

# Development of a German Instrument for Self-Perceived Data Literacy

An Algorithm-based Approach to Scale Development

Leonie Hagitte

2024-03-14

## Contents

<b>Abstract</b>	<b>3</b>
<b>Acknowledgements</b>	<b>3</b>
<b>1 Intro</b>	<b>5</b>
<b>2 Background</b>	<b>5</b>
<b>3 Methods</b>	<b>6</b>
3.1 Sample . . . . .	6
<b>4 Analysis</b>	<b>7</b>
<b>5 Results</b>	<b>7</b>
<b>6 Discussion</b>	<b>7</b>
<b>References</b>	<b>7</b>



## **Abstract**

The increasing relevance of competent and critical handling of data in society not only makes it possible to record this competence, but also makes self-perception with regard to this competence increasingly clear. Previous approaches consider this competence primarily against the specific background of individual target groups, jobs or roles (Cui et al., 2023). In addition, only a few explicitly refer to the general population (Carmi et al., 2020; Cui et al., 2023). In view of the various theoretical approaches, there is a need for a uniform definition of data literacy in order to create comparability. Our aim is therefore to derive a holistic definition based on these approaches and to develop a questionnaire for self-perception of one's own data literacy. To this end, the decisive factors for the construct from previous definitions and operationalizations in various disciplines are brought together. Cognitive interviews are conducted iteratively to create and refine the items. The items are then selected using algorithm-based item selection. The facets of data literacy are comprehensively tested for factorial, discriminant, convergent and congruent incremental validity in order to promote a differentiated understanding of the construct. Construct and criterion validity are tested using correlations and hierarchical regression analyses, while cross-validation checks the robustness of the instrument. Based on a cross-sectional online questionnaire study, we first examine a representative sample of people from the general population. Limitations arise from the cross-sectional design and the heuristic item reduction, which limit predictions of predictive validity. The heterogeneous nature of the construct makes global instrument development and understanding of all participants difficult. The self-assessment questionnaire promotes a holistic assessment of competence and its perception for further research, for example by comparing self-assessment and actual performance.

## **Acknowledgements**

I dedicate this thesis to

I want to thank my advisers, Prof. Martin Schultze, Prof. Timo Lorenz, and Prof. Manuel Völkle for their time and patience, and my friends for their resourceful advice:



# 1 Intro

## 2 Background

The relevance of data literacy in today's society becomes evident as it serves as a potent tool in navigating the complex data-driven environment. In a world characterized by information overload and rapid technological advancements, individuals equipped with strong data literacy skills can discern patterns, critically evaluate information, and make informed decisions

The exploration of citizens' interaction with media and the cultivation of their agency has traditionally centered around concepts such as written literacy, media literacy, information literacy, and digital literacy. In more recent discussions, Data Literacy has been approaching relevance among discussed competencies regarding what is necessary for agency in the current society (carmi2020?). Deficiency in data literacy not only exposes individuals to various risks and harms on personal, social, physical, and financial levels but also constrains their capacity to actively engage as informed citizens within an evolving, data-driven society {(carmi2020?)}. Thus, Data Literacy is a competency that is becoming increasingly important to everyone. And research has acknowledged this in recent years, as more and more research is being done in that direction (Cui et al., 2023); And in praxis there are Training programmes and Workshops sprouting to enhance ones Data Literacy as well (QUELLE).

In "Thinking, Fast and Slow," Kahneman (2011) introduces the two systems of thinking: System 1 and System 2. System 1 operates swiftly, relying on intuition and often succumbing to biases, while System 2 operates more deliberately, engaging in analytical thinking. Kahneman's exploration extends to the realization that even experts across diverse domains can fall prey to cognitive biases, leading to errors in judgment (Kahneman, 2011). Applying Kahneman's insights to the realm of data literacy and research questions sheds light on a critical aspect. It underscores the notion that individuals, irrespective of their statistical expertise, remain vulnerable to cognitive biases when interpreting and analyzing data. In the current societal landscape, where data plays an increasingly pivotal role in shaping decisions and policies, acknowledging and addressing these cognitive biases becomes paramount.

We noticed several gaps in the current research that we want to try addressing with this study. For once the target groups for Data Literacy are currently each having their own definition of the construct it seems (Cui et al., 2023). This not only makes comparisons impossible but also makes a general understanding of

the topic as well as communication in the community and science communication to the public harder.

data vs. information

## 3 Methods

### 3.1 Sample

The sample for this study comprised XXX participants ( $M=$ ,  $SD=$ ). Within the sample, XXX% identified as female, XXX% as male, and xxx% did not identify with binary gender categories. All participants were aged 18 and above. Regarding education, all participants exhibited a [insert educational level- specifying the range or types of educational levels observed in the sample]. Among the participants,  $n=$  reported higher knowledge on items x, x, x, leading to their selection for an additional set of items as a preliminary survey for factors four and five. The study encompassed every sector within the occupational classification (Bundesagentur für Arbeit, 2020), ensuring comprehensive representation. Conducted in German, the participation in the study was entirely voluntary, with no external incentives provided. The recruitment of participants was carried out through a combination of personal and professional networks, along with outreach on various online social media platforms.

Our study sample serves as a focal point for comparison against the demographic landscape of the general public in Germany. In 2022, the mean age of the German population was 44.6 years, with 45,457,000 individuals engaged in employment. Educational backgrounds varied (XXX), and for gender distribution, the split was nearly 50/50 (41,616,473 males and 42,816,197 females) according to the Statistisches Bundesamt (source: <https://www.destatis.de/DE/Themen/Gesellschaft-Umwelt/Bevoelkerung/Bevoelkerungsstand/Tabellen/liste-zensus-geschlecht-staatsangehoerigkeit.html#651186>).

Analyzing our sample against these benchmarks provides a comprehensive understanding of any distinctions or parallels in age, employment, education, and gender. This comparison enhances the applicability of our findings to the broader German population.

## 4 Analysis

We employed an automated item selection algorithm to craft the [insert name] scale. The process of scale development, involving the strategic selection of items to ensure psychometric soundness, is conceptualized as a combinatorial problem (Kerber et al., 2022). Combinatorial problems, exemplified by the knapsack problem (Schroeders et al., 2016), entail identifying a discrete and finite solution within predefined constraints (Hoos and Stützle, 2005).

Contemporary approaches to address these combinatorial problems leverage automatic optimization algorithms, such as Genetic Algorithms (GA; Holland, 1975), Ant Colonization Algorithms (ACO; XXX), brute force (XXX), or random sampling(XXX). (Schultze, 2022). Unlike classical approaches that consider items based on their individual merits, heuristic item selection algorithms aim to enhance the psychometric properties of a set of items within predetermined constraints (Schultze, 2017). Noteworthy is the inherent approximate, rather than deterministic, nature of metaheuristics (Schultze & Lorenz ,2023; Blum and Roli, 2003). Wich makes brute force approaches the preferred choice, if applicable (Schultze & Lorenz ,2023). Nevertheless, as brute force often isnt fesable because of timely and computational costs, approximate algorithms are indispensable for obtaining near-optimal solutions to complex combinatorial problems in a timely or computationally efficient manner (Schultze & Lorenz ,2023; Dorigo and Stützle, 2010).

## 5 Results

## 6 Discussion

## References

- Carmi, E., Yates, S. J., Lockley, E., & Pawluczuk, A. (2020). Data citizenship: Rethinking data literacy in the age of disinformation, misinformation, and malinformation. *Internet Policy Review*, 9(2). <https://doi.org/10.14763/2020.2.1481>
- Cui, Y., Chen, F., Lutsyk, A., Leighton, J., & Cutumisu, M. (2023). Data literacy assessments: A systematic literature review. *Assessment in Education: Principles, Policy & Practice*, 30, 1–21. <https://doi.org/10.1080/0969594X.2023.2182737>
- Kahneman, D. (2011). *Thinking, fast and slow*. Farrar, Straus; Giroux.