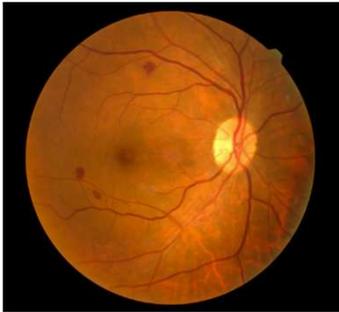


# **The METRIC-Framework: Vertrauenswürdige KI-Daten für die Medizin**

Katinka Becker, Physikalisch-Technische Bundesanstalt (PTB), Berlin



# Einsatz von KI im Medizinbereich



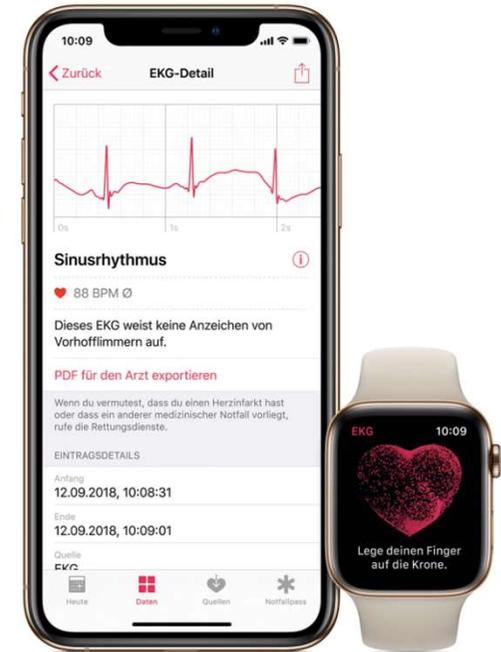
Diabetische Retinopathie

Graphic: aerzteblatt.de



Kardiovaskuläres Ultraschallsystem

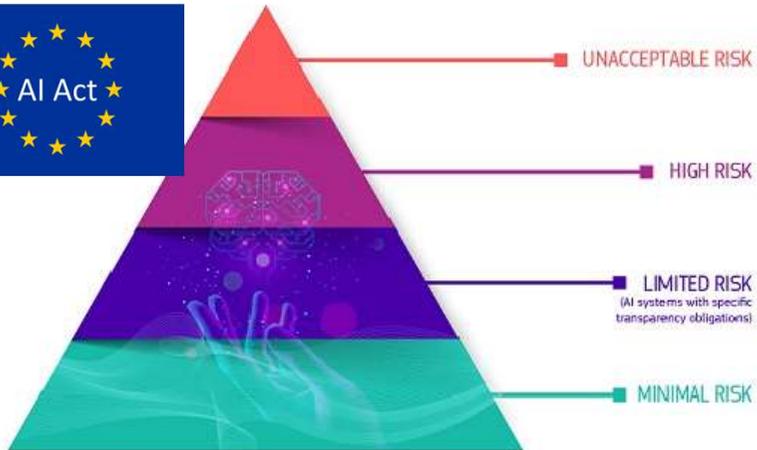
Graphic: philips.de



Apple Watch und Health App

Graphic: apple.com

# Konformitätsbewertung für KI Medizinprodukte

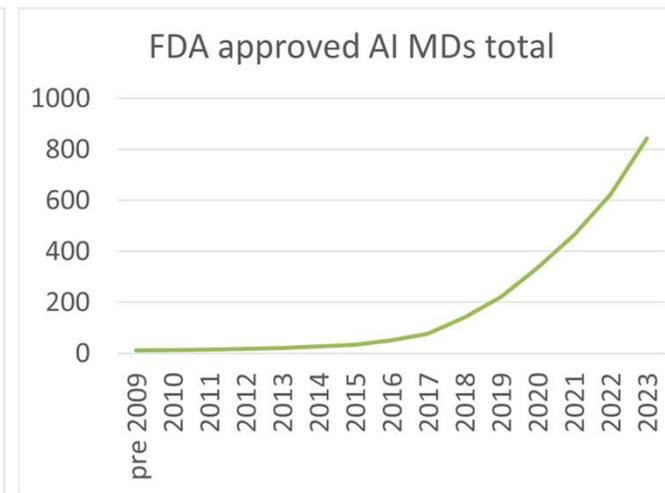
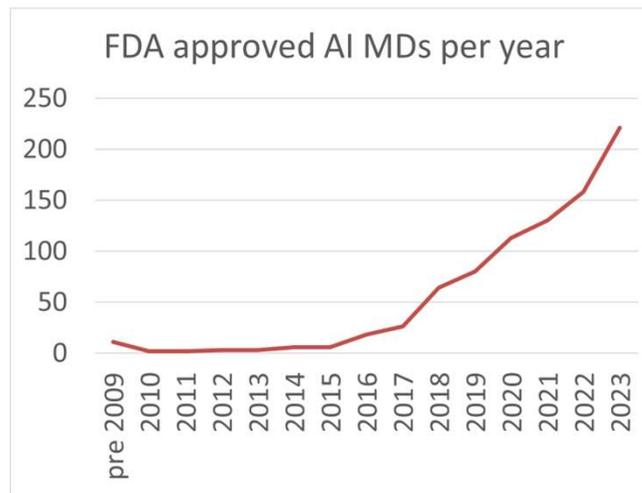


KI Medizinprodukte

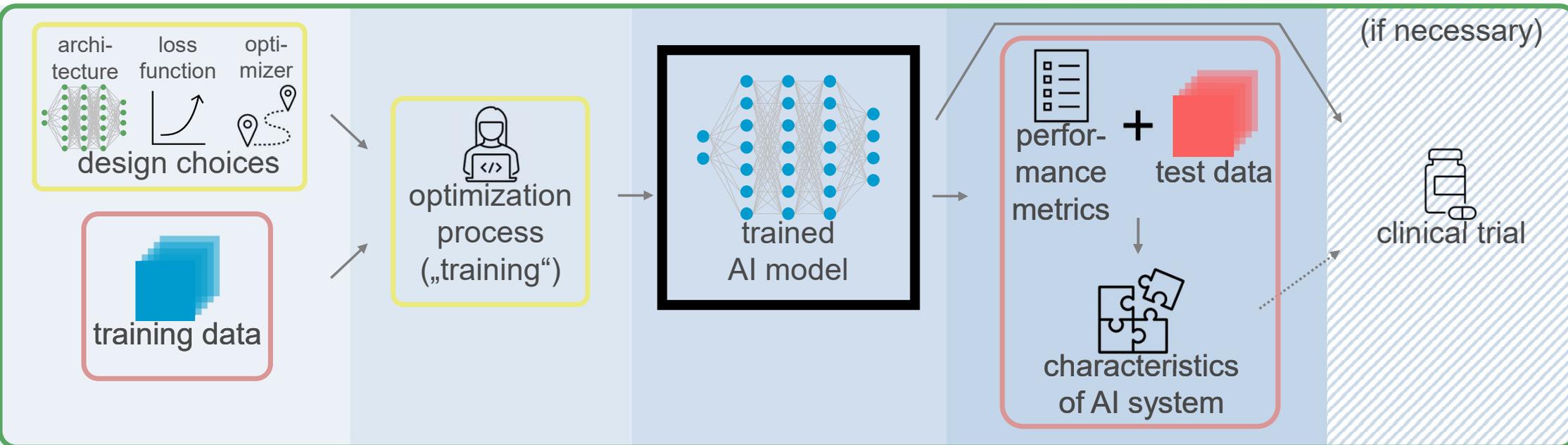


KI Konformitätsbewertung

© Daniel Schwabe



# Quantitatives Testen von KI im Medizinbereich



## AI conformity assessment

*general*

QMS

RMS

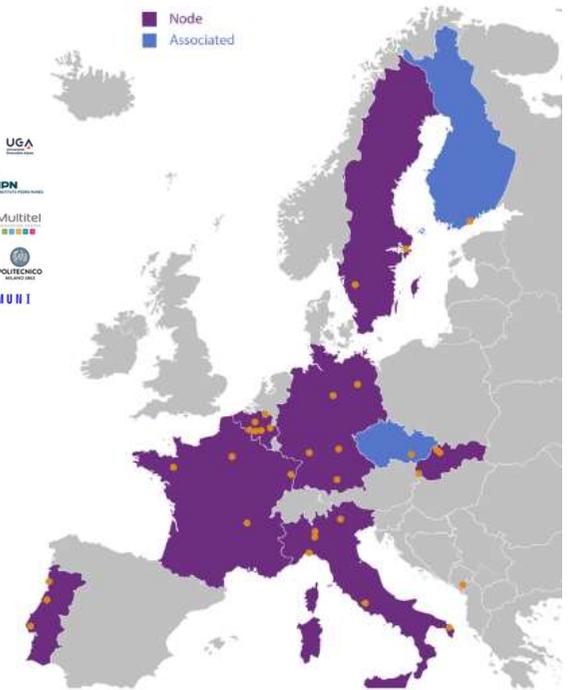
*device specific*

technical documentation

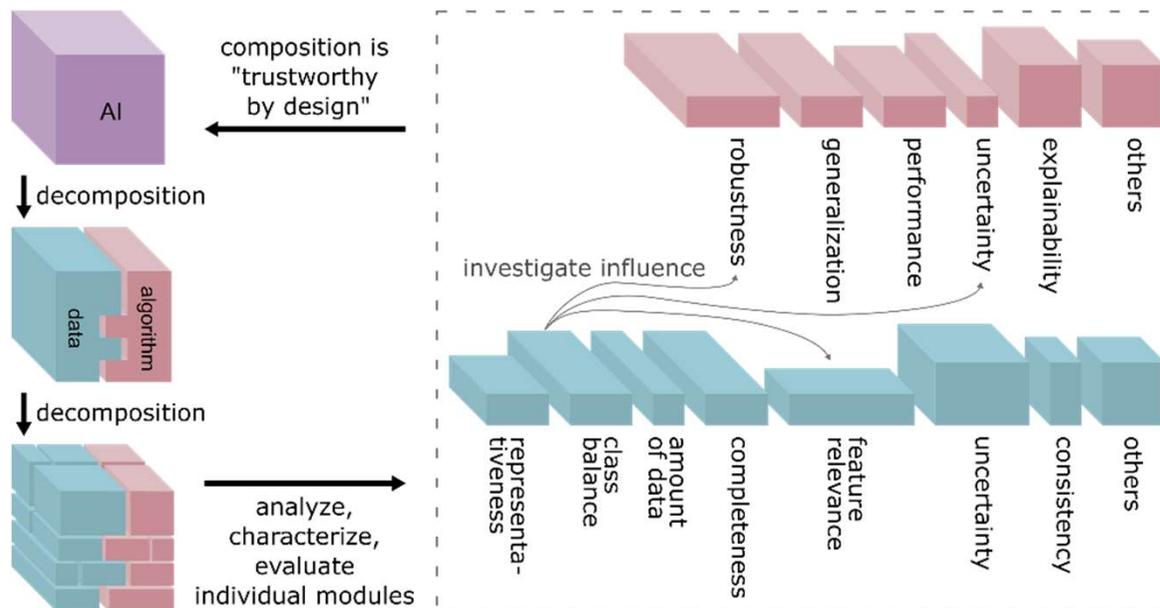
quantitative technical tests

# Das EU Projekt TEF- Health

- „Testing and Experimentation Facility for Health AI and Robotics“
- 52 Partner in der EU
- Leitung: Prof. Petra Ritter, Charité
- PTB leitet das Workpackage „Standards and Quality“



# Vertrauenswürdigkeit per Design



© Daniel Schwabe

- KI System unterteilt sich in zwei entscheidende Komponenten:

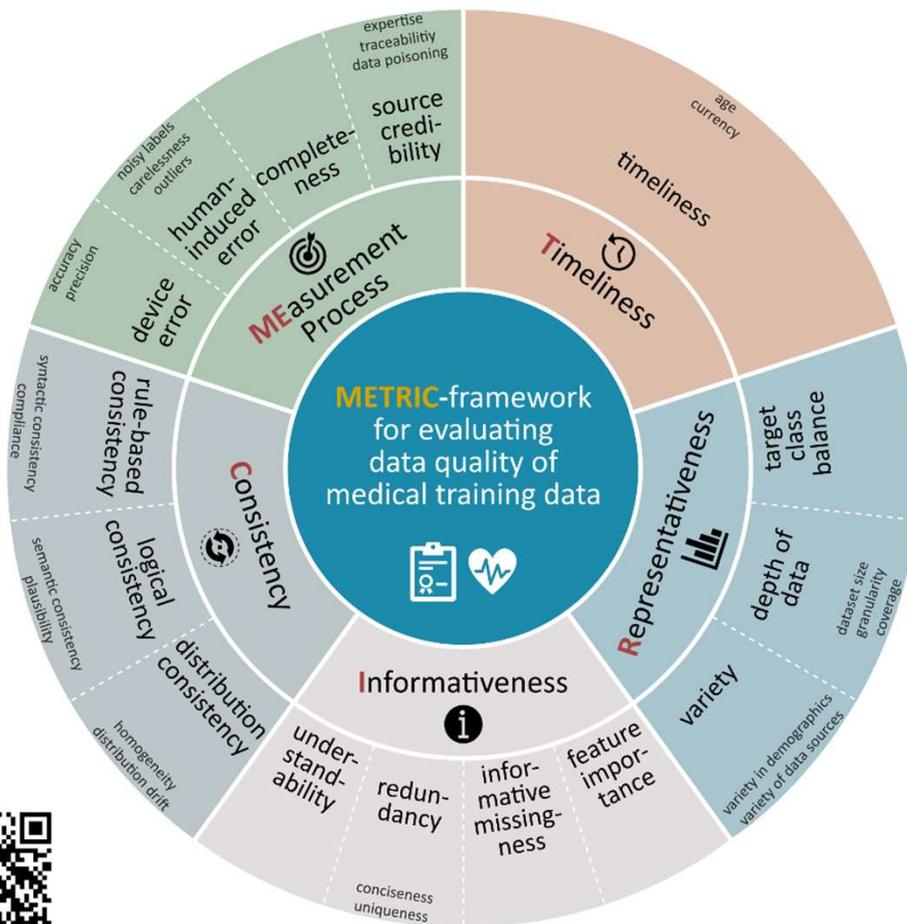
- **Daten und Algorithmus**

Datenqualität

Idee

Qualitätssicherung des Systems durch Qualitätssicherung der Komponenten

# Das METRIC-framework



Systematic Review von ~5400  
Artikeln (120 included)



Das „Rad der Datenqualität“

- 5 Cluster mit 15 „Awareness Dimensions“
- Grundlage für **systematische Evaluierung von Qualität von medizinischen Daten** für z.B.:

- Referenzdatensätze
- Design von Testdatensätzen
- Validierung von Trainingsdatensätzen

© D. Schwabe, K. Becker, M. Seyferth, A. Klaub, T. Schäfer. The METRIC-framework for assessing data quality for trustworthy AI in medicine: a systematic review. *npj Digit. Med.* 7, 203 (2024)

# AI Act Mapping



AI Act Article 10, Paragraph 3, fordert:

„...data sets shall be...“:

- „relevant“
- „sufficiently representative“
- „to the best extent possible, free of errors“

# Von der Theorie in die Praxis

- Quantitative Dimensionen „messbar machen“ mithilfe einer Metrikensammlung
- Website in Planung

### Concordance Correlation Coefficient

Synonyms: Lin's Concordance Correlation Coefficient

$$\rho_c = \frac{2\rho\sigma_x\sigma_y}{\sigma_x^2 + \sigma_y^2 + (\mu_x - \mu_y)^2}$$

with  $\rho$  being the Pearson's correlation coefficient between two variables,  $\mu$  the means and  $\sigma^2$  the variances for the two measurements  $x$  and  $y$ . A statistical metric for assessing the agreement between two methods to evaluate reproducibility and inter-rater reliability.

**Value Range:**  $[-1, 1]$  ↑  
A value of +1 indicates a perfect positive agreement, 0 no agreement and -1 a perfect negative agreement.  
For interpreting different scales, the following references could be used: Altman, 1990; Landis & Koch, 1977.

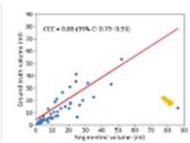
Use in METRIC-framework



Noisy Labels

References  
(Altman, 1990; Akoglu, 2018; Fahrmeir et al., 2016; Lin L., 1989)

**Example**  
Evaluating the consistency between two radiologists in assigning tumor stage labels based on MRI scans. Inter-rater variability may occur when radiologist A assigns stage T2 (encoded as 2.0), while radiologist B assigns stage T3 (encoded as 3.0) for the same lesion. This leads to noisy labels, e.g. important in training data for AI models.



Link to standards and norms

**Relation to other metrics**

- Cohen's Kappa
- Fleiss' Kappa
- Krippendorff's Alpha
- Pearson correlation coefficient

**Applicability**

Data: Images, Time series, Tabular, Text, Multi-modal

Variable: Numerical, Categorical, Ordinal

**Pitfalls and limitations**

- Sensitive to outliers that can significantly influence the results.
- Not interpretable in isolation, needs comparisons.

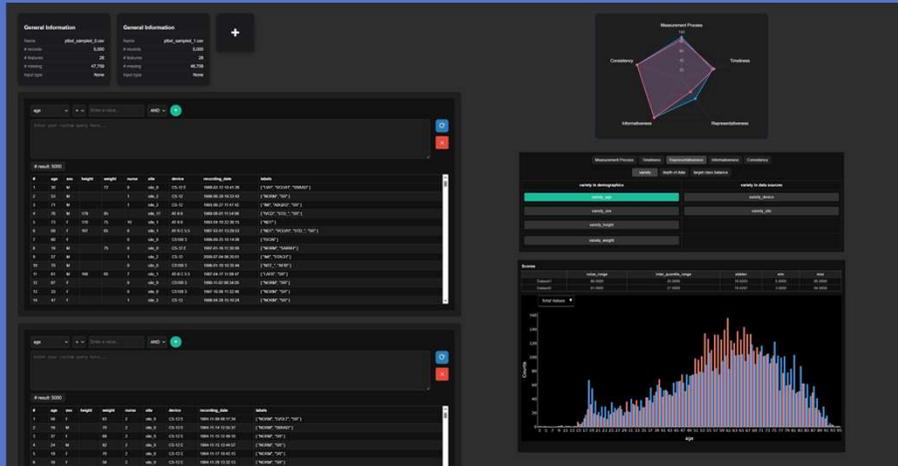
**Prerequisites and recommendations**

- Variables must be continuous and on the same scale.
- The data should be approximately normally distributed.
- Pearson's correlation coefficient should be calculated first.
- Best used for paired data with repeated measurements, e.g. two observers.

- Qualitative Dimensionen „messbar machen“ mithilfe von Fragebögen

Plausibility				
<b>1. Relevance</b>				
1.1 Is the dimension "Plausibility" relevant for your application?	<input type="checkbox"/> Yes [-]	<input type="checkbox"/> No [-]		
2.1 Are units documented that were used for measurements? [QUANTUM, ISO 8000-6:2015]				
2.1	<input type="checkbox"/> Yes [1]	<input type="checkbox"/> No [0]		
2.2 Are acceptable ranges of variables documented?				
2.2	<input type="checkbox"/> Yes [1]	<input type="checkbox"/> No [0]		
2.3 If yes, do all variables contain plausible value ranges (measured with the same units)?				
2.3	<input type="checkbox"/> Yes [1]	<input type="checkbox"/> No [0]		
<b>3. Methods and measures</b>				
<b>3.1 Internal implausibilities</b>				
3.1.1 Was the data checked (by you or someone else) for internal implausibilities/relational contradictions? [Quantum, ISHIEC 25012]	<input type="checkbox"/> Yes, by outlier detection method [1]	<input type="checkbox"/> Yes, by visualizing/plotting [1]	<input type="checkbox"/> Yes, by a person looking over at the data. [0.5]	<input type="checkbox"/> No [0]
3.1.2 If yes, is the data free of internal implausibilities/relational contradictions?				
3.1.2	<input type="checkbox"/> Yes [1]	<input type="checkbox"/> No [0]		
3.1.3 Are date and time format consistent and valid?				
3.1.3	<input type="checkbox"/> Yes [1]	<input type="checkbox"/> No [0]		
3.1.4 Are there any implausible sequences (e.g. discharge before admission)?				
3.1.4	<input type="checkbox"/> Yes [1]	<input type="checkbox"/> No [0]		
<b>3.2 External implausibilities</b>				
3.2.1 Was the data checked (by you or someone else) for implausibilities regarding external knowledge? [Quantum, ISHIEC 25012]	<input type="checkbox"/> Yes, by an expert. [1]	<input type="checkbox"/> Yes, in comparison to relevant literature. [1]	<input type="checkbox"/> Yes, by any other person. [0.5]	<input type="checkbox"/> No [0]
3.2.2 If yes, is the data free of implausibilities regarding external knowledge?				
3.2.2	<input type="checkbox"/> Yes [1]	<input type="checkbox"/> No [0]		
<b>4. Definition-specifics</b>				
<b>4.1 Availability</b>				
4.1.1 Can people involved in measurements and processing of the data be reached to confirm extreme outliers or unusual values?	<input type="checkbox"/> Yes [1]	<input type="checkbox"/> No [0]		
<b>4.2 Human expertise</b>				
4.2.1 Did humans interacting with the data have high expertise? [SPIRIT-AI]	<input type="checkbox"/> Yes, they are experts. [1]	<input type="checkbox"/> They had some general expertise but not specifically qualified. [0.5]	<input type="checkbox"/> No expertise at all. [0]	<input type="checkbox"/> Don't know [0]
<b>4.3 Statistical plausibility</b>				
4.3.1 Do all frequencies match known epidemiological patterns? (E.g. the prevalence of diabetes, hypertension, etc., within expected national or regional ranges)	<input type="checkbox"/> Yes [1]	<input type="checkbox"/> No [0]	<input type="checkbox"/> Don't know [0]	
<b>4.4 Labels/Annotation</b>				
4.4.1 Were the labels/annotations assigned using the <b>medical gold standard</b> ? (E.g. diagnostic test)	<input type="checkbox"/> Yes [1]	<input type="checkbox"/> No [0]	<input type="checkbox"/> Don't know [0]	
4.4.2 Were the labels/annotations <b>derived by humans</b> ?	<input type="checkbox"/> Yes [-]	<input type="checkbox"/> No [-]	<input type="checkbox"/> Don't know [-]	
4.4.3 Was the data <b>labeled by multiple raters</b> ?	<input type="checkbox"/> Yes [1]	<input type="checkbox"/> No [0]		
<b>4.5 Metadata</b>				
4.5.1 Does the dataset contain variables with <b>interdependencies</b> ?	<input type="checkbox"/> Yes [-]	<input type="checkbox"/> No [-]	<input type="checkbox"/> Don't know [-]	
4.5.2 Is all <b>metadata correctly linked to other datatypes</b> ? (E.g., are	<input type="checkbox"/> Yes [1]	<input type="checkbox"/> No [0]	<input type="checkbox"/> Don't know [0]	
<b>4.6 Noise/error</b>				
4.6.1 Does the dataset contain any <b>sources of errors or noise</b> ? [Fehr et al. 2022, Datasheets, FactSheets]	<input type="checkbox"/> Yes [0]	<input type="checkbox"/> No [1]		
Plausibility score:				

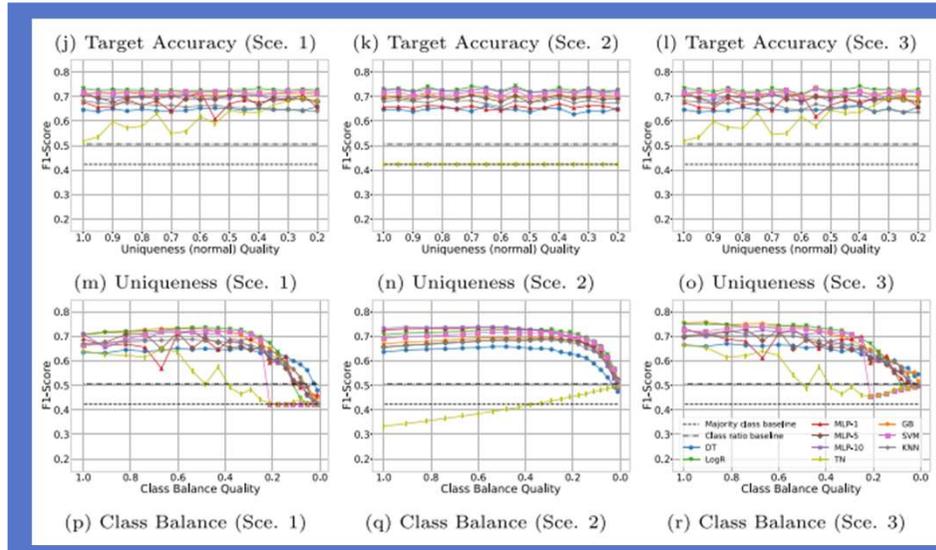
# Von der Theorie in die Praxis



- Tool:
- Evaluierung und Vergleich von Datensätzen
  - Perspektivisch: Kuratieren von Datensätzen

© Martin Seyferth

© S.Mohammed et al. The effects of data quality on machine learning performance on tabular data. *Information Systems*, 132 (2025)



- Wie stark wirken sich **einzelne Datenqualitätsdimensionen** auf die Qualität des Systems aus (Performance, Robustheit, Fairness, Erklärbarkeit,...) ?

# Danke

