

Consequences of Measurement Instability on the Mapping of Real Quantities to Rating Scale Scores

Maas van Steenbergen^{1,*}

¹ Faculty of Behavioural and Social Sciences, Methodology & Statistics, Utrecht University, the Netherlands

Correspondence*:
Corresponding Author
m.vansteenbergen@uu.nl

“*Eh bien!*”—exclaimed Walras characteristically—“this difficulty is not insurmountable. Let us suppose that this measure exists, and we shall be able to give an exact and mathematical account of it”. [...] In view of the fact that theoretical science is a living organism, it would not be exaggerating to say that this attitude is tantamount to planning a fish hatchery in a moist flower bed”.

— *Nicholas Georgescu-Roegen*

“I think the foundations of measurements [...] have a lot of implications about the way you do and actually think about measurement. There is a great deal of feedback from work on foundations on the actual practice, unlike a lot of other fields of mathematics, where work on foundations is divorced from the practice of other disciplines”.

— *Amon Tversky*

1 INTRODUCTION

Variables in psychology often aim to quantify strength of sensation, personality traits, or attitudes. Because of the inherently subjective nature of these variables, we cannot measure or validate those concepts without research participants communicating it to us. In other words, we cannot look in their heads and experience what they experience. This makes it more challenging to assess such variables and their underlying structure. Over the course of time, psychologists developed tools to quantitatively estimate the value of such variables. The rating scale is the most popular technique by far. This family of instruments always involves a series of ordered options that represent an increasing magnitude on the variable of interest. There is a great variety in the implementation of these instruments, but the central logic remains the same.

Yet, what the numbers that are generated through this method mean has always been unclear. Over the course of time, researchers have become accustomed to a series of prescriptions to ensure that the result of the instrument is valid. Some say you should treat these variables as continuous interval-level measurements. Others claim that this is the wrong way to study these phenomena, and prescribe non-parametric methods instead (?). For some, parametric methods are only recommended after meeting certain criteria. A popular prescription, for example, is that you can treat an ordinal variable as continuous as long as you have more than a certain number of answering options (?). Another popular prescription is that ordinal variables become continuous after summing a sufficient number of them if they measure the same construct.

We intend to study how different assumptions about the mapping of the quantitative structure of a variable on a rating scale influence the number of false positives. We do this by imagining a perfect (ratio-level) continuous homomorphic mapping of a *quantifiable* psychological variable. Then, we specify a certain structural relationship between the perfectly measured variable and an ordinal equivalent with a limited number of equivalence classes. We use a series of increasingly loose restrictions to the mapping, studying how the original variables are deformed when looser restrictions are assumed. Then, we will administer some simple statistical tests to see the degree to which the conclusions based on the ‘perfect’ homomorphic mapping of the score correspond to ratings on the rating scale. We use this to argue that one clearly needs to define how the test relates to a hypothetical zero-point and to the stability of the mapping. We also study the influence and make a case for the importance of stability, scaling, and zero point for psychological tests. Finally, we make clear that this is necessarily a question of *degree*. We show that stronger or weaker violations of these assumptions are possible.

For the current project, we focus our attention on cases where an attribute of interest is principally continuously measurable. Therefore homomorphic mapping of the qualitative structure towards a ratio-scale quantitative representation can be constructed. While there might be practical limitations in the construction of such a mapping, this is of no consequence to our experimental set-up: we posit to know the true value, and use this as a starting point for a simulation of the consequences. Then, we do an extra mapping: from the perfect measure to the scale. We map this quantitative representation again for each individual rating scale representation, adding some difference. We call the differences in mapping *measurement instability* and make a distinction between four kinds. We add *zero point instability*, which adds noise to the zero-point. We add *scaling instability* which relates to a changing range of the instrument. We add and the relative position of the thresholds. Finally, we perform some commonly, and test the robustness of these measures

1.1 Representational Theory of Measurement

We assume, in alignment with Krantz, Luce, Suppes, and Tversky (henceforth referred to as KLST), that an attribute can be quantified only if it is proven to be a quantity (). We will use the criteria formalized in the representational theory of measurement (RTM). KLST is at present the most fully developed account of RTM. In line with this theory, we see measurement from the perspective of a representation theorem and a uniqueness theorem. The representation theorem asserts that a set, together with one or more relations on that set, follows certain axioms so that it can truthfully preserve the structure of qualitative attributes of an object to a numerical representation of that quality. In other words, a *homomorphic* mapping of the qualitative relational structure into a quantitative representation can be constructed.

Suppose that the transformation function $\phi(a)$ is a homomorphism and the representation theorem holds for the mapping of qualitative ordering \succ to quantitative ordering $>$. If $\phi(a)$ then maps a qualitative attribute of set A into \mathbb{R} , it should maintain the structure of the qualitative attribute in its numerical representation: if and only if a is qualitatively greater than b , it is numerically greater than b . The uniqueness theorem specifies the permissible homomorphic transformations of A that lead to the same structure of numerical relations. E.g., both Fahrenheit or Celsius measurements are legitimate homomorphic mappings that maintain the qualitative structure of temperature, differing in zero-point and scale. Note that these aspects of measurement do not concern modelling or statistical analysis, but are a *necessary condition* for those.

We align ourselves with “representational minimalism” that rids it of most of its epistemological basis so that it can serve as a common ground for discussion about measurement, recognizing that most of the critical literature about measurement takes place either within or in discussion of this framework (?).

The acceptance of RTM is not universal. We do not make firm claims whether a realist, an operationalist, or a representationalist epistemology should be matched to the formal aspects of the framework. We do, however, choose to use representationalist language, as this is the most well-known and fully realized. It should be noted that choosing a different framework means that the formalism might either change slightly or majorly. If this change is subtle, then the current study might still be interpretable. In some frameworks, however, especially those that doubt the hypothetical quantifiability of psychological constructs, the current discussion might be uninterpretable.

1.1.1 Ordinal measurement

An ordinal score would imply that the score is only dependent on rank order, and does not support concatenation operations. A score in P can only indicate that a score is higher or lower than another score in P . Based on the measure only, beyond that ordering, we know nothing about P . In other words, any mapping of which we can only assume that these characteristics hold is an ordinal mapping of the relational structure of Q . Note that an ordinal mapping of the structure of Q can be made if a ratio-mapping can also be made, but not vice versa.

This can also be written down axiomatically. Let x, y & z by any values of quantity Q . Then the result of the projection of the quantity to an ordinal scale, which we will refer to as P , holds to the following conditions. Let A be a set and \succeq be a binary relation on A .

- if $X \succeq Y$ & $Y \succeq Z$, then $X \succeq Z$ (transitivity);
- Either $X \succeq Y \parallel Y \succeq X$ (connectedness);
- For strongly tied order: If $X \succeq Y$ & $Y \succeq X$, then $X = Z$ (antisymmetry).

Two elements have a weak order iff, for all a, b , axiom 1 and 2 are met. A set of all tied elements with the same value is called an equivalence set. If the third axiom is also met, then the ordering is strong and elements cannot be tied. Because the likert scale obviously allows for ties (multiple people can be ranked on the same score), and is thus a weak ordering, we assume transitivity and connectedness for the mapping. Transitivity implies that all order-relations need to be consistent. Connectedness implies that a connection between two elements is either larger, smaller, or equal. If antisymmetry is met, equal scores are only possible if they refer to the same object.

1.1.2 Ratio measurement

As for the ‘perfect’ ratio-mapping of the principally measurable attribute Q we will assume that the following additional characteristics will hold above axiom 1 and 2 presented above (?):

- $X \oplus (Y \oplus Z) = (X \oplus Y) \oplus Z$ (associativity);
- $X \oplus Y = Y \oplus X$ (commutativity);
- $X \succeq Y$ iff $X \oplus Z \succeq Y \oplus Z$ (monotonicity);
- if $X \succ Y$ then there exists a Z such that $X = Y \oplus Z$ (solvability);
- $X \oplus Y \succ X$ (positivity);
- there exists a number n such that $nX \succeq Y$ (where $1X = 1$ and $(n \oplus 1)X = nX \oplus X$) (Archimedean condition).

This essentially means that a value q can always be put in terms of another value r in Q . Every ratio-scale is homomorphic to the qualitative attribute that is being measured, and ratio-scales are thus homomorphic

110 to each other (?). The last axiom (the archimedean condition) is added to ensure that the set of possible
 111 scores is finite. E.g., say that we have developed a standard measure for happiness: H . Then we can say
 112 that any value x is written in terms of H iff the finite ratio $\frac{x}{H}$ holds¹.

113 1.1.3 Rating scales and ordinality

114 Likert scale ratings are normally understood to be ordinal, but we argue that this is not unambiguously
 115 true if there is between-person or within-person scaling instability. If the ordering of attributes or their
 116 relationship to the hypothetical perfect measure is not invariant to time and is not invariant between people,
 117 then the relationship of the attribute to the rating-scale representation of that attribute is not consistently
 118 weakly ordered. To show this, we first need to formalize the connection between the hypothetical perfect
 119 measure and a measure that is ordinal for one specific person at one specific time.

120 Assume that the empirical structure of an attribute is in correspondence with the axioms for ratio-level
 121 measurement scaling. This means, by extension, that all axioms for ordinal measurement are in correspon-
 122 dence with the empirical structure of that attribute. Let us assume we have a hypothetical homomorphic
 123 mapping of the empirical structure of this attribute expressed quantitatively (with ratio-scaling). Then, each
 124 possible weakly² ordinal representation of quantitative values of this perfect measurement measure can
 125 result in an unambiguous weakly ordered classification of scores in Q to scores in terms of P ³. Assume
 126 that instrument P has n elements. Then we can define a series of $n - 1$ ‘thresholds’ R as $\{a, b, c \dots z\}$.
 127 These thresholds are the point of q where a score p_1 in P jumps to another classification p_2 in P on the
 128 basis of a score q in Q . We then define a score p in P for each q in Q as follows, assuming we have an
 129 n -level measurement scale for P :

$$\begin{cases} 1 & q \leq a \\ 2 & a \leq q < b \\ 3 & b \leq q < c \\ \dots & \dots \\ n & z < q \end{cases}$$

130 Note that if the mapping is inconsistent with this formalism, then the first axiom for ordinality is broken.
 131 Thus, this step function is a logical consequence of the empirical structure of the variable. It needs to be
 132 met. Otherwise, the two representations of the attribute are mutually inconsistent. This representation
 133 follows from the assumed structure of our hypothetical perfect representation of the variable and the ordinal
 134 rating scale.

135 While a strong ordering of an ordinal representation of the variable is a unique representation of the
 136 variable up to monotonic transformations, a representation that groups values into a certain number of
 137 equal segments can be inconsistent with another representation of that variable that splits the range into
 138 an identical number of equivalence sets. In fact, many different constructions can be made that are not

¹ We agree with the point by Franz that it is misleading to give physical magnitudes as examples when discussing psychological phenomena (?). Therefore, we choose to use psychological examples when we are thinking about psychological phenomena. This adds an extra face validity check: does discussing psychological attributes in this manner make sense?

² Meaning elements can be tied, see above.

³ This lack of ambiguity is only a result of this formalization. We may find that the relationship of the ‘real’ variable to its ordinal estimate is not quite so straightforward. We will discuss this later on.

139 compatible with each other. Say we want to map a continuous variable into five categories. The values of
140 our continuous variable fall in range [0, 100] and we would like to map this range to five equivalence sets.

141 If so, then both

$$\left\{ \begin{array}{l} 1 \quad q \leq 22.4 \\ 2 \quad 22.4 \leq q < 48.9 \\ 3 \quad 48.9 \leq q < 59.4 \\ 4 \quad 59.4 \leq q < 82.2 \\ 5 \quad 82.2 < q \end{array} \right. \& \left\{ \begin{array}{l} 1 \quad q \leq 20.0 \\ 2 \quad 20.0 \leq q < 40.0 \\ 3 \quad 40.0 \leq q < 60.0 \\ 4 \quad 60.0 \leq q < 80.0 \\ 5 \quad 80.0 < q \end{array} \right. .$$

142 are valid representations of the variable with 5 equivalence sets. Yet, they are inconsistent, because a
143 quantity can be sent by the mapping to one category that would be sent to another category in another
144 mapping. E.g., the value 21.1 is seen as a 1 in the first representation while a value of 21 is seen as a
145 2 in the second representation. This breaks the axiom of transitivity. If we were to sum scores gained
146 through multiple ordinal representations of the same score, the mapping can become inconsistent between
147 these representations. Note also that this is not measurement error per sé, but another type of error. The
148 measurement itself can be seen as correct, but the issues come to be through the incompatibility of one
149 measurement with another.

150 1.2 Between-person mapping inconsistencies

151 It is often assumed that aggregations of individual measures rescind the limitations of different scaling
152 because these measurement differences cancel out when enough samples are drawn. The origins of this
153 assertion are from Knapp. This idea can be seen as follows: if the limit of the number of bins becomes
154 infinite, and the underlying variable is continuous, then the measuring instrument will approach the normal
155 distribution.

156 1.3 Within-person mapping inconsistencies

157 When studying within-person mappings, we can assume that those psychological processes part of Q are
158 time-dependent. These measures are shaped by different forces and the previous states of the construct (?
159). Their values are continuous and are related to each other in a structured manner (?). We also assume
160 that their values are differentiable over time, changing smoothly. Their values can increase or increase
161 very quickly, but not instantaneously. This, they are modeled best using the rate of change of the variable,
162 through differential equations (?). We assume that the mapping function of a variable develops in the same
163 way: continuously, changing smoothly, and differentiable over time.

164 Making a, b, c, and d functions of time, step function F(x, t) becomes:

$$\left\{ \begin{array}{l} 1 \quad q \leq a(t) \\ 2 \quad a(t) \leq q \leq b(t) \\ 3 \quad b(t) \leq q \leq c(t) \\ \dots \quad \dots \\ n \quad z(t) \leq q \end{array} \right. .$$

165 This means that each threshold becomes a function, with time as input. That is, for each time t , a .

166 1.4 Measurement instability

167 We introduce the *measurement instability framework* to study the consequences of this type of instabi-
168 lity. *Measurement instability* can be used to make a distinction between the influence of different parameters
169 that affect the scaling of the mapping function. Note that none of the measurement instability parameters
170 influence the quantity of the real variable.

171 1.4.1 Zero-point instability

172 *Zero-point instability* means that the zero-point of the mapping is unstable. Because the absolute zero-
173 point of the real quantity is stable, the thresholds move away or towards the absolute zero point of the real
174 quantity by a constant amount in the same direction. The range of the thresholds will stay identical.

175 1.4.2 Scaling instability

176 *Scaling instability* is instability that comes from a changing distance between the highest threshold and
177 the lowest threshold (the zero-point). The distance of the range of the measurement instrument is unstable.
178 In other words, the segment of the real quantity that is mapped to the representation can be wider or
179 more narrow in the real quantity depending on the person or the timepoint on which a person is measured.
180 Scaling instability is independent of zero-point instability.

181 1.4.3 Interval instability

182 *Interval instability* means that the distance between consecutive pairs of thresholds is not identical.
183 Interval instability can have a structural or a random appearance. An example of structural interval
184 instability would be a growing distance between each subsequent neighbouring pair of thresholds. In
185 random interval instability, there is no relationship between the order of the pairs and the distance of
186 neighbouring pairs of thresholds. Interval instability can simultaneously have a structural and a random
187 component.

188 1.4.4 Combinations

189 These three instabilities operate mostly independently but a few things should be noted. One is that real
190 scores can fall outside of the range of the thresholds. These scores will naturally be mapped to the lowest
191 or highest option. Second, if the zero-point of the mapping is stable then scaling instability can only affect
192 the right bound of the range. Further, thresholds cannot be outside of the range of the measuring device by
193 definition. Therefore, interval instability is limited by scaling and zero-point instability. We also reiterate
194 once more that the values of the actually existing quantity are not affected by the mapping. We consider
195 these absolute up to the point of the uniqueness theorem. Moreover, the chosen scaling of the real quantity
196 is of no practical concern. If the original attribute would be scaled differently, it would affect the thresholds
197 proportionately. Statistical conclusions would be identical after these scaling differences are propagated.
198 Lastly, all forms of instability can be either systematic or unsystematic. If the instability is systematic, then
199 the mapping of the variable from an existing quantity to the rating score is impacted by other variables
200 of relevance to the test. If the instability is unsystematic, there is no relationship between other relevant
201 variables and the mapping. Systematism can be consequential in the context of applying a statistical test,
202 because if the mapping is inconsistent between test conditions it might lead to biased results.

2 METHOD

If you want to make claims about the consequences of different structural deficiencies of a framework, it can be helpful to run a simulation study to the effect of those deficiencies to show the consequences of those effects to researchers.

2.1 Aims

The aim of this simulation is to illustrate and clarify the potential consequences of measurement instability. This helps us to formalize and . Given these aims, we initially have decided to use the simplest test and data-generating mechanisms we could think of.

2.2 Data-generating mechanism

Our data-generation pipeline consists of multiple stages. First, we generate the data on the data itself.

We generate data for a simple independent-sample t-test with the aim that this test performed on the unprocessed data has a statistical power of 0.80 for a one standard deviation difference based on an α of 0.05. Our sample size was chosen to ensure an empirical power of 0.807: 18 participants per group.

In the first test, the two distributions should be equivalent. Two samples are drawn from the same distribution. Therefore, our test will be rejected in 0.05 percent of cases, as should be expected. We draw both samples from a standardized normal distribution:

$$\text{Normal } (\mu = 0; \sigma = 1)$$

In the second test, there is a one standard deviation difference between the null distribution and the alternative distribution. The parameters are as follows:

$$0: \text{Normal } (\mu = 0; \sigma = 1)$$

$$A: \text{Normal } (\mu = 1; \sigma = 1)$$

Random numbers are generated using the Xoshiro256++ algorithm. The seed is 8508845.

2.3 Threshold transformation pipeline

We have implemented a full-factorial design by planning a range of parameters for each sub-type of measurement instability. These factors include the distribution of the zero-point, the number of bins the data should be divided into, the distribution of the scaling, and both random and structural threshold instability. We emulate differences in the degree of measurement stability by running the thresholds through a transformation pipeline. Each stage of the pipeline draws on the results of the last, and the stages result in one set of thresholds. Thresholds are generated for each iteration of the simulation, right before we perform the analysis.

2.3.1 Setting the number of response categories

We set the number of response categories as one of the parameters for the simulation. This is fed into the function that generates the thresholds. The set of values includes each integer from four to fifteen, 20, 50, 80, and 100.

2.3.2 Generating the zero point

First we set a zero-point. This is based on a three standard deviation difference from the mean. If it is stable, the left bound is constant. If it is unstable, we use a normal distribution to generate a zero-point with μ is three standard deviations lower than μ is variable. In the last stage, this actual bound is cut off

because the left-most category extends theoretically to lower values up to the absolute zero-point. Numbers are either a constant or drawn from a distribution. The constant is always -3. If drawn from a distribution, μ is always -3. σ is set at each interval when incrementing by 0.1 from 0.1 to 1.0.

2.3.3 Generating the scaling

We placed the maximum at the right bound of the 99% confidence interval on the basis of the null-data so that the mapping is equivalent for the two groups. The generation mechanism is similar to that of the zero-point: it is either a constant, or we use a normal distribution to generate a max-point where μ is equal to three standard deviations to the right of μ . The constant is thus always 6. If drawn from a distribution, μ is always 6. σ is set at each interval when incrementing by 0.1 from 0.1 to 1.0.

2.3.4 Generating the thresholds

We generated the intervals (all values of the thresholds aside from the outer bounds) by dividing the distance of the low bound to the right bound by the number of thresholds that are not at the edges. Unequal thresholds were generated using the cumulative distribution function of the beta-distribution. We chose this set-up because it results in either increasing or decreasing distance between, controllable by the β -parameter. If it is stable, the same parameters are chosen for each . If it is unstable, α and β were

2.4 Mapping the real scores to the bins

Afterwards, we cut the data using the thresholds we generated. This works using the simple step function described in the introduction: each value we've generated

2.5 Analysis

We ran two independent two-sample t-tests to perform the analysis. Each We save the t-value and the p-value of the test for each mapping, and save this

2.6 Software

We used the Julia language, and in particular the 'Statistics.jl' packages, to set up the simulation and to run the statistical analysis (?). Analyses were run on a personal computer. Full information about dependencies and version numbers can be found in a machine-readable format in the Manifest.toml file in the Github-repository. Instructions for running the analysis through a sandboxed project environment identical to our system can be found on the main page of this repository. Both the distributions and the test were chosen to use simple, commonly known statistics with the aim to show the potential consequences of zero-point, scaling, and interval instability. The experiment can be repeated using any possible test: the values are mapped *before* the analysis is run, so the full study is independent of the implementation of the test.

3 DISCUSSION

This study needs to be seen as an attempt at a bridge between scientists who study psychometrical testing and theoretical psychologists with an interest in formal measurement theory. We believe that the dots yell out to be connected, but it can be difficult to see the structural equivalence in . Either the abstract philosophical nature of

Another reason for the current treatment is the frustration about the limitations of modern treatments of this topic. Many modern simulation studies implicitly assume the thing that ought to be proven: that the

mapping of the continuous variable to the measurement scales result in equal intervals. We believe that part of the reason for this confusion is that psychologists are allergic to making claims about the quantitativity of variables. Yet, most of the statistical tests we commonly use to test data rely on that claim. It can be difficult to understand why we rely on those assumptions for testing, but are reluctant when it comes to justification. This is an attempt to show that it can bring a great deal of clarity if we extend our lenience to allow for these assumptions in justification.

Our choice to give the hypothesized perfect measurement an absolute zero is based on both a defence and a justification. The defence concerns the claim that one never really lacks a psychological variable is no challenge. Each quantity has to have a theoretical zero-point, because otherwise the quantity is meaningless. Temperature, for example, does have an absolute zero-point, at 0° K, yet it cannot ever be reached. Nobody would argue because it cannot ever be reached that we should stop attempting to measure temperature in Kelvin. As for the justification: using the absolute zero-point (or the absence of the posited attribute) allows for an unambiguous relationship of the perfect quantity measurement

3.1 Ceiling Floor effect

3.2 Distributions: not as expected

3.3 Operational ideal

3.4 Relationship of this framework to other theories

We would like to

3.5 Relationship to Classical Test Theory

In classical test theory,

3.6 Relationship to Latent Variable Theory and Measurement Invariance

The current framework has the strongest resemblance to measurement invariance models Simulation studies comparing measurement invariance approaches

3.7 Relationship to the Rasch Model

3.8 Validity

We subscribe to the point of . However, we do not agree that it should be so casually treated as a question that can be answered with a simple yes or no. Not only the entity

3.9 Measurability

The experiment above aims to show that the factors that go into scale instability, discretization, can influence.

Some concepts are easier to measure than others.

In the social sciences, there is often something left when you try to quantify kitchen-and-sink concepts that is not captured in the measurement (). The economist Georgescu-Roegen has called this the qualitative residual (?). It is not far-fetched to assume that this qualitative residual plays a role in the difficulty of achieving scale stability or defining the relative zero-point. Georgescu-Roegen asserts the importance of quantifying that residual regardless of the challenges: ‘We buy and sell land by acres, because land is

often homogeneous over large areas; and because this homogeneity is not general, we have differential rent. How unimaginably complicated economic life would be if we adopted an ordinal measure of land chosen so as to eliminate differential rent, let alone applying the same idea to all economic variables involving qualitative variations' (?)!

What is problematic, however, is that we do not know how our variables relate to

If there is no one-to-one relationship from the measurement device to the measure, it is wise not to mistake the measurement scores for the empirical relationships that are described. This problem is called the 'map-territory' problem in the philosophy of science. Rating scales were invented to make an estimate of the level of psychological variables in certain contexts, but they make a whole lot of assumptions about the nature of what those variables are. Mistaking the scores on that rating scale for the quantity of .

Going o

It must also be said that we do feel some sympathy towards the viewpoint by Sijtsma(?).

We agree with Borsboom, Mellenbergh and van Heerden that what validity should be defined as a scale measuring what it purports to measure. However, their conception is probabl. A scale might very well be able to measure what it purports to measure, but if this is only an attempt

3.10 Utility Measures

An hereto unnoticed (as far as I know) strong resemblance to the discussion from economists about the status of utility measures. This problem is in many ways identical to the problems that we see in psychology. What are the limits of

There is a distinction between unobservable and observable properties made in psychology that pops up often, such as in the discussion of latent variables. This distinction is not without criticism. As Burgos notes, the distinction is a (?)

3.11 Takeaway points

1. When thinking about measurement, it can be useful to separate *measurement error* from *measurement instability*. The former is error of a measure already assuming a certain scaling context. The latter is caused by inconsistencies in the mapping process from the real value to its rating scale equivalents.
2. The *operational ideal* of the measurement instability framework is to minimize the issues that inconsistent mapping between participants or timepoints can cause. This, for example, can be through rescaling or the development of statistical measures. Moreover, uncertainty calculations can be made on the basis of estimated measurement instability.
3. If it were easy to measure the relationship of the scaling of a rating scale and a 'real' magnitude of a quantity, it would have been as common as

NOTATION

CONFLICT OF INTEREST STATEMENT

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

FUNDING

344 No external funding was used for this project.

ACKNOWLEDGMENTS

345 I acknowledge the work of my thesis supervisors, who introduced me to the method and left me free to
 346 pursue the project as I imagined it in all its weird and shifting shapes. The great help of the Julia community
 347 was also appreciated, as they have been pushing me forward where I got stuck and took the time to respond
 348 to my stupid questions. Finally, I would like to acknowledge the feedback and conversations between me
 349 and my thesis group, who have been working through my text and made sure that it is easy to follow and
 350 well-written. Special thanks to Giuliana Orizzonte, Daniel Anadria, dr. Hessen, dr. Derksen, dr. Grelli
 351 and dr. Bringmann for fruitful discussions and feedback about my topic. Lastly, I would like to thank my
 352 girlfriend, family, and friends for the mental support throughout.

DATA AVAILABILITY STATEMENT

353 The code, additional material, and generated data for this study can be found on GitHub.

354 This relationship can also be visualized on a number line. Consider a five-point instrument is used to
 355 measure a quantitative psychological variable. At a particular point in time t , we assume that the projection
 356 on the real number line is divided into four segments of equal length and one element which goes on
 357 infinitely. It should be noted that the scaling is arbitrary, but it has been set here to multiples of two for
 358 convenience.

