



Können Algorithmen diskriminieren?



Wie funktioniert sie, wo ist sie hilfreich,
wo sollte sie nicht eingesetzt werden?

Prof. Dr. Katharina A. Zweig,
TU Kaiserslautern

Konstituierende Sitzung der
Enquete-Kommission
„Künstliche Intelligenz“ am 27.9.

Aus der Rede von Bundestagspräsidenten
Dr. Schäuble:

- „Die künstliche Intelligenz gilt
Vielen als neue Zauberformel des
technischen Fortschritts, ...
- ... sie wird dichten, ...
- ... sie wird belohnen und bestrafen ...“



Die zwei Ängste

Sie wird dichten

Sie wird richten





Aber was ist überhaupt KI?

SNV, „Eckpunkte einer nationalen Strategie für Künstliche Intelligenz“:

„(...) Verfahren (...),
welche die Übertragung
von bislang
menschlich getroffenen
Entscheidungen,
Bewertungen und
Handlungen
**auf Computer und
Maschinen erlauben.**“

Nun, das deckt ganz schön viel ab....

Superintelligenz



KI

Data
Science

ML

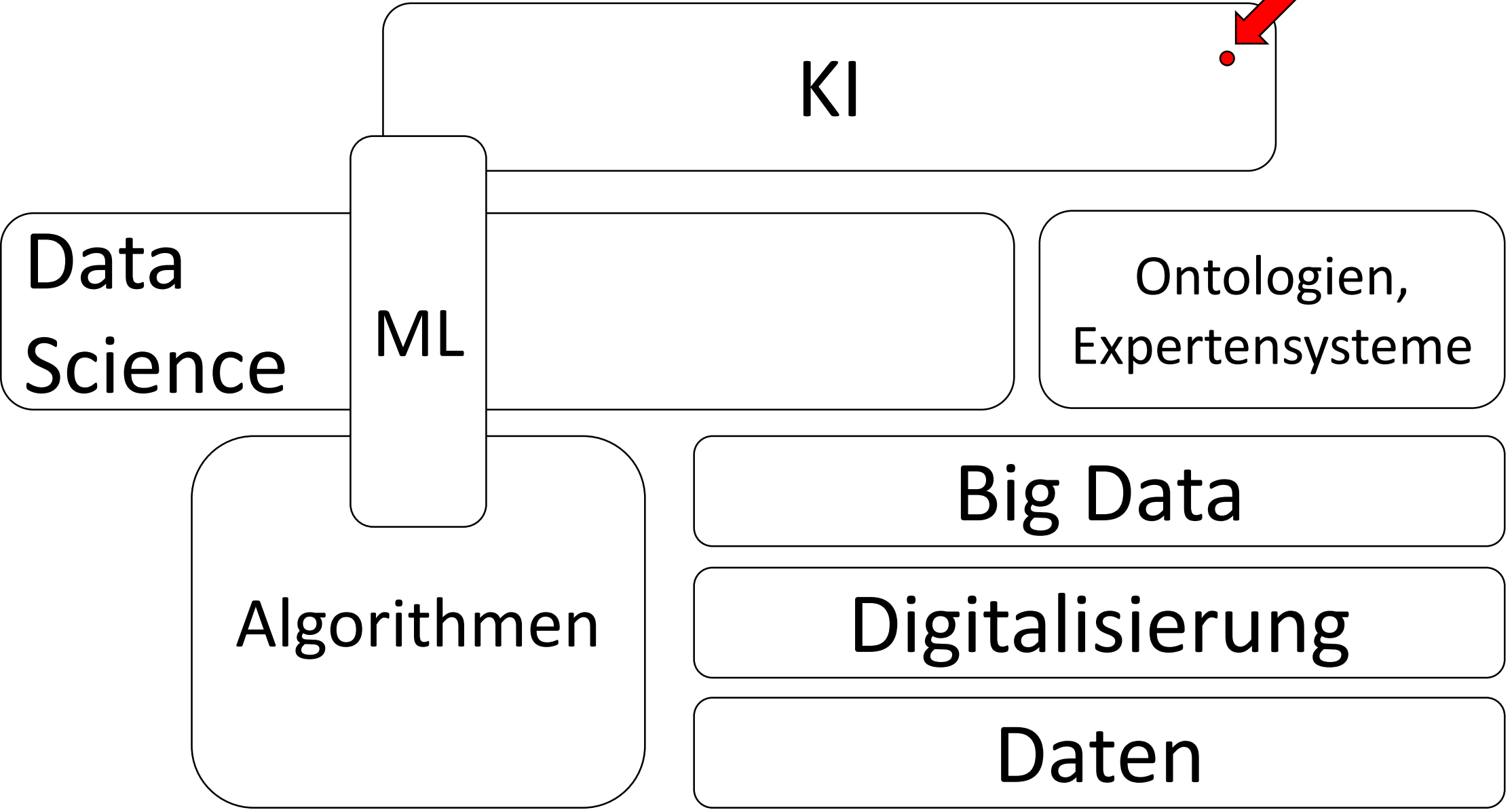
Ontologien,
Expertensysteme

Algorithmen

Big Data

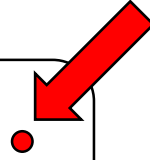
Digitalisierung

Daten



Nun, das deckt ganz schön viel ab....

Superintelligenz



Data
Science

„KI“ im
öffentlichen
Diskurs

Ontologien,
Expertensysteme

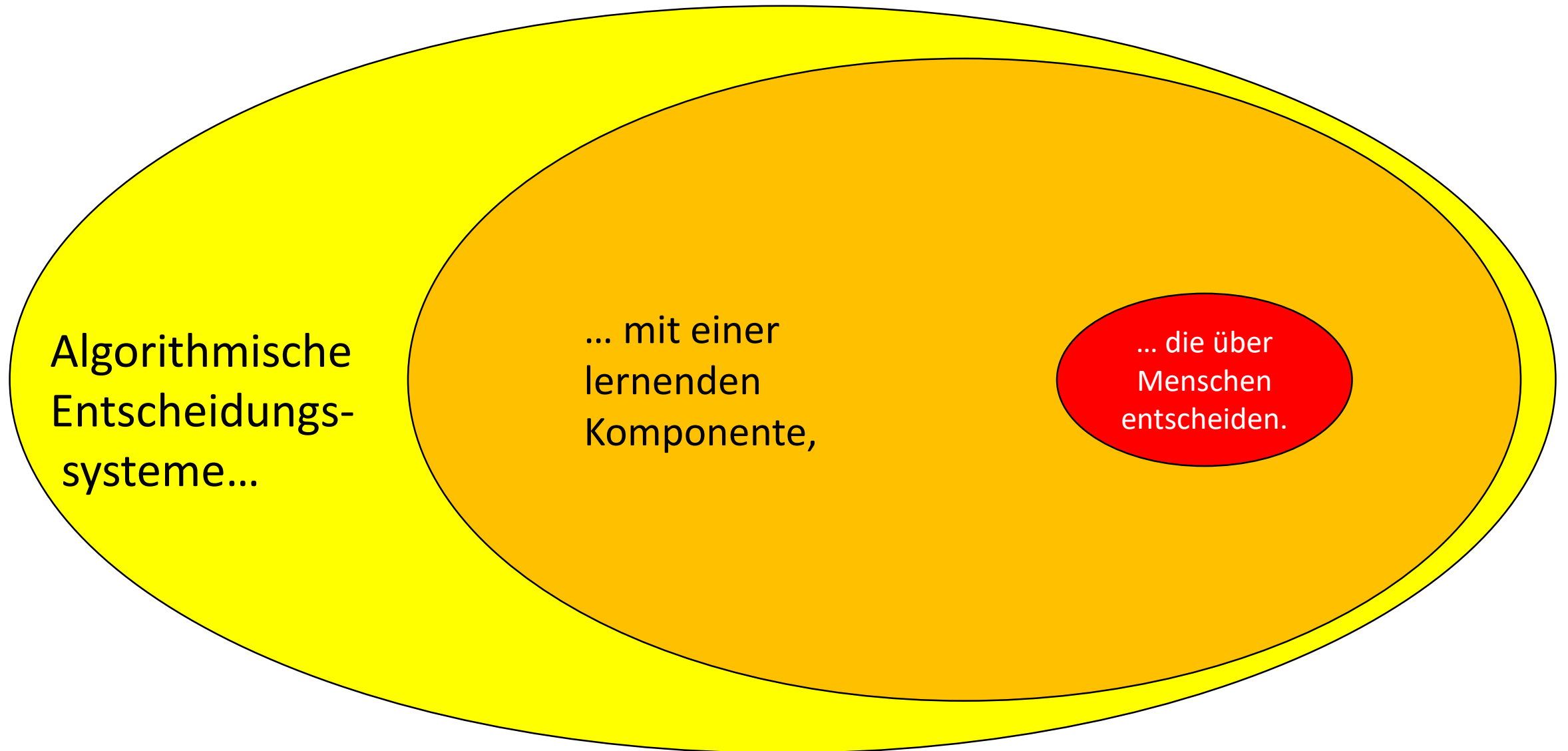
Big Data

Algorithmen

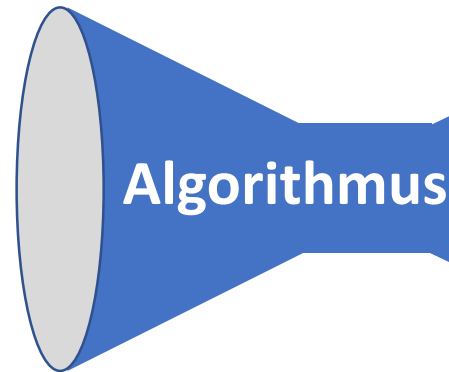
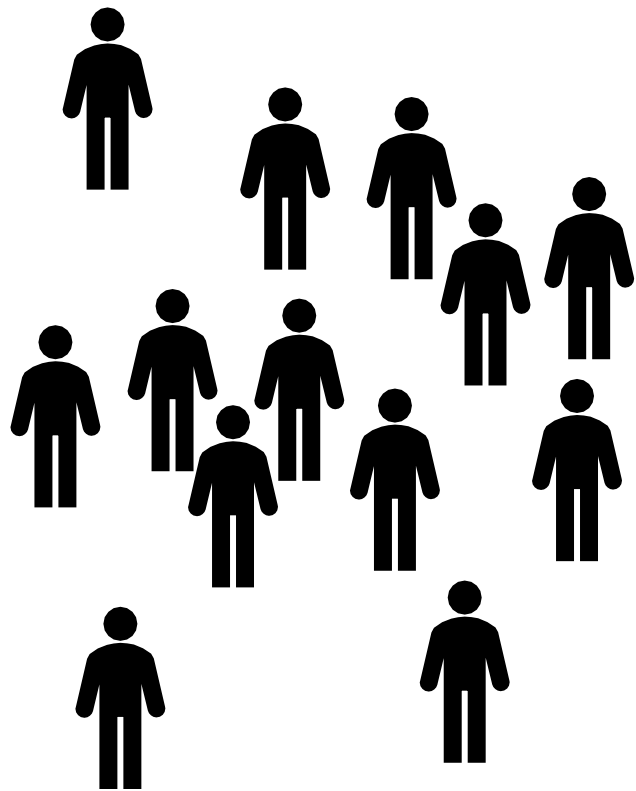
Digitalisierung

Daten

Welche algorithmischen Entscheidungssysteme (ADM Systeme) sind problematisch?



Algorithmische Entscheidungssysteme (ADM Systeme)



oder

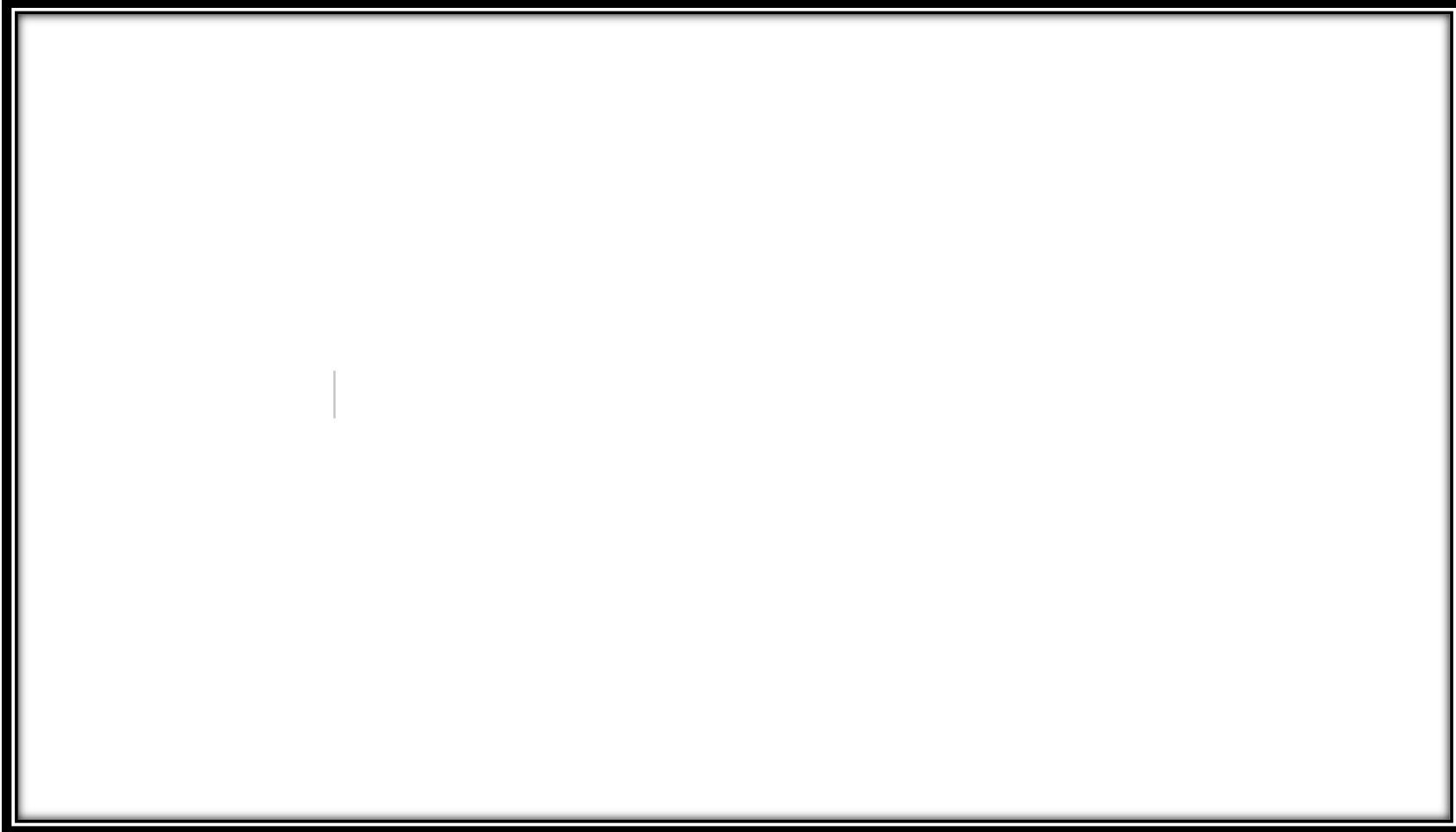


Scoring-Verfahren



Klassifikation

Kann KI HR?



Diskriminierung bei Bewerbungen

- Lebensläufe mit „deutschen“ Namen bekommen 14% Vorstellungsangebote 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 Ghumman, S. & Ryan, A. M.: "Not welcome here: Discrimination towards women who wear the Muslim headscarf , human relations, 2013 , 66(5) , 671-698

Könnten Computer das besser?

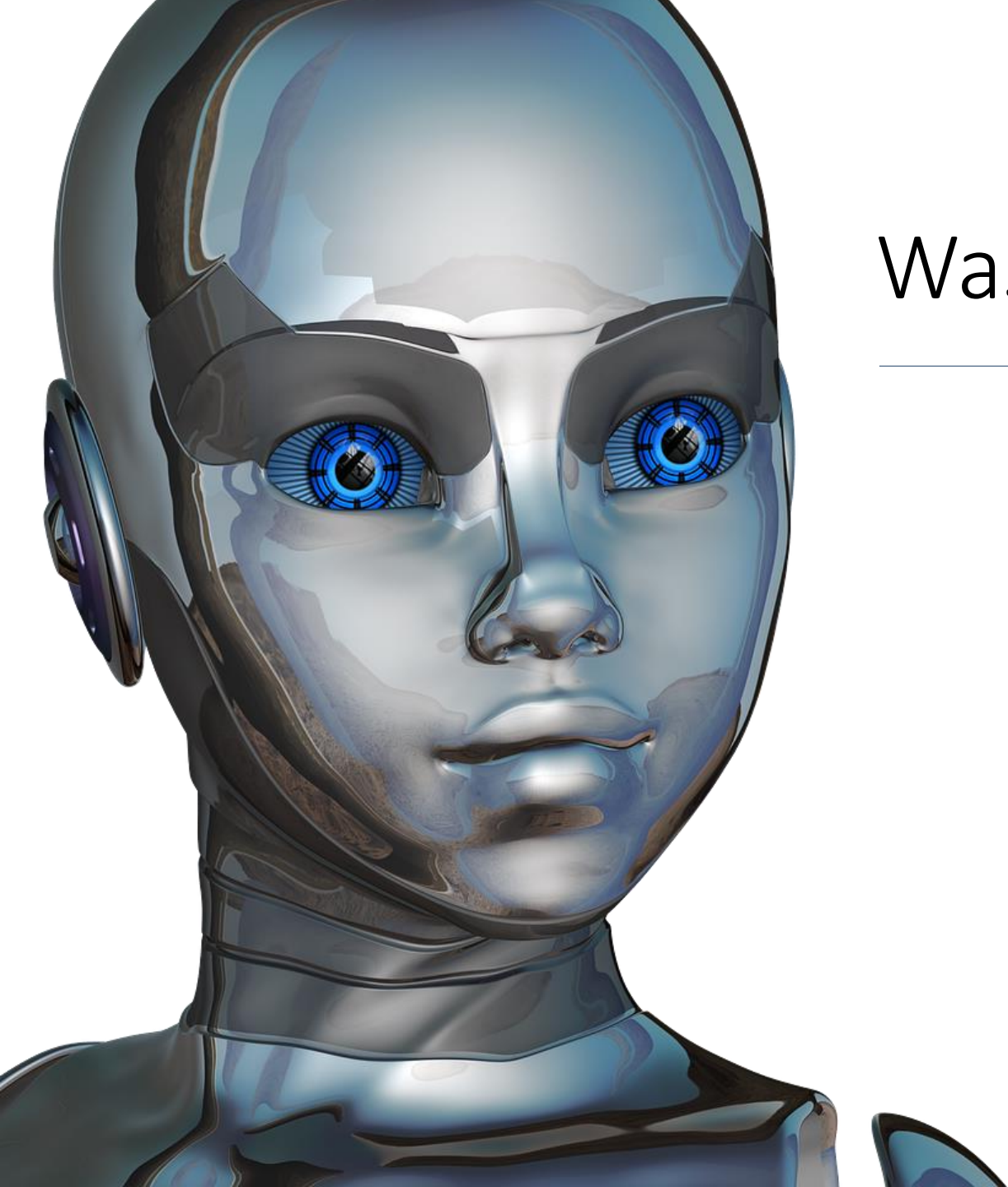
- Die ersten Firmen testen *algorithmische Entscheidungssysteme* (oder Entscheidungsunterstützungssysteme)¹.
- Eigenschaften, nach denen nicht diskriminiert werden darf, können vor ihnen besser verborgen werden.
- Sie sind objektiv und arbeiten nahezu fehlerfrei.
- (objektiv := „reproduzierbar dieselbe Entscheidung bei derselben Eingabe von Daten“)



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>



Können Computer lernen?



Was heißt Lernen?

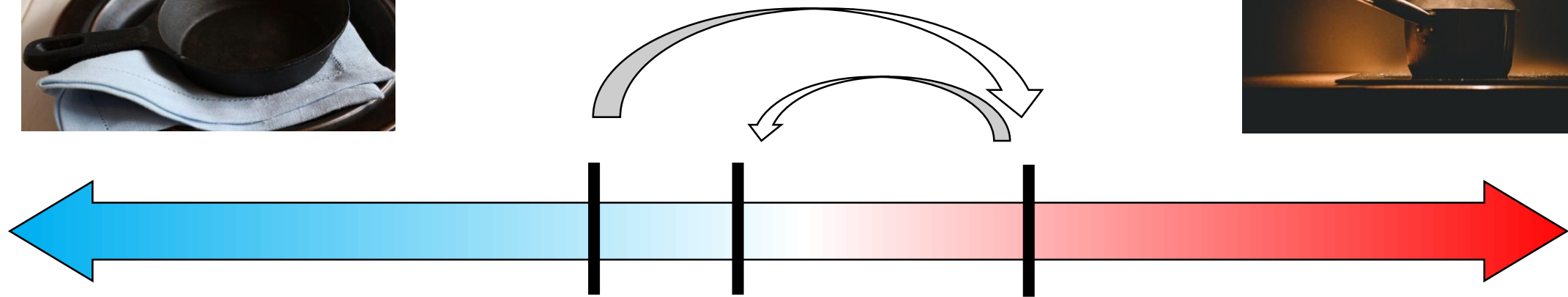
Einfach:

In derselben Situation ein vorher gezeigtes Verhalten wiederholen.

Generalisiert:

In derselben Art von Situation das richtige Verhalten aus einer Reihe von Möglichkeiten auswählen.

Sebastian lernt „heiss“ und „warm“

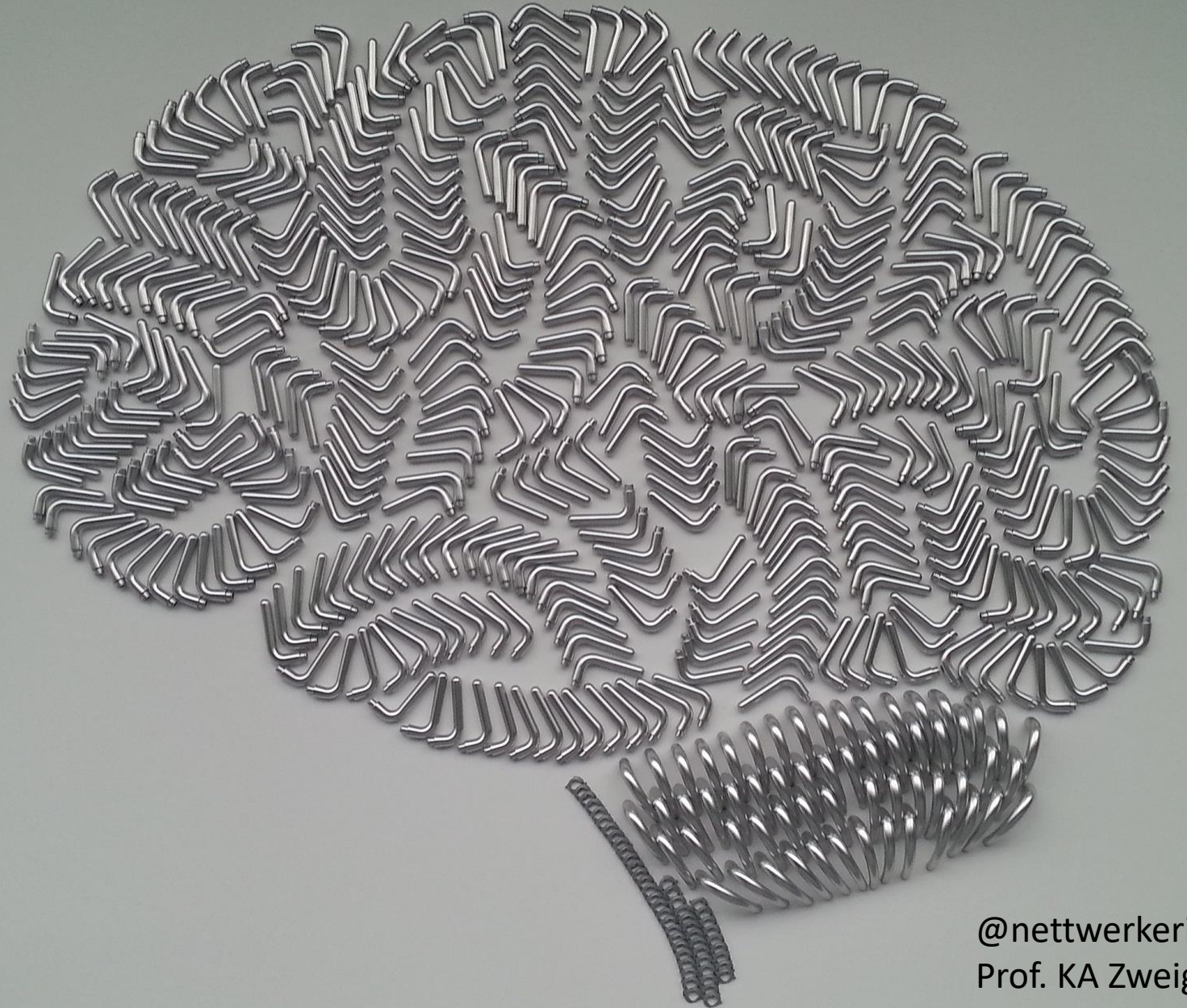


Juli März September

**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.**

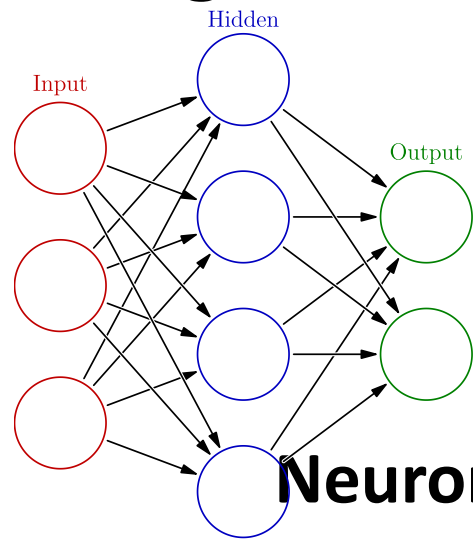


Computer lernen

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

Optimal auch **Rückkopplung**.

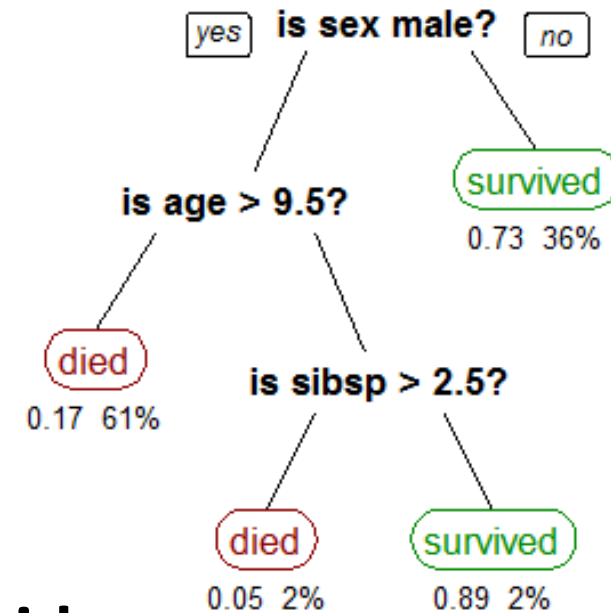
Er lernt **generelle Regeln**.



Neuronales Netz

Formel

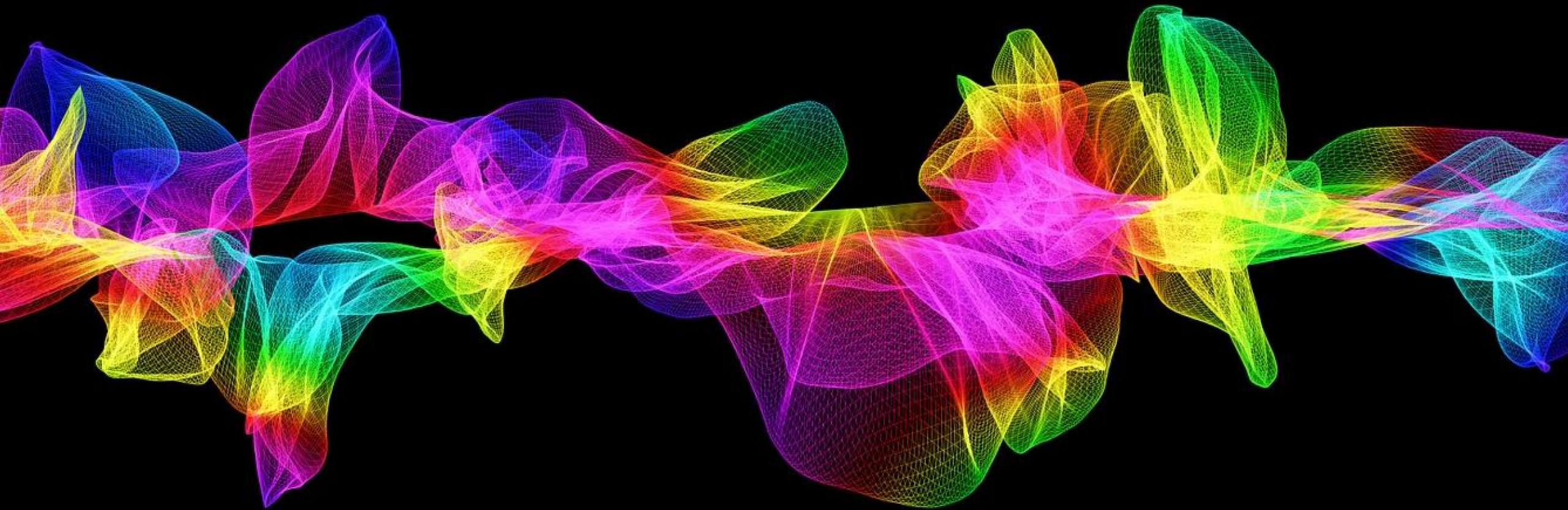
$$w_1 * \#V_h - w_2 * \#dayI V_h + w_3 * I[g = male] * 1 + w_4 * I[T = R] * 1.0 + \dots$$



**Entscheidungs-
bäume**

@nettwwerkerin
Prof. KA Zweig
TU Kaiserslautern

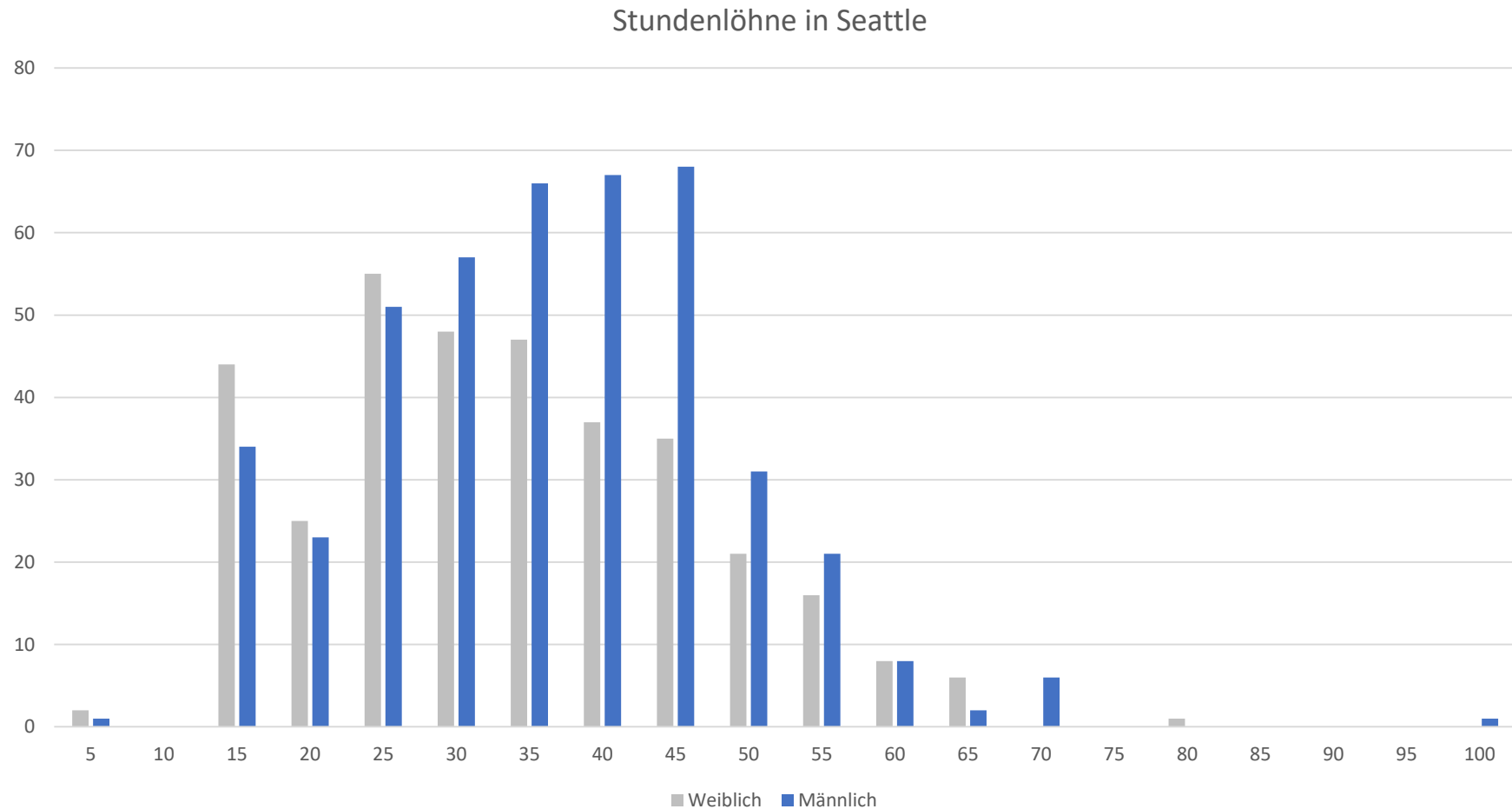
"CART tree titanic survivors" by Stephen Milborrow - Own work. Licensed under CC BY-SA 3.0 via Wikimedia Commons - https://commons.wikimedia.org/wiki/File:CART_tree_titanic_survivors.png#/media/File:CART_tree_titanic_survivors.png
By Glosser.ca - Own work, Derivative of File:Artificial neural network.svg, CC BY-SA 3.0, <https://commons.wikimedia.org/w/index.php?curid=24913461>



“Lernen” mit Korrelationen

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

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





“Lernen” mit SVMs

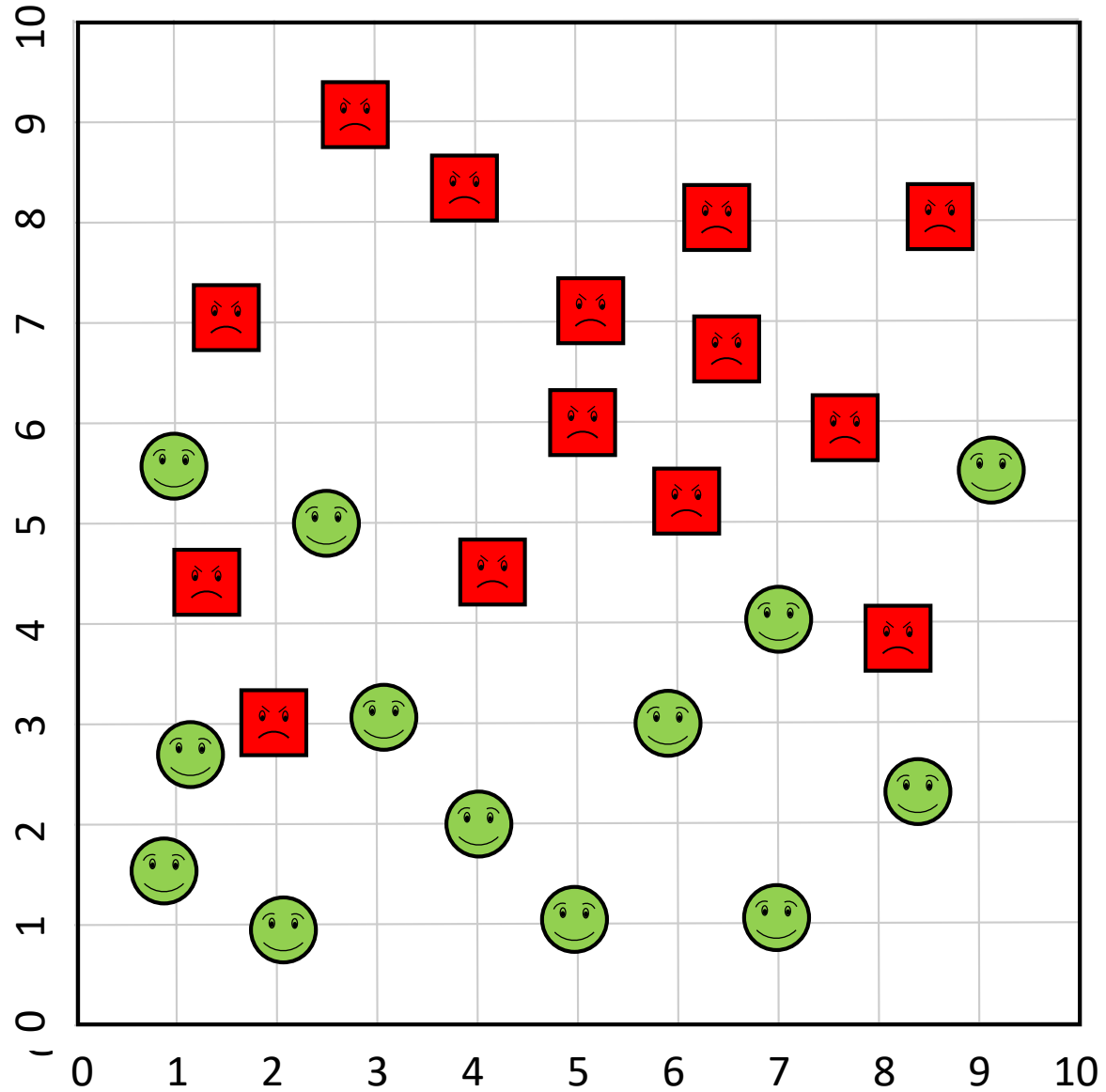


Weniger erfolgreiche
Arbeitnehmer:innen



Erfolgreiche Arbeit-
nehmer:innen

Jahre arbeitslos



Expertise



Weniger erfolgreiche Arbeitnehmer:innen



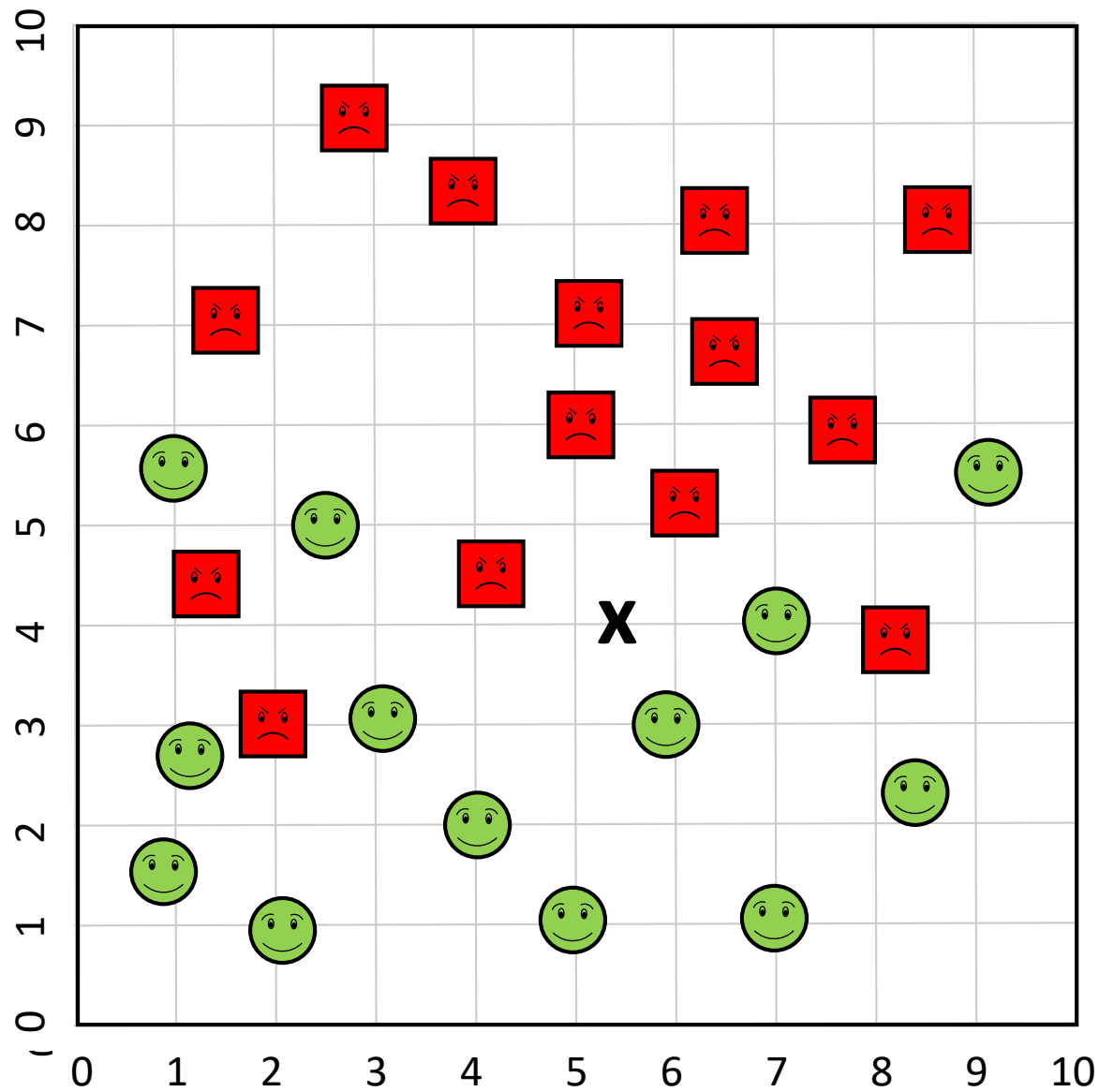
Erfolgreiche Arbeitnehmer:innen

Bewerten Sie Frau Müller:

5.5 Jahre Erfahrung

4 Jahre arbeitslos

Jahre arbeitslos



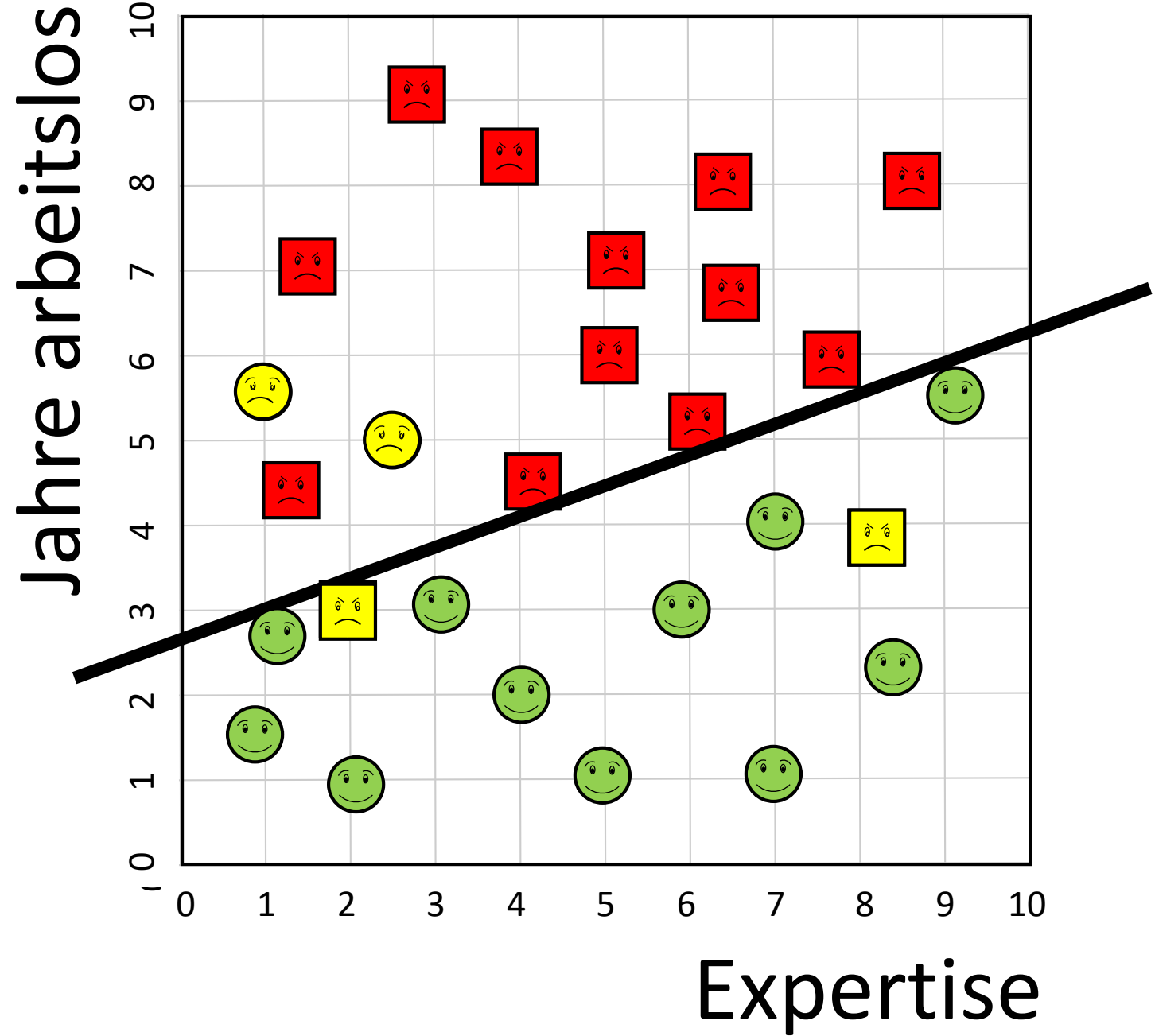
Expertise

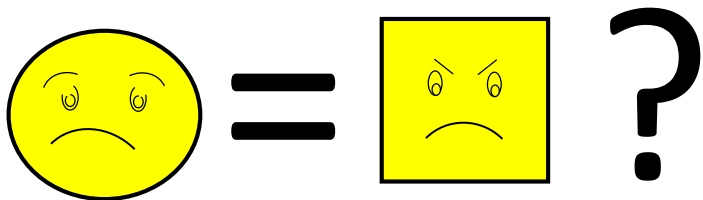




Weniger erfolgreiche Arbeitnehmer:innen

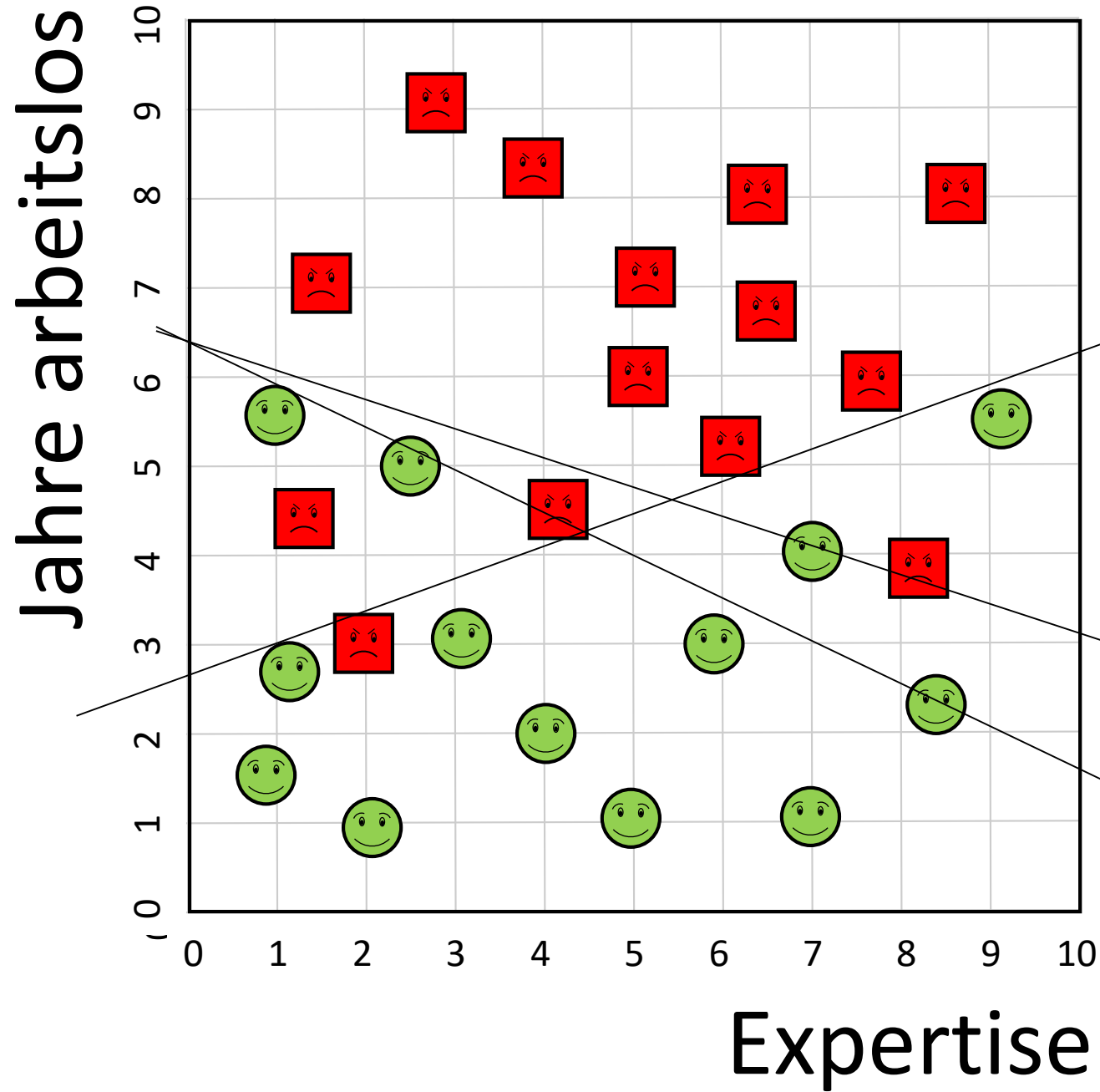


Erfolgreiche Arbeitnehmer:innen





-  Weniger erfolgreiche Arbeitnehmer:innen
-  Erfolgreiche Arbeitnehmer:innen



Datengrundlagen

- Meistens mehr als zwei Eigenschaften.
- Am wichtigsten:
 - **War Einstellung erfolgreich?**

Ausbildung

Jahre der
Arbeitslosigkeit

Alter

Arbeitgeber
-wechsel

Bewerbungs-
schreiben

Rechtschreibung

Wortvielfalt

Ton

Social Media?



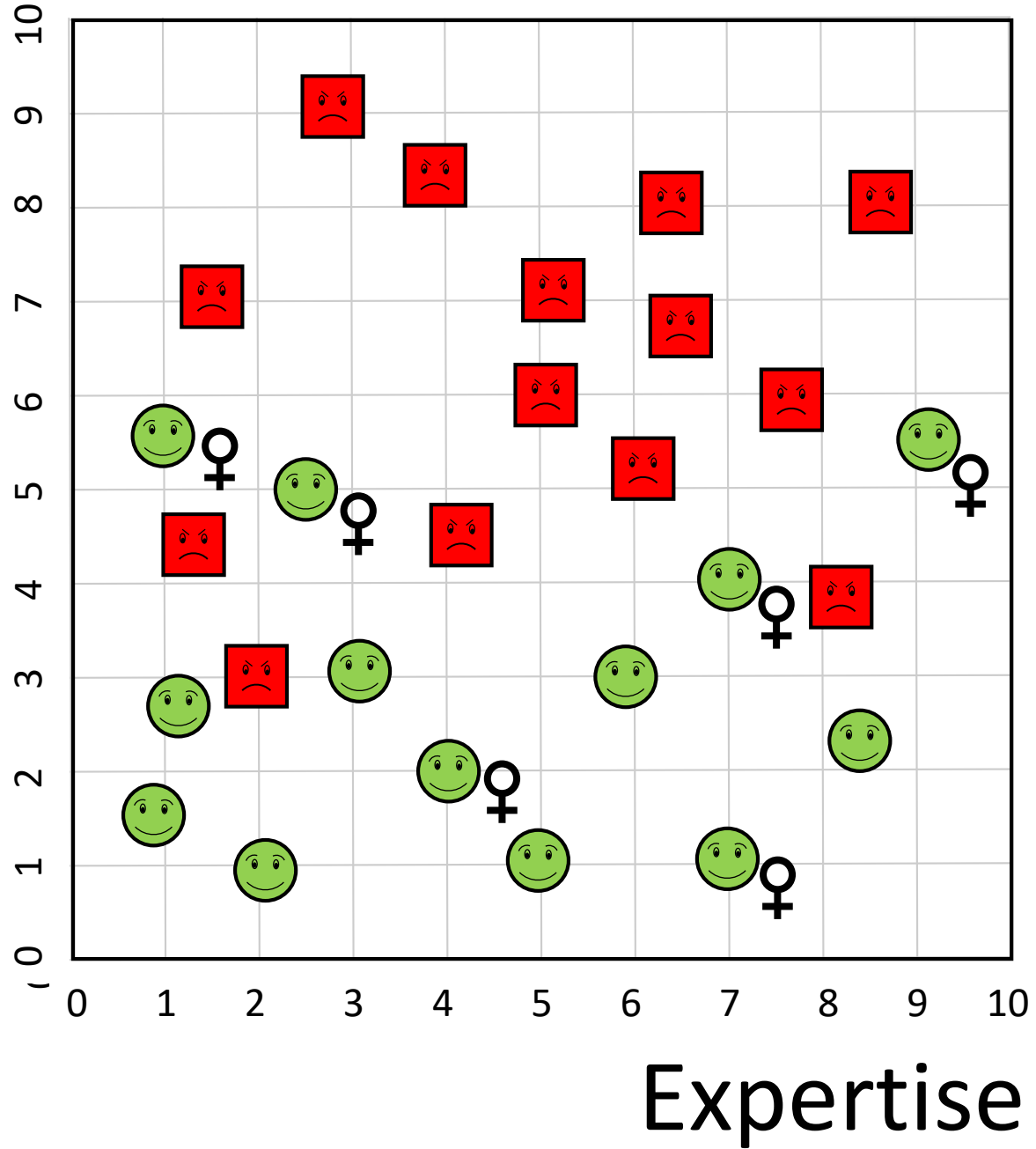
Weniger erfolgreiche
Arbeitnehmer:innen



Erfolgreiche Arbeit-
nehmer:innen

Daten- grundlage

Jahre arbeitslos





Ausbildung

~~Leerzeiten~~

~~Alter~~

Arbeitgeber
-wechsel

Bewerbungs-
schreiben

Rechtschreibung

Wortvielfalt

Ton

Social Media?

Systemen:
Einstellung
gleich?



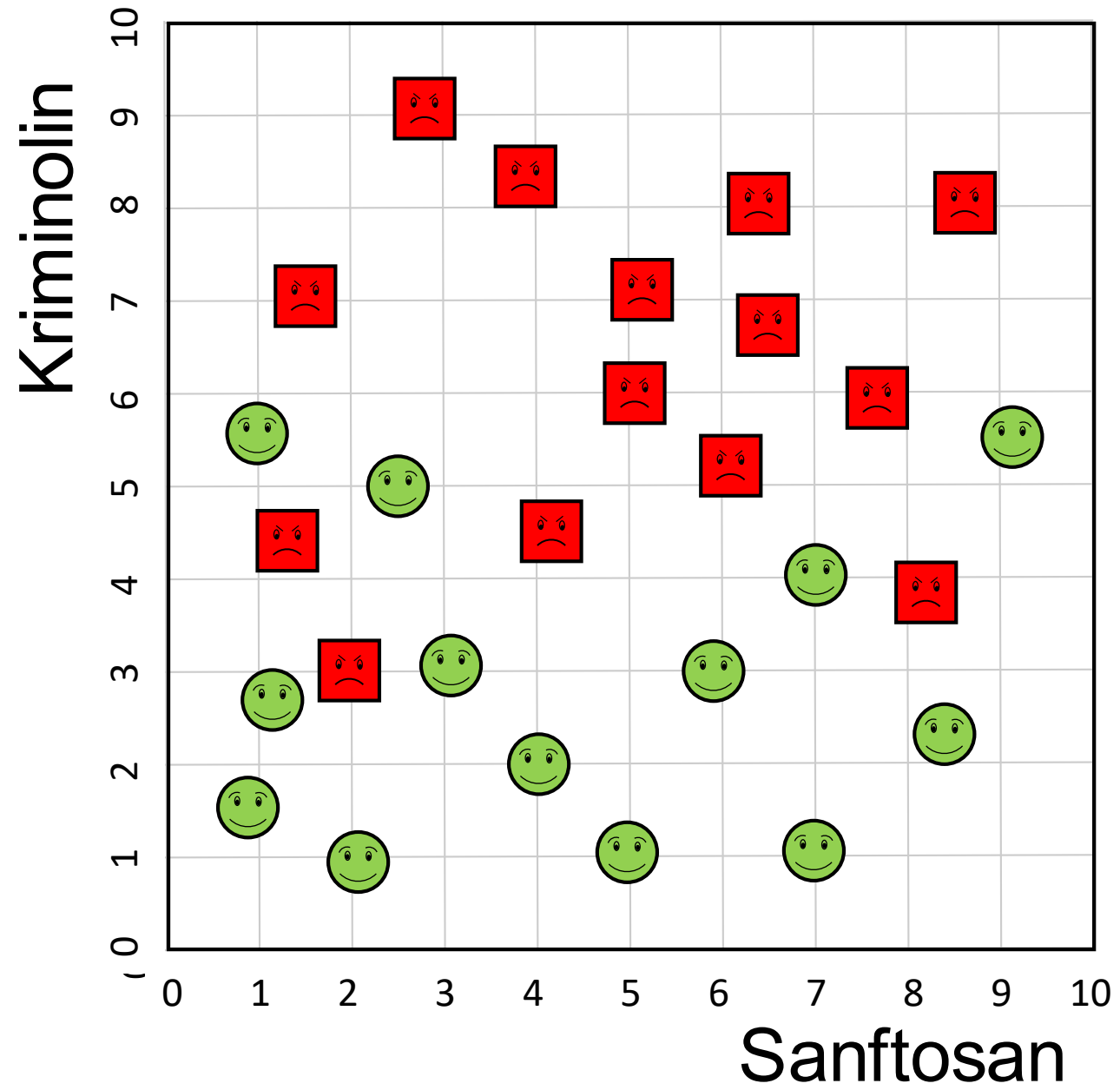
Qualität eines Algorithmus



Bösartige Kriminelle

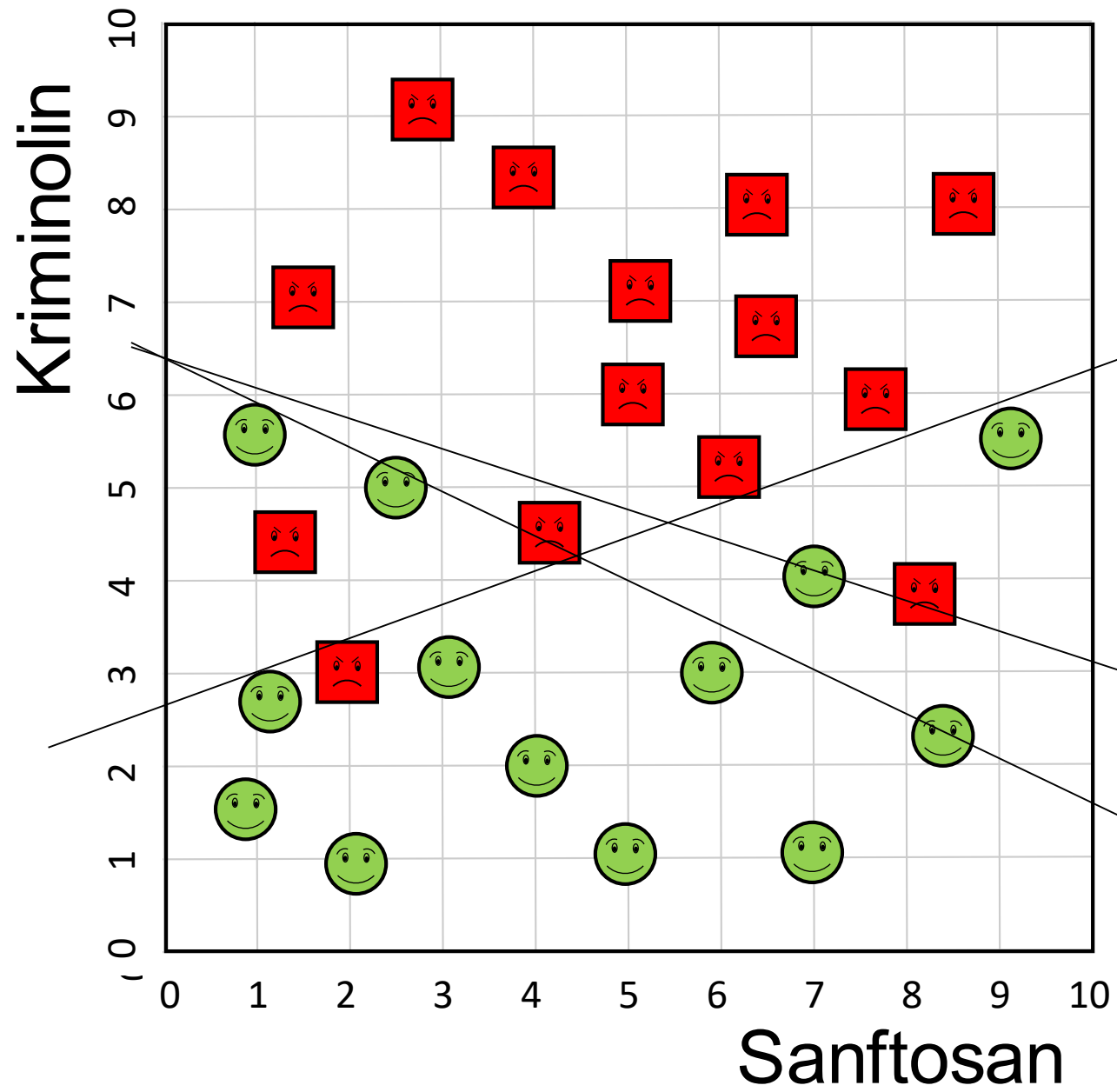
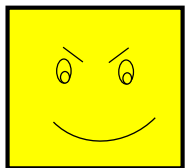


Unschuldige Bürger



 Böartige Kriminelle

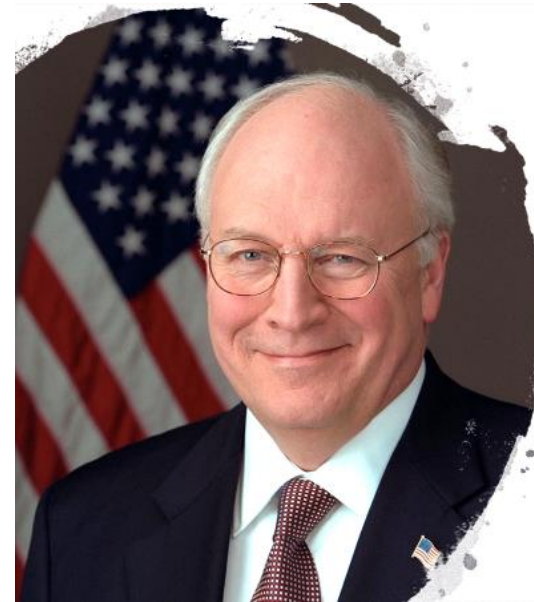
 Unschuldige Bürger





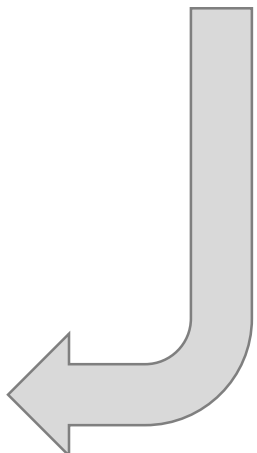
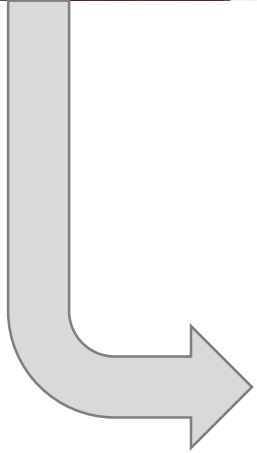
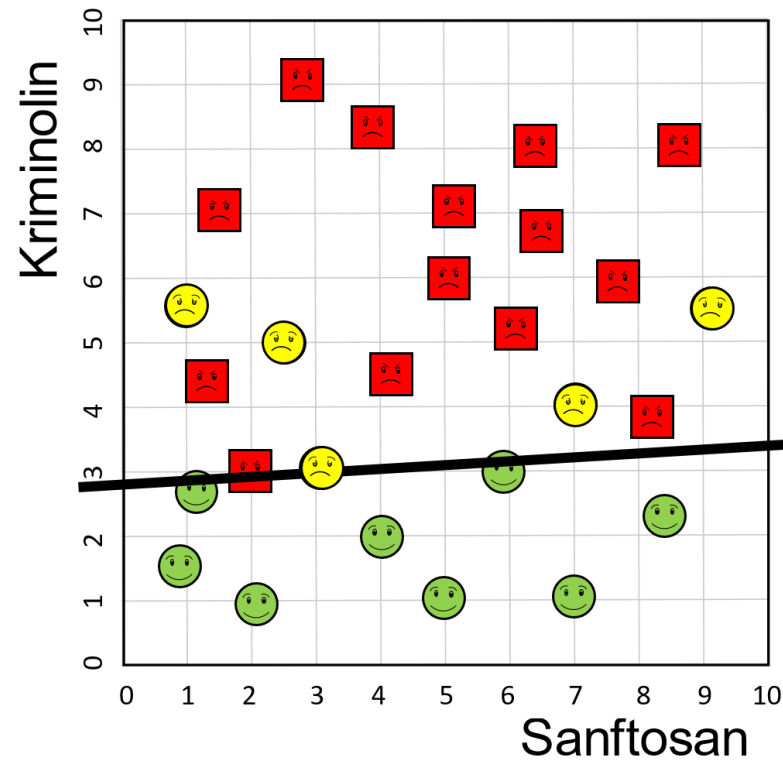
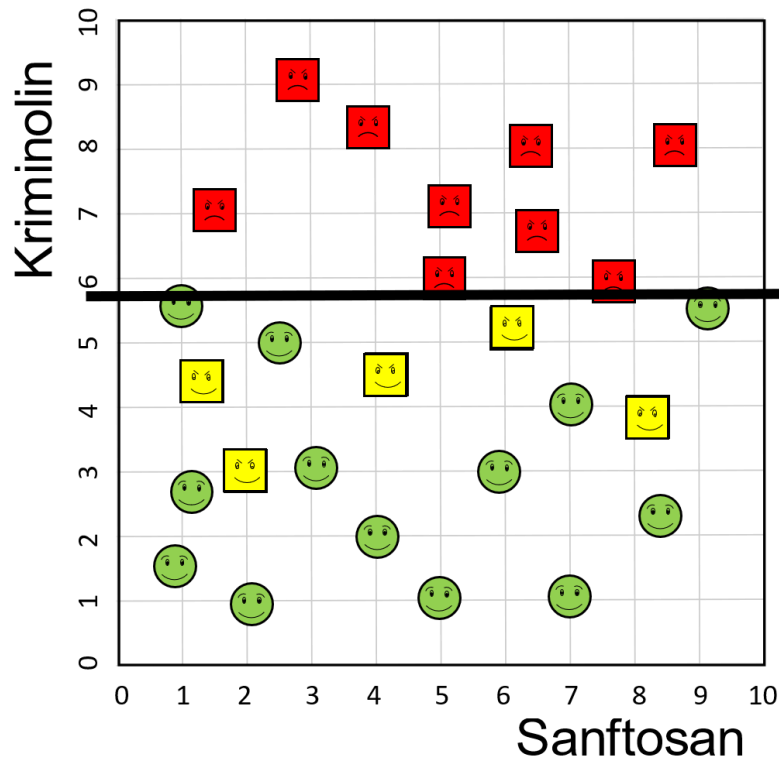
„It is better that ten guilty persons escape than that **one** innocent suffer.“

William Blackstone, Rechtsphilosoph, 1760



"I am more concerned with bad guys who got out and released than I am with a few that, in fact, were innocent."

Dick Cheney, ehemaliger Vizepräsident der USA,



Qualität von ADM Systemen

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

```
is},a(window).on( load...
e strict";function b(b){return this.each(function(){var
ction(b){this.element=a(b)};c.VERSION="3.3.7",c.TRANSITION_D
.data("target");if(d||(d=b.attr("href"),d=d&&d.replace(/.*(?
ide.bs.tab",{relatedTarget:b[0]}),g=a.Event("show.bs.tab",{r
ar h=a(d);this.activate(b.closest("li"),c),this.activate(h,h
.bs.tab",relatedTarget:e[0]}))}}},c.prototype.activate=fun
Class("active").end().find('[data-toggle="tab"]').attr("ar
b[0].offsetWidth,b.addClass("in"):b.removeClass("fade"),l
="tab"]').attr("aria-expanded",!0),e&&e()}var g=d.find(">
e").length);g.length&&h?g.one("bsTransitionEnd",f).emula
tab=b,a.fn.tab.Constructor=c,a.fn.tab.noConflict=functionio
"click.bs.tab.data-api",[data-toggle="tab"],e).on("
return this.each(function(){var d=a(this),e=d.data(
function(b,d){this.options=a.extend({},c.DEFAULTS,c
,this)).on("click.bs.affix.data-api",a.proxy(thi
is.checkPosition());c.VERSION="3.3.7",c.RESET=
his.$target.scrollTop(),f=this.$element.offse
l=c?!(e+this.unpin<=f.top)&&"bottom":!(e+
bottom"},c.prototype.getPinnedOffset=fu
get.scrollTop(),b=this.$element.of
this.checkPosition,this) 1))
```

Qualität von ADM Systemen

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



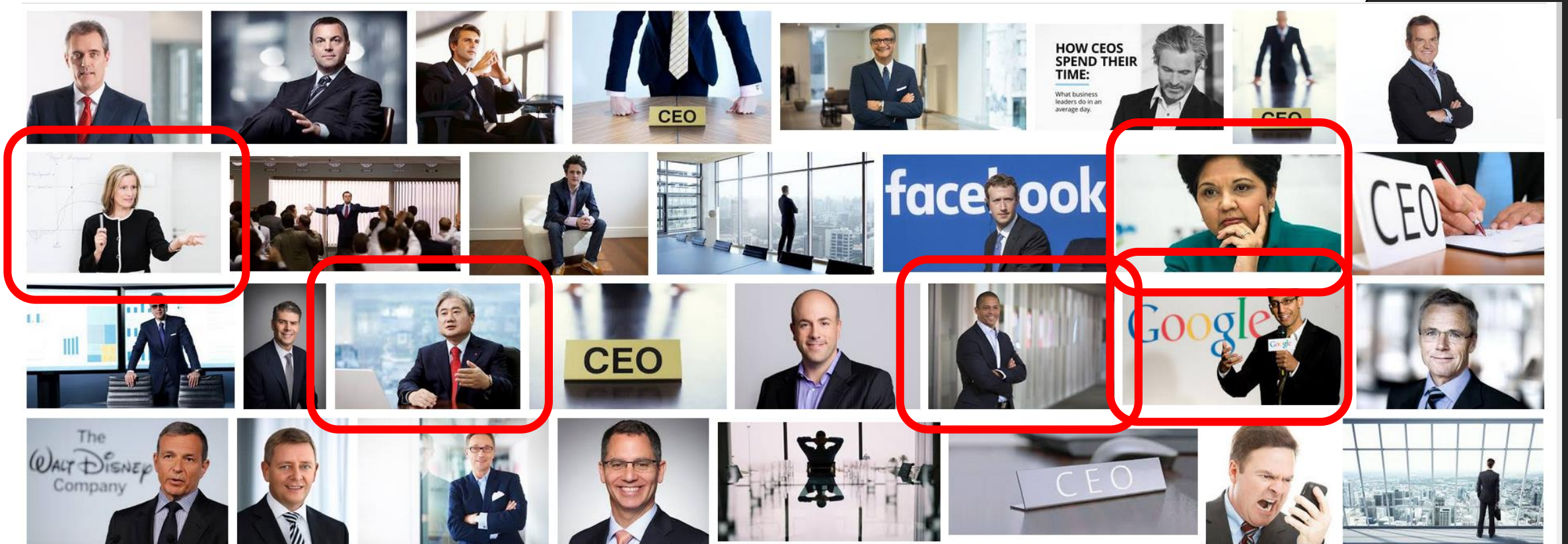


Computer: Objektiv und fehlerfrei?

Gleichberechtigung




Wenn man auf Google nach „CEO“ sucht...

Passt das einigermaßen: ca. 8%
der (deutschen) CEOs sind Frauen



Übersetzungen

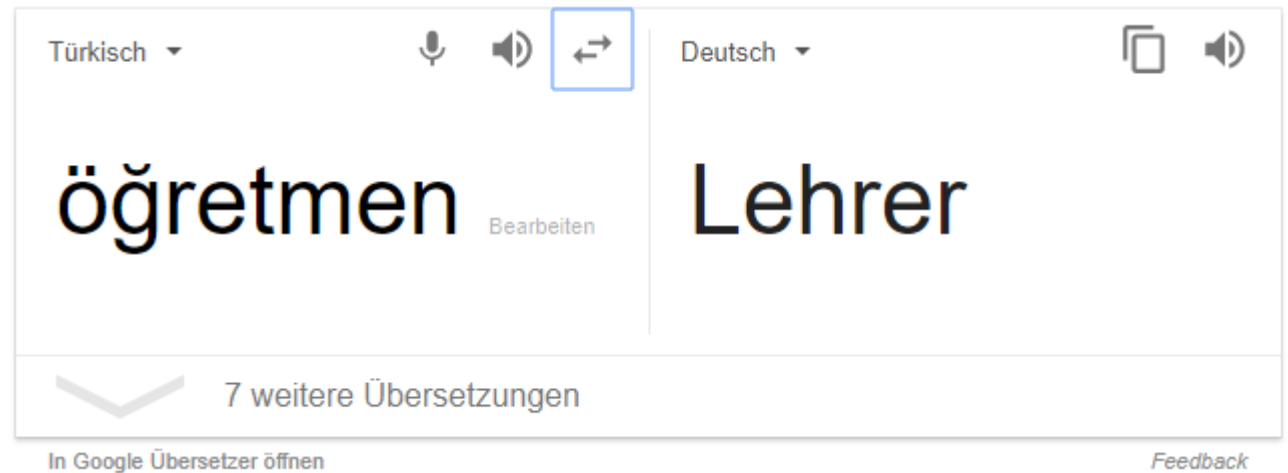






Deutsch ▾    Türkisch ▾  

lehrerin Bearbeiten **öğretmen**

In Google Übersetzer öffnen Feedback

A red arrow points to the swap languages icon (two arrows pointing in opposite directions) between the source and target language dropdowns.



Türkisch ▾    Deutsch ▾  

öğretmen Bearbeiten **Lehrer**

 7 weitere Übersetzungen

In Google Übersetzer öffnen Feedback



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

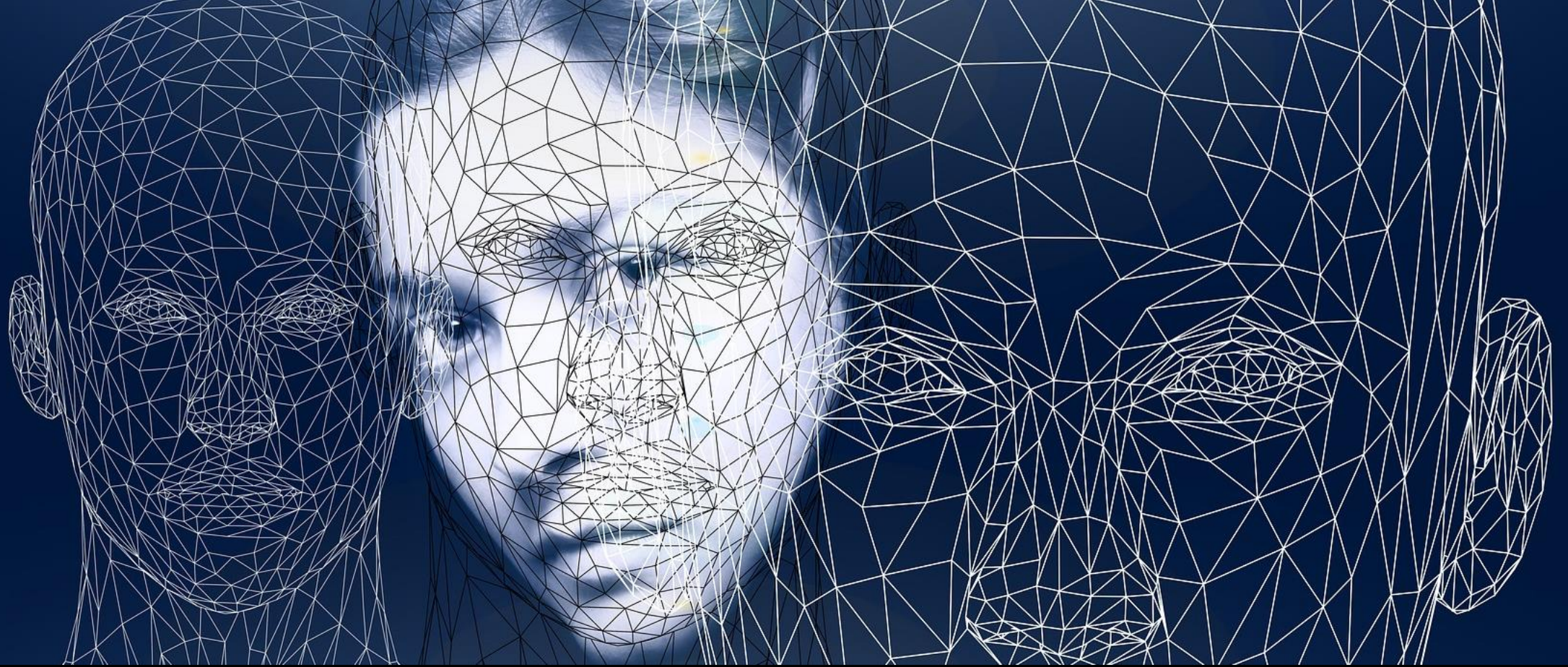
Diskriminierung

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

Qualität von ADM Systemen

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





Sozio-informatische Gesamtbetrachtung



Sozio-informatische
Betrachtungsweise

- KI schafft neue Anreize für menschliche Akteure.
- Diese reagieren auf die Anreize und wirken auf die KI ein.

Probleme der Einbettung der ADM in den sozialen Prozess

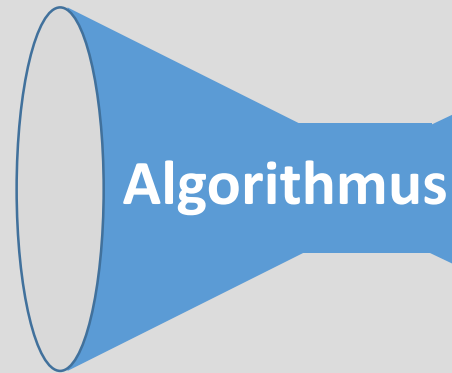
- **Aufmerksamkeitsökonomie** von Entscheiderinnen und Entscheidern.
- „**Best practice**“ erfordert Nutzung der Software.
- **Delegation von Verantwortung!**
- Manchmal kann ein falsch Beurteilter **die Vorhersage prinzipiell nicht entkräften!**
 - Z.B. abgelehnte Bewerberin, inhaftierter Krimineller

Algorithmische
Entscheidungssysteme
(ADM Systeme)

Bewertete

Nutzer des
ADM Systems

Daten



Soziales System

Scoring-Verfahren

Klasse 1

oder

Klasse 2

Klasse 3

Qualität von ADM Systemen

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 bedürfen einer sozio-informatischen Gesamtanalyse.**

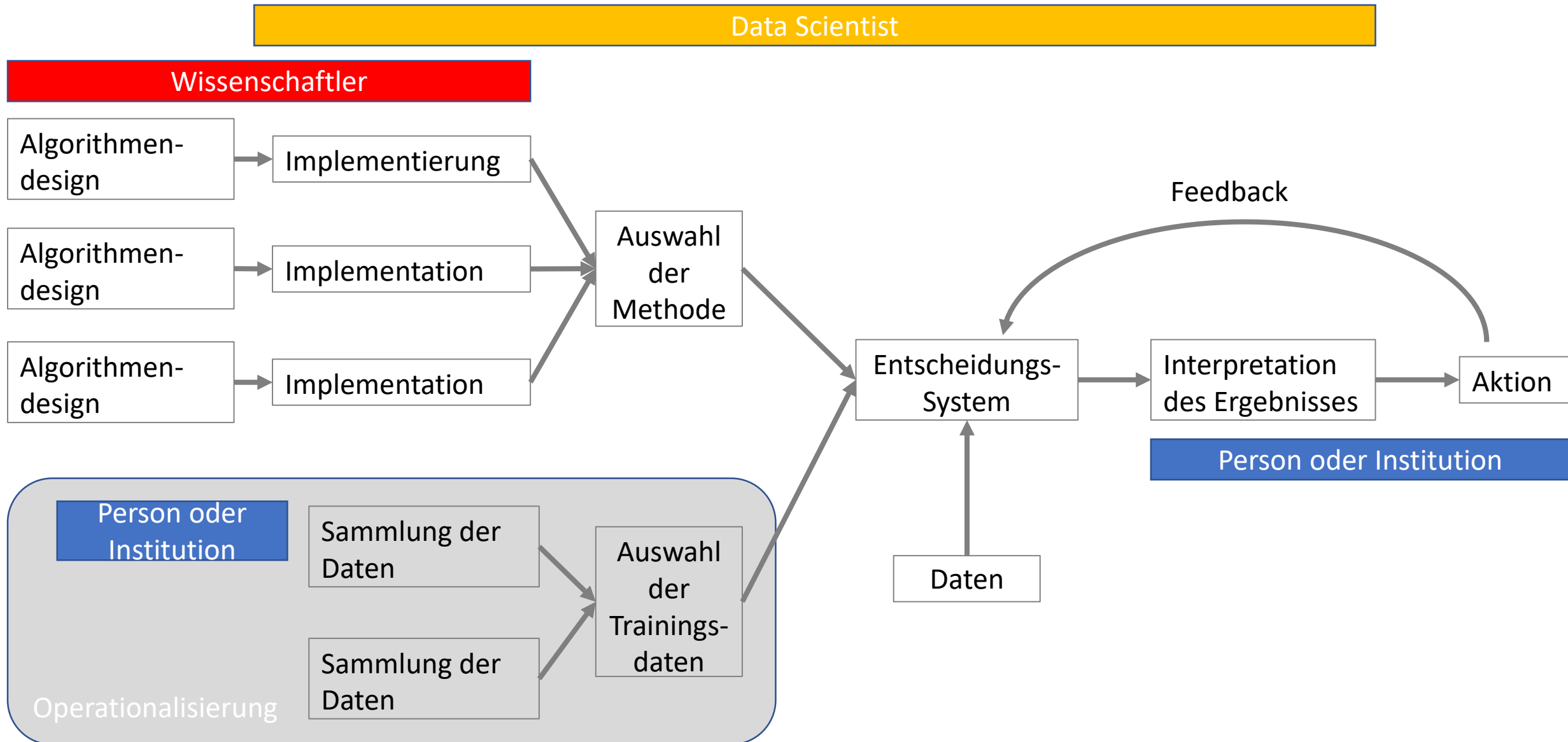



Wie gut sind die Robo-Richter?

- Ganz schön schlecht: COMPAS
 - Hochrisiko-Kategorie:
 - Gewöhnliche Kriminaltaten: nur zu 50% richtig!
 - Schwere Straftaten: nur zu 20% richtig!
- Ein amerikanisches Terroristenidentifikationssystem tönt:
 - „Nur 0.008% falsch Positive!“
 - Bei 55 Millionen Einwohner sind das 4.400 Unschuldige, um wenige Hundert zu identifizieren.
 - Von den „Hochrisikopersonen“ also vermutlich unter 20%!
- Im medizinischen Bereich teilweise besser als Doktoren!



Lange Kette der Verantwortlichkeiten



- 
- Algorithmen können diskriminieren
 - Algorithmen können trotzdem Entscheidungsprozesse verbessern
 - Es gibt Situationen, die keine algorithmischen Entscheidungen erlauben.

Zusammenfassung

Weitere Literatur



- Katharina Zweig: „Auch Maschinen können diskriminieren“, 16.1.2018, MERTON Magazin, <https://merton-magazin.de/auch-algorithmen-koennen-diskriminieren>
- Studie für die Bertelsmann-Stiftung:
Zweig, Fischer & Lischka: „[Wo Maschinen irren können](#)“ (Serie AlgoEthik, No. 4, 2018)
- [Broschüre der Bayerischen Landesmedienanstalt](#) (Zweig, Krafft & Hauer, 2016): „Dein Algorithmus - meine Meinung“