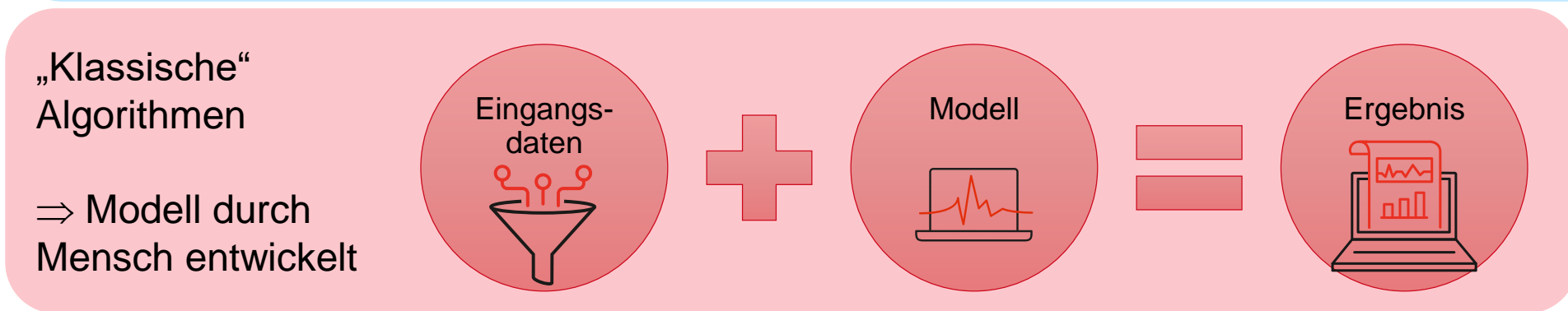
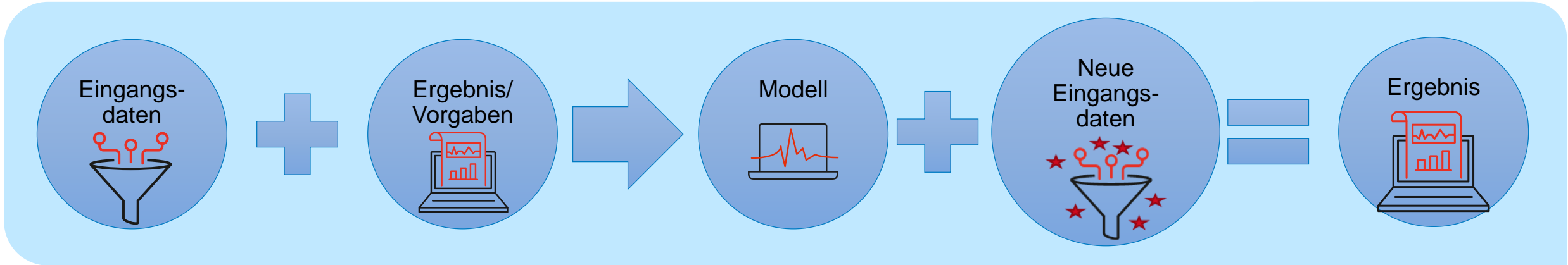
- bis zu 5 Monate Projektdauer
- Ist-Zustand gemeinsam analysieren
- Lösungskonzept (und Prototyp) gemeinsam entwickeln
- Ergebnisse dokumentieren / Lastenheft erstellen
- öffentliche Berichterstattung

Agenda



- Vorstellung des Mittelstand-Digitalzentrum Chemnitz
- **Einführung in erklärbare KI**
- Technische Möglichkeiten der XAI mit Fokus auf virtuelle Sensorik
- Beleuchtung und Diskussion über rechtliche Auswirkungen beim Einsatz von KI

Vergleich ML zu klassischen Algorithmen



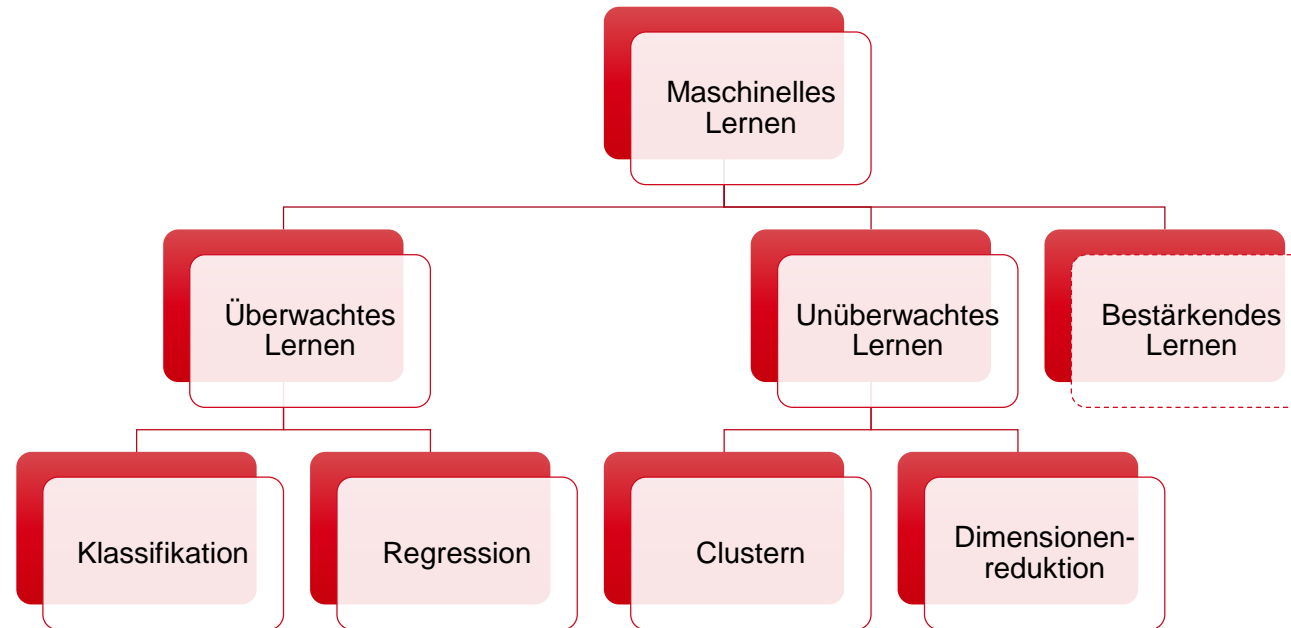
„Klassische“
Algorithmen

⇒ Modell durch
Mensch entwickelt

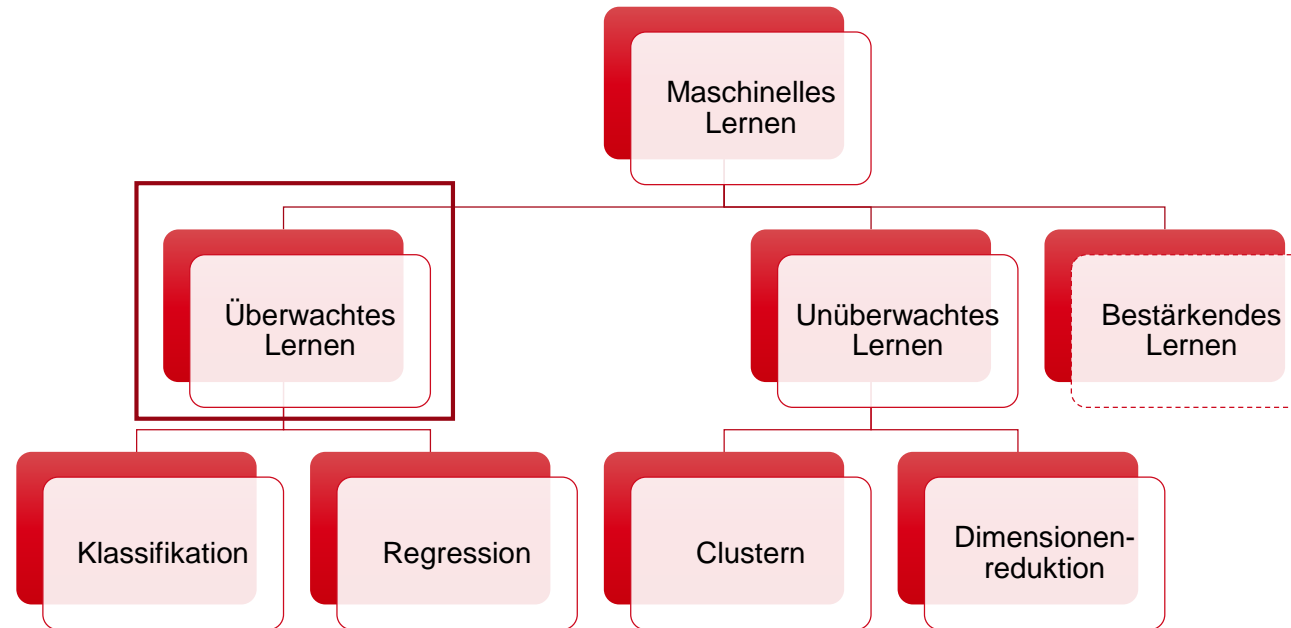
Maschinelles Lernen

⇒ Modell durch
Maschine „erlernt“

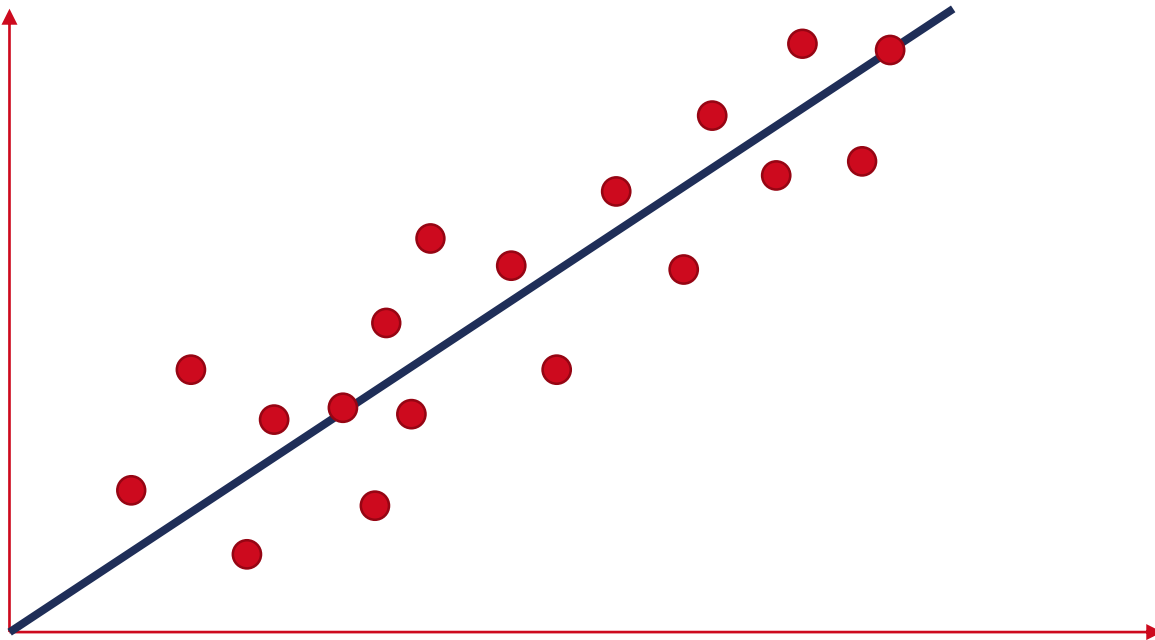
Ansätze des maschinellen Lernens



Ansätze des maschinellen Lernens



Lineare Regression

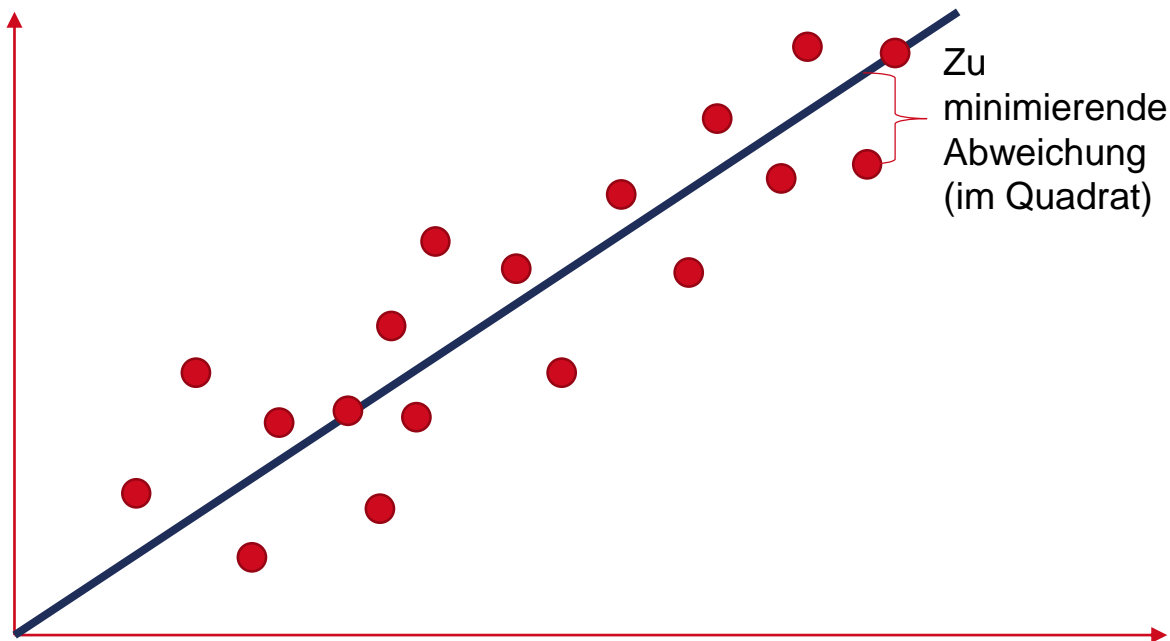


→ Minimiert die Fehlerquadrate

→ $y = m * x + b$

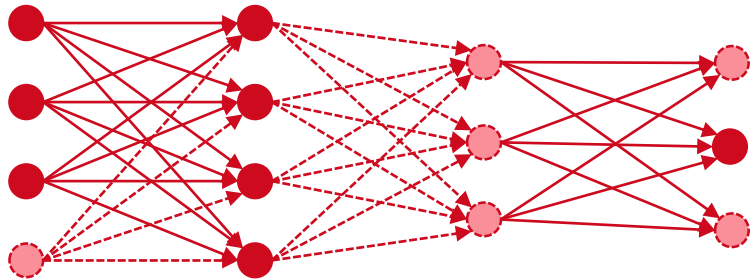
wobei y die abhängige Variable (Zielwerte),
 x die unabhängige Variable(n) (Eingangswerte),
 m die Steigung der Geraden und
 b der Achsenabschnitt ist

Lineare Regression

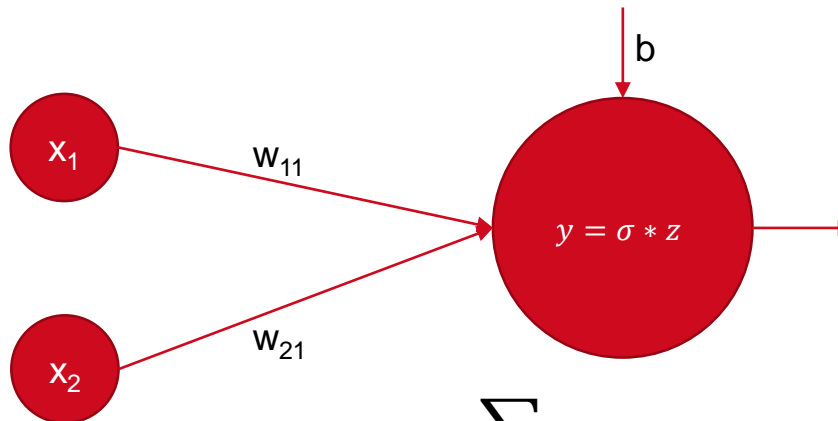


- Minimiert die Fehlerquadrate
- $y = m * x + b$
wobei y die abhängige Variable (Zielwerte),
 x die unabhängige Variable(n) (Eingangswerte),
 m die Steigung der Geraden und
 b der Achsenabschnitt ist
- Vorhersage des Modells ist bei neuen Werten nachvollziehbar

Neuronale Netze



Input-Schicht **versteckte Schicht** **versteckte Schicht** **Output-Schicht**

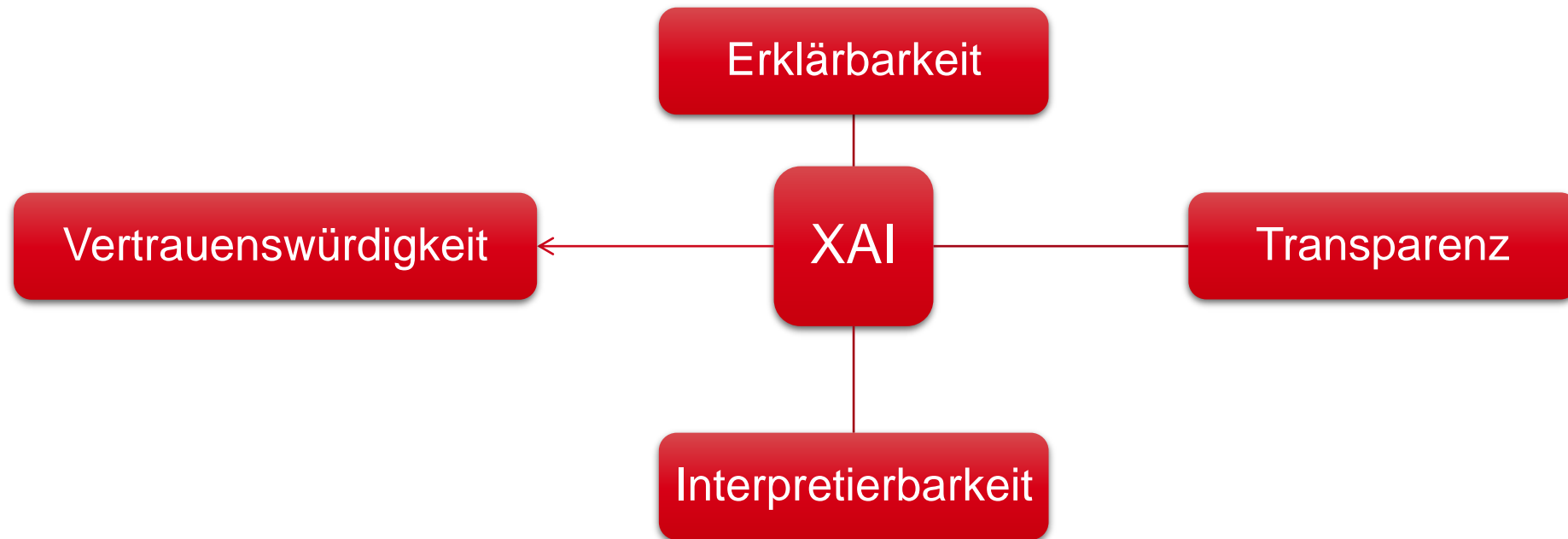


$$z = \sum_i w_{ij} * x_i + b$$

σ : Aktivierungsfunktion,
bspw: Sigmoid, ReLU, tanh

Vorhersage des Modells ist bei neuen Werten schwer intuitiv nachvollziehbar

Begrifflichkeiten im Zusammenhang mit XAI



Fehlermöglichkeiten bei der Klassifikation

Reales Ergebnis		Ergebnis der KI	
		In Ordnung (H_0)	Nicht in Ordnung (H_1)
Ergebnis der KI	In Ordnung	Spezifität	Fehler 2. Art Kunden Risiko
	Nicht in Ordnung	Fehler 1. Art Produzenten Risiko	Sensitivität

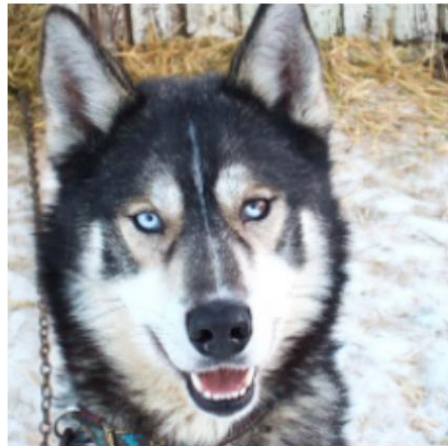
Fehlermöglichkeiten bei der Klassifikation

Reales Ergebnis		Ergebnis der KI	
		In Ordnung (H_0)	Nicht in Ordnung (H_1)
Ergebnis der KI	In Ordnung	Spezifität	Fehler 2. Art Kunden Risiko
	Nicht in Ordnung	Fehler 1. Art Produzenten Risiko	Sensitivität

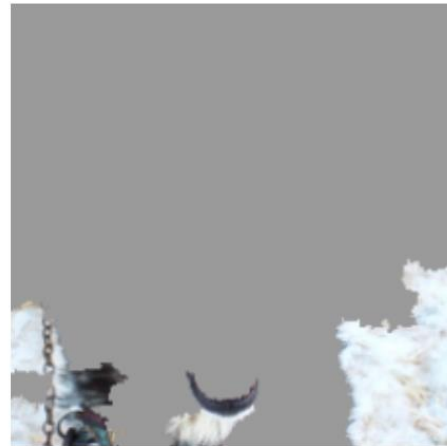
↔ Risikoabwägung ↔

Wie kommt die KI zu einer Entscheidung?

Richtige Entscheidung, falscher Grund: Beispiel Wolf/Husky-Klassifikator



(a) Husky classified as wolf



(b) Explanation

Figure 11: Raw data and explanation of a bad model's prediction in the "Husky vs Wolf" task.

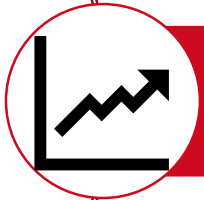
- Klassifikator soll Bilder von Huskies und Wölfen den entsprechend abgebildeten Tieren zuweisen.
- Auf dem rechten Bild ist zu sehen, dass der Klassifikator Schnee als relevantes Merkmal erlernt hatte
- Fehler wurde bewusst durch Bias in den Trainingsdaten erzeugt

Quelle: Ribeiro et al, 2016: "Why Should I Trust You?" Explaining the Predictions of Any Classifier <https://arxiv.org/pdf/1602.04938.pdf>

Gründe für eine erklärbare KI



Vertrauen in die Ergebnisse einer KI schaffen



Erleichterte Optimierung von KI-Algorithmen



Einhaltung regulatorischer und ethischer Vorgaben

Nutzung domänenbezogener Merkmale

- Durch ausschließliche Verwendung domänenbezogener Merkmale kann die Vertrauenswürdigkeit erhöht werden
- Domänenexperte benennt relevante Merkmale, die eine kausale Auswirkung auf das Ergebnis haben können
- Beispiel 3D-Druck: Ist das Wetter ein relevantes Merkmal?
- In der Banken Branche sind ökonomisch motivierte Faktoren regulatorisch verpflichtend

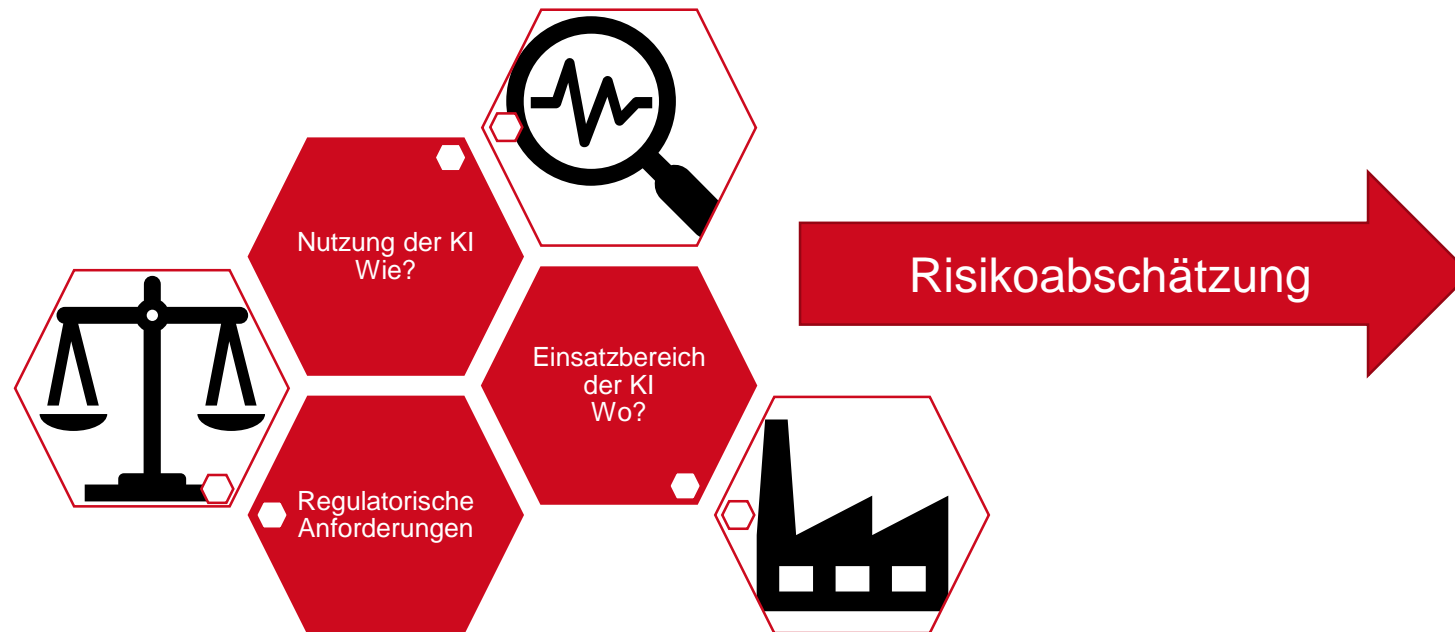
Dokumentationsmöglichkeiten bei KI

Welche Bestandteile der KI sollten dokumentiert bzw. versioniert werden?



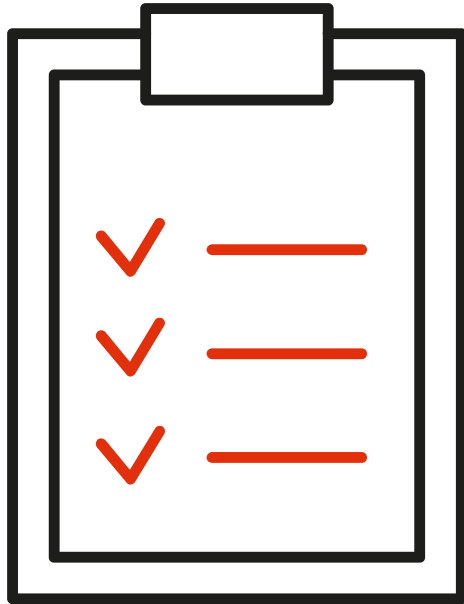
Wichtigkeit der Nachvollziehbarkeit beim KI-Einsatz

Wie viel Aufwand sollte für eine Nachvollziehbarkeit erbracht werden?



- Bedeutung des Fehlers
- Auftretenswahrscheinlichkeit des Fehlers
- Entdeckungswahrscheinlichkeit des Fehlers

Agenda



- Vorstellung des Mittelstand-Digitalzentrum Chemnitz
- Einführung in erklärbare KI
- Technische Möglichkeiten der XAI mit Fokus auf virtuelle Sensorik
- Beleuchtung und Diskussion über rechtliche Auswirkungen beim Einsatz von KI

Technische Möglichkeiten der XAI mit Fokus auf virtuelle Sensorik