



Entseelte Entscheidungen

Wenn Computer über Menschen entscheiden

IG-Metall Betriebsrätinnen Automobilbranche,
Kassel 6.6.2018

Prof. Dr. Katharina A. Zweig
Algorithm Accountability Lab
TU Kaiserslautern

@netwerkerin

Diskriminierung bei Bewerbungen

- Lebensläufe mit „deutschen“ Namen bekommen 14% mehr Vorstellungangebote als solche mit „türkischen“ Namen¹.
- US-amerik. Studie: Frauen mit Kopftuch erhalten weniger Jobangebote als solche ohne².



¹ Kaas, L. & Manger, C.: "Ethnic Discrimination in Germany's Labour Market: A Field Experiment", German Economic Review, 2011 , 13 , 1-20

² Ghuman, S. & Ryan, A. M.: "Not welcome here: Discrimination towards women who wear the Muslim headscarf , human relations, 2013 , 66(5) , 671-698

Richter

- Richter müssen vorzeitige Haftentlassungsanträge begutachten.
- Studie: je weiter von der letzten Pause weg, desto weniger risikoreiche Entscheidungen¹.
- Eine Vielzahl solcher Studien scheint zu beweisen:

¹ Danziger, S.; Levav, J. & Avnaim-Pesso, L.: "Extraneous factors in judicial decisions", Proceedings of the National Academy of the Sciences, 2011, 108, 6889-6892



Menschen – so irrational!

- Richter müssen vorzeitige Haftentlassungsanträge begutachten
- Studie: In den letzten 10 Jahren wurden weniger Haftentlassungsanträge bewilligt als in den vorherigen 10 Jahren
- Eine Vielzahl solcher Studien scheint zu beweisen:

Menschen sind irrational und vorurteilsbeladen.

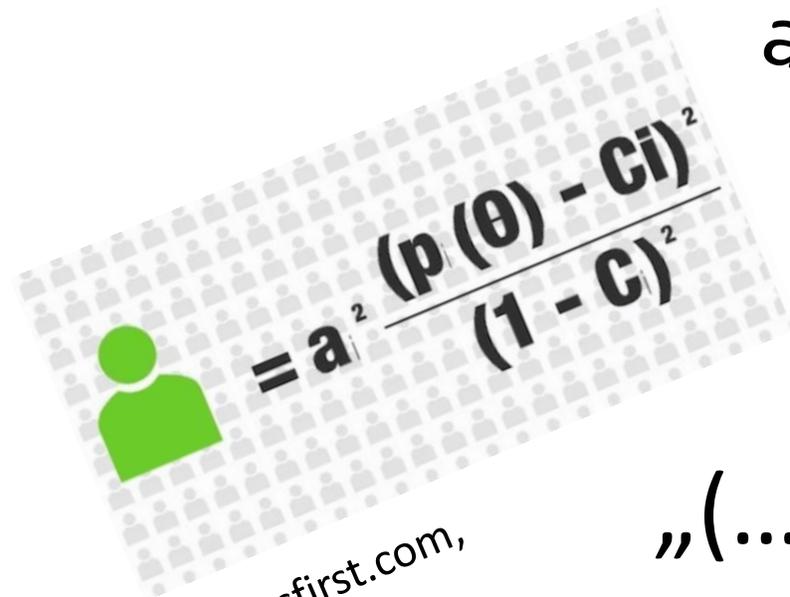
Könnten Computer das besser?

- Die ersten Firmen testen *algorithmische Entscheidungssysteme*¹.
- Eigenschaften, nach denen nicht diskriminiert werden darf, können vor ihnen besser verborgen werden.
- Sie entscheiden konsistent.

¹ Claire Miller: "Can an Algorithm hire Better than a Human?", The New York Times, June 25, 2015, <https://www.nytimes.com/2015/06/26/upshot/can-an-algorithm-hire-better-than-a-human.html>



„Employment assessment software“


$$= a^2 \frac{(p(\theta) - ci)^2}{(1 - c)^2}$$

Assessfirst.com,
16.11.2017

„(...) with the availability of good data, the predictive possibilities are virtually unlimited (...)“

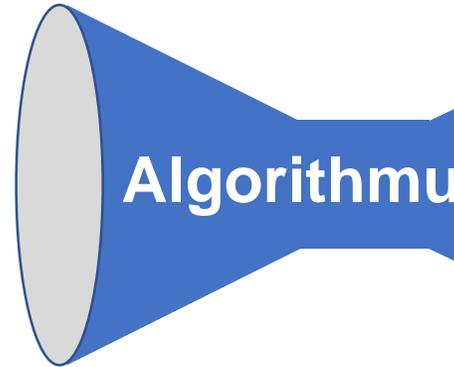
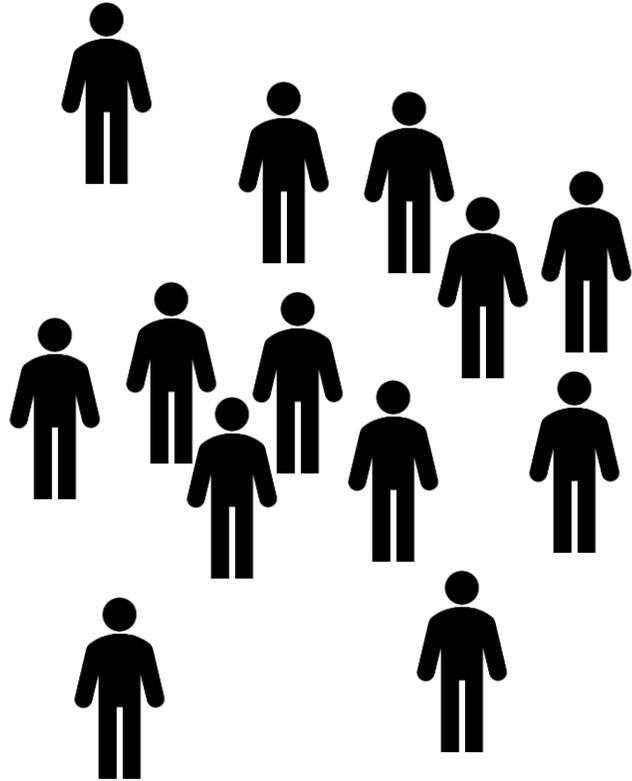
<https://www.inostix.com/predict-hiring-success/>
16.11.2017

Let's take the emotion out of the process and replace it with a data-driven approach...“

iNostix (by Deloitte),
16.11.2017

@netwerkerin
Prof. KA Zweig
TU Kaiserslautern

Algorithmische Entscheidungssysteme



Scoring-Verfahren

oder



Klassifikation

A close-up photograph of a person's hands gripping vertical metal bars, likely in a prison cell. The lighting is dramatic, with strong highlights and deep shadows, emphasizing the texture of the skin and the metallic surface of the bars. The background is dark, making the hands and bars stand out.

Forschung

—

Vorhersage des
Rückfallrisiko
von Kriminellen

Das kleine ABC der Informatik

Können

Algorithmen,

Big Data und

Computerintelligenz

Menschen besser bewerten und richten als
Menschen?



A wie Algorithmus

Ein Algorithmus ist ein Problemlöser

Mathematisches Problem



INPUT

**Der OUTPUT
der uns sagt,
wie Input
mit Output
zusammenhängt.**



OUTPUT



Beispiel für ein Problem: Navigation

Navigation

Gegeben das Kartenmaterial und weitere Daten, berechne die kürzeste Route zwischen Start und Ziel

Das **Problem** sagt nicht, wie man die Lösung **findet**.



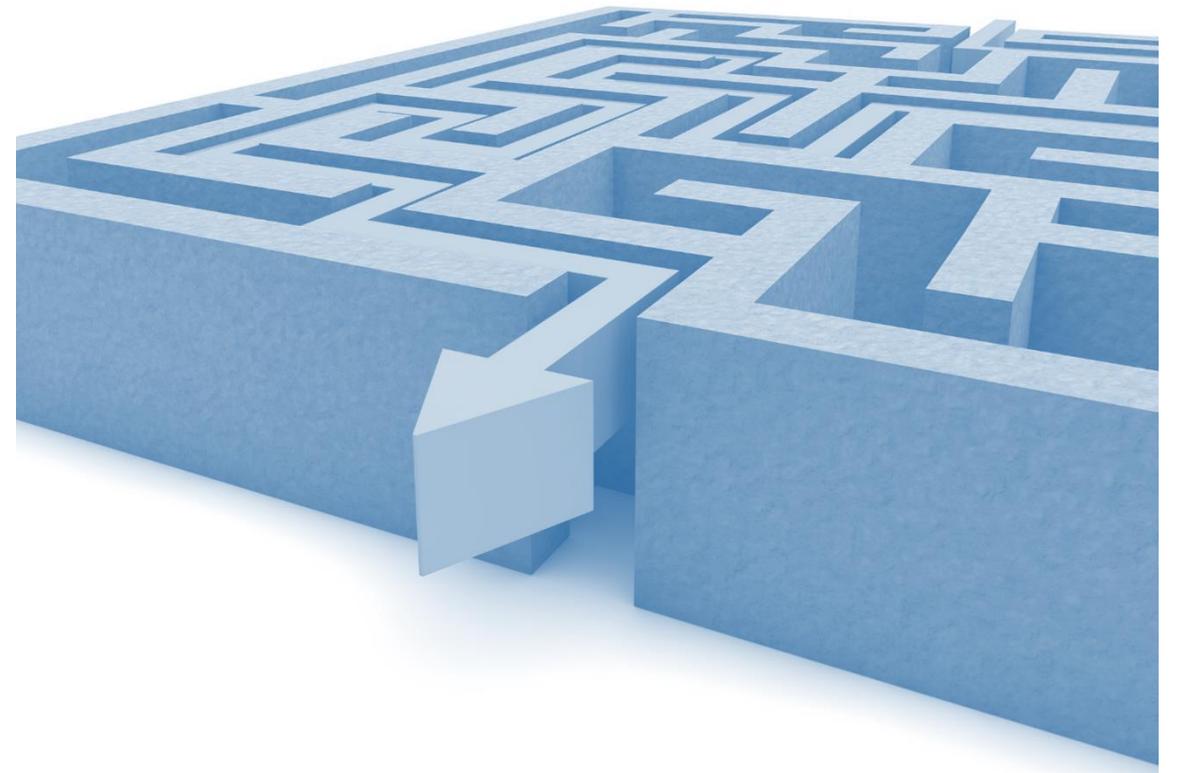
Input: Straßen, Länge, Staus, ...
Start und Ziel



Output: optimale Route

Ein Algorithmus ist...

...eine für jede **erfahrene Programmiererin** ausreichend **detaillierte Lösungsvorschrift**, so dass bei **korrekter Implementierung** der Computer **für jede korrekte Inputmenge den korrekten Output** berechnet – in endlicher Zeit.



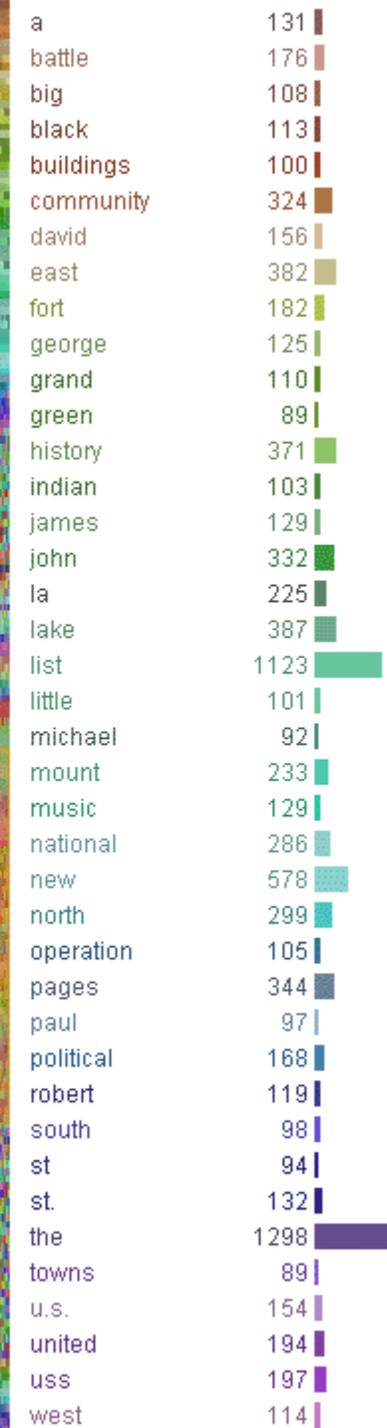


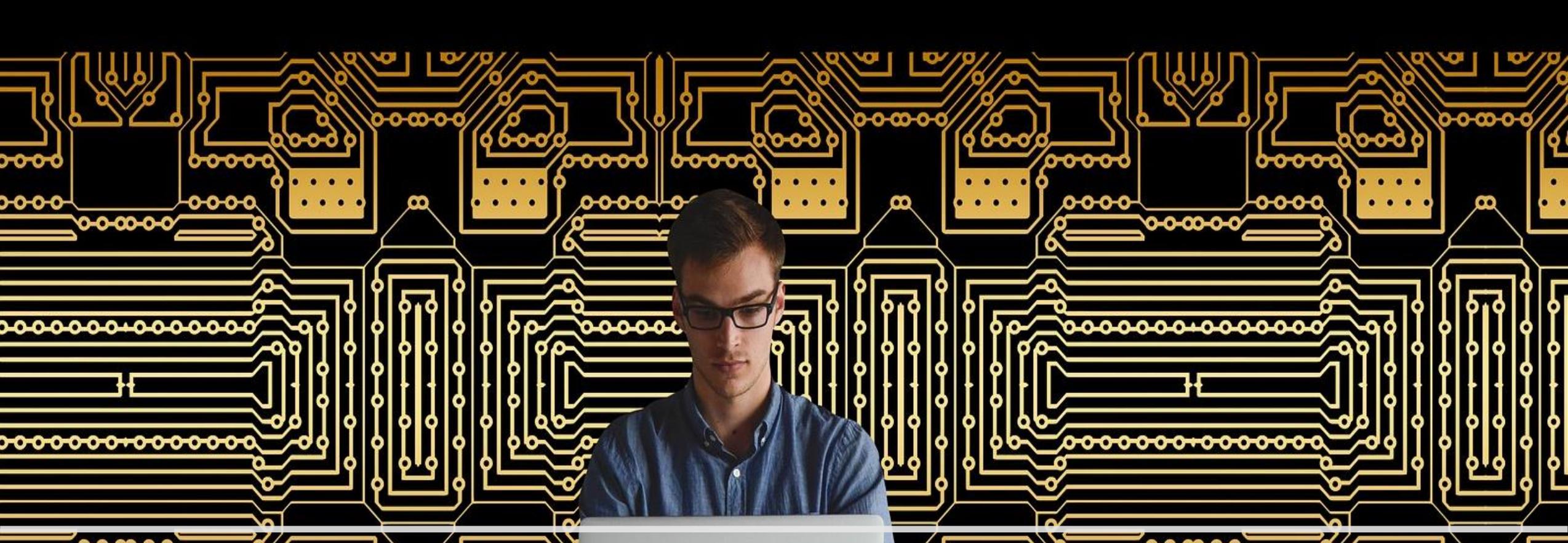
B wie Big Data

Daten als Grundlage

Big Data

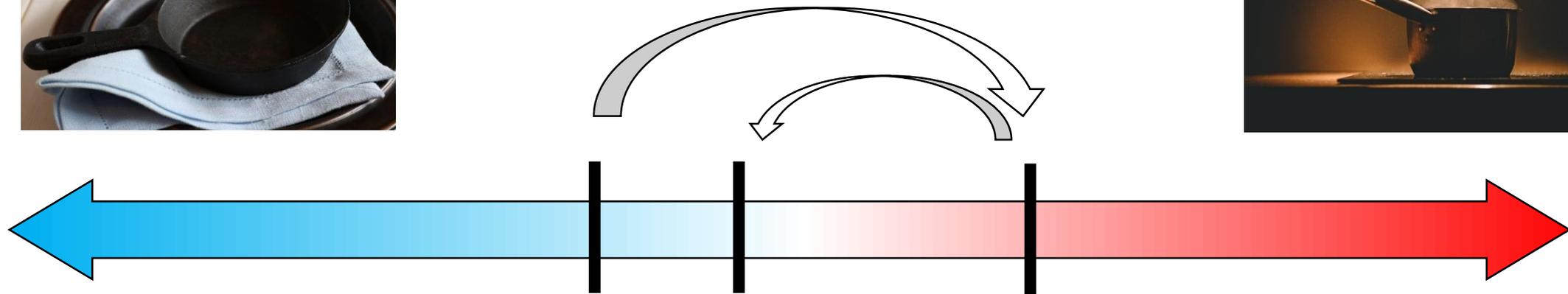
- Große Datenmengen.
- Außerhalb ihres spezifischen Zwecks genutzt.
- Daher im Einzelnen vermutlich fehlerbehaftet.
- Dank großer Masse und wenig individualisiertem Verhalten statistisch nutzbar.
- Hier werden Methoden des maschinellen Lernens benötigt.





C wie Computerintelligenz

Sebastian lernt „heiss“ und „warm“

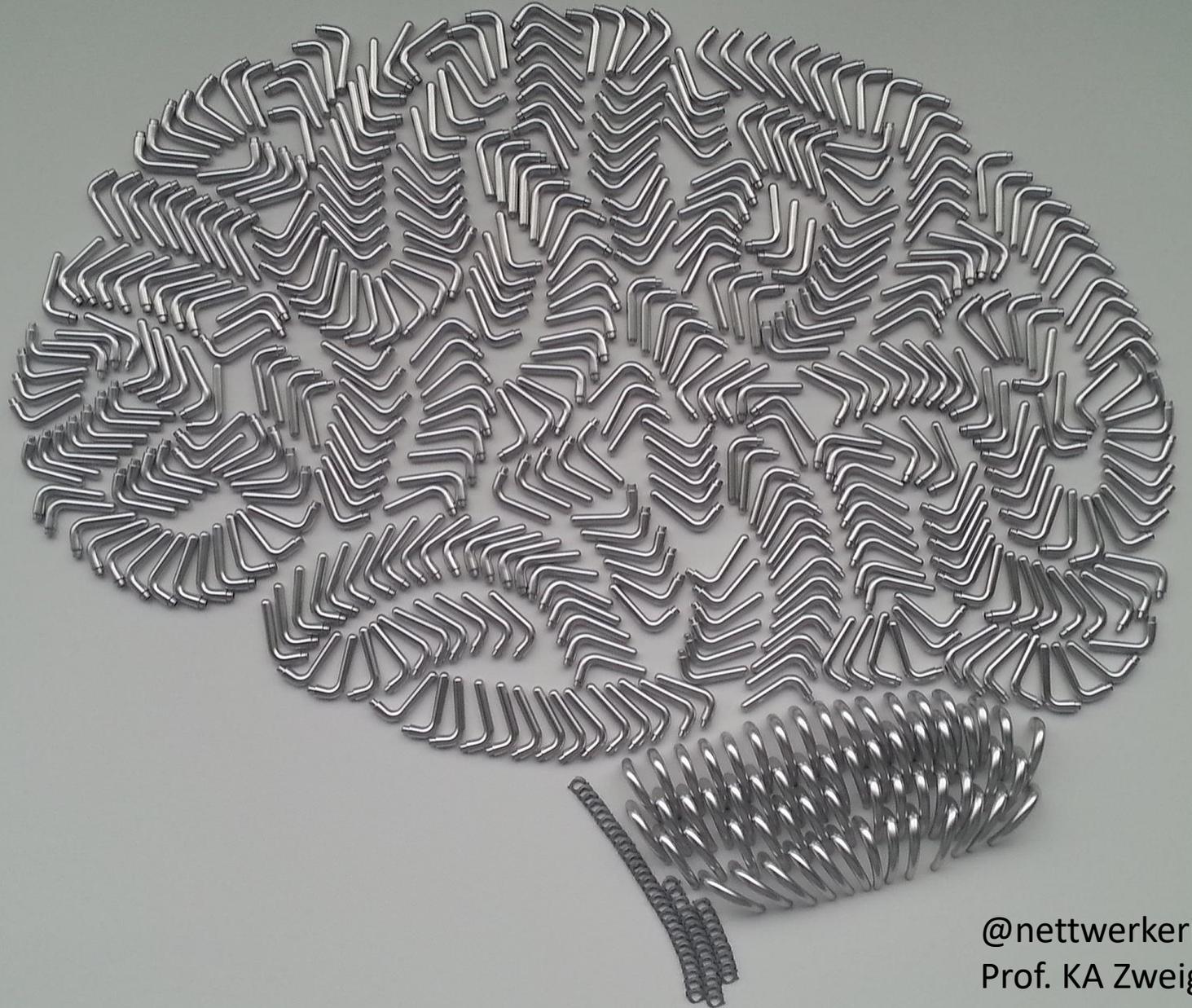


Juli März September

**Zu vorsichtig: Darf nicht dampfen Zu mutig geworden
Alles muss kalt sein**

Sebastian lernt...

- Durch **Rückkopplung:** unerwartet heiß, unerwartet kalt
- Durch **Speicherung in einer Struktur:** in Neuronen und deren Verknüpfung.
- Durch **Generalisierung des Gelernten.**

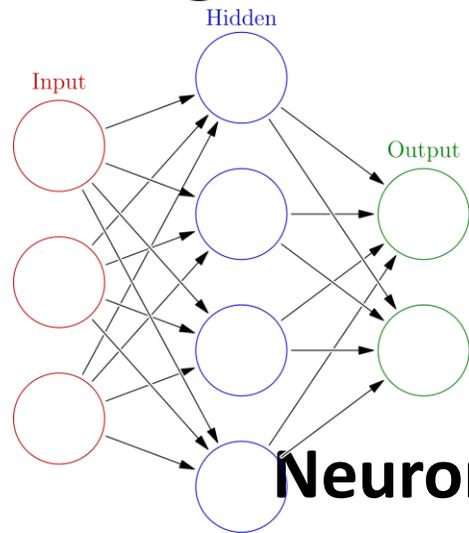


Computer lernen

Damit ein Computer lernen kann, benötigt er ebenfalls eine **Struktur**, um Gelerntes abzuspeichern.

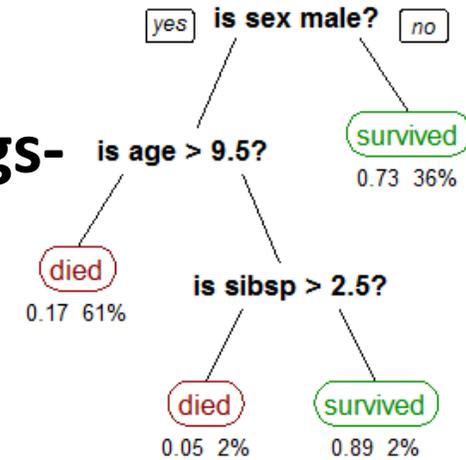
Optimal auch **Rückkopplung**.

Er lernt **generelle Regeln**.



Neuronales Netz

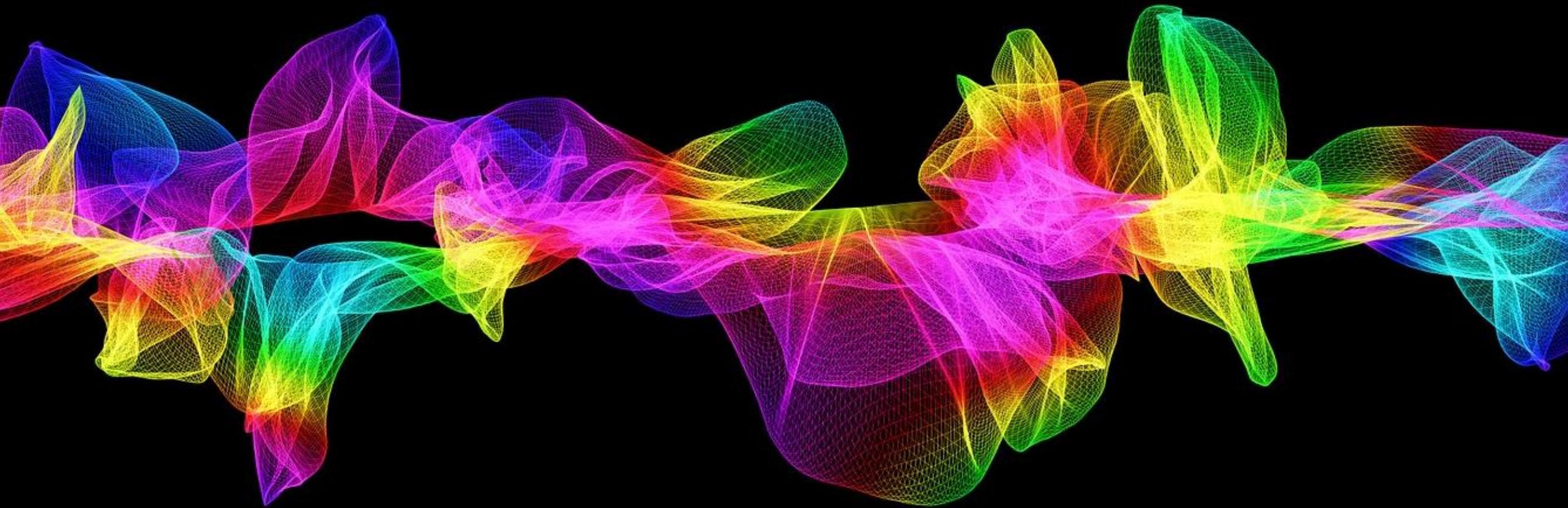
Entscheidungsbäume



Formel

$$w_1 * \#Vh - w_2 * \#day_1Vh + w_3 * I[g = male] * 1 + w_4 * I[T = R] * 1.0 + \dots$$

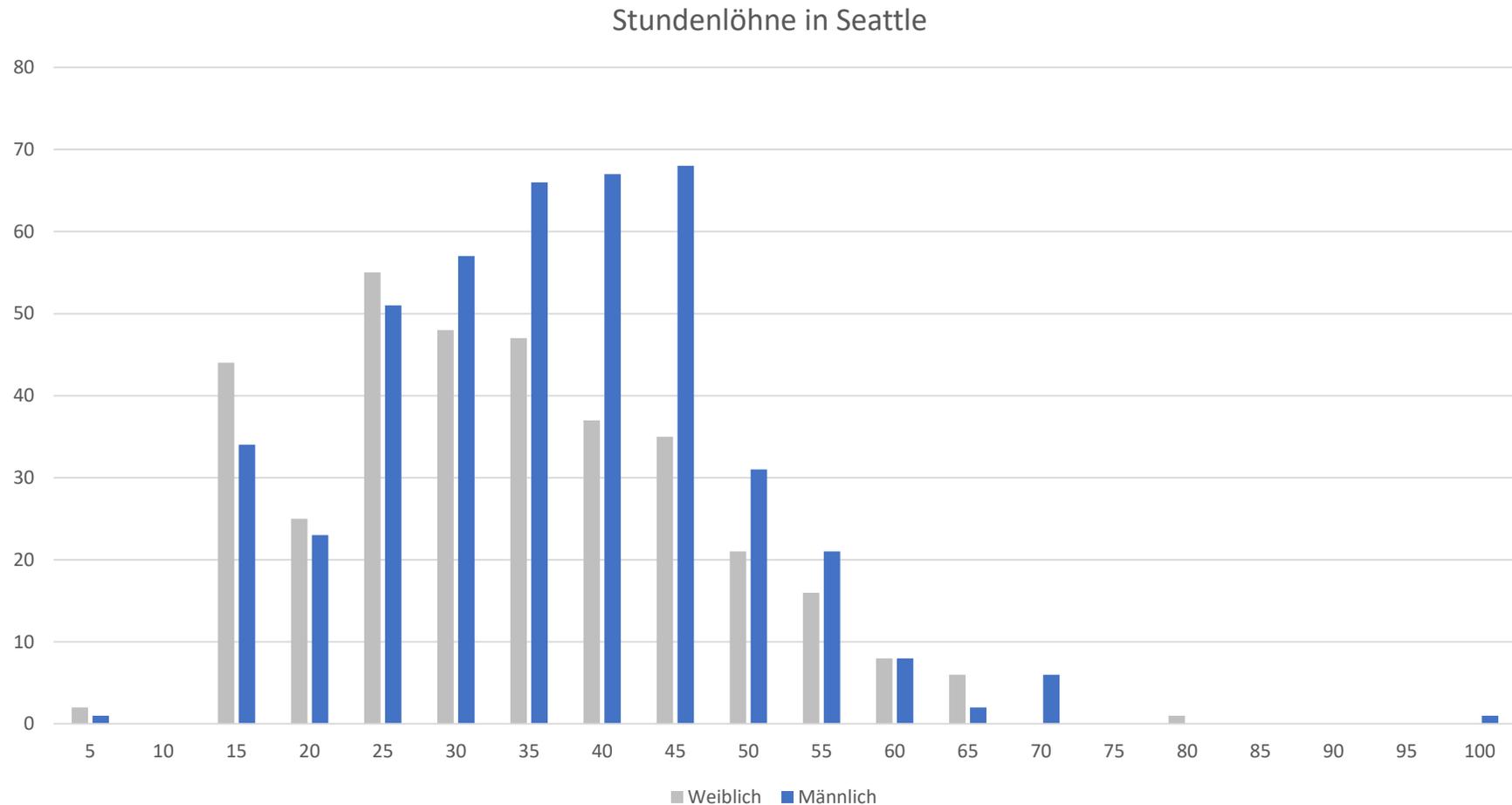
@netwerkerin
Prof. KA Zweig
TU Kaiserslautern



“Lernen” mit Korrelationen

Heißen Sie unsere(n) neue(n) Mitarbeiter(in) willkommen!

- Anteil weiblicher Angestellter?
 - 44%
- Anteil weiblicher Angestellter mit Lohn unter \$25?
 - 55%



$$X_{1/2} = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$



$$X^2 + px + q = 0$$



$$X_{1/2} = -\frac{p}{2} \pm \sqrt{\left(\frac{p}{2}\right)^2 - q}$$

$$x = b - 2v$$

Lernen mit Formeln

Datengrundlagen

- Data Mining Methoden nutzen verschiedene Informationen
- Am wichtigsten:
 - **War Einstellung erfolgreich?**

Ausbildung

Leerzeiten

Arbeitgeber
-wechsel

Alter

Bewerbungs-
schreiben

Rechtschreibung

Wortvielfalt

Ton

Social Media?

Regressionsansätze

- Die Algorithmen-designer entscheiden, welche Daten vermutlich mit „erfolgreicher Einstellung“ korrelieren.
- Die Software sollte eine einzige Zahl ausgeben.
- Je höher die Zahl, desto höher die Erfolgswahrscheinlichkeit.
- Beispiel Formel:

$$\begin{aligned} & 3 * \text{Jahre im Job} \\ - & 2 * \text{Anzahl Arbeitgeber} \\ + & 3 * (\text{Wenn Auslandserfahrung,} \\ & \text{dann 1, sonst 0}) \\ + & 2,5 * (\text{Wenn Fortbildung,} \\ & \text{dann 1, sonst 0}) + \dots \end{aligned}$$

Allgemein

Der Computer bestimmt die Gewichte und bekommt ein Feedback (Rückkopplung), inwieweit die resultierende Bewertung mit dem (beobachteten) Verhalten übereinstimmt.

$$\begin{aligned} & w_1 * \text{Jahre im Job} \\ + & w_2 * \text{Anzahl Arbeitgeber} \\ + & w_3 * (\text{Wenn Auslandserfahrung,} \\ & \text{dann 1, sonst 0)} \\ + & w_4 * (\text{Wenn Fortbildung,} \\ & \text{dann 1, sonst 0)} + \dots \end{aligned}$$



Qualität eines Algorithmus

ROC AUC

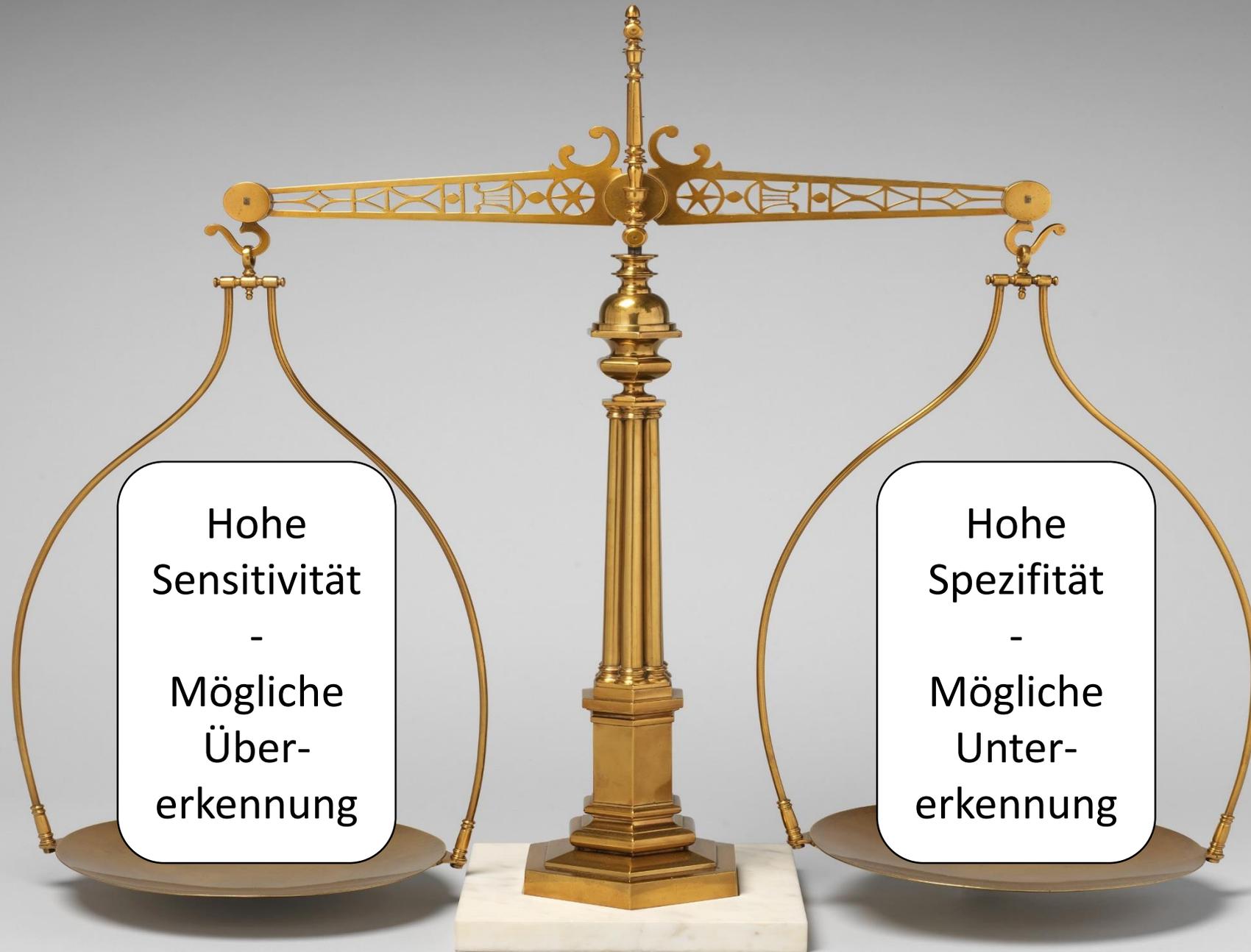
...und 20 mehr

Positive Predictive
Value

Sensitivität

Accuracy

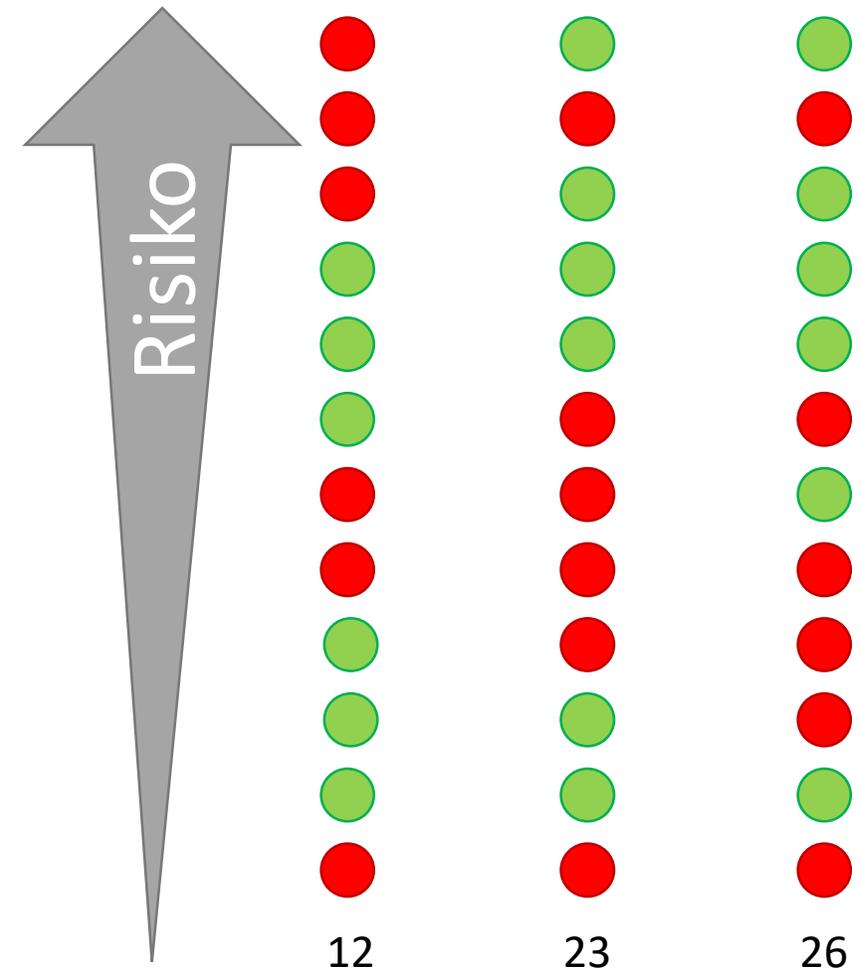
Qualitäts-
Maße



Hohe
Sensitivität
-
Mögliche
Über-
erkennung

Hohe
Spezifität
-
Mögliche
Unter-
erkennung

- Grüne Kugeln symbolisieren erfolgreiche, rote nicht erfolgreiche BewerberInnen.
- Optimale Sortierung: Alle grünen oben, alle roten darunter.
- Qualitätsmaß: Paare von rot und grün, bei denen die grüne Kugel über der roten einsortiert ist. (**ROC AUC**)



Ist das Qualitätsmaß sinnvoll?

- Wenn die Stelle **sofort** besetzt werden muss, und nur 5 Bewerber da sind: **ja**
- Wenn es langfristig um die Identifikation der besten Talente geht: **nein**
- Hier müssen andere Qualitätsmaße benutzt werden.





einen Jagdhund zu kaufen,



um Schafe zu hüten.

Das ist wie...

Probleme von algorithmischen Entscheidungssystemen (ADM Systemen) im People Assessment

- 1. Wer entscheidet, wann ein
ADM System „gut“ ist?**





Wahrscheinlichkeit & Wahrheit

Regel

Algorithmen der künstlichen Intelligenz werden da eingesetzt, wo es **keine einfachen Regeln** gibt.

Sie suchen **Muster** in hoch-verrauschten Datensätzen.

Die Muster sind daher grundsätzlich **statistischer Natur**.

Versuchen fast immer, eine **kleine Gruppe** von Menschen zu identifizieren (Problem der **Unbalanciertheit**)

Algorithmen...

- ... basieren auf Korrelationen von Eigenschaften mit gewünschtem Verhalten.
- **Quasi algorithmisch legitimierte Vorurteile:**
 - Zu 70% erfolgreich heißt:
 - Von 100 Personen, die „genau so sind wie dieser Mensch“, sind 70 nachher erfolgreich.



Probleme von algorithmischen Entscheidungssystemen (ADM Systemen) im People Assessment

1. Wer entscheidet, wann ein ADM System „gut“ ist?
2. **ADM Systeme ergeben nur Wahrscheinlichkeiten, keine Wahrheiten.**





Können Algorithmen diskriminieren?



Und das, wenn ich auf Pixabay nach „Chef“ suche...

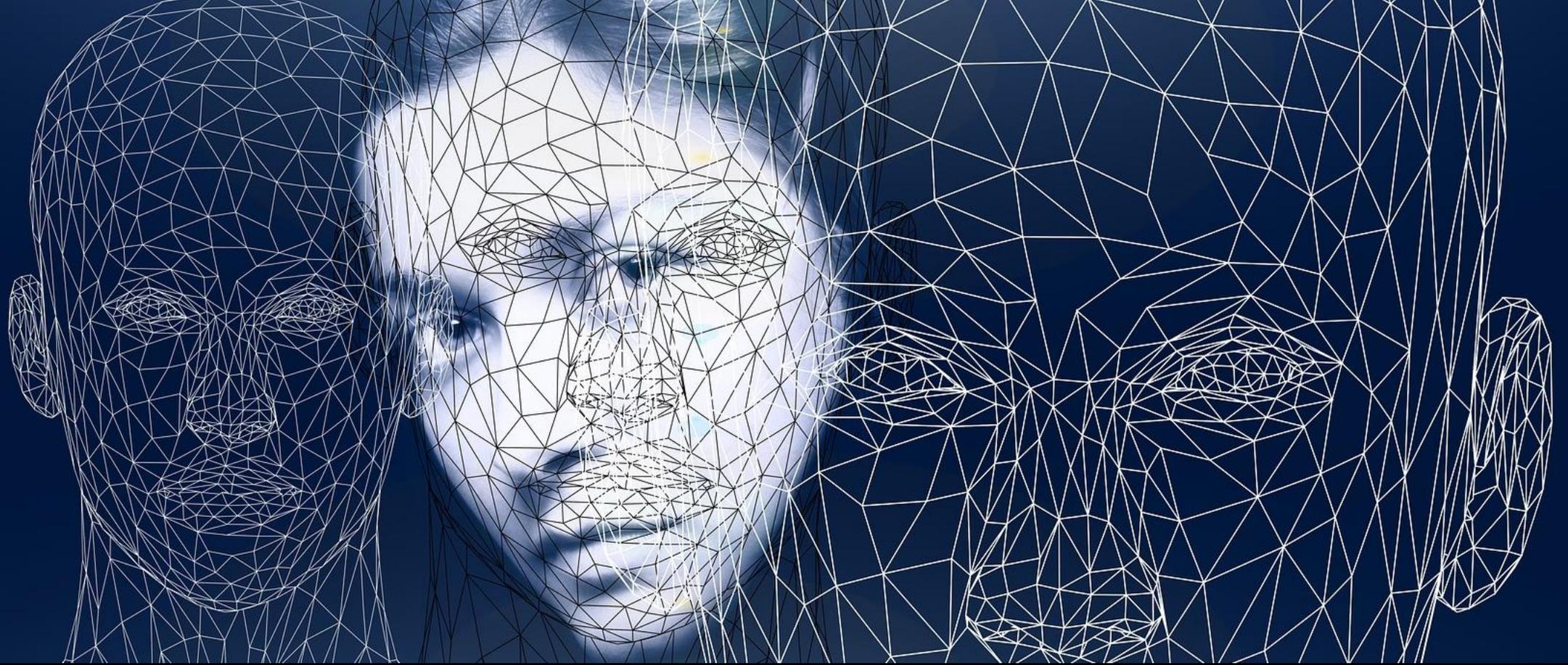
Diskriminierung

- Google zeigt weiblichen Surfern schlechtere Jobs an.
- Rückfälligkeitsvorhersagealgorithmen sind rassistisch.
- Diskriminierungen in Trainingsdaten werden „mitgelernt“.
- Wenn Trainingsdaten zu wenig Daten über Minderheiten enthalten, werden deren Eigenschaften nicht „mitgelernt“.

Probleme von algorithmischen Entscheidungssystemen (ADM Systemen) im People Assessment

1. Wer entscheidet, wann ein ADM System „gut“ ist?
2. ADM Systeme ergeben nur Wahrscheinlichkeiten, keine Wahrheiten.
3. **ADM Systeme können diskriminieren.**





Sozio-informatische Gesamtbetrachtung

Probleme der Einbettung der ADM in den sozialen Prozess

- **Aufmerksamkeitsökonomie** von Entscheiderinnen und Entscheidern.
- „**Best practice**“ erfordert Nutzung der Software.
- **Delegation von Verantwortung!**
- Manchmal kann ein falsch Beurteiler **die Vorhersage prinzipiell nicht entkräften!**
 - Z.B. abgelehnte Bewerberin

Probleme von algorithmischen Entscheidungssystemen (ADM Systemen) im People Assessment

1. Wer entscheidet, wann ein ADM System „gut“ ist?
2. ADM Systeme ergeben nur Wahrscheinlichkeiten, keine Wahrheiten.
3. ADM Systeme können diskriminieren.
4. **ADM Systeme können soziale Prozesse verändern.**



Probleme von algorithmischen Entscheidungssystemen (ADM Systemen) im People Assessment

- 1. Wer entscheidet, wann ein ADM System „gut“ ist?**
- 2. ADM Systeme ergeben nur Wahrscheinlichkeiten, keine Wahrheiten.**
- 3. ADM Systeme können diskriminieren.**
- 4. ADM Systeme können soziale Prozesse verändern.**



„Employment assessment software“

~~Let's take the emotion out of the process
and replace it with a data-driven approach...“~~

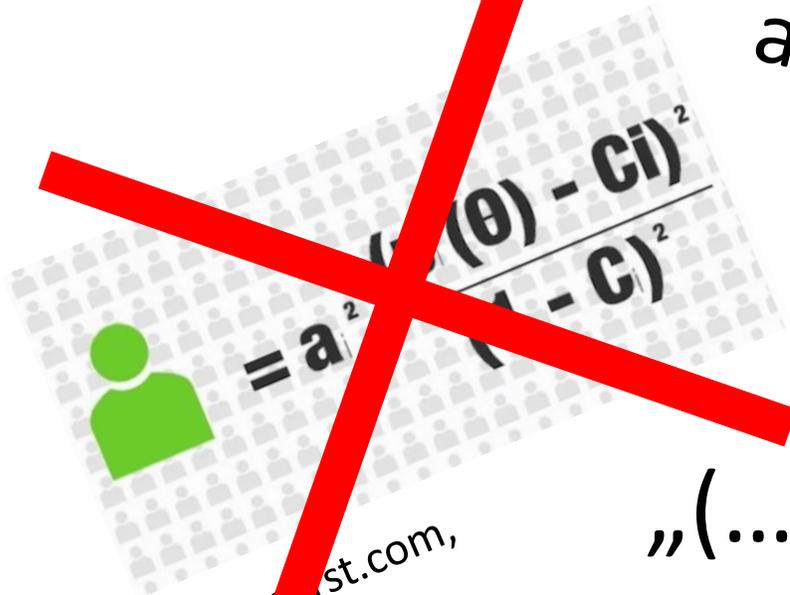
~~„(...) with the availability of
good data, the predictive
possibilities are virtually
unlimited (...)“~~

iNostix (by Deloitte),
16.11.2017

@netwerkerin
Prof. KA Zweig
TU Kaiserslautern

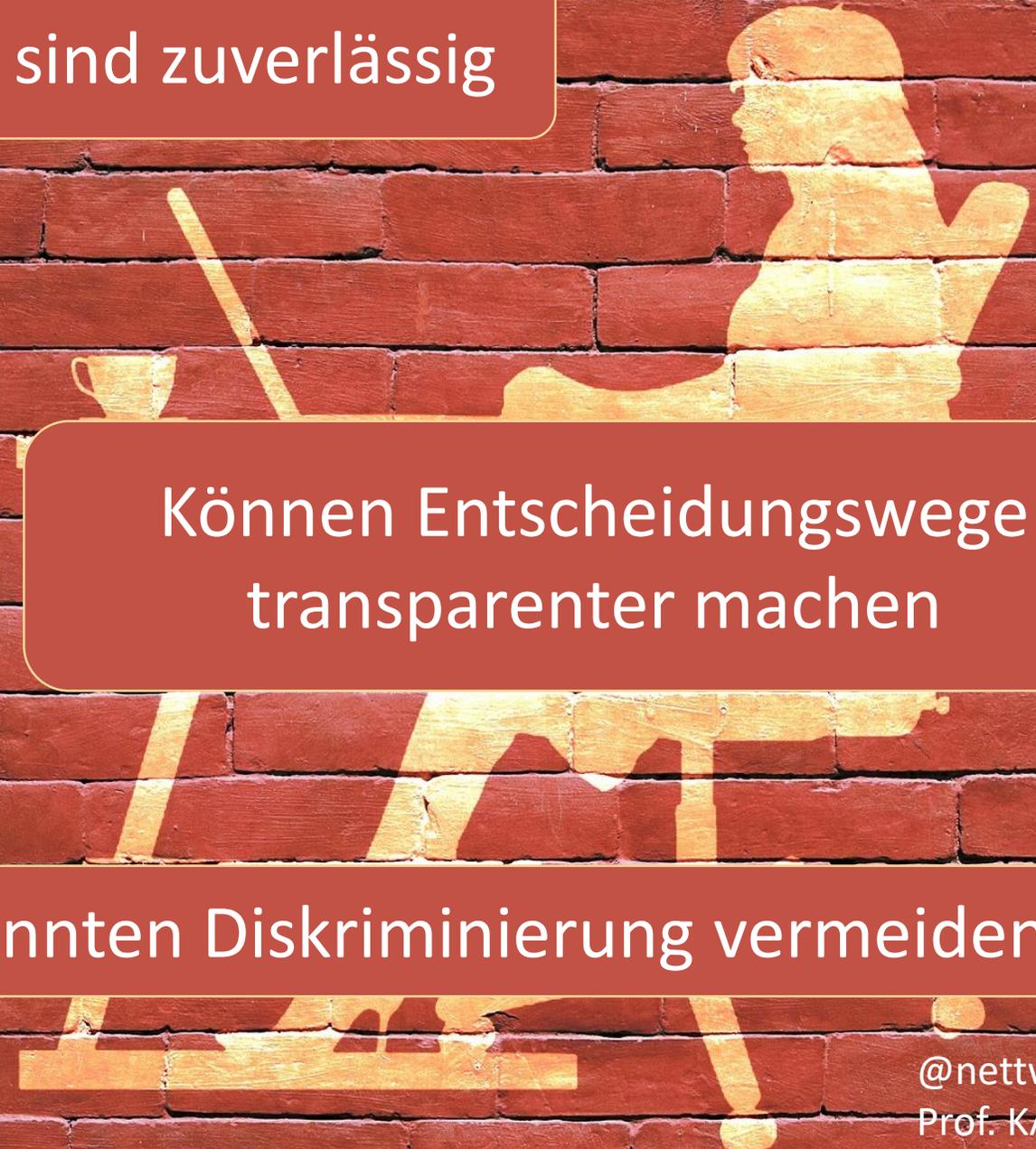
<https://www.inostix.com/predict-hiring-success/>
16.11.2017

Asses...st.com,
16.11.2017



Einschätzung

- People Assessment Systeme könnten dabei helfen, bessere Entscheidungen zu treffen.
- Allerdings ist es schwierig, sie transparent, fair und nachvollziehbar zu gestalten.



Sie sind zuverlässig

Können Entscheidungswege transparenter machen

Könnten Diskriminierung vermeiden

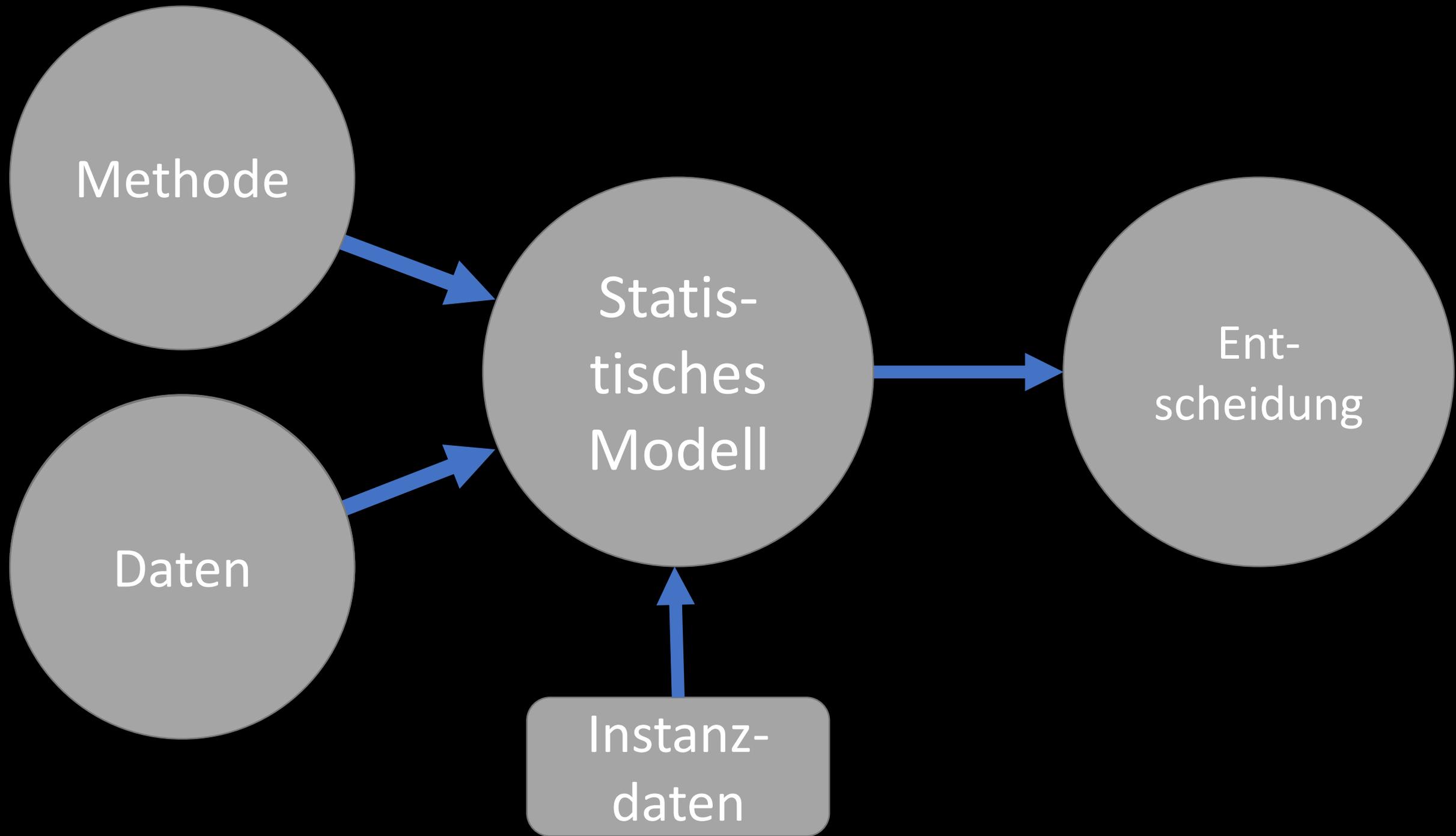
Ethik in der Softwareentwicklung



...in der
Automobilbranche



Lassen Sie uns einen
Aufmerksamkeits-
Sensor bauen!



Daten → Datenauswahl

Sensoren

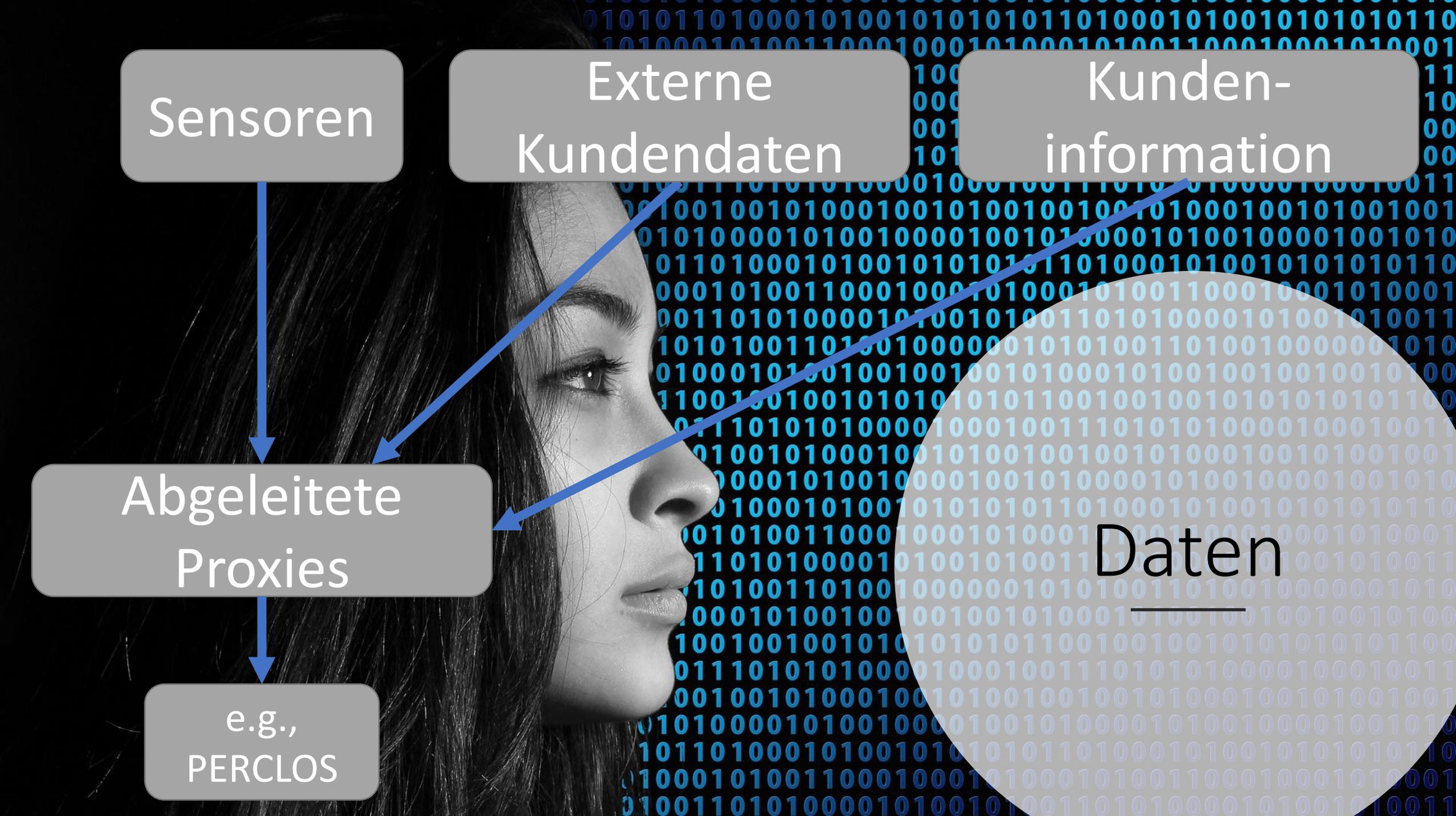
Externe
Kundendaten

Kunden-
information

Abgeleitete
Proxies

e.g.,
PERCLOS

Daten



Drowsiness detection system

Abstract

This invention describes a non-intrusive system used to detect and at risk of falling asleep at the wheel due to drowsiness. The system includes drowsiness detection systems and a control unit. This reduces drowsiness assessment. The first subsystem consists of an interior headliner and seat, which detects head movements that are characteristic of a driver. The second subsystem consists of heart rate monitoring at the wheel. The control unit is used to analyze the sensory data to determine the state and therefore corresponding risk of falling asleep while driving using intelligent software algorithms, and the data provided by the sensors. Characteristics that may indicate a drowsy driver are outputted which may be used to activate a response system in automobiles; this system may be used in any type of vehicle.

Kopfbewegungen
Herzschlag
Augenlider
(PERCLOS)

Images (3)



Classifications

G08B21/06 Alarms for ensuring the safety of persons indicating a condition of sleep, e.g. anti-dozing alarms

View 1 more classifications

B2

Find Prior Art Similar

Adam Basir, Jean Pierre Bhavnani, Fakhreddine Desrochers

Intelligent Mechatronic Systems Inc

Intelligent Mechatronic Systems Inc

2-01-18

Family: US (1)

Date	App/Pub Number	Status
2003-01-21	US10348037	Active
2003-08-14	US20030151516A1	Application
2004-11-23	US6822573B2	Grant

Info: Patent citations (14), Cited by (47), Legal events, Similar documents, Priority and Related Applications

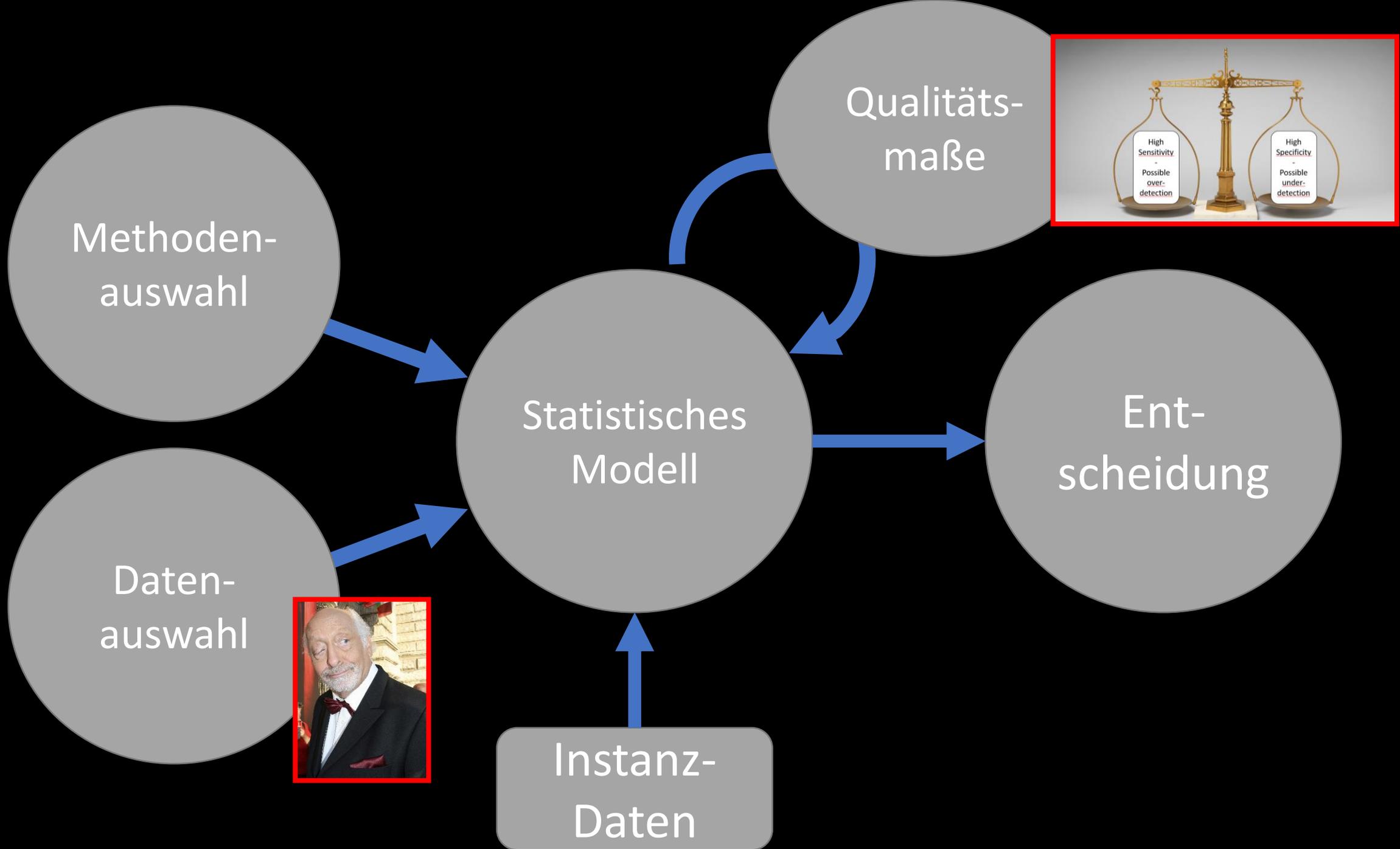
External links: USPTO, USPTO Assignment, Espacenet, Global Dossier, Discuss

Ethische Überlegungen bei der Datenauswahl

Durch die Videosensorik diskriminieren wir vielleicht...

- ...Personen mit trockenen Augen oder Kontaktlinsen?
- ...Personen mit einer Ptosis wie Karl Dall?

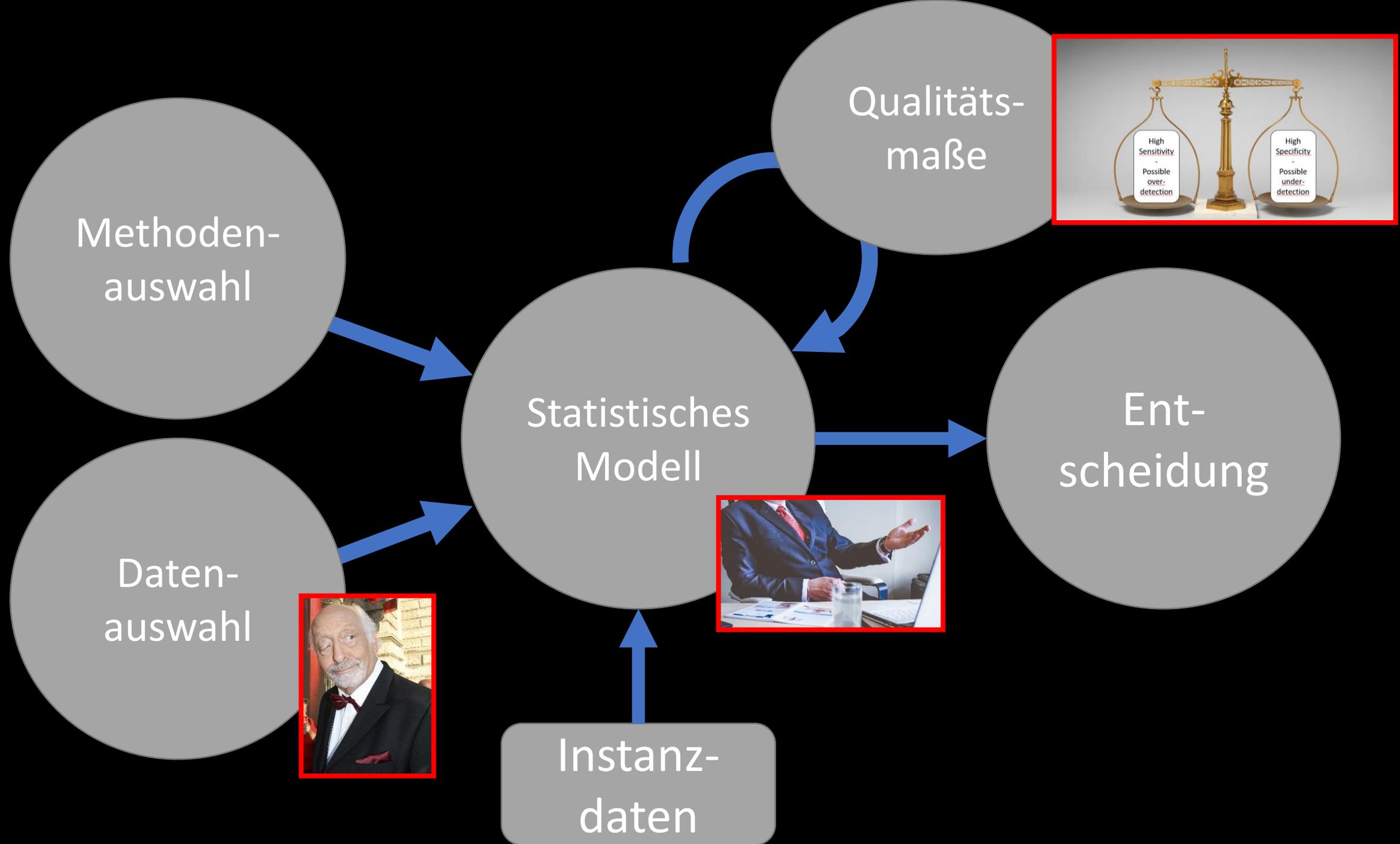


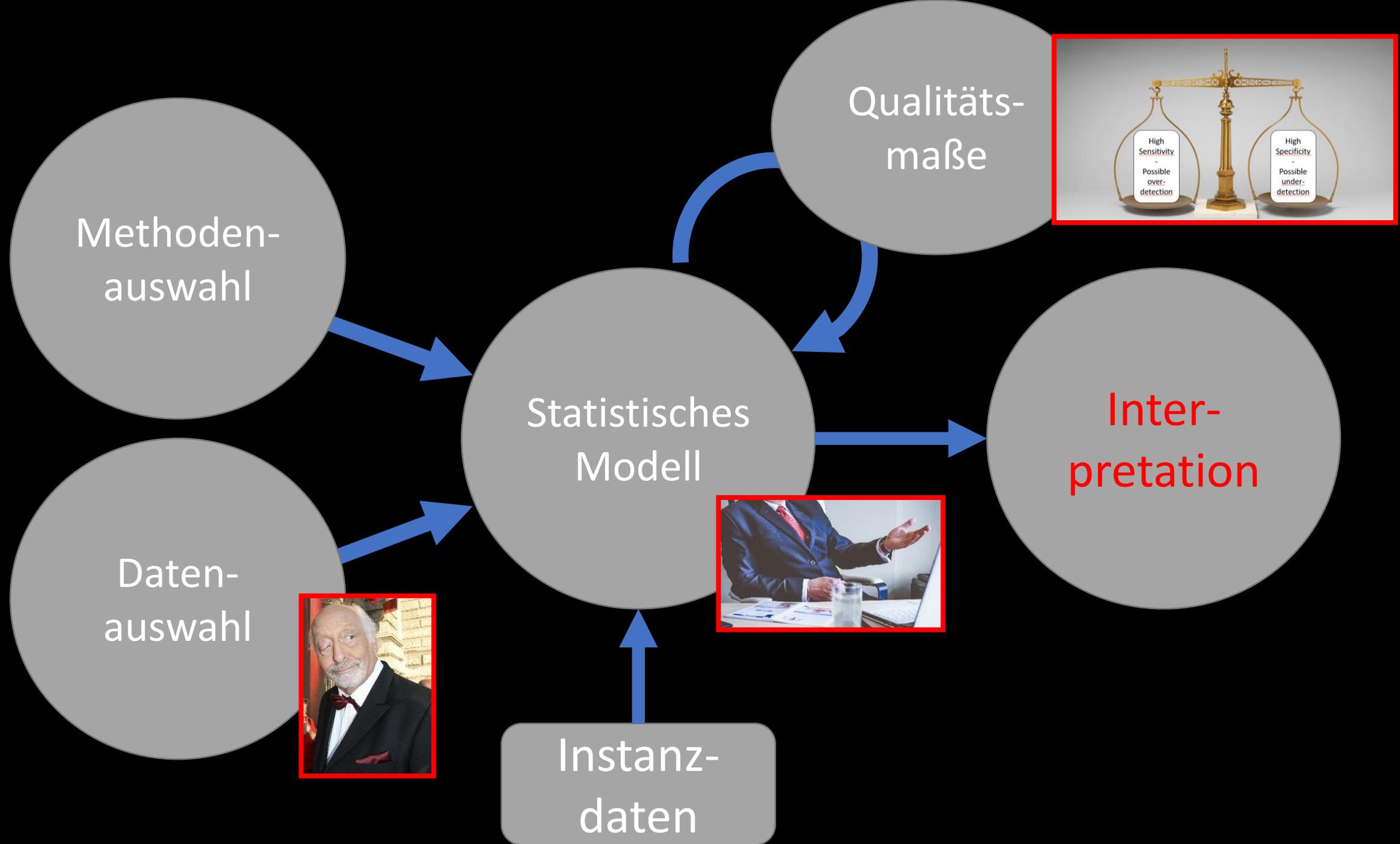


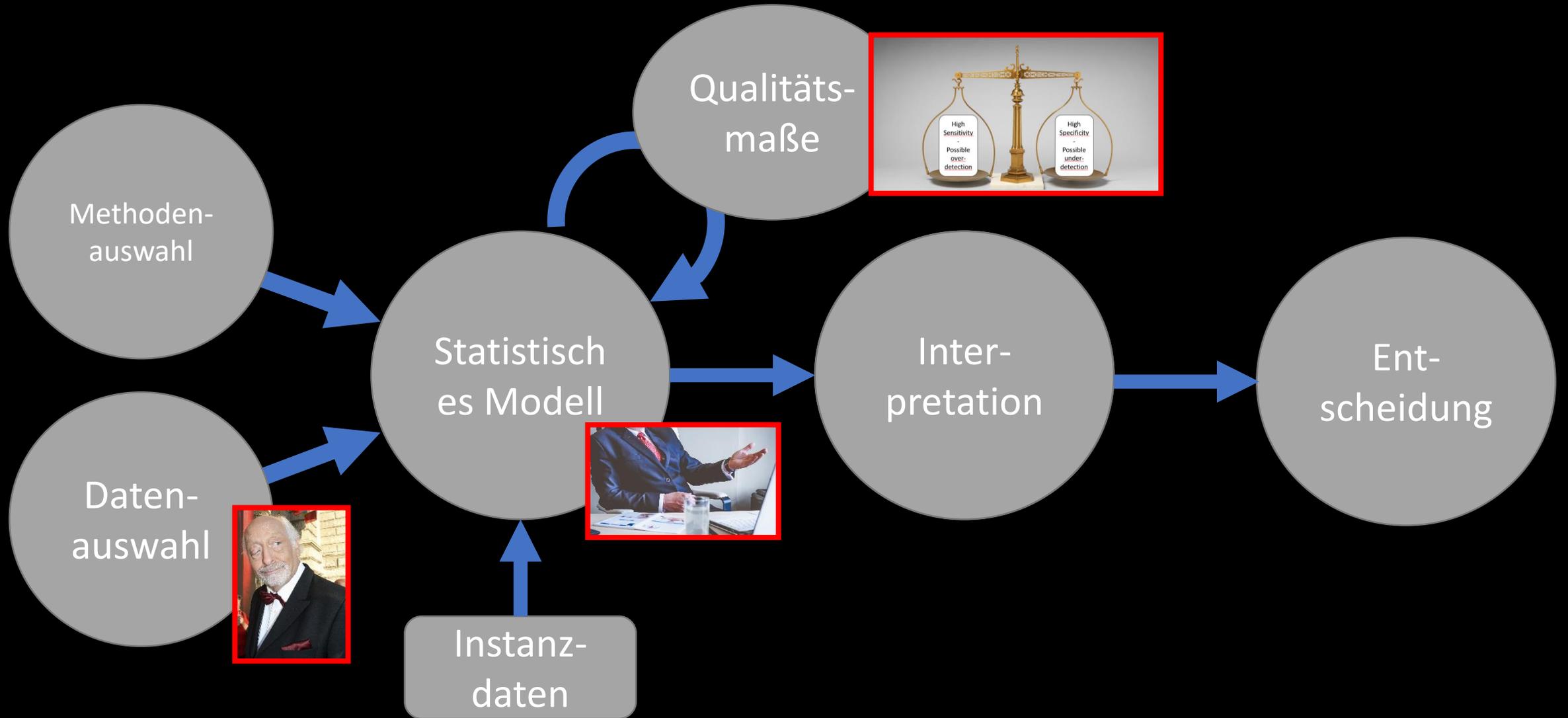


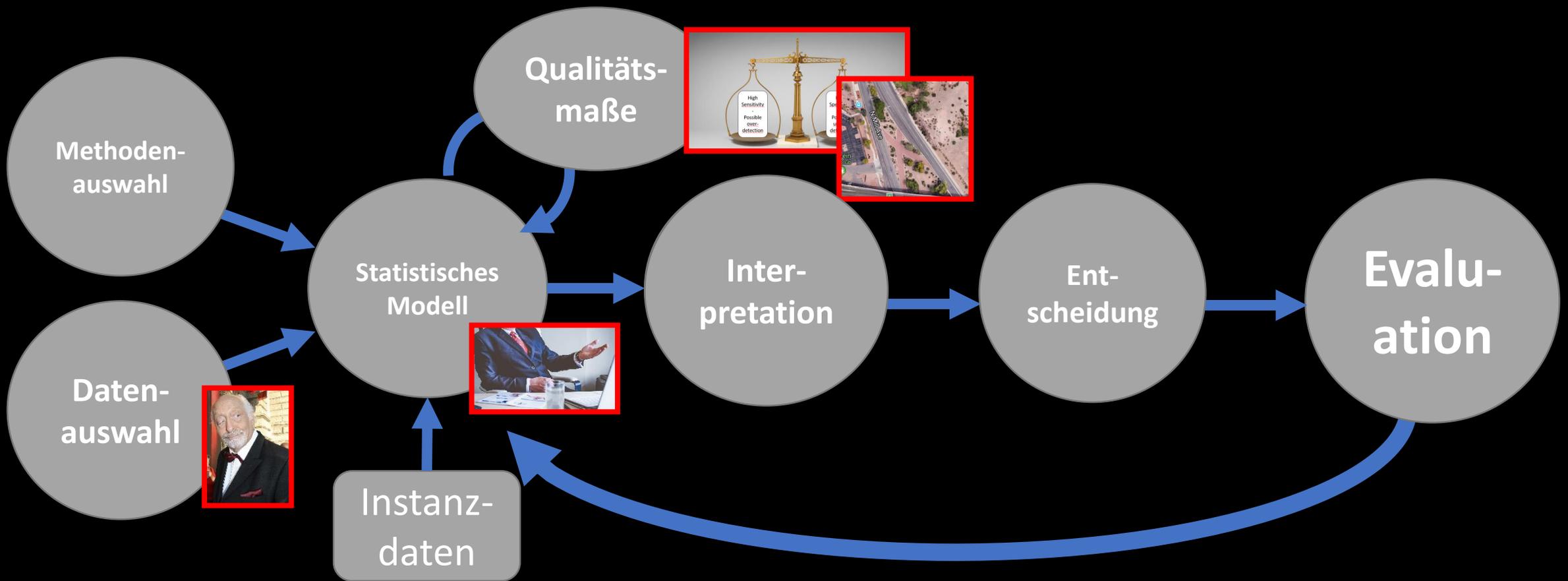
Soziale Einbettung des Systems

- Wird sich der Fahrer oder die Fahrerin auf das System verlassen?
- Wollen Speditionen Zugriff auf Fahrerdaten?
- Könnte das System auch für Büros nützlich sein?
 - Welche Verantwortlichkeiten erwachsen daraus?





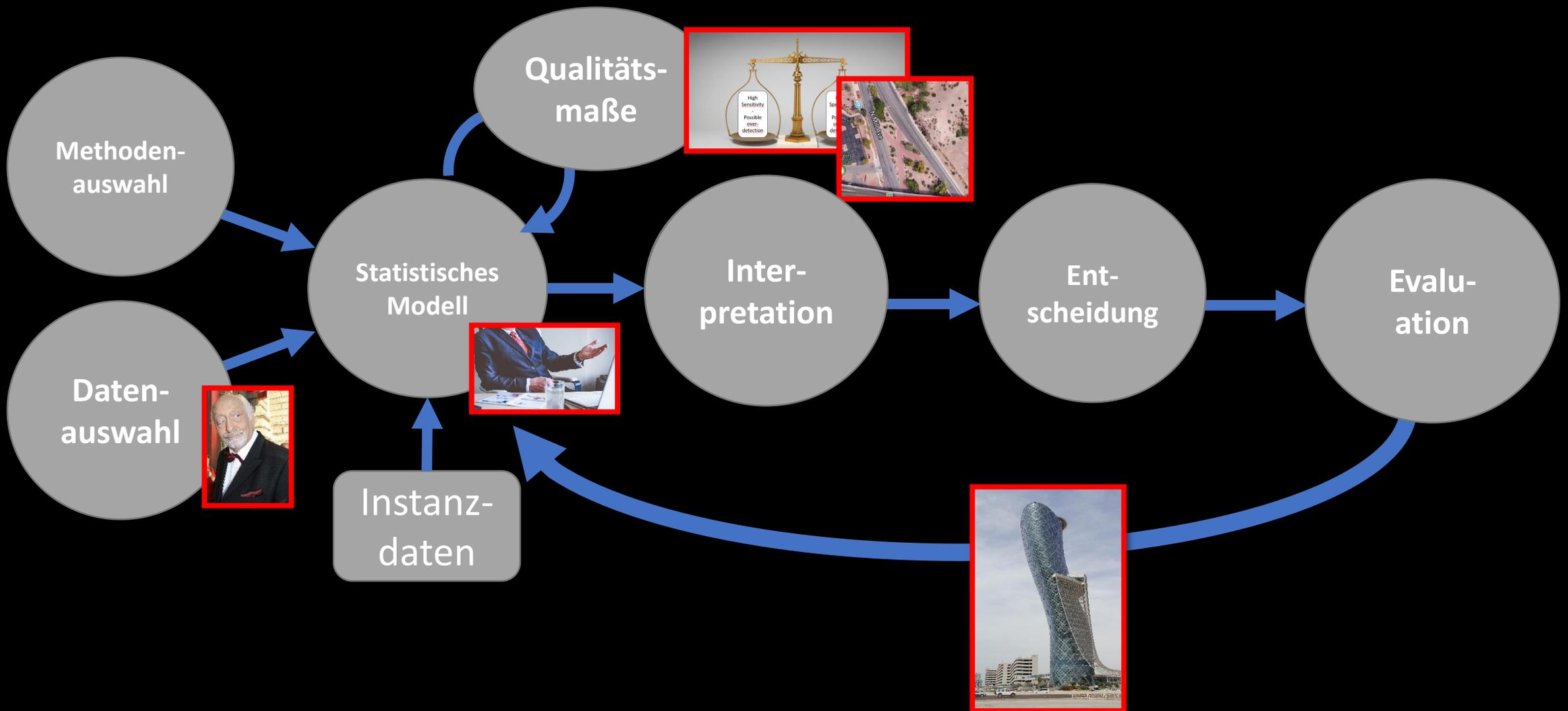




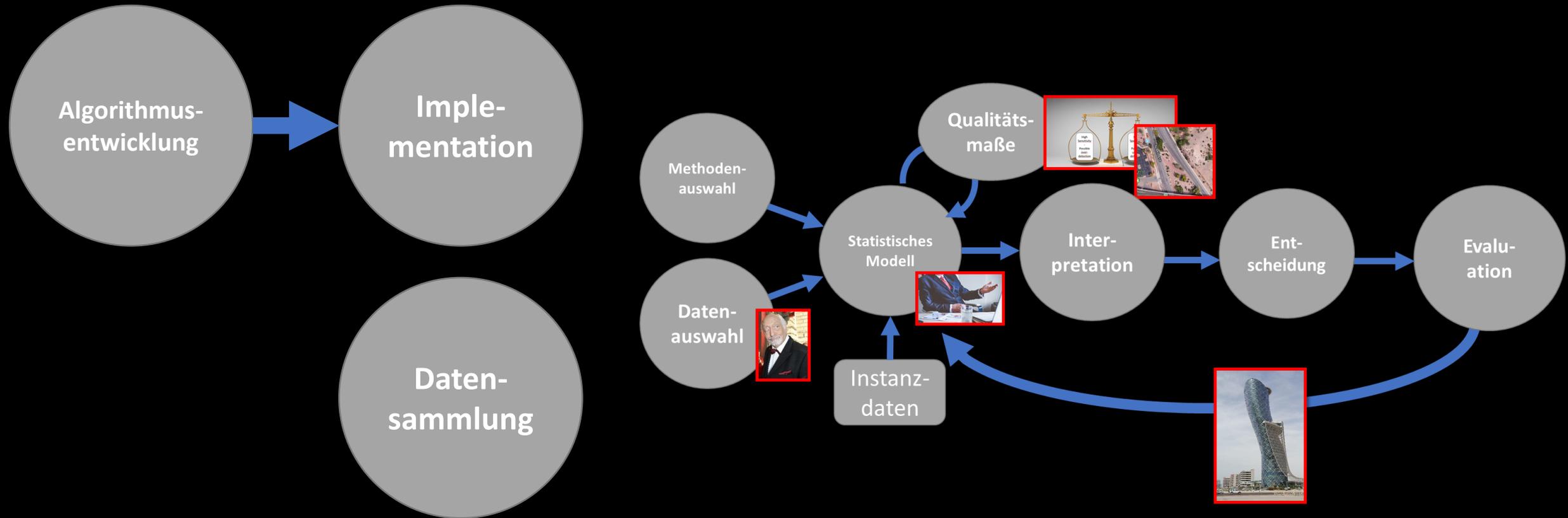


Feedback-Probleme

Personen, die wegen der Entscheidung des Sensors nicht fahren dürfen, können nicht nachweisen, dass sie sicher gefahren wären.

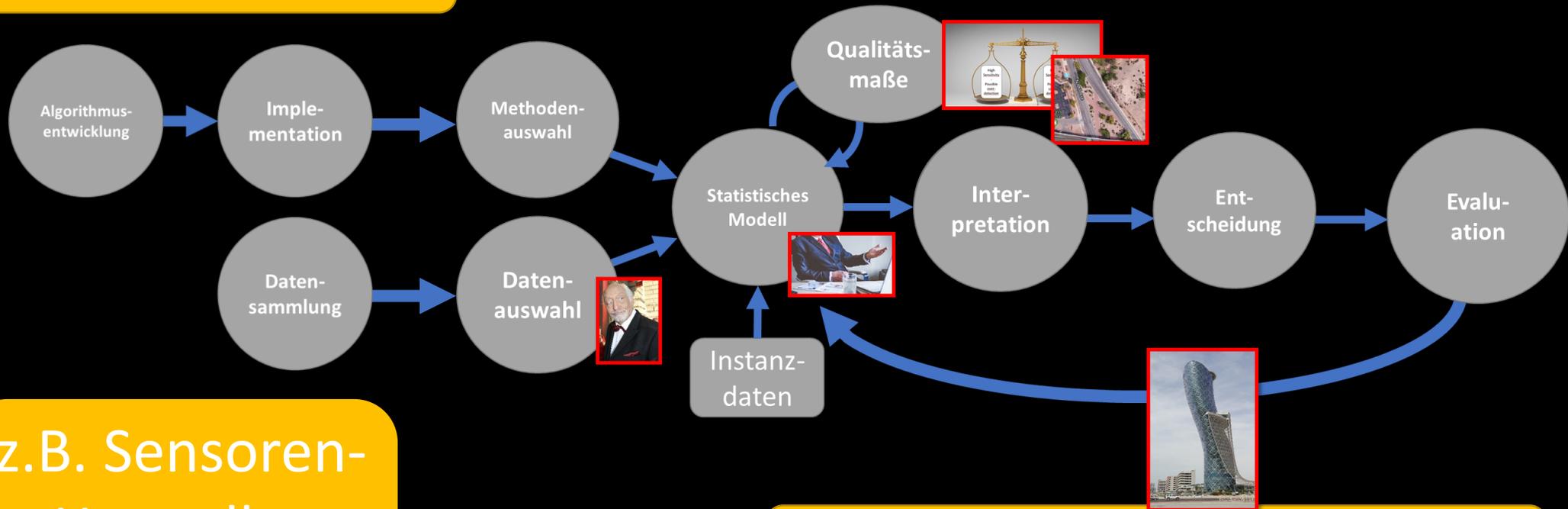


Wer ist verantwortlich?



Wer ist verantwortlich?

Informatiker



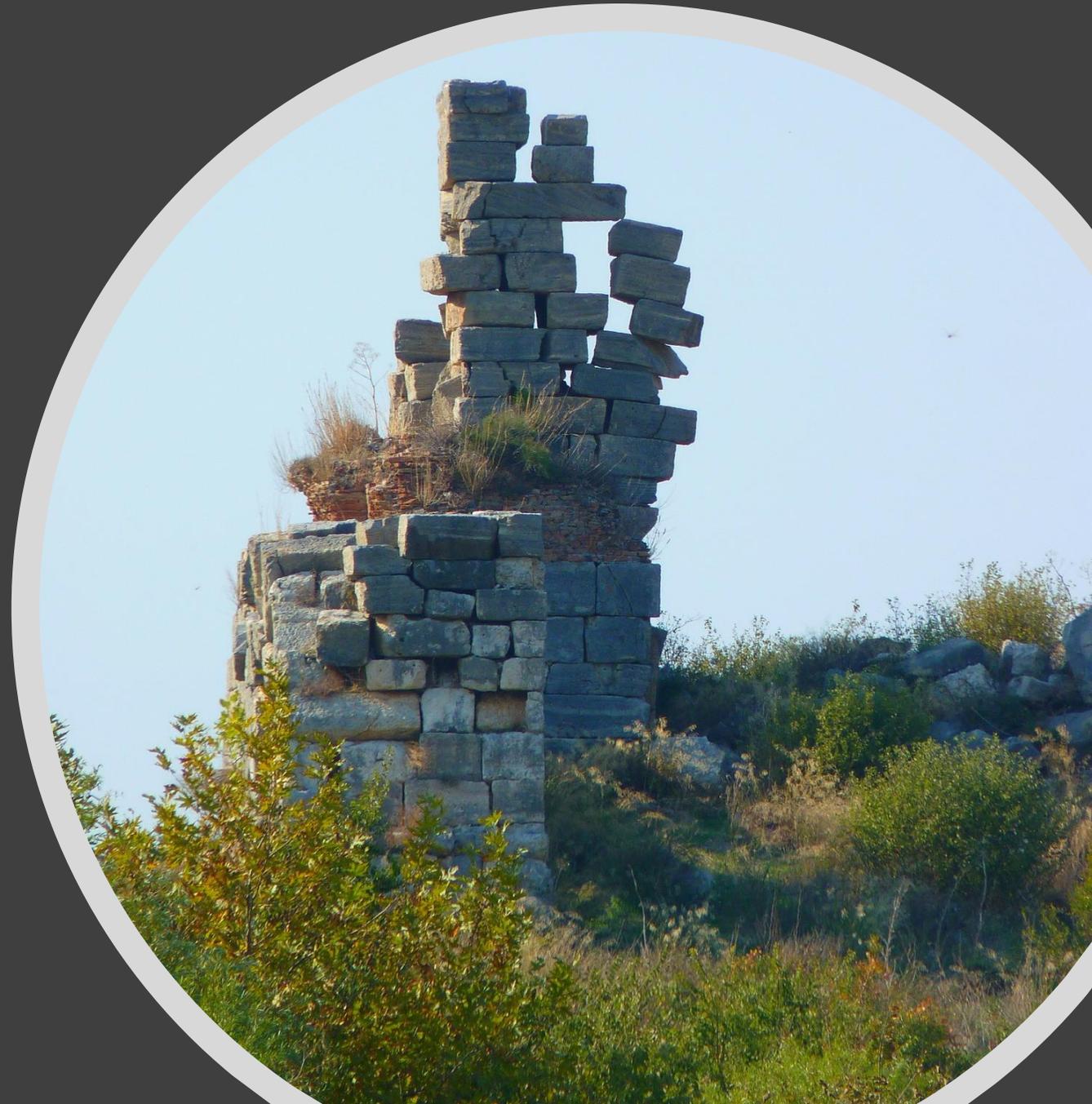
z.B. Sensoren-Hersteller

Kunden Ihres ADM Systems

Data Scientists / Ingenieure

Kann ein fehlendes ADM System
unethisch sein?

Und darüber
hinaus?



Der Uber-Unfall in Tempe

Warning

Some viewers may find the following footage distressing

**The
Guardian**

Zwischen- bericht

“According to Uber, emergency braking maneuvers are not enabled while the vehicle is under computer control, to reduce the potential for erratic vehicle behavior.”

“The vehicle operator is relied on to intervene and take action.”

“The system is not designed to alert the operators.”

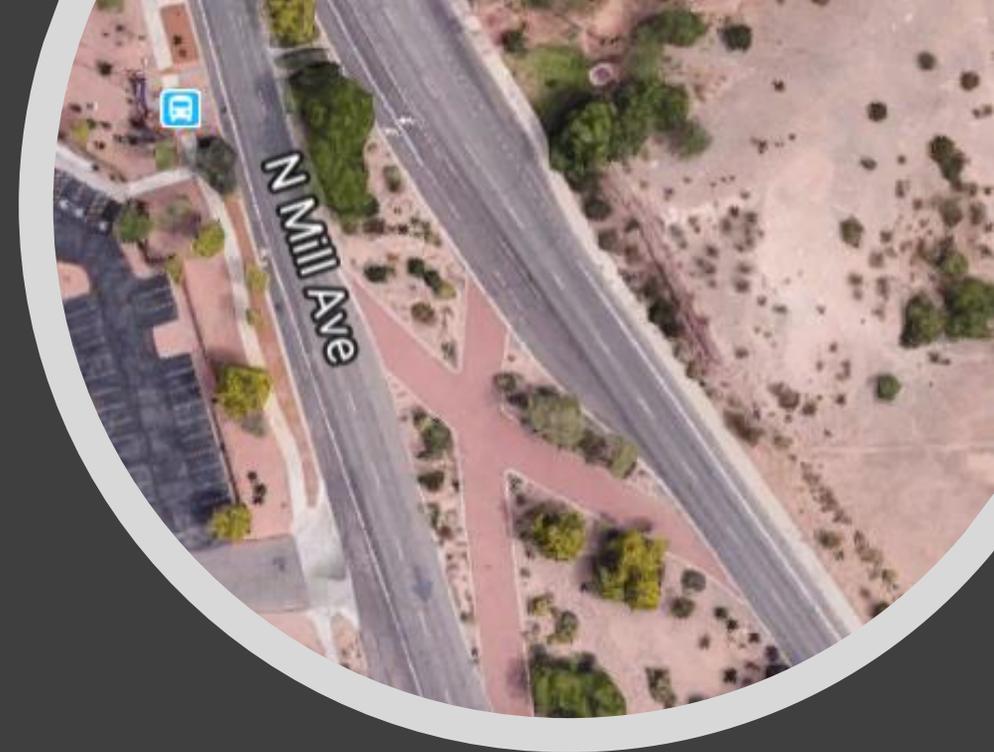
<https://www.nts.gov/investigations/AccidentReports/Reports/HWY18MH010-prelim.pdf>



<https://www.google.de/maps/place/N+Mill+Ave,+Tempe,+AZ+85281,+USA/@33.4364084,-111.9436953,335m/data=!3m1!1e3!4m5!3m4!1s0x872b09338c84a6a3:0x4eb48ca97885c3f7!8m2!3d33.4379952!4d-111.9435544>

Welche ethischen Entscheidungen traf Uber?

- Wurde die Wahrscheinlichkeit für Fußgängerübergänge von Straßenkarten und Satellitenbildern abgeleitet?
- Balance zwischen “Sicherheit (anderer)” und “Komfort der Passagiere”.
- Es fehlte ein Alarmsystem!
- Warum kein Alarmsystem, wenn der “operator” unaufmerksam ist?



Wie können wir
bessere Systeme
bauen?



Best Practice für die Entwicklung und Nutzung von ADM Systemen, die Menschen bewerten

- „Normale“ Algorithmen sind gut überprüfbar.
- Fokus auf **lernende Systeme**, deren Entscheidungen Menschen betreffen.
- Wichtig sind Entwicklungsteams **mit hoher Diversität**, sowohl was die Personen als auch Ausbildungen angeht.
- Weiterbildung im Bereich „Data Science Literacy“ und „Ethik für Data Scientists“.
- **Klare Kommunikationsprozesse** entlang der „langen Kette der Verantwortlichkeiten“:
 - **Z.B. Fehlerraten der Daten von Sensoren.**
 - Über alle Designentscheidungen, die getroffen wurden, und in welchem Kontext sie gültig sind.
- **Klare Abgrenzungen der jeweiligen Verantwortlichkeiten.**

„There are no simple solutions for complex problems.
Whoever promises them, is fooling himself and others.“

Volkswagen at the IAA 2017