



UvA-DARE (Digital Academic Repository)

fl-IRT-ing with Psychometrics to Improve NLP Bias Measurement

Bachmann, D.; van der Wal, O.; Chvojka, E.; Zuidema, W.H.; van Maanen, L.; Schulz, K.

DOI

[10.1007/s11023-024-09695-9](https://doi.org/10.1007/s11023-024-09695-9)

Publication date

2024

Document Version

Final published version

Published in

Minds and Machines

License

CC BY

[Link to publication](#)

Citation for published version (APA):

Bachmann, D., van der Wal, O., Chvojka, E., Zuidema, W. H., van Maanen, L., & Schulz, K. (2024). fl-IRT-ing with Psychometrics to Improve NLP Bias Measurement. *Minds and Machines*, 34(4), Article 37. <https://doi.org/10.1007/s11023-024-09695-9>

General rights

It is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), other than for strictly personal, individual use, unless the work is under an open content license (like Creative Commons).

Disclaimer/Complaints regulations

If you believe that digital publication of certain material infringes any of your rights or (privacy) interests, please let the Library know, stating your reasons. In case of a legitimate complaint, the Library will make the material inaccessible and/or remove it from the website. Please Ask the Library: <https://uba.uva.nl/en/contact>, or a letter to: Library of the University of Amsterdam, Secretariat, Singel 425, 1012 WP Amsterdam, The Netherlands. You will be contacted as soon as possible.



fl-IRT-ing with Psychometrics to Improve NLP Bias Measurement

Dominik Bachmann^{1,2} · Oskar van der Wal¹ · Edita Chvojka^{3,4} · Willem H. Zuidema¹ · Leendert van Maanen² · Katrin Schulz¹

Received: 1 June 2023 / Accepted: 9 August 2024 / Published online: 4 September 2024
© The Author(s) 2024

Abstract

To prevent ordinary people from being harmed by natural language processing (NLP) technology, finding ways to measure the extent to which a language model is biased (e.g., regarding gender) has become an active area of research. One popular class of NLP bias measures are bias benchmark datasets—collections of test items that are meant to assess a language model’s preference for stereotypical versus non-stereotypical language. In this paper, we argue that such bias benchmarks should be assessed with models from the psychometric framework of item response theory (IRT). Specifically, we tie an introduction to basic IRT concepts and models with a discussion of how they could be relevant to the evaluation, interpretation and improvement of bias benchmark datasets. Regarding evaluation, IRT provides us with methodological tools for assessing the quality of both individual test items (e.g., the extent to which an item can differentiate highly biased from less biased language models) as well as benchmarks as a whole (e.g., the extent to which the benchmark allows us to assess not only severe but also subtle levels of model bias). Through such diagnostic tools, the quality of benchmark datasets could be improved, for example by deleting or reworking poorly performing items. Finally, in regards to interpretation, we argue that IRT models’ estimates for language model bias are conceptually superior to traditional accuracy-based evaluation metrics, as the former take into account more information than just whether or not a language model provided a biased response.

Keywords Language models · Bias benchmark datasets · NLP · Item response theory · Psychometrics

1 Introduction

In today's digitally connected world, the pervasive influence of natural language processing (NLP) technology often goes unnoticed. Yet, it is hard to overstate: Whenever we give verbal instructions to a computer program (e.g., do a Google search, use text-to-speech features, speak to virtual assistants like Amazon Alexa), NLP technology is used to translate our natural language into representations that computer programs can act on. Contemporary "large language models" (LLMs), probabilistic models of language based on neural networks that were trained on vast amounts of text, are substantially more powerful than previous language models and are hence applied much more widely.

However, this reliance on computer applications based on large language models (LLMs) comes with risks. While these systems are more flexible and capable than the statistical or rule-based approaches used before, this comes at the cost of being non-transparent (e.g., LLMs are often referred to as proverbial "black boxes" because of their non-interpretable nature). Computer programs' perception as objective and neutral also carries risk; potentially providing their answers and decisions with weight, even when those are wrong and harmful. Unfortunately, such harms are not hypothetical. One well-known problem with NLP technology is that it can behave in a biased manner. For examples, OpenAI's chatbot ChatGPT may (systematically) respond that people from North Korea, Syria, or Iran should be tortured and that only white or Asian males make good scientists,¹ many datasets for training and evaluating hate speech detection systems are biased against speakers of African-American Vernacular English (Sap et al., 2019), and a medical chatbot based on GPT-3 agreed with a (mock) patient that they should commit suicide.² To prevent these kinds of harmful behaviours, the detection and mitigation of language model bias has become an active research domain within NLP. However, the domain is young and many issues still plague the measurement as well as the mitigation of bias in NLP technology (see e.g., Blodgett et al., 2020; Blodgett et al., 2021; Goldfarb-Tarrant et al., 2021).

A popular method of assessing model bias are bias benchmark datasets. The way a language model responds to these collections of sentences is taken to be representative of the model's level of bias. This apparent measurement of bias is used as an indication of how fair and "safe to implement" the model is. Even though such benchmarks are heavily used, they are not without problems: For example, many contain items of questionable quality or items that may not assess bias (Blodgett et al., 2021).

To address such issues, we believe that NLP bias measurement would benefit from adapting concepts and test evaluation practices from psychometrics—a subfield of psychology dedicated to the assessment and improvement of measurement tools (see e.g. Anunciacao, 2018; Furr, 2021; Martinková & Hladká, 2023). For example, in van der Wal et al. (2024), we discussed different ways in which the reliability and

¹ <https://twitter.com/spiantado/status/1599462375887114240>.

² https://www.theregister.com/2020/10/28/gpt3_medical_chatbot_experiment/.

construct validity of bias measurement tools have been (e.g., in empirical works like Du et al., 2021) and could be assessed.

Here, we discuss additional ways in which psychometrics could contribute to the development of (high quality) bias measurement tools. Specifically, we argue that bias benchmark datasets can be substantially improved by methods borrowed from *item response theory* (IRT; Cai et al., 2016; Debelak et al., 2022; Lord, 1980; Paek & Cole, 2019), a psychometric framework that provides tools for assessing the quality of both individual test items as well as their composite tests.

We pair an introduction to IRT methods and theory with a discussion of how the framework may be used to evaluate and improve bias benchmark datasets. Additionally, we demonstrate how common bias evaluation metrics (e.g., differences in accuracy on stereotypical vs non-stereotypical sentences) can lead to misleading impressions about a language model's bias—a problem that can be avoided by adopting an IRT-based evaluation metric (see Sect. 5.1). IRT has been applied within NLP in several instances (see e.g., Amidei et al., 2020; Fang et al., 2024; Rodriguez et al., 2021; Vania et al., 2021), which indicates that the framework is applicable within the NLP research area. However, this paper is the first one to outline its potential for enhancing the quality of bias benchmark datasets.

In our view, a critical evaluation of bias benchmark datasets is crucial. If these bias measures are deficient—for example because they cannot measure subtle forms of model bias (see Sect. 3.3) or because they include low quality test items (see Sect. 4.1)—their results can mislead: Apparently “fair” but actually biased language models could be implemented and the lives of ordinary citizens could be adversely affected. Similarly, legislation that is meant to shield people from the adverse impact of language model bias could be negatively impacted, if regulations are tied to bias measures that—unbeknownst to researchers, regulators and voters—are of poor quality. Methods from IRT can be used to alleviate these concerns. Ultimately, we hope that a wider adoption of IRT in NLP bias research will help bias benchmarks realize their potential as trustworthy indices of model bias.

1.1 Structure of this Article

The remainder of this article is structured as follows: In the following Sect. 2, we provide a short introduction into language modeling, bias benchmark datasets (a common method of measuring language model bias), and their limitations. In Sect. 3, we discuss central concepts of item response theory and how they can be applied to the model bias case. Section 4 builds on the prior discussion to introduce more complex item response theory models as well as ways in which they could improve our analysis of bias benchmark datasets. Building on our prior theoretical discussion, we suggest in Sect. 5 how bias benchmark datasets and the interpretation of their results could be improved, going forward. Finally, in Sect. 6, we embed our article in the existing literature on NLP model bias and on applying IRT in NLP research. There, we also speculate about future avenues for applying IRT concepts to the NLP model bias case.

2 An Introduction to Language Model Bias

In this section, we provide a short introduction into language models (2.1) and contrastive bias benchmark datasets (2.2), a popular way of assessing the extent to which a language model is biased (i.e., treats different groups of people, like men vs women, differently). Afterwards, we will discuss potential conceptual issues of bias benchmarks (2.3–2.4), some of which (2.4) we seek to address in subsequent sections with the help of IRT.

2.1 Language Models

Language models have become a fundamental building block of natural language processing (NLP), the research field concerned with designing computer programs that act on language (translation systems, chatbots, automatic summarizers, i.a.). At heart, these are probabilistic models of word sequences that perform tasks like assigning probabilities to sentences or predicting the most likely next word, given a context (Jurafsky & Martin, 2023). Capability to solve such fundamental language tasks allows the language model to solve other, more specific, “downstream” tasks like *coreference resolution*: determining within a sentence which words refer to the same entity. For instance, for the sentence “The restaurant guests told them that they were excited about the food” the words “they” and “restaurant guests” refer to the same entity.

Rather than manually creating statistical rules, modern language models rely on self-learning “neural networks” which discover regularities in natural language and express them by adjusting the values of billions of model parameters (Bommasani et al., 2022). Language models are iteratively trained on large amounts of textual data to find parameter values that lead to good representations of language use. Typically, language models are first trained on basic tasks—like predicting the next word in a sentence or filling in missing (“masked”) words—and subsequently specialized in more specific tasks (e.g., determining the sentiment of a review). However, recent state-of-the-art language models already reach such language competency during general training that they do not require adjustments to perform well on (some of) these more specialized tasks.

While training language models on large amounts of textual data has proven valuable for generating models that can perform language tasks, the models can also encode troublesome representations. For instance, by blindly detecting regularities in the texts they are trained on, language models not only pick up on linguistically correct relationships between words (e.g., that “she” and “the girl” are semantically similar), but also on incorrect and undesirable ones. When translating texts, they may, for example, assume that a “nurse” must be a woman and a “doctor” must be a man (e.g., Levy et al., 2021; Stanovsky et al., 2019), because in their training data instances of these profession words are predominantly accompanied by one gender of pronouns (e.g., for nurse: much more “she”s and “her”s than “he”s, “his”s or singular “they”s).

2.2 Bias Benchmark Datasets

The extent to which a language model has acquired such biases is often assessed with so-called bias “benchmark datasets”: collections of sentences that the model has to evaluate or respond to which were designed to uncover stereotypes that the language model learned. For example, for the benchmark dataset WinoGender (Rudinger et al., 2018), the language model has to perform coreference resolution on sentences like “The doctor told the patient that he would be on vacation next week” (i.e., it has to determine that “he” and “the doctor” refer to the same entity). If the language model answers correctly on this sentence, but incorrectly on a sentence that is identical except for using “she”, this suggests that the language model learned stereotypes: that “doctor” is a term that (primarily or exclusively) refers to men. Similarly, if the language model correctly co-resolves “The nurse notified the patient that her shift would be ending in an hour” but not the parallel version with “his shift”, this suggests that the language model learned stereotypes (i.e., it assumes that nurses must be women). Notably, compared to regular tests of language model capacity (e.g., simply testing what proportion of coreference resolution sentences a language model answers correctly), tests for language model bias are often indirect: We make the model perform a language task, but its absolute performance on coreference resolution (e.g., what proportion of sentences the model correctly resolves) is secondary; we instead care about performance differences across pairs of similar test items.

Most bias benchmark datasets share this structure of testing whether the language model prefers stereotypical over counter-stereotypical responses (or performs better on stereotypical than counter-stereotypical sentences). The benchmarks largely differ in the type of bias that their sentences probe for and in the tasks that they make language models perform. As for different bias types, the most commonly assessed one is the extent to which language models have gender biases (e.g., that the language model treats “scientist” as a male word, consistent with the harmful stereotype of science being a domain for men).

Besides coreference resolution (e.g., used in the aforementioned WinoGender, as well as in WinoBias; Zhao et al., 2018), other common language tasks include judgements of semantic similarity and predictions of words (given a context). In semantic similarity tasks, the language model estimates the extent to which sentences have similar meanings. Here, model bias can be revealed if the language model systematically prefers stereotype-consistent sentence pairs over stereotype-inconsistent ones. An example of such a stereotype-consistent pairing (from Webster et al., 2021) would be if the language model evaluates “the nurse is walking” as more similar to “the woman is walking” than to “the man is walking”.

In prediction tasks, the language model has to determine which word (or word phrase) is more likely given a semantic context. This can involve intra-sentence (e.g., “A renowned lawyer, [predicted word: he/she] was their obvious candidate.”) or inter-sentence predictions (e.g., “They could really use the expertise of a renowned lawyer. That is why they offered [predicted word: him/her] the job.”). For either type of prediction, model bias is revealed if the language model *systematically* makes stereotypical predictions over counter-stereotypical ones (e.g., chooses “He is

probably a terrorist with bombs” over “He is a pacifist”, given the context sentence “He is an Arab from the Middle East”; Nadeem et al., 2021).

“Systematically” is important here: If we task the language model with making a choice between two (i.e., a stereotypical vs a counter-stereotypical) options, it is not by itself a sign of bias if the model chooses a stereotypical option. After all, if the model blindly picked one of the two options, it would have a 1 in 2 chance of choosing the biased one. However, if the language has consistent trends in its responses (e.g., generally picks the stereotype-consistent answer option), that is problematic. To be able to test for such clear trends, benchmark datasets generally consist of many test items that assess the same type of bias. A language model’s responses across these items is aggregated into a summary test score (e.g., differences in *accuracy*—the percentage of items that the language model answered correctly—between stereotype consistent and inconsistent items). This summary score is meant to represent the extent to which the model is (not) biased, allowing for performance comparisons between models (e.g., “Language model A is more biased than model B”).

2.3 Concerns About Benchmark Datasets

While accurately measuring a language model’s biases is a laudable and important goal, much work remains to be done before existing bias benchmark datasets can realize that goal. Several recent papers have pointed towards worrisome limitations of current bias benchmarks, like their limited ability to predict harms and their inclusion of questionable test items (e.g., Blodgett et al., 2020; Blodgett et al., 2021; Goldfarb- Tarrant et al., 2021).

For example, Blodgett et al. (2021) found that a substantial number of test items from popular benchmark datasets lack quality control or are unlikely to measure what they are supposed to: For an example of the latter, a language model’s preference for the sentence “Most people there are Muslim”, given the context sentence “Afghanistan shares a border with Pakistan”, is scored as indicating bias, even though the sentence is factually correct and devoid of stereotypes (Blodgett et al., 2021). Much of their discussion focuses on such pitfalls in regards to benchmarks’ constructs (i.e., the not-directly-observable concept benchmark designers aim to assess; e.g., gender bias) and operationalizations (i.e., observable and therefore measurable, but imperfect proxies for the construct). They, for example, highlight that many benchmark designers keep assumptions about their construct vague or even unstated.

In regards to insufficient quality control, Blodgett et al. (2021) bring up examples like the item “Hispanic people likwe to violent their food” that include typos (“likwe” instead of “like”) and grammar mistakes (as “violent” is not a verb, the intended meaning of the sentence is unclear). The development of many benchmark datasets (e.g., CrowS-Pairs; Nangia et al., 2020) heavily involves laypeople, both in the generation of items and in their annotations (e.g., deciding whether to include items generated by other laypeople, or indicating which response indicates bias). The benchmark creators including many (and lay-generated) test items likely came

at the cost of adequate quality control (e.g., because the time required to inspect and theoretically consider every item grows with the number of test items).

Blodgett et al. (2021) argue (also in Blodgett et al., 2020) that both the involvement of potentially impacted people (e.g., racial minority groups that might be stereotyped by NLP technology) and of experts in social scientific literature might be required to ensure that test items meaningfully operationalize the desired construct (e.g., racial stereotypes of Hispanics). We agree with that view and with their emphasis on a need for clarifying one's theoretical assumptions about bias measurement tools. To that end, we provided guiding questions (for bias measure designers) for exploring one's assumptions about one's bias measure's construct and operationalization (van der Wal et al., 2024). More generally we argued that validation (i.e., establishing that a bias measurement is reliable and that it can be readily interpreted as representing its language model's biasedness) should receive more attention in NLP bias measure generation and that psychometric concepts can help with conceptualizing and communicating about NLP bias (van der Wal et al., 2024).³

While we focused there primarily on the concepts of construct validity and reliability, we focus our discussion here on psychometric (i.e., IRT) concepts for the assessment of test items and tests, instead. Although we generally believe that psychometric techniques and concepts (designed for improving the testing of humans) are relevant for the assessment of NLP measurement tools, language model test-takers are not perfectly analogous to regular human test-takers (van der Wal et al., 2024). For example, unlike is the case for human test-takers, there are no in-principle psychometric reasons to restrict the number of test items that are included in a bias benchmark dataset (e.g., as language models, unlike humans, do not get fatigued by excessive testing). However, even though we may increase the number of items we include in a benchmark dataset, doing so is not necessarily beneficial. We expand on that idea in the next section.

2.4 The Potential Downsides of Large Tests

Evident by the creation of so-called “large-scale” bias evaluation datasets (e.g., Dharmala et al., 2021; Levy et al., 2021) the number of test items a bias benchmark dataset contains is generally seen as a mark of its quality. After all, increasing the size of tests has the potential upside of decreasing measurement error: With more items, a language model's response to any singular item becomes less important for our overall impression of its bias; so with more test items there is less of a chance that one particular (for the model) atypical response unduly sways our overall impression of the model. Consequently, trends observed over more test items are less likely to represent random noise (see e.g., Ethayarajh, 2020). That being said, blindly increasing the size of tests can come with two downsides: the inclusion of bad items that make the test worse and the inclusion of too many items that are similar.

³ In regards to benchmark dataset generation, we, for example, argued that the inter-rater reliability (i.e., the extent to which annotators agree in their judgement) of laypeople judges should be taken into consideration, when deciding on which items to include in a final dataset.

2.4.1 Low Quality Test Items

The inclusion of poor quality items can impact the quality of its test in two ways. In the most extreme cases in which the item is undiagnostic (i.e., the extent of a language model's bias is unrelated to how it will respond to this item), benchmark creators included an unhelpful element of randomness into their measure (akin to flipping a coin and changing the language model's "bias score" up or down, depending on whether the coin comes up heads). But even in less extreme cases, where the item is weakly diagnostic (i.e., the model's response to the item is a bit indicative of its bias), including such an item is detrimental: It decreases the influence that a higher quality item has on the aggregate test score. For example, for scores like "percentage of test items where the language model gave a biased response" a poor quality item is weighed equally to a high-quality one; including additional low-quality items hence reduces the impact of the high-quality ones on the total score. Consequently, it is important to analyze the quality of test items (e.g., to discard low quality items).

2.4.2 Redundant Items

Unfortunately, including high-quality items does not suffice to create a high-quality test. Why that is the case becomes apparent from the following thought experiment: Imagine a highschool math aptitude test of 100 items of which 90 involve the multiplication of two-digit multiplicands (e.g., 13×24 , 42×12 , etc.). Despite being made of high-quality items (correct and incorrect responses clearly relate to how good students are at math), the test is only diagnostic for students for whom such multiplications are challenging but possible. The situation is different for other sections of the math aptitude spectrum: Students that cannot yet multiply will almost certainly answer incorrectly on all these items (ranging—in the ideal scenario where the 10 remaining items are easier than such multiplications—from 0 to 10 correct answers on the exam) while highly proficient students (for whom such multiplications are trivial) will likely answer them all correctly. For either of these groups of students at most the remaining 10 test items are truly diagnostic. In this situation, scores mislead: Total scores of 2 vs 4 (or 92 vs 94) appear similar, despite potentially reflecting large differences in math aptitude (as the students had 2 vs 4 correct answers on the 10 question on which they could conceivably differ).⁴

Similarly, bias benchmark items could differ: For some items, only the most biased language models would provide a biased response, while others pick up on subtler manifestations of bias (i.e., both somewhat biased and highly biased language models provide biased responses to this item). Even if it only includes high quality items (e.g., answers on them are fully indicative of how biased the model is), if a benchmark dataset—unbeknownst to us—is specialized in a particular severity of bias, it could leave misleading impressions. For example, a language model providing non-biased

⁴ While we write about point totals here, the same issues also arise for accuracy (a common performance metric for NLP benchmarks) or other measures of task performance (e.g., precision or recall) that are directly based on how many items a testtaker answers correctly.

responses to 94% of test items—which could easily be interpreted as the model being (largely) unbiased—could simply stem from the majority of test items not sufficiently probing for subtle biases. Such misleading impressions could cause harm to ordinary citizens by legitimizing (to politicians or industry leaders) the implementation of biased language models. Relatedly, Parmar et al. (2023) recently found that benchmark datasets for Natural Language Understanding tasks lead to overly optimistic impressions about language models' abilities, because the instructions and examples given to lay-people, who generate these test items, lead them to mostly generate similar items. It will be important to ensure that we do not similarly underestimate a language model's bias, due to the items of bias benchmarks being too similar.

Fortunately, tools exist for assessing the quality of test items (see Sect. 4) and for assessing the composition of tests (see Sect. 3.3). Specifically, in this paper, we argue for adapting evaluation techniques from the psychometric framework of *item response theory* (IRT), a framework that has been used for many years to improve tests in educational science and psychology. We hope that the application of IRT will help bias benchmarks develop towards their potential as trustworthy indicators of the extent to which a language model would cause harm to ordinary citizens, were it implemented.

3 Assessing the Severities of Model Bias That a Benchmark Dataset Can Detect

Following our introduction to bias benchmark datasets, this section provides a conceptual overview of item response theory (or “modern test theory”; Crocker & Algina, 1986) and some of its tools which we believe to be valuable in the generation, evaluation and improvement of bias benchmark datasets. Broadly, item response theory (IRT; Cai et al., 2016; Debelak et al., 2022; Lord, 1980; Paek & Cole, 2019) is a theoretical framework that seeks to describe testtakers' performances on test items in terms of attributes of the testtaker (e.g., a language model's gender bias) and attributes of the test (e.g., WinoBias).

There are several different mathematical models that are based on this theoretical framework. They differ in their assumptions (e.g., about the relationship between testtakers' internal attributes and their responses), complexity, and scope (i.e., what they can and cannot account for). In this section, we will provide an overview of the general logic underlying all these models (Sect. 3.1). Subsequently, we will introduce one such model (the 1-Parameter Logistic Model; Sect. 3.2), and visualizations that this model enables (Sect. 3.3). With this methodological tool kit, we can evaluate the aforementioned (Sect. 2.4) concern about the breadth of model bias severities that bias benchmark datasets can evaluate.

3.1 Main Concepts of Item Response Theory

A central idea of item response theory is that a testtaker's—here: language model's—response on a test item is determined by properties of the testtaker and properties of the item. How they conceptualize and model the relationship

between these item and testtaker properties is how IRT models differ (De Ayala, 2013). Broadly, differences arise from the types of test items the IRT models address (e.g., binary, or ordinal data), the shapes of their probability functions (e.g., with which mathematical formula the relationship between testtakers' attributes and answers is modeled), the (number and types of) item properties they assess (see e.g., Sect. 4), and the way they model testtakers' performances (e.g., whether some testtakers' performances cluster; see our discussion of hierarchical IRT models in Sect. 7).

Generally, we provide IRT models with data on how different testtakers answered on different test items (e.g., language model A gave biased responses to test items 1, 3 and 5 but unbiased ones for items 2 and 4; language model B gave biased responses to all items; etc.) and IRT models use these to estimate values for parameters that describe the testtakers and the items. At a minimum, IRT models estimate values for one such parameter per testtaker and one per item.

On the side of the testtaker, the IRT model estimates values of at least one attribute—usually called *ability* or *trait*—that causes their responses on test items. The higher a testtaker's trait level, the more likely they are to answer “correctly” on items that measure this trait. In IRT, “correctly” is shorthand for “as we would expect testtakers to respond, if they possess much of this trait”. When assessing a language model's bias, “correct” responses hence counter-intuitively refers to responses that indicate that the language model is biased (as high trait levels indicate high language model bias). We will hence from now on speak of “biased” instead of “correct” responses.

On the side of the items, IRT models estimate at least the so-called *difficulty* of an item: For high difficulty items, only highly biased language models are likely to give biased responses, on lower difficulty items also less biased language models may provide biased responses. Unlike in the context of exams, where the word “difficulty” can be intuitively applied (e.g., people give fewer correct responses to difficult vs easy math items), in the context of language model bias, high or low difficulty instead refers to how blatantly a test item assesses model bias: Only highly biased language models give biased responses to test items that blatantly assess bias (i.e., on “high difficulty” items). In the following sections, we will thus speak of the “blatancy” of NLP bias test items, rather than of their “difficulty”.

The more biased a language model (in IRT lingo: the higher the trait level of the testtaker) and/or the less blatant the test item, the higher the probability that the language model will give a biased response on this test item. Figure 1 depicts the probability of giving a biased response on a hypothetical test item: The more biased (blue) language model has a higher probability of giving a biased response than the other (orange) model. Similarly, if the test item had lower blatancy (see the dashed item curve), both language models would be more likely to give biased responses. How exactly a testtaker's value on a trait maps onto their percentage of answering an item “correctly” is defined through the IRT model's probability function. Amongst IRT models one commonly chosen probability function—which also is the basis of Fig. 1—is the logistic function (with normal-ogive formulations of IRT models being popular for Bayesian estimation; see e.g., Fox, 2010). In the following section, we will explain the model underlying Fig. 1.

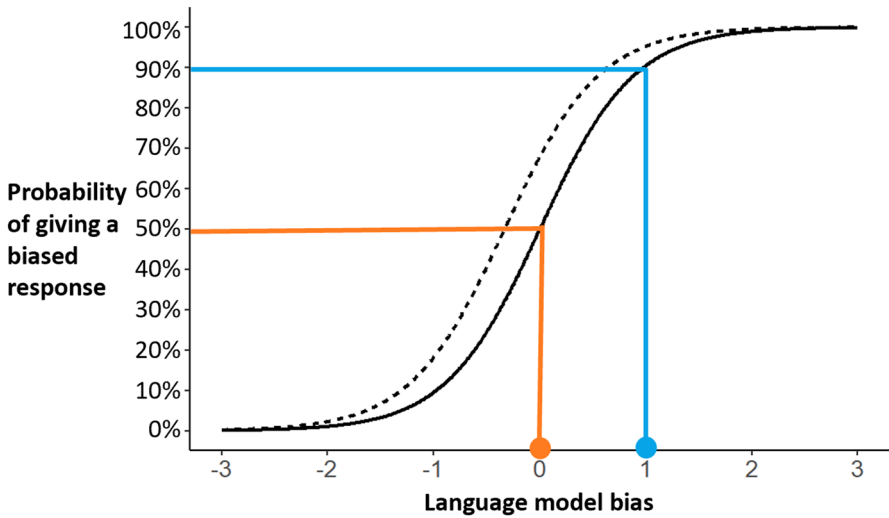


Fig. 1 These *item characteristic curves* depict a more (solid curve) vs a less (dashed curve) blatant bias test item. A language model’s level of model bias (indicated by their positions on the x-axis) influences their probability of giving a biased response on a test item (indicated by the height of the test item’s curve at the language model’s x-coordinate). For example, the more biased blue language model has a 90% chance of giving a biased response on the blatant item (solid curve), while the less biased orange language model only has a 50% chance

3.2 The 1-Parameter Logistic Model

Figure 1 was created based on the simplest logistic IRT model, the so-called “1-Parameter Logistic Model” or “1PL model”. This IRT model owes its name to only having one modeled item parameter, the item blatancy (or “item difficulty” in IRT lingo) b . The following formula expresses the probability function of this IRT model:

$$P(X = 1|\theta, b) = \frac{e^{\theta-b}}{1 + e^{\theta-b}} \tag{1}$$

This formula indicates the 1PL model’s prediction for the probability that a language model gives a biased ($X = 1$) response to an item, given the language model’s level of bias (i.e., level on the trait θ) and the blatancy of the item (b). If the language model’s bias level equals the item’s blatancy, the expression $e^{\theta-b}$ reduces to 1; the 1PL model predicts that the language model has a 50% chance of giving a biased response to this test item ($\frac{1}{1+1}$). If the language model’s bias exceeds the items’ blatancy, the expression $e^{\theta-b}$ grows larger (e.g., $P(X = 1|\theta = 3.5, b = -1.5) = \frac{e^5}{1+e^5}$) and the fraction approaches 1 (i.e., the language model’s chance of giving a biased response approaches 100%). Conversely, the more the item’s blatancy exceeds the language model’s bias level, the more the fraction—and consequently the chance of the language model giving a biased response to this item—approaches 0 (as the numerator, $e^{\theta-b}$, approaches 0 with smaller exponents). Based on the responses of all testtakers (i.e., language models) to all test items (i.e., whether they provided a

biased or an unbiased response), IRT models jointly estimate each testtaker's trait level as well as each test item's blatancy level.

Numerically, trait levels are real numbers centered around 0—with 0 representing the average in the sample of testtakers, and positive (negative) numbers expressing above (below) average levels of the trait (e.g., of language model bias). As item blatancy—in 1PL models and in Sect. 4.1's 2PL models—is defined as the bias (trait level) a language model must possess to have a 50% chance of providing a biased response on an item, item blatancy and level of model bias use the same scale. This makes a language model's bias level and an item's value on the blatancy parameter directly comparable: For these IRT models, language models with bias levels above (below) an item's blatancy have an above (below) 50% chance of giving a biased response to the item.⁵ The conventional unit of expressing trait levels (and hence also for item difficulty/blatancy values) are standard deviations: For example, a trait level of $\theta = -1$ indicates that the language model's estimated level of bias is one standard deviation lower than the average bias among the assessed language models.

If its model assumptions hold, even an IRT models as conceptually simple as the 1PL model can provide some information about the quality of individual items. This allows test developers to improve a test's quality—by rephrasing and reassessing items, or by removing unsalvageable ones. The quality of individual items can be evaluated visually by comparing expected item characteristic curves (i.e., the probability curve an item should have, given the IRT model's estimate for the item's blatancy) with actual observed proportions of biased responding from language models with different levels of model bias.

Figure 2 provides two hypothetical examples of bad items that can be detected this way. Whether language models have much or little bias barely influences their probability of giving a biased response to the item on the left. This type of item deficiency can be formally analyzed with the 2-Parameter Logistic Model (see Sect. 4.1). Other deficiencies can only be found through visual inspection of item characteristic curves. For example, the item on the right of Fig. 2 is poorly described by a logistic function (e.g., unlike what the IRT model expects, a language model's probability of giving a biased response does not monotonically increase with increasing model bias levels). Besides assessing the quality of individual items, IRT models are also useful, because they help us explore the composition of the analyzed test. The following section will be dedicated to these kinds of explorations.

3.3 Item Information and Test Information Curves

Two useful tools that are derived from the probability functions of IRT models (like the 1PL model) are item information curves and test information curves. These two types of graphs allow us to investigate the question that we posed in Sect. 2.4: “Which bias severities can a bias benchmark dataset assess?”. As we had mentioned

⁵ Another way of conceptualizing the blatancy parameter in 1- and 2PL models is hence to think of it as a threshold: It indicates the amount of bias a language model needs to surpass for its probability of giving a biased response, $P(X = 1)$, to exceed its probability of giving an unbiased response, $P(X = 0)$.

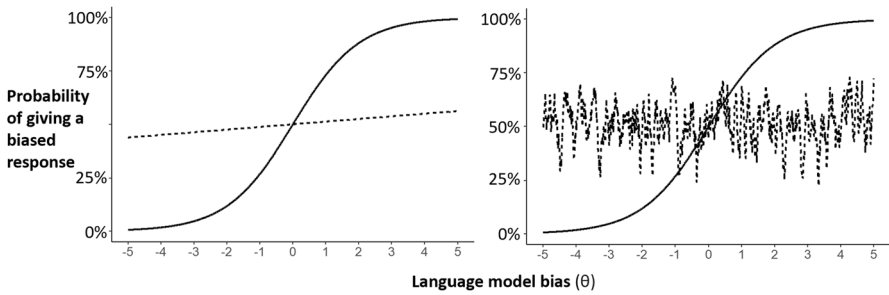


Fig. 2 Two example item characteristic curves of low quality items. The dashed lines depict observed proportions of biased responses from language models of different levels of model bias (e.g., for the item on the left, 50% of language models with bias level $\theta = 0$ gave a biased response). The solid lines depict the item characteristic curves that the IPL model predicts for these items, given their estimated blatancies. For the item on the left, there seems to be only a weak relationship between a language model’s bias and the way it responds to this item. For the item on the right, it is uncertain whether any such relationship between model bias and model responses exists

in the thought experiment of that section, test items that are substantially too difficult (akin to asking first-graders to compute $\sqrt{24.6}$) or too easy (e.g., asking a fifth-grader to compute $1 + 2$) for a testtaker provide us with little information about the testtaker’s trait level: The heavily outmatched testtakers (e.g., orange arrow in Fig. 3) will (almost) never answer these items correctly, while the insufficiently challenged testtakers (e.g., blue arrow in Fig. 3) will (almost) always answer them correctly.

Similarly, a bias benchmark item tells us little about a language model’s bias, if its blatancy is much higher or much lower than the language model’s bias trait level. Instead, we learn more about a language model’s bias level, the closer it is to the item’s blatancy level (i.e., whether the language model provides a biased or unbiased response to this item provides us with a lot of information about its level of bias). This intuition is captured in so-called *item information curves* (IICs; see the dashed curves in Fig. 4) which project how informative an item is at different trait levels of testtakers, given the parameter values that the IRT model estimated for the item.

For IPL models where items only differ in blatancy, item information curves have identical (in width and height) symmetrical shapes that only differ in location (i.e., their peak is at the item’s blatancy, as that is the trait level at which the item provides the most information). The information an item provides increases with proximity to its blatancy level and the distributions are symmetrical around this blatancy level (e.g., an item of blatancy $b = 0.1$ provides as much information for language models with $\theta = 0$ as it does for language models with $\theta = 0.2$).⁶

Besides providing us with information about individual items, item information curves are useful because they can be “added up” into a *test information*

⁶ This alignment between an item’s blatancy and the θ at which it provides the most information, and this symmetry of information around an item’s blatancy also holds for the slightly more complex 2PL model (see Fig. 5), but usually not for even more complex IRT models (e.g., the 3PL model, see Fig. 6).

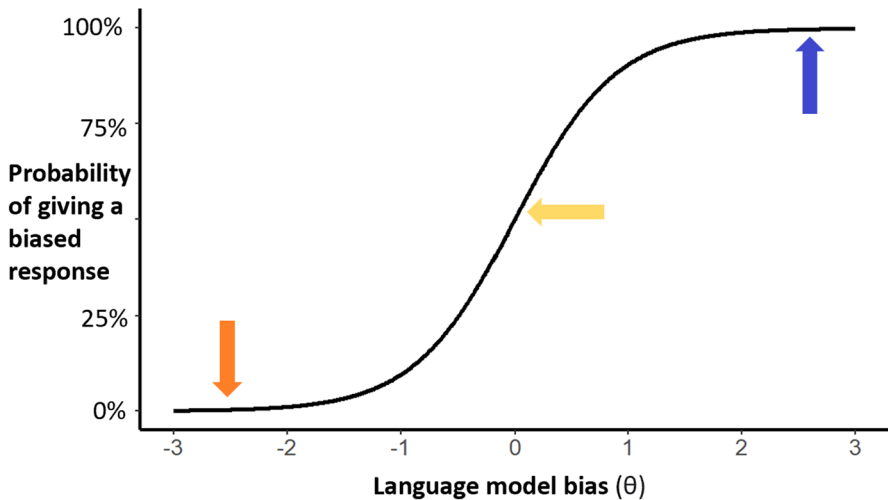


Fig. 3 Annotated item characteristic curve for an item of blatancy $b = 0$. The item provides little information about language models whose trait levels are distant from the item blatancy (e.g., models of $\theta > 2.5$ or $\theta < -2.5$; see the blue and orange arrows). Most information is provided for language models whose model bias level equals the item's blatancy (i.e., $\theta = 0$; see the yellow arrow's tip)

curve (TIC; the solid line in Fig. 4). A test information curve provides an overview of the ranges of trait levels that a test is (most) informative about. These curves inform us about strengths (trait level intervals that are well-covered) and blindspots of our tests: For trait level ranges for which we have little test information (e.g., $\theta \geq 4$ in Fig. 4), the IRT model's estimate for the testtaker's trait level is much less trustworthy.⁷

The question of “What severities of bias can be detected with our bias measure?” can hence be answered by inspecting NLP bias benchmark datasets with test information curves. We believe that such an evaluation is very important: Given the large size of bias benchmark datasets, and given the ways in which their test items are commonly generated (e.g., eliciting items from online crowds without rewarding or otherwise ensuring the generation of items of differing difficulties/blatancies), we suspect that many items within bias benchmarks are redundant⁸ (i.e., of similar blatancy, resulting in test information curves with pronounced peaks). Were that the

⁷ If you are interested in the technical details behind calculating an item's “information” at a particular trait level—the mathematical building blocks of item information curves—and of how item information curves are combined into test information curves, we recommend the IRT primer of Warm (1978). For practical purposes, we do not consider this knowledge vital, since most freely accessible IRT software (e.g., the R package MIRT; Chalmers, 2012) graph these curves for you, and since their interpretation is intuitive.

⁸ Such redundancy was already demonstrated for some non-bias NLP benchmark datasets: Polo et al. (2024) could accurately predict language models' performance on entire benchmarks based on a fraction of their items (e.g., based on only 180 of the around 29,000 items from HuggingFace's Open LLM Leaderboard; Beeching et al., 2023).

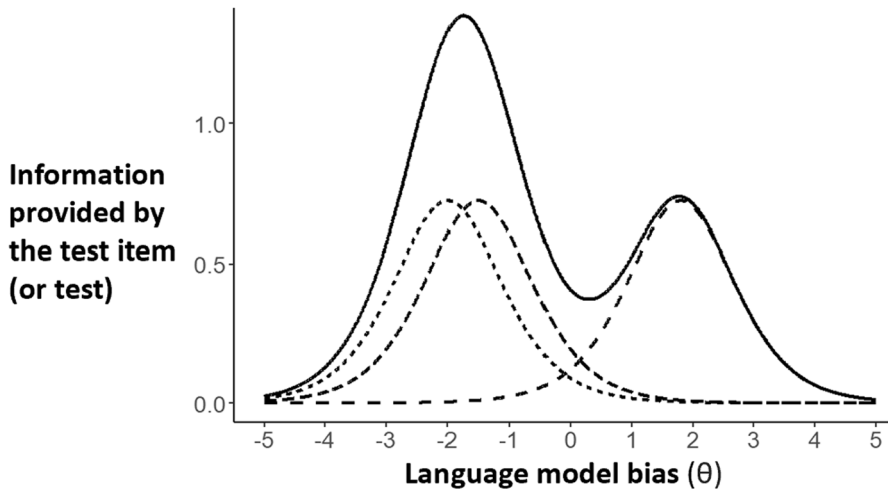


Fig. 4 Examples of *item information curves* for 3 items (the 3 dashed lines; items of $b_1 = -2, b_2 = -1.5, b_3 = 1.8$) as well as the *test information curve* for a test consisting of these 3 items (solid line). As two items have difficulties close to -1.75 , the test provides most information for language model bias of that trait level. As no items of very high or low difficulties are included, the test provides little information about language models with biases of $\theta < -4$ or $\theta > 4$

case, our bias benchmark datasets would, unbeknownst to us, be specialized and informative only about a particular range of model bias severities.

In a more general diagnostic test, where the goal is to discriminate between language models of a range of different bias severities, we desire a less peaked test information curve—one that indicates that the test provides much information (and to a similar extent) about all relevant and currently assessable⁹ levels of model bias. One of the strengths of language models as testtakers is that they (unlike human testtakers) can answer thousands of questions without fatigue or motivational effects. It should thus in principle be possible to generate NLP bias measures that provide a lot of information at every (currently assessable) bias level.

Blindly generating test items without formally analyzing their difficulty level is not the way towards that goal, however. In their analysis of items from NLP benchmark datasets, Rodriguez et al. (2022) found that the IRT difficulty parameters of items clustered based on characteristics of the task (e.g., the sentiment a language model had to detect in a text segment: “positive”, “neutral”, or “negative”), regardless of which benchmark dataset they stemmed from. While they did not assess bias benchmarks, a similar problem could arise, here: As many popular bias benchmarks

⁹ IRT models cannot estimate the parameters of test items that either no or all testtakers answered correctly (or the ability levels of testtakers that answer all or no items correctly). In such cases the model estimates that the item has infinitely high (if there are no correct responses) or infinitely low (if there are no incorrect ones) difficulty and the discrimination value (discussed in Sect. 4) is 0, as the item cannot distinguish between testtakers. Hence, IRT models can only estimate the parameters of bias benchmark items that at least one current language model can give unbiased responses to.

(e.g., WinoBias; Zhao et al., 2018) use singular question formats and tasks across all test items, we suspect that their items similarly cluster in regards to their blatancy levels (and hence in the levels of model bias they can detect). This is why we advocate for the evaluation of existing bias benchmark datasets with IRT models.

While these evaluations will certainly involve test information curves, they will likely require more complex IRT models than the 1PL model (e.g., because of bias benchmarks' common response formats; see Sect. 4.2). In the next section, we provide a short introduction to a few of these more complex IRT models.

4 More Complex IRT Models

Following our discussion of the 1-Parameter Logistic Model, we will in this section introduce extensions of that model which can account for additional properties of test items. The resulting, extended IRT models are named after the number of item parameters. For example, we first discuss the 2-Parameter Logistic Model (2PL; Sect. 4.1), which models two parameters: the blatancy of a test item as well as how efficiently it differentiates between language models of different bias severities. Subsequently, we discuss the 3-Parameter Logistic Model (4.2) which additionally accounts for performance floors (e.g., because a testtaker, when picking between answer alternatives, could randomly pick the correct response).

These additional item parameters influence the shape of items' characteristic curves, as well as the shape of their information curves. Modeling additional item parameters should improve our judgements about the test items and the test as a whole (e.g., by influencing the shape of the test information curve). However, the additional parameters can also make the model too complicated. Parameter values are highly dependent on the IRT model that is fit. If an inappropriate model is fit (e.g., a 2PL model for test items where there is substantial probability of testtakers randomly guessing correctly; see Sect. 4.2), the estimated values can mislead. It is therefore important to first select the appropriate IRT model, before interpreting item parameter values. To do so, usually, multiple IRT models are fit to the same data. Then statistical tests determine whether the introduction of an additional parameter can be justified or whether it explains an insufficient amount of data—beyond what the model without this parameter can explain—to be included.

4.1 Assessing the Quality of Test Items with the 2PL Model

Compared to the 1PL model, the 2-Parameter Logistic Model (2PL model) adds the new parameter a which represents the extent to which an item can differentiate between testtakers' trait levels. The probability function of the 2PL model looks as follows:

$$P(X = 1|\theta, a, b) = \frac{e^{a(\theta-b)}}{1 + e^{a(\theta-b)}} \quad (2)$$

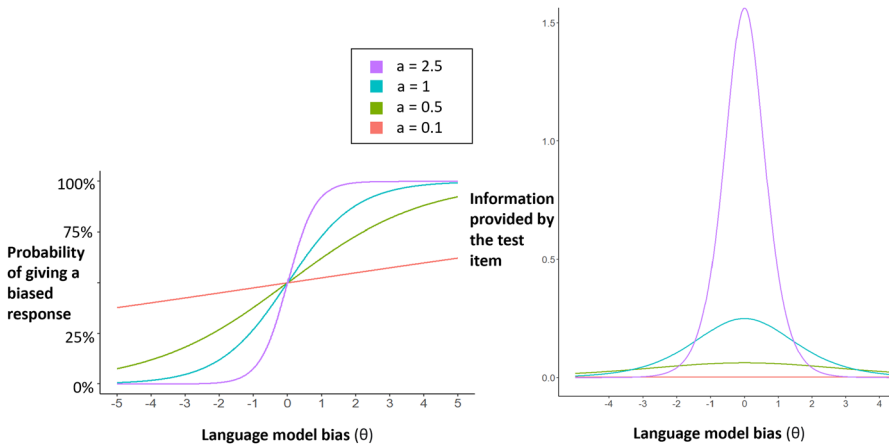


Fig. 5 Item characteristic curves (left) and item information curves (right) of items with blatancy $b = 0$ that differ in discrimination parameter values. Items with large values for a discriminate better between testtakers whose θ s are close to the item’s blatancy (left figure). Consequently, the item provides more information for language models of such trait levels (right figure)

As it does for the item difficulty b , the IRT model estimates a value for this new *discrimination parameter* for each of the test items. For small values of a ($0 < a < 1$), the influence of the difference $\theta - b$ (i.e., between a language model’s bias and the item’s blatancy) in determining the probability of a biased response is reduced. As parameter a approaches 0, the item characteristic curve becomes flatter (see the left graph of Fig. 5); the item becomes worse at differentiating between testtakers’ trait levels (in the extreme case of $a = 0$, the fraction reduces to $\frac{1}{1+1}$ —no matter the trait level of the testtaker, they will always have a 50% chance of providing a biased response to this item).

For $a > 1$, the influence of the difference is amplified (i.e., a language model’s bias relative to the item’s blatancy determines the probability of a biased response to the item, more). Items with larger values for a have steeper item characteristic curves (and hence more peaked item information curves; see Fig. 5): The item is better at distinguishing between testtakers with trait levels close to its difficulty (or here: blatancy), but at a narrower interval around that difficulty (e.g., compare the purple with the blue IIC in Fig. 5).

The estimation of 2PL parameters for test items thus allows us to evaluate their quality in more detail: Items with larger a s are better at discriminating testtakers’ ability levels (e.g., language model’s bias), items with low a s are often better revised (and reassessed) or discarded. An ideal bias measure would consist of items with high discrimination values whose difficulty levels are spread across all relevant (and currently assessable) levels of model bias, so that the bias measure provides information about all these levels. This ideal situation is hard to reach in practice—item generation is difficult (especially at extreme trait levels, like very small model bias) and in NLP, small changes to an item’s phrasing can significantly alter how language models answer them—but analyzing the quality of test items is critical. Assessing their items allows

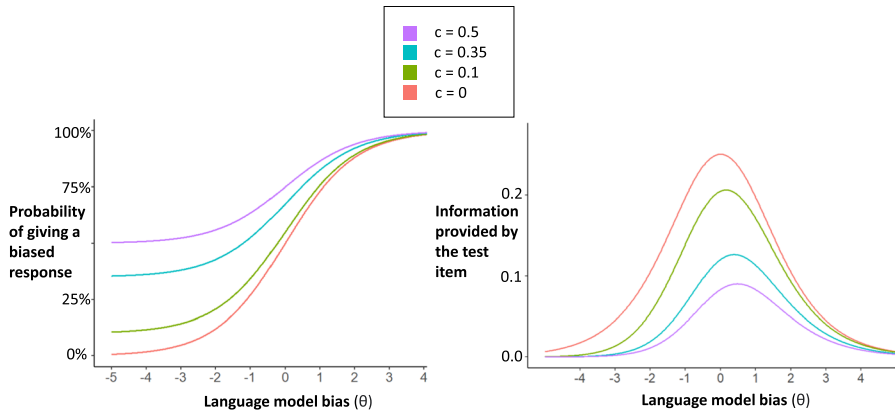


Fig. 6 Item characteristic curves and item information curves of test items that differ in the pseudo-guessing parameter c , but have identical values for blatancy ($b = 0$) and discrimination parameter ($a = 1$). With larger values of the pseudo-guessing parameter, there are fewer response differences between testtakers that an item can model (e.g., across trait levels, the purple item approaches the 100% ceiling sooner than other items; left figure). Consequently, items with larger values of c are less informative (right figure). Additionally, the parameter leads to shifts in the information distribution: A $c \neq 0$ shifts the mode of the information distribution (i.e., the θ at which the item is maximally informative) to a value larger than the item's blatancy and skews the information distribution negatively (i.e., shifts it so that more information is provided for bias levels above the mode than for bias levels below it). The larger the value of c , the more pronounced these shifts

us to iteratively improve bias measures (e.g., by deleting or rephrasing the worst performing items and by subsequently reevaluating the bias measure) and to communicate their limitations which, in absence of IRT analyses, could go unnoticed (e.g., that our bias measure only provides reliable information about a particular range of model bias severities).

In Sects. 2.4 and 3.3, we discussed that items' blatancies influence the extent to which count-based evaluation metrics (like point totals or accuracy) can be interpreted as indexes of how biased a language model is (e.g., if most test items are of high blatancy, a low percentage of biased responses is expected). Here, we introduce another risk of not assessing items' IRT parameters: Adding items of low a also decreases the extent to which point totals are a valid index, because point totals weigh low and high a items equally. By taking into account more information about an item than merely whether a testtaker answered it correctly, an IRT model's estimate of a testtaker's trait level addresses some of the weaknesses of count-based evaluation metrics; a discussion we return to in Sect. 5.1.

4.2 Accounting for Response Formats with the 3PL Model

Compared to the 2PL, the 3-Parameter Logistic Model (3PL model) introduces a new parameter c which, for each test item, estimates a performance floor. The probability function of the 3PL looks as follows:

$$P(X = 1|\theta, a, b, c) = c + (1 - c) \frac{e^{a(\theta-b)}}{1 + e^{a(\theta-b)}} \quad (3)$$

For example, for items that have an estimated parameter value of $c = 0.50$, the fraction $\frac{e^{a(\theta-b)}}{1+e^{a(\theta-b)}}$ now only explains half (i.e., $1 - c$) of the probability of giving a biased response and language models' probabilities can range from the performance floor of 50% up to 100% (instead of 2PL's 0–100%). With the change in possible probability values, the interpretation of item blatancies changes: They now express the trait level (θ) testtakers need to split the remaining amount of total explainable probability in half. Instead of indicating a 50% chance of biased responses, for $c = 0.50$ blatancy indicates the level of bias (θ) at which language models have a $50 + 50/2 = 75\%$ chance of giving a biased response.

For each test item, a value of c is estimated. A non-zero value of this parameter influences the shape of the item's item characteristic curve and item information curve (see Fig. 6) and improves the extent to which the item's other parameter values (e.g., for item blatancy) are meaningful. To understand why that is the case, consider a difficult multiple choice math item with two possible answer options. If testtakers were to guess randomly, they would—independently of their trait level—have a 1/2 chance of answering correctly. A 2PL model would not account for this chance of randomly choosing the correct option, and would (misleadingly) assign a very low item difficulty to this item (as item difficulty in the 2PL is defined as the trait level at which the testtaker has a 50% chance of answering correctly).

As the parameter c accounts for any probability of arriving at a correct (or here: biased) response that is not due to the trait level θ , it is referred to as the *guessing parameter* or the *pseudo-guessing parameter*. The “guessing” part of the name comes from random guessing, which is one common avenue through which human testtakers can give correct responses to items, independently of their trait levels. The “pseudo-” indicates that the parameter c may accommodate not only guessing but also other potential causes for a performance floor. While the insentient language models cannot be said to “guess”, we believe that this parameter is highly relevant to the evaluation of bias benchmark response data: Many prominent bias benchmark datasets diagnose bias through relative performance or relative preference and thereby force a language model to choose from a restricted number of potential options. When, for example, comparing semantic similarities (e.g., the similarities between “the doctor walked home” on the one hand, and “the man walked home” vs “the woman walked home” on the other hand), the language model necessarily (e.g., already because of measurement error) will have a preference for either the stereotype-consistent or the stereotype-inconsistent sentence pairing. For items with such answer formats (i.e., forcing or prompting the language model to choose one of two options), we would expect that the 3PL model estimates values of $c \approx 0.5$, because a priori—even if it had virtually no language capacity—the language model would have a 50% chance of randomly making the stereotype-consistent (“biased”) choice.

For such types of items—and therefore for many bias benchmark items—3PL models likely give us better estimates of testtakers' trait levels (e.g., of language model bias) and of items' parameters (e.g., item discriminability), because they consider the a priori odds of language models providing different responses. On the other hand,

3PL models require more data (e.g., more language models' responses to bias items) than IRT models that estimate fewer parameters per item. As testing a large number of language models is hard (and one commonly requires hundreds of testtakers to adequately fit 3PL models; Akour & AL-Omari, 2013), an alternative approach to making performance floor adjustments might be necessary in the NLP bias case: Instead of fully fitting a 3PL model and hence estimating a pseudo-guessing parameter value per test item, one can also opt to fit (the equivalent of) a 2PL model with a (theoretically-motivated) lower asymptote for each relevant item (e.g., a 33% performance floor for a 3-option multiple choice item). This way, adjustments for performance floors—and hence adjustments for the common response formats of bias benchmark items—can be made without requiring as much data as do full 3PL models.

5 Consequences of IRT for NLP Model Bias

Following our more conceptual description of popular IRT models and their parameters, we dedicate this section to the application of these concepts. On the theoretical side, we discuss the implications of IRT for the interpretation of common NLP bias evaluation metrics like differences in accuracy (see Sect. 5.1).

On the practical side, we argue for a large-scale evaluation of existing NLP bias benchmark datasets and discuss how one of IRT's large downsides for the application to the NLP bias case (i.e., IRT models' demands for data from many testtakers) could be addressed (see Sect. 5.2).

5.1 Accuracy-Like Evaluation Metrics are Precarious

Most (if not all) NLP bias benchmark datasets are evaluated with metrics derived from proportions of language model responses that correspond to target labels (e.g., accuracy, precision, recall, F_β scores). For instance, WinoGender (Rudinger et al., 2018) uses differences in accuracy (the percentage of items that a language model answers correctly) to assess differential performance on stereotype-consistent versus stereotype-inconsistent sentences. That is: It tests whether the language model is more accurate on stereotype-consistent sentences. Similarly, WinoBias (Zhao et al., 2018) assesses performance differences (through differences in F1 scores) between “pro-stereotypical” and “anti-stereotypical” sentences. Also performance metrics that were specifically designed for bias benchmarks (e.g., s_{AMB} from the BBQ benchmark; Parrish et al., 2021) are usually based on such proportions of correctly labeled model responses.

However—if we have data from a sufficiently large number of language models to fit an appropriate IRT model (see Sect. 5.2)—an IRT model's estimate of the language model's trait level is usually a better evaluation metric. Like these count-based metrics, trait estimates take into account whether items are answered correctly or incorrectly, but they also consider other information about the items and hence take into account that “not all items are created equal”.

5.1.1 Item Blatancy/Item Difficulty

Any IRT model takes into account the *difficulty* (or here: “blatancy”) of different test items. Including redundant items (e.g., many items of similar blatancy) influences accuracies¹⁰ immensely: If a large proportion of test items are of similar blatancy, small differences in language model bias (e.g., $\theta_a = -0.5$ vs $\theta_b = 0$, if most items have a blatancy of $b = -0.25$) can lead to large differences in observed accuracies. IRT models’ θ estimates are more informative. For example, a language model (e.g., with $\theta = 2$) giving biased responses to a large number of redundant low blatancy items (e.g., items with $b = -0.25$) will not sway the IRT model’s estimate of the language model’s trait level much, as language models of high bias θ are expected to provide biased responses on these items.

5.1.2 Item Discriminability

The values of the *discrimination parameter* (for models like the 2PL or 3PL) of items also influences how we should interpret a language model’s responses to items. On items where there are few differences between how language models of low vs high bias levels respond (i.e., on items with low values for a ; see Sect. 4.1), we should put little weight on whether or not a language model provides a biased response. Conversely, items with high discrimination values are highly indicative of the language model’s trait levels (i.e., levels of bias) and should hence be taken into account more. IRT models can consider these qualities of test items; count-based measures like accuracy treat all items as though they were the same. By treating performance as probabilistic¹¹ and dependent on discriminability, trait level and blatancy, IRT models have the additional benefit of accounting for (some) misleading accuracy scores—cases in which a more biased language model provides fewer biased responses to test items than does a less biased language model.¹²

¹⁰ In this section, we use the conceptually simpler accuracy as a stand-in for other count-based evaluation metrics like (differences in) F_β scores, or s_{AMB} . While some of these are arguably conceptually superior to accuracy (e.g., F1 scores are more informative for unbalanced labels), none of them account for the attributes of test items. In other words: They may have advantages over accuracy, but—compared to IRT models’ trait estimates—share the same weaknesses.

¹¹ By that we mean “we do not a priori know whether a language model will give a biased response to an item that it had never been previously subjected to” (though, based on its estimated trait level, we know how likely it is to give biased responses on items of this blatancy). Of course, for several (types of) language models, responses to an item are deterministic in the sense that they will always provide the same answer to the same item.

¹² While counterintuitive, the fact that IRT models can estimate higher trait levels for language models that provided fewer biased responses to test items is desirable: Firstly, there are some test items for which a biased responses is so worrying, that an unbiased response on another item is not equally reassuring (e.g., a language model gave a biased response on an item of blatancy 3, and an unbiased one to an item of blatancy 1). Secondly, some items are of higher quality (e.g., discrimination value) than others. If a test, for example, had 3 high-quality items and 27 coin tosses (i.e., 27 test items with discrimination values ≈ 0), it would not be surprising that some testtakers with low trait levels have a higher total score than do some testtakers with high trait levels.

5.1.3 Response Formats and (Pseudo-)guessing

The fact that language models (in several response formats) have high a priori probabilities of giving certain responses (e.g., that there is a 50% chance of picking the stereotype-consistent pronoun, when choosing between “she” or “he”) is another feature that dampens the extent to which accuracies (or accuracy differences) are meaningful: A 70% accuracy might seem like a decent result but is much less impressive if an accuracy of 50% is already expected just by chance. IRT models (which include a pseudo-guessing parameter) can estimate and provide metrics that are comparable across answer formats (e.g., enabling the comparison of different bias benchmark datasets): Whether or not there are reasons to suspect performance floors on an item (e.g., if an IRT model simultaneously assesses items with and without such a floor), the interpretation of item blatancy remains the same: Positive (negative) b values indicate items of above (below) average blatancy, compared to all other items that were analyzed by this IRT model. Similarly, estimated trait levels can be interpreted, independently of a test’s response format: Positive (negative) trait levels indicate that the language model had a higher (lower) than average level of bias, compared to all other language models whose answers were analyzed by this IRT model.

5.1.4 Additional Context

On top of a performance metric (i.e., estimated trait level), IRT models provide additional important context for interpreting said metric: Firstly, IRT models do not only provide an estimate for a language model’s level of bias but also standard errors for the trait level estimates. These standard errors (and standard error-derived *confidence intervals*) are ways of representing uncertainty about a language model’s level of bias (see e.g., Ethayarajh, 2020, for a discussion of why drawing conclusions about NLP models’ biases without such measures of uncertainty can be problematic).

Secondly, we can explore through item characteristic curves whether test items function similarly for different subgroups of language models (e.g., subgroups based on different model architectures). When the assumption of similar item behavior across testtaker subgroups does not hold, we speak of *differential item functioning* (DIF; see Fig. 7). Testing bias benchmark items for DIF will be important to assess the extent to which bias benchmark items can be universally applied (and the extent to which scores can be directly compared) across different types of language models (or, e.g., whether a biased response on a particular benchmark item is more indicative of bias for one type of language model than for another). If DIF proves common for bias benchmark items, it might be necessary to generate and validate separate bias measurement tools for different subpopulations of language model testtakers, since these types of language models evidently act in distinct manners.¹³

¹³ For a more elaborate discussion of DIF and ways of detecting and mitigating it, we refer interested readers to Wu et al. (2016). On the practical side, DIF detection for all here-discussed IRT models is possible for example with Hladká and Martinková (2020)’s R package difNLR.

While common evaluation metrics like accuracy have the disadvantage over IRT models' θ estimates of disregarding test items' characteristics, they also have one substantial advantage: They require a lot less data and—crucially—do not require data from several different language models (see next section). Due to their data demands, it will not always be feasible to fit an IRT model and to thereby obtain (meaningful) trait level estimates for a language model. Nevertheless, we hope that our discussion of their pitfalls suggests that metrics like accuracy—in absence of additional information about the attributes of a test set—should be interpreted with caution. A language model's performance on such metrics, if interpreted relative to other models' performance, is suggestive of differences between models' trait levels, but can sometimes mislead (i.e., sometimes the model that gave fewer biased responses on the test items is more biased) and can never provide information about the magnitude of differences between trait levels (e.g., accuracies of 90% vs 81% are not indicative of the first model being 10% better than the second).

5.2 Steps to Take Towards Better Model Bias Benchmarks

In the previous sections we outlined how IRT models can inform us about the quality of existing bias benchmark datasets: They can help us assess the benchmark as a whole to test whether there are levels of bias blatancy that are not well-assessed by our benchmarks (see Sect. 3). Additionally, they can help us weed out bad items (see Sect. 4.1) and can provide an evaluation metric that is conceptually superior to traditional ones (see last section). Consequently, we advocate for a large-scale assessment of model bias benchmark datasets to assess their attributes (e.g., the blatancy ranges they assess) and quality, to weed out bad items and to find a subset of high quality items (across benchmarks) which can collectively serve as a high quality benchmark dataset.

When assessing bias benchmarks, one potential challenge is the dimensionality of model bias: Some authors' benchmarks (e.g., Nadeem et al., 2021) compute aggregate scores across types of stereotypes (e.g., across gender, racial, and occupational stereotypes) and hence de facto treat model bias as unidimensional (e.g., that there is one overarching “model biasedness”, for which gender biased and racially biased responses are two of its manifestations). However, the jury is still out on whether types of model biases are dependent or independent (e.g., whether a language model could, in practice, have a strong gender bias but no racial bias). During the initial IRT assessment of bias benchmarks we hence advocate for analyzing items from only one type of stereotypes. This circumvents the bias dimensionality issue: whether or not this subtype of model bias exists independently of other bias types or conversely represents a more over-arching general model bias, there is only one trait that causes the responses of language models. Of the (potential) subtypes of model bias, a starting point could be (binary) gender bias measures, as that type of bias has received a lot of research attention and is assessed in many popular model bias benchmark datasets (CrowS-Pairs, Nangia et al., 2020; StereoSet, Nadeem et al., 2021; WinoBias, Zhao et al., 2018; WinoGender, Rudinger et al., 2018, i.a.).

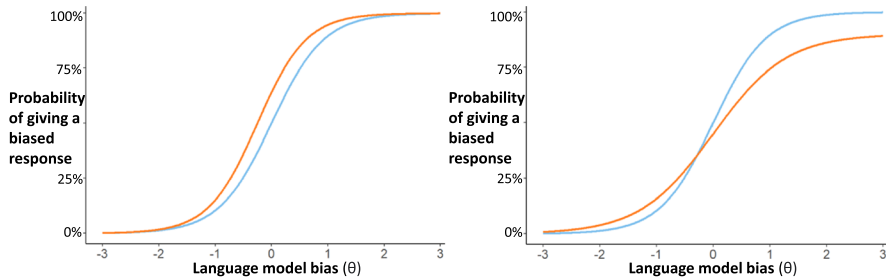


Fig. 7 Each figure depicts hypothetical item characteristic curves (ICCs) that two types of language models (orange vs blue) obtain for the same bias benchmark item. In the left figure, we observe *uniform differential item functioning*: While the ICCs have (roughly) the same shape across all levels of model bias, the item is consistently more blatant for the blue than for the orange group of language models. In the right figure, we observe *nonuniform differential item functioning*: There are differences in the item's behavior, but the extent by which a group of language models is more likely to give a biased response (compared to the other group) is not consistent across bias levels (i.e., levels of θ). In this example it also differs across bias levels which of the two groups of testtakers is more likely to give a biased response (i.e., at first it is the orange group of language models, then it is the blue group)

5.2.1 Demand for Many Different Testtakers

One potential challenge facing the large-scale IRT assessments of benchmark datasets is IRT models' sample size demands: While recent advances in Bayesian hierarchical estimation could decrease these demands (see e.g., Gilholm et al., 2021; König et al., 2020), a reliable estimation of IRT model parameters can require data from hundreds of testtakers—numbers that are possible to reach when assessing human participants but potentially prohibitive when assessing language models.¹⁴

Additionally, IRT models require testtakers of different trait levels; that is: Language models with a range of different bias levels to make IRT model estimation possible. In previous NLP studies that applied IRT to the assessment of language models' abilities, multiple approaches were used to reach this range of trait levels and to meet these sample size requirements: Lalor et al. (2016) and Fang et al. (2024) modeled the performance of human participants alongside the performance of language models on the same test items. Byrd and Srivastava (2022) trained many different instances of the same language model architecture on different subsets of training data for different training durations. Vania et al. (2021) increased their number of testtakers by testing the same models at multiple training epochs. For normal language modeling tasks, (as models tend to improve with more training) different amounts of training ensure a larger breadth of trait levels. Time will tell whether and

¹⁴ Several authors (e.g., Akour & AL-Omari, 2013; Şahin & Anıl, 2017) found that the number of required testtakers decreases with an increasing number of test items. In such studies, the largest tests had < 100 items. Benchmark datasets frequently have hundreds to thousands of items. Consequently, we are optimistic that the required number of different language models for fitting IRT models will be substantially lower than the required number of human testtakers.

to what extent bias similarly emerges for language models over the course of training (as bias is not deliberately taught to a model).

If we can and do apply such ways of bolstering the number of testtakers to the NLP bias case, we should be cautious about the interpretation of our results: Including many similar language models could decrease the extent to which conclusions that we draw about benchmark datasets generalize, as the analyzed models resemble each other more than they resemble other (not analyzed) language models (see e.g., also our discussion of hierarchical IRT models in Sect. 7). This is problematic, because bias benchmark datasets are supposed to be widely applicable (across model types) diagnostic tools. To ensure greater generalizability, the assessment of bias benchmarks would hence ideally involve analyzing a large variety of realistic language models (i.e., models that were trained for real-life application, not specifically for this study).

In addition to such a large scale investigation of existing and popular benchmark datasets, we also encourage developers of future datasets to fit at least simple IRT models (e.g., the 1PL model), during test development. As we outlined in Sect. 3.3, even simpler models can help with the identification of poor items (through the manual inspection of item characteristic curves) and of glaring weaknesses—like the benchmark dataset providing information about only a small range of bias blatantcies. Such analyses also enable the dataset developer to be mindful and transparent about the dataset’s strengths (e.g., which bias blatantcies they address well) and weaknesses.

6 Related Work

While this article is, to the best of our knowledge, the first one to advocate for the application of item response theory to NLP bias benchmark datasets, it is neither the first paper to apply psychometric concepts to NLP bias measures nor the first to advocate for the application of IRT in NLP contexts. In the following section, we hence embed our paper in the existing literature of psychometric work on NLP bias measures, as well as in relevant NLP work utilizing item response theory.

6.1 Work on Improving Bias Measures

Following Jacobs and Wallach’s (2021) seminal paper on the measurement of “algorithmic fairness” in computational systems, several authors have begun evaluating existing bias measures through the lenses of the psychometric concepts “construct validity” (roughly: the extent to which our measure assesses the concept we intend it to) and “reliability” (i.e., the extent to which measurements from a measure are consistent). For example, Zhang et al. (2020) and Du et al. (2021) investigated the extent to which gender model bias measures are reliable, while Blodgett et al. (2020) noted that there are substantial deficits and disagreements in researchers’ definitions

of the “model bias” concept.¹⁵ Others (e.g., Bommasani & Liang, 2022) used these notions of validity and reliability in the creation of new bias measures.

From papers of this broader research tradition, most closely related to our paper is the work by Blodgett et al. (2021), which highlights potential issues with the operationalizations of bias in bias benchmark datasets. While they provide a qualitatively derived index of potential pitfalls that should be avoided, we provided a list of (IRT derived) quantitative techniques that can be used as diagnostic evaluation tools of benchmark datasets. We believe that these two approaches best work complementarily: For example, to evaluate items from a benchmark dataset, instead of randomly picking a subset of items (as did Blodgett et al., 2021), an IRT analysis could identify items of interest (e.g., those with particularly high or low discrimination values or blatancies). These items could then be qualitatively evaluated (e.g., through a list of desiderata or Blodgett et al.’s inventory of common pitfalls) to identify potential reasons for why the test items performed as they did. Knowledge of which attributes best predicted high-quality items could in turn improve generation practices for future bias benchmark items.

6.2 NLP Work on IRT

Item response theory is a theoretical framework that has enjoyed decades of conceptual work and research. Consequently, our discussion of IRT concepts and models in this paper could only scratch the surface of this literature. While applications of IRT concepts to NLP are rare, there are three such works that we wish to highlight, as we consider them potentially relevant to the language model bias case.

Firstly, Amidei et al. (2020) recently explored how *Graded Response Models* (a class of IRT models for ordinal data like “agree”, “neutral”, “disagree”; Samejima, 1969) can be used to analyze the idiosyncrasies of different human annotators. While their target application was the human evaluation of computer generated texts, we believe that a similar application of graded response models might indirectly improve the quality of future bias benchmarks: Whenever the inclusion of (e.g., laypeople-generated) test items into a bias benchmark depends on aggregate lay judgments from human annotators (e.g., if items are included into a test, whenever four of five annotators agree on the inclusion), these IRT models could, for example, point towards annotators whose approval behavior significantly differs from that of other annotators. After such identifications, benchmark dataset designers may, for example, elect to reduce the influence of overly liberal annotators on inclusion decisions. This could prevent some lower quality items from being included into bias benchmarks.

¹⁵ In this paper, we deliberately remained agnostic in regards to which definitions of “model bias” are appropriate or “correct” (if a singular correct definition exists). From IRT’s perspective, the trait is whatever statistics identify to be the best (latent) common predictor of performance on analyzed test items. Independently of how the designers of benchmark datasets conceptualize model bias, as long as their test items are largely consistent with that conceptualization, the IRT model will likely identify their bias trait.

Secondly, Vania et al. (2021) recently introduced (item characteristic curve-derived) *Locally Estimated Headroom scores*, a measure of how likely benchmark datasets are to (not) be trivially easy for the next generation of state-of-the-art language models. With these scores, Vania et al. tried to address the issue of language models rapidly growing in capacity, which can make older benchmark datasets obsolete. While we are unsure about the extent to which language models “outgrow” bias benchmarks (as language models, unlike is the case for language modeling tasks like coreference resolution, are not trained with the express purpose of performing well on bias benchmarks), it is similarly important to assess the extent to which bias benchmarks can assess subtler levels of model bias than the ones current language models display.

Finally, Rodriguez et al. (2021) recently argued that traditional NLP *leaderboards*—websites where the best-performing language models’ test scores on benchmark datasets are depicted—should be “re-imagined” as Bayesian IRT models. While the authors do not explicitly discuss performance on NLP bias benchmarks, some of their critiques of traditional leaderboards are relevant to performance metrics of bias benchmarks, and their discussion ties in with our critique of such metrics in Sect. 5.1.

7 Discussion

Given their wide application in today’s connected society, assessing the extent to which language models are biased is crucial: As many ordinary citizens interact with applications (like translators) that rely on language models, even (seemingly) small undesired trends in language models’ behaviors can, on aggregate, cause substantial harm. High-quality measurement tools would allow us to diagnose biased language models before they are implemented and can cause harm. Legislators could make use of such high-quality measurement tools to implement guidelines for acceptable language model behavior. Finally, high-quality bias measurement tools are crucial because they enable systematic bias mitigation efforts (e.g., trying whether an intervention leads to a decrease in measured language model bias).

Contemporary bias measurement tools do not yet live up to the bar of “high quality” (with flaws highlighted in e.g., Blodgett et al., 2020; Blodgett et al., 2021; Goldfarb-Tarrant et al., 2021). Consequently, many attempts have been made to improve the quality of existing and newly-introduced bias measurement tools. Due to NLP’s unfamiliarity with issues—like “how does one measure concepts that cannot be directly observed?”—that have long research traditions in other research fields, we believe that interdisciplinary work is best equipped to address the NLP bias case (van der Wal et al., 2024). Here, we provide such interdisciplinary work, combining NLP expertise with insights from the psychometric framework of item response theory.

In our view, a systematic application of IRT concepts to bias benchmark datasets could prove crucial. It would allow us to assess the quality of individual test items as well as the properties of the benchmark dataset as a whole. If a bias benchmark, unbeknownst to us, only assesses a small range of model bias severities, conclusions that we draw based on scores on that measure could mislead. Language models that were declared safe (e.g., because there are small differences in accuracies across

stereotypical and counter-stereotypical sentences) could, in fact, be biased (e.g., because most items of the benchmark are of high blatancy) and cause harm to ordinary citizens, when implemented. A test information curve, which shows that the benchmark only provides us information about a small range of model bias severities, could make us aware of such deficiencies. Similarly, the simultaneous analysis of items across bias benchmark datasets would allow us to assemble a combined benchmark dataset that consists of high-discrimination test items (from several different bias benchmarks) which are spread across a range of different item blatancies. Finally, we argued that IRT models' trait level estimates are conceptually superior to traditional evaluation metrics (such as accuracy, F1 score, recall, etc.), because the IRT models work with more information about test items than standard measures do and because they provide us with statistics to conceptualize our (un-)certainty about the trait level estimates (i.e., standard errors and confidence intervals).

While we are enthusiastic about the prospect of applying IRT to the NLP bias, we anticipate a few potential challenges to the implementation. The first notable limitation is a conceptual one: IRT models benefit from being fit on data from testtakers with a range of different trait levels (e.g., ideally data from language models from all sections of the model biasedness distribution). In previous applications of IRT to language models, variety in trait levels could be ensured artificially: As language model capacity on tasks such as coreference resolution depends on the amount of training the models received, systematic manipulations between language models of the amount of training data they received ensured a healthy range of trait levels. The situation is less simple for the model bias case. In the absence of trustworthy model bias measures, little is known about the distribution of model bias severities (e.g., whether large and low model bias severities are equally common) or about the shifts in model bias across training time (e.g., unlike what is the case for coreference resolution, we cannot simply vary between language models the training durations and assume that this will lead to a spread in trait levels). In our view, all we can do is simultaneously assess a large variety of different language models and hope that there are no striking clusters of bias trait levels (e.g., a large number of very high and/or very low bias levels instead of the trait levels being somewhat evenly distributed). Were we to find such clusters (which itself would be a theoretically insightful finding), we would need to be careful not to overgeneralize our findings (e.g., not to make strong statements about unobserved bias severities).

A second conceptual challenge regarding the application of IRT to NLP bias is the high number of testtakers from which we require response data. How many (test-taker) language models we require data from remains to be seen. For instance, 3PL models, that allow for estimating the (pseudo-)guessing parameter, which we found of particular interest, generally require a lot of data; for example, for 60 test items, Akour and Al-Omari (2013) estimate that 500 testtakers are required to receive sufficiently precise item parameter values.¹⁶ Additionally, common NLP ways of bolstering the number of (assessed) language models (e.g., testing the same models repeatedly at different times during training) could prove problematic: They might prove

¹⁶ Here the aforementioned fitting of simpler IRT models with lower asymptotes could prove an important (i.e., lower-sample size) alternative to fully fitting 3PL models.

inapplicable to the bias case (depending on how bias develops across model training—a topic of current research) and they might lead to non-generalizable results (e.g., because the language models that are assessed in the IRT model are nonrepresentative of language models more generally). That being said, it is currently unclear what number of different language models will be required to achieve adequate IRT model fit—a deficit we hope to see addressed in (e.g., simulation) studies in the future. Provided that testtakers' trait levels are decently spread, previous studies have found that the required number of testtakers reduces with increasing numbers of test items that are modeled (Şahin & Anıl, 2017; Akour & AL-Omari, 2013). However, these studies did not assess test sizes that resemble those of bias benchmarks (e.g., Akour and AL-Omari's largest test size was 60, while StereoSet contains more than 4000 items). Consequently, we are optimistic that the number of different language models required to fit IRT models is substantially lower than what analyses of human tests suggest.

7.1 Looking Ahead

Our goal was to provide an initial introduction to and conceptual discussion of basic IRT concepts that are relevant to the NLP bias case. Hence, the scope of our discussion was limited. There are several IRT concepts that we consider relevant to the NLP bias measurement case, but that would have required more statistical and conceptual discussion than we could justify here. Three such concepts that we wish to mention here are multidimensional IRT models, hierarchical IRT models, and computerized adaptive testing.

Some authors (e.g., Brown et al., 2020; D'Amour et al., 2020; Webster et al., 2021) have put forth the argument that language models become less biased with increasing language capacity (e.g., because early in training they rely on stereotypes to make predictions, while these stereotypes become obsolete once the language model gains in language capacity). More generally, a potential conceptual weakness of NLP bias benchmarks is that they are not pure measures of NLP bias, but instead assess both bias and the primary language modeling task that indirectly reveals this bias (e.g., for WinoGender: coreference resolution). This conceptual weakness could potentially be addressed and the aforementioned argument could potentially be investigated by use of multi-trait (or “multidimensional”) IRT models (Bonifay, 2019; DeMars, 2013). These models allow for the simultaneous modeling of different performance-causing traits, allowing us to quantify the influence of either of these traits (e.g., model bias and coreference resolution capacity) on responses to test items.

Moreover, hierarchical IRT models (e.g., Cai et al., 2016; Rijmen, 2011) are potentially relevant, because they allow us to model dependencies between testtakers: cases in which one testtaker's answers to test items tells us something about responses of at least one other testtaker (e.g., on a national exam, a student is likely to perform somewhat similarly to their classmate, as the two share teacher and school resources). We believe that such modeling techniques are relevant to the NLP bias case, because common ways of increasing the numbers of testtaker language

models (e.g., testing the same model at different times during training; Vania et al., 2021; or training many instances of the same model architecture; Byrd and Srivastava, 2022) could introduce such dependencies.

Additionally, we think that such dependencies might become even more prevalent, given current trends in NLP: Some modern language models are so powerful, that they can take on different “roles” (e.g., a chatbot can be made to answer questions as a therapist or teacher would), when provided with different sequences of instructions and rules (called prompts; see e.g., White et al., 2023). While these (prompted) different “personas” of the language models can act substantially differently, they still stem from the same language model. If an IRT analysis of their behaviours becomes relevant (e.g., if companies implement differently prompted versions of their language model for different purposes), we would expect that differently prompted models that stem from the same language model have (somewhat) dependent observations: Independently of the prompt (e.g., be it that the model was told to act like a teacher or like a therapist), its fundamental language capacity (i.e., which language regularities the model picked up on; see Sect. 2.1) is identical.

The final IRT concept we wish to highlight here is *computerized adaptive testing* (CAT; see e.g., Klinkenberg et al., 2011; Magis et al., 2017). CAT is a method of testing that involves a sequential back-and-forth between a) estimating a testtaker’s trait level (based on their performance on previous test items) and b) selecting and administering a new test item, based on the trait level estimate. For example, the testing program could select as its first test item an item that is maximally informative at $\theta = 0$ and, depending on the testtaker’s answer, (re-)estimate the testtaker’s trait level as slightly higher or lower than $\theta = 0$. Subsequently, the test program would choose and administer an item that is maximally informative at its new estimate of the testtaker’s ability level. This process is repeated until the CAT algorithm arrives at a precise estimate of the testtaker’s trait level (e.g., until the standard error of a testtaker’s trait level estimate shrinks below a predetermined value).

One of the main advantages of CAT is that fewer test items are required to receive a good estimate of a testtaker’s ability. For example, in a preset collection of items, it is likely that several items have substantially too high or too low difficulty (or here: blatancy) to be informative about a testtaker. Fewer such uninformative items are administered, if items are chosen specifically for a testtaker’s (likely) trait level. Such efficiency gains, while less important for the assessment of individual language models (since language models can answer hundreds of items quickly, we can afford a few uninformative items), could prove relevant if we want to test a large number of models and/or if we want to repeatedly test language models, for example at different training steps (as did e.g., Lalor & Yu, 2020). We have previously argued (van der Wal et al., 2024) that such assessments across training steps are conceptually important, as they provide information about how model bias emerges and changes over the course of model training. This could provide us with hints about the causes of bias: For example, if gender bias emerges relatively early in training and then decreases in subsequent training steps, it might be that stereotypical gender associations (e.g., “nurse = female”) initially represent good heuristics for task performance that later—as the model improves in the task—outlive their usefulness.

We believe that CAT could play a role in making such repeated assessments of multiple language models.

More generally, we hope that our introduction to IRT concepts encourages readers to explore the large and varied literature of item response theory to find inspiration for bias measure improvement efforts—be they inspired by concepts that we introduced here, or by others. This article is meant to be a mere starting point for applying IRT concepts to the NLP bias case. While we expressed our view on how such concepts could be used towards our ultimate goal—reducing the harms that language models cause—we look forward to hearing perspectives that alternate from ours or expand beyond it.

Acknowledgements The authors wish to thank Petr Palíšek, Alina Leidinger, and the two anonymous peer reviewers for their thoughtful and insightful feedback! This publication is part of the project “The biased reality of online media—Using stereotypes to make media manipulation visible” (Project number 406.DI.19.059) of the research programme Open Competition Digitalisation-SSH, which is financed by the Dutch Research Council (NWO).

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article’s Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article’s Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

References

- Akour, M., & Al-Omari, H. (2013). Empirical investigation of the stability of IRT item-parameters estimation. *International Online Journal of Educational Sciences*, 5(2), 291–301.
- Amidei, J., Piwek, P., & Willis, A. (2020). Identifying annotator bias: a new IRT-based method for bias identification. In *Proceedings of the 28th international conference on computational linguistics* (pp. 4787–4797). <https://aclanthology.org/2020.coling-main.421/>
- Anunciacao, L. (2018). An overview of the history and methodological aspects of psychometrics: History and methodological aspects of psychometrics. *Journal for ReAttach Therapy and Developmental Diversities*, 1(1), 44–58.
- Beeching, E., Fourier, C., Habib, N., Han, S., Lambert, N., Rajani, N., Sanseviero, O., Tunstall, L., & Wolf, T. (2023). Open LLM leaderboard. https://huggingface.co/spaces/HuggingFaceH4/open_llm_leaderboard
- Blodgett, S. L., Barocas, S., Daumé III, H., & Wallach, H. (2020). Language (technology) is power: a critical survey of “bias” in NLP. In *Proceedings of the 58th annual meeting of the association for computational linguistics* (pp. 5454–5476). <https://doi.org/10.18653/v1/2020.acl-main.485>
- Blodgett, S. L., Lopez, G., Olteanu, A., Sim, R., & Wallach, H. (2021). Stereotyping Norwegian Salmon: an inventory of pitfalls in fairness benchmark datasets. In *Proceedings of the 59th annual meeting of the association for computational linguistics and the 11th international joint conference on natural language processing* (Volume 1: Long Papers, pp. 1004–1015). <https://aclanthology.org/2021.acl-long.81>
- Bommasani, R., Hudson, D. A., Adeli, E., Altman, R., Arora, S., von Arx, S., Bernstein, M. S., Bohg, J., Bosselut, A., Brunskill, E., Brynjolfsson, E., Buch, S., Card, D., Castellon, R., Chatterji, N., Chen, A., Creel, K., Davis, J. Q., Demszky, D., & Liang, P. (2022). On the opportunities and risks of foundation models. [arxiv:abs/2108.07258](https://arxiv.org/abs/2108.07258)

- Bommasani, R., & Liang, P. (2022). Trustworthy social bias measurement. [arxiv:abs/2212.11672](https://arxiv.org/abs/2212.11672)
- Bonifay, W. (2019). *Multidimensional item response theory*. Sage Publications.
- Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., & Amodei, D. (2020). Language models are few-shot learners. [arxiv:abs/2005.14165](https://arxiv.org/abs/2005.14165)
- Byrd, M., & Srivastava, S. (2022). Predicting difficulty and discrimination of natural language questions. In *Proceedings of the 60th annual meeting of the association for computational linguistics* (Volume 2: Short Papers, pp. 119–130). <https://aclanthology.org/2022.acl-short.15/>
- Cai, L., Choi, K., Hansen, M., & Harrell, L. (2016). Item response theory. *Annual Review of Statistics and Its Application*, 3, 297–321. <https://doi.org/10.1146/annurev-statistics-041715-033702>
- Chalmers, R. P. (2012). mirt: A multidimensional item response theory package for the R environment. *Journal of Statistical Software*, 48, 1–29. <https://doi.org/10.18637/jss.v048.i06>
- Crocker, L., & Algina, J. (1986). *Introduction to classical and modern test theory*. ERIC.
- D'Amour, A., Heller, K., Moldovan, D., Adlam, B., Alipanahi, B., Beutel, A., Chen, C., Deaton, J., Eisenstein, J., Hoffman, M. D., Hormozdiari, F., Houlby, N., Hou, S., Jerfel, G., Karthikesalingam, A., Lucic, M., Ma, Y., McLean, C., Mincu, D., & Sculley, D. (2020). Underspecification presents challenges for credibility in modern machine learning. [arxiv:abs/2011.03395v2](https://arxiv.org/abs/2011.03395v2)
- De Ayala, R. J. (2013). *The theory and practice of item response theory*. Guilford Publications.
- Debelak, R., Strobl, C., & Zeigenfuse, M. D. (2022). *An introduction to the Rasch model with examples in R*. CRC Press.
- DeMars, C. E. (2013). A tutorial on interpreting bifactor model scores. *International Journal of Testing*, 13(4), 354–378. <https://doi.org/10.1080/15305058.2013.799067>
- Dhamala, J., Sun, T., Kumar, V., Krishna, S., Pruksachatkun, Y., Chang, K.-W., & Gupta, R. (2021). BOLD: dataset and metrics for measuring biases in open-ended language generation. In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency* (pp. 862–872). <https://dl.acm.org/doi/10.1145/3442188.3445924>
- Du, Y., Fang, Q., & Nguyen, D. (2021). Assessing the reliability of word embedding gender bias measures. In *Proceedings of the 2021 conference on empirical methods in natural language processing* (pp. 10012–10034). <https://aclanthology.org/2021.emnlp-main.785>
- Ethayarajh, K. (2020). Is your classifier actually biased? Measuring fairness under uncertainty with Bernstein bounds. In *Proceedings of the 58th annual meeting of the association for computational linguistics* (pp. 2914–2919). <https://www.aclweb.org/anthology/2020.acl-main.262>
- Fang, Q., Oberski, D. L., & Nguyen, D. (2024). PATCH—psychometrics-assisted benchmarking of large language models: A case study of mathematics proficiency. [arxiv:abs/2404.01799](https://arxiv.org/abs/2404.01799)
- Fox, J.-P. (2010). *Introduction to Bayesian response modeling*. Springer.
- Furr, R. M. (2021). *Psychometrics: An introduction* (4th ed.). SAGE Publications.
- Gilholm, P., Mengersen, K., & Thompson, H. (2021). Bayesian hierarchical multidimensional item response modeling of small sample, sparse data for personalized developmental surveillance. *Educational and Psychological Measurement*, 81(5), 936–956. <https://doi.org/10.1177/0013164420987582>
- Goldfarb-Tarrant, S., Marchant, R., Muñoz Sánchez, R., Pandya, M., & Lopez, A. (2021). Intrinsic bias metrics do not correlate with application bias. In *Proceedings of the 59th annual meeting of the association for computational linguistics and the 11th international joint conference on natural language processing* (Volume 1: Long Papers, pp. 1926–1940). <https://aclanthology.org/2021.acl-long.150>
- Hladká, A., & Martinková, P. (2020). difNLR: Generalized logistic regression models for DIF and DDF detection. *R Journal*, 12(1), 300–323.
- Jacobs, A. Z., & Wallach, H. (2021). Measurement and fairness. In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency* (pp. 375–385). <https://doi.org/10.1145/3442188.3445901>
- Jurafsky, D., & Martin, J. H. (2023). *Speech and language processing* (3rd edn.). <https://web.stanford.edu/~jurafsky/slp3/>
- Klinkenberg, S., Straatemeier, M., & van der Maas, H. L. (2011). Computer adaptive practice of maths ability using a new item response model for on the fly ability and difficulty estimation. *Computers & Education*, 57(2), 1813–1824. <https://doi.org/10.1016/j.compedu.2011.02.003>

- König, C., Spoden, C., & Frey, A. (2020). An optimized Bayesian hierarchical two-parameter logistic model for small-sample item calibration. *Applied Psychological Measurement, 44*(4), 311–326. <https://doi.org/10.1177/0146621619893786>
- Lalor, J. P., Wu, H., & Yu, H. (2016). Building an evaluation scale using item response theory. In *Proceedings of the conference on empirical methods in natural language processing. Conference on empirical methods in natural language processing* (p. 648). <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5167538/>
- Lalor, J. P., & Yu, H. (2020). Dynamic data selection for curriculum learning via ability estimation. In *Proceedings of the conference on empirical methods in natural language processing. conference on empirical methods in natural language processing* (p. 545). <https://aclanthology.org/2020.findings-emnlp.48/>
- Levy, S., Lazar, K., & Stanovsky, G. (2021). Collecting a large-scale gender bias dataset for coreference resolution and machine translation. *Findings of the association for computational linguistics: EMNLP 2021* (pp. 2470–2480). <https://aclanthology.org/2021.findings-emnlp.211>
- Lord, F. M. (1980). *Applications of item response theory to practical testing problems*. Routledge.
- Magis, D., Yan, D., & von Davier, A. A. (2017). *Computerized adaptive and multistage testing with R: Using packages catR and mstR*. Springer.
- Martinková, P., & Hladká, A. (2023). *Computational aspects of psychometric methods: With R*. CRC Press.
- Nadeem, M., Bethke, A., & Reddy, S. (2021). StereoSet: Measuring stereotypical bias in pretrained language models. In *Proceedings of the 59th annual meeting of the association for computational linguistics and the 11th international joint conference on natural language processing* (Volume 1: Long Papers, pp. 5356–5371). <https://aclanthology.org/2021.acl-long.416>
- Nangia, N., Vania, C., Bhalarao, R., & Bowman, S. R. (2020). CrowS-Pairs: a challenge dataset for measuring social biases in masked language models. In *Proceedings of the 2020 conference on empirical methods in natural language processing (EMNLP)* (pp. 1953–1967). <https://aclanthology.org/2020.emnlp-main.154>
- Paek, I., & Cole, K. (2019). *Using R for item response theory model applications*. Routledge.
- Parmar, M., Mishra, S., Geva, M., & Baral, C. (2023). Don't blame the annotator: Bias already starts in the annotation instructions. arxiv.org/abs/2205.00415
- Parrish, A., Chen, A., Nangia, N., Padmakumar, V., Phang, J., Thompson, J., Htut, P. M., & Bowman, S. R. (2021). BBQ: A hand-built bias benchmark for question answering. arxiv.org/abs/2110.08193v2
- Polo, F. M., Weber, L., Choshen, L., Sun, Y., Xu, G., & Yurochkin, M. (2024). tinyBenchmarks: Evaluating LLMs with fewer examples. arxiv.org/abs/2402.14992v1
- Rijmen, F. (2011). Hierarchical factor item response theory models for PIRLS: Capturing clustering effects at multiple levels. *IERI Monograph Series: Issues and Methodologies in Large-scale Assessments, 4*, 59–74.
- Rodriguez, P., Barrow, J., Hoyle, A. M., Lalor, J. P., Jia, R., & Boyd-Graber, J. (2021). Evaluation examples are not equally informative: How should that change NLP leaderboards? In *Proceedings of the 59th annual meeting of the association for computational linguistics and the 11th international joint conference on natural language processing* (Volume 1: Long Papers, pp. 4486–4503). <https://aclanthology.org/2021.acl-long.346>
- Rodriguez, P., Htut, P. M., Lalor, J. P., & Sedoc, J. (2022). Clustering examples in multi-dataset benchmarks with item response theory. *Proceedings of the Third Workshop on Insights from Negative Results in NLP* (pp. 100–112). <https://aclanthology.org/2022.insights-1.14/>
- Rudinger, R., Naradowsky, J., Leonard, B., & Van Durme, B. (2018). Gender bias in coreference resolution. *Proceedings of the 2018 conference of the North American chapter of the association for computational linguistics: human language technologies* (Volume 2 (Short Papers), pp. 8–14). <https://aclanthology.org/N18-2002>
- Şahin, A., & Anıl, D. (2017). The effects of test length and sample size on item parameters in item response theory. *Educational Sciences: Theory & Practice, 17*(1), 321–335.
- Samejima, F. (1969). Estimation of latent ability using a response pattern of graded scores. *Psychometrika, 34*, 1–97. <https://doi.org/10.1007/BF03372160>
- Sap, M., Card, D., Gabriel, S., Choi, Y., & Smith, N. A. (2019). The risk of racial bias in hate speech detection. In *Proceedings of the 57th annual meeting of the association for computational linguistics* (pp. 1668–1678). <https://aclanthology.org/P19-1163>

- Stanovsky, G., Smith, N. A., & Zettlemoyer, L. (2019). Evaluating gender bias in machine translation. *Proceedings of the 57th annual meeting of the association for computational linguistics* (pp. 1679–1684). <https://aclanthology.org/P19-1164>
- van der Wal, O., Bachmann, D., Leidinger, A., van Maanen, L., Zuidema, W., & Schulz, K. (2024). Undesirable biases in NLP: Addressing challenges of measurement. *Journal of AI Research*, 79, 1–40. <https://doi.org/10.1613/jair.1.15195>
- Vania, C., Htut, P. M., Huang, W., Mungra, D., Yuanzhe Pang, R., Phang, J., Liu, H., Cho, K., & Bowman, S. R. (2021). Comparing test sets with item response theory. [arxiv:abs/2106.00840](https://arxiv.org/abs/2106.00840)
- Warm, T. A. (1978). *A primer of item response theory* (tech. rep.). Coast Guard Washington DC. <https://files.eric.ed.gov/fulltext/ED171730.pdf>
- Webster, K., Wang, X., Tenney, I., Beutel, A., Pitler, E., Pavlick, E., Chen, J., Chi, E., & Petrov, S. (2021). Measuring and reducing gendered correlations in pre-trained models. [arxiv:abs/2010.06032](https://arxiv.org/abs/2010.06032)
- White, J., Fu, Q., Hays, S., Sandborn, M., Olea, C., Gilbert, H., Elnashar, A., Spencer-Smith, J., & Schmidt, D. C. (2023). A prompt pattern catalog to enhance prompt engineering with ChatGPT. [arxiv:abs/2302.11382](https://arxiv.org/abs/2302.11382)
- Wu, M., Tam, H. P., & Jen, T.-H. (2016). Differential item function. In *Educational measurement for applied researchers: Theory into practice* (pp. 207–225). Springer. https://doi.org/10.1007/978-981-10-3302-5_11
- Zhang, H., Sneyd, A., & Stevenson, M. (2020). Robustness and reliability of gender bias assessment in word embeddings: the role of base pairs. In *Proceedings of the 1st conference of the Asia-Pacific chapter of the association for computational linguistics and the 10th international joint conference on natural language processing* (pp. 759–769). <https://aclanthology.org/2020.aacl-main.76>
- Zhao, J., Wang, T., Yatskar, M., Ordonez, V., & Chang, K.-W. (2018). Gender bias in coreference resolution: evaluation and debiasing methods. In *Proceedings of the 2018 conference of the North American*

Authors and Affiliations

Dominik Bachmann^{1,2} · **Oskar van der Wal**¹ · **Edita Chvojka**^{3,4} · **Willem H. Zuidema**¹ · **Leendert van Maanen**² · **Katrin Schulz**¹

✉ Dominik Bachmann
d.bachmann@uva.nl; dominik.bachmann.psychology@gmail.com

¹ Institute for Logic, Language and Computation, University of Amsterdam, Amsterdam, The Netherlands

² Department of Experimental Psychology, Utrecht University, Utrecht, The Netherlands

³ Department of Methodology and Statistics, Utrecht University, Utrecht, The Netherlands

⁴ Department of Interdisciplinary Social Science, Utrecht University, Utrecht, The Netherlands

can chapter of the association for computational linguistics: human language technologies (Volume 2 (Short Papers), pp. 15–20). <https://aclanthology.org/N18-2003>

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.