

# Measured Inference: Scales, Statistics, and Scientific Inference

Conor Mayo-Wilson

University of Washington, Seattle, USA

## Abstract

Despite the recent “epistemic turn” in the philosophy of measurement, philosophers have ignored a nearly 80 year controversy about the relationship between statistical inference and measurement theory. Some scholars maintain that measurement theory places no constraints on statistics, whereas others argue that the measurement scale (e.g., ordinal or interval) of one’s data determines which statistical methods are “permissible.” I defend an intermediate position: even if existing measurement theory were irrelevant to statistical inference, it would be critical for *scientific* inference, which requires connecting statistical hypotheses to broader research hypotheses.

## 1 Introduction

Despite the recent “epistemic turn” in the philosophy of measurement, philosophers have ignored a nearly 80 year controversy about the relationship between statistical inference and measurement theory.<sup>1</sup> Statistical libertarians, as I will call them, maintain that measurement theory places essentially no constraints on statistics.<sup>2</sup> In contrast, measurement bureaucrats (again, my term) endorse Stevens’ doctrine of permissible statistics, according to which parametric methods (e.g., t-tests) should be applied only to interval or ratio-scaled data, whereas ordinal data requires the use of non-parametric tests (e.g., Mann-Whitney U).<sup>3</sup>

To see what is at stake, imagine a CEO is worried about sexual harassment in her company. She issues a two-question survey to 50 randomly selected employees. The first question asks employees to identify their gender, and the second asks, “On a 1-7 scale – with 1 representing ‘completely dissatisfied’ and

---

<sup>1</sup>See [Tal, 2015] for a discussion of the “epistemic turn”. To my knowledge, the only philosophical work that engages with this controversy is [Larroulet Philippi, 2021] and [Larroulet Philippi, 2022]. As is common (e.g., [Tal, 2015]), I use the phrase “measurement theory” to refer to the mathematical work that culminates in the three-volume *Foundations of Measurements* texts [Krantz et al., 1971]. I avoid using the term “representational measurement theory” (RTM) to describe those mathematical results because RTM is often used to denote several further epistemological theses [Tal, 2021].

<sup>2</sup>See Atkinson [1988], Gaito [1980], Lord [1953, 1954], Velleman and Wilkinson [1993].

<sup>3</sup>Bureaucrats include Blalock [1960], Senders [1958], Siegel and Jr [1988], Thomas [2006], Wilson [1971]. An intermediate position is defended by [Marcus-Roberts and Roberts, 1987], who argue that only “meaningful” statistical hypotheses are of scientific interest (see Section 3.2) but that there are no restrictions on what statistics it is appropriate to calculate.

7 representing ‘completely satisfied’ – how satisfied are you with the company’s sexual harassment policies?” When the survey has been completed by all 50 selected employees, the CEO divides responses according to gender and calculates the averages/means of men’s and women’s responses (3.2 and 2.1 respectively). She then performs a t-test to assess whether those averages differ. She finds a statistically significant difference. May the CEO conclude that men and women in the company are satisfied to different degrees with the company’s sexual harassment policies?

Libertarians maintain “yes”; bureaucrats say “no.” According to Libertarians, if the *averages* of the men’s and women’s responses differ, then so must the two distributions of responses. End of story.

For bureaucrats, however, the CEO’s data are merely *ordinal*, and averages of ordinal data should not be invoked in statistical reasoning. To motivate that prohibition, suppose the survey had omitted a numerical scale and instead asked respondents to choose from seven categories describing their degree of satisfaction in English words. Just as a researcher might be hesitant to “average” non-numerical responses like “somewhat dissatisfied” and “very satisfied”, one should be reluctant to calculate the means of the responses from the original survey.

I defend an intermediate position: even if existing measurement theory were irrelevant to statistical inference, it is critical for *scientific* inference which requires connecting statistical hypotheses to broader research hypotheses.<sup>4</sup> In the CEO’s case, the statistical hypotheses concern the relationship between two (probability) distributions over the numbers 1-7, which represent employee responses on a fixed numerical scale. In contrast, the research hypothesis of interest likely concerns whether *attitudes* about sexual harassment differ, or whether men and women’s *behavior* differ in ways that matter to the CEO (e.g., whether productive women are more likely to leave the company within the year). Libertarians are correct that the CEO’s statistical inferences require no measurement theory, but bureaucrats are correct that further assumptions are necessary to draw research conclusions from the CEO’s statistical analysis.

To show how the distinction between statistical and research hypotheses arises in scientific practice, I summarize the controversy on “interpretable effects” in psychology in [Section 2](#).<sup>5</sup> I argue that the controversy amounts to the following: statistical conclusions reached in memory experiments do not *mathematically entail*<sup>6</sup> research hypotheses about some purported latent attributes, specifically, an attribute that might be “memory strength.” Moreover, many psychologists believe that the data from some memory experiments are of questionable scientific interest unless the statistical hypotheses that the data

<sup>4</sup>The distinction between statistical and research hypotheses is standard in medical science. See [\[Lawler and Zimmermann, 2021\]](#) for a discussion of cases in which the two types of hypotheses are misaligned.

<sup>5</sup>The controversy originated with [\[Loftus, 1978\]](#). See [\[Wagenmakers et al., 2012\]](#) for the history.

<sup>6</sup>Henceforth, I say a set of premises *mathematically entail* a conclusion if the premises of the argument and axioms of set theory together *logically* entail the conclusion. I say an argument is *mathematically valid* if its premises mathematically entail its conclusion.

support mathematically entail the relevant research hypotheses.

I remain agnostic what it is of “scientific interest” in psychology, but I investigate the consequences of the skeptical position about memory experiments. In [Section 3](#), I argue that the theory of “meaningfulness” developed in measurement theory can help one identify general conditions under which which inferences from statistical hypotheses to research ones are mathematically valid.<sup>7</sup>

Before beginning, it will be helpful to characterize the distinction between statistical and research hypotheses more precisely. Statistical hypotheses are (sets of) *probability distributions* that specify how likely various data are. Such hypotheses always concern a *particular* experimental setup (which might be repeatable). Statistical methods (e.g., hypothesis tests and estimators) allow one only to evaluate how well different statistical hypotheses are supported by data.

In contrast, research hypotheses have implications beyond a given experimental context, and they may not specify any precise probabilities whatsoever. A research hypothesis might, for example, concern (i) latent or unmeasured attributes or (ii) the outcomes of different measurement procedures in other experimental contexts. In the CEO’s case, the latent attributes are attitudes or behavioral dispositions, which are not measured in the survey.

## 2 Conflating latent attributes with measured ones

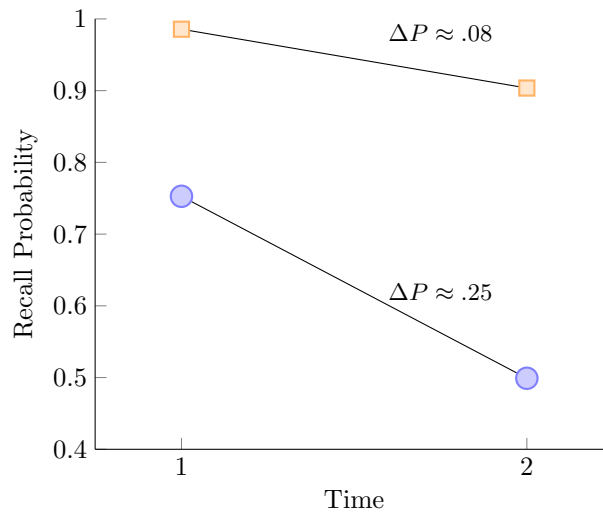
Nearly 50 years ago, [Loftus \[1978\]](#) famously argued that many memory experiments in psychology suffer from a serious methodological problem: a latent attribute – call it *memory strength* – is conflated with what is directly measured in experiments, e.g, the probability of correctly recalling a stimulus.

Imagine experimental subjects are divided into two groups; call them A and B. Participants in both groups are presented with a sequence of five “random” letters, which they will be asked to recall at two different later times (e.g., after 5 and 20 seconds respectively). But prior to the recall phase of the experiment, different groups are subject to different conditions. Groups A and B might receive different instructions, for example.

Suppose the results of the experiment are shown in Figure 1. Group A’s average recall rates are represented by the two circular blue endpoints of the bottom line, and Group B’s average recall rates are represented by the two square orange endpoints of the top line.

---

<sup>7</sup>I use the theory of *semantic* meaningfulness in [\[Adams et al., 1965\]](#). This theory was inspired by the theory of “empirical” or “scientific” meaningfulness that was developed by [\[Suppes and Zinnes, 1962\]](#) and later defended by [\[Roberts, 2009\]](#) and [\[Narens, 2012\]](#), among others. I do not endorse the latter theories.

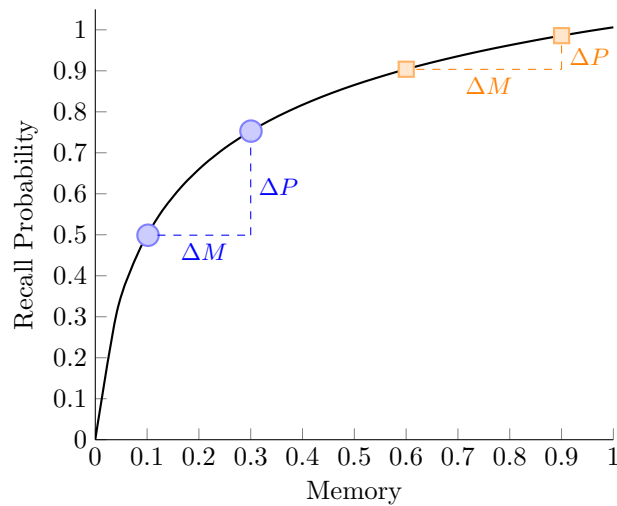


**Figure 1**

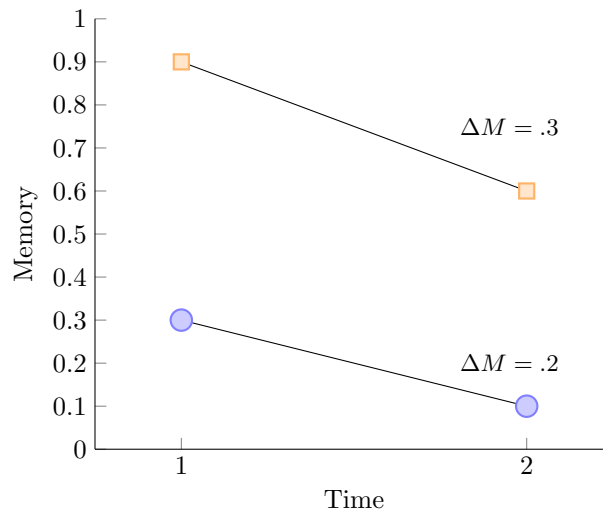
The lines in the diagram are merely heuristic. Both groups are tested for recall at only two discrete times, and so the lines do **not** indicate that, in the experiment, the probability of correct recall decreases linearly over time. However, the lines help one see an important fact: the slope of the A-group line is steeper than that of the B-group. One might hypothesize, therefore, that participants in condition A *forget* at a rate faster than do participants in condition B. That hypothesis, Loftus argues, is underdetermined by the experiment.

Why? Suppose memory strength – call it  $q$  – is a quantifiable, and suppose observed recall rate at time  $t$  is a function  $r_t$  of  $q$ . Loftus shows that, even if recall rate increases with  $q$ , it is possible for  $q$  to decrease *at the same rate* or even *faster* in condition B than in condition A unless one makes further assumptions about the mathematical form of the function  $r_t$ . Which assumptions? Loftus proves that, if the recall rate is a *linear* function of  $q$ , then the desired inference about memory strength is valid.

Figures 1 through 3 illustrate Loftus' critique. Suppose the function  $r_t$  is the one shown in Figure 2. Then *memory* decreases faster in Condition B than in Condition A (as shown in Figure 3), even though *recall rates* in Condition B decrease more slowly than in Condition A (as shown in Figure 1).



**Figure 2**



**Figure 3**

Loftus never distinguishes between statistical and research hypotheses, but his critique beautifully illustrates the distinction. Even if statistical methods establish that the distribution of *observed recall rates* in condition A differs from that in B, Loftus' critique challenges the inference from those statistical conclusions to research hypotheses about memory.

One might object to Loftus' critique by arguing that, in the memory experiments in question, psychologists are not trying to draw inferences about a latent attribute: the research hypotheses and statistical hypotheses alike concern the measured probability of recall. Memory is "operationalized" in recall rates. That critique would be legitimate if such recall rates were known to be

of independent scientific interest. For example, perhaps those recall rates can predict performance in many other important “memory” tasks. But crucially, the *differences* in recall rates must be predictive, for otherwise the conclusion that subjects’ recall rates *decrease* faster in one condition than in another is irrelevant.

### 3 Inference, Meaningfulness, and Scales

Loftus shows that, in important scientific settings, there may be an inferential gap between research hypotheses about latent attributes and statistical hypotheses about measured outcomes. I now argue that, in many of those settings, the theory of meaningfulness developed by measurement theorists specifies the assumptions necessary to bridge the inferential gap. To do so, I first argue that the theory of meaningfulness provides a plausible answer to the question, “Under what conditions does an attribute (e.g., memory strength) have the type of quantitative structure for which differences (e.g., between memory strength at two different times) are meaningful?” The answer, I claim, is that the attribute admits an *interval-scale*.<sup>8</sup> Similarly, claims about ratios of an attribute are meaningful if and only if the attribute admits a ratio-scale. I discuss scale types in [Section 3.1](#).

My main contribution is to show that scale-classifications also play an important *epistemic* role as they can be used to identify *mathematically valid* inferences from statistical hypotheses to a research ones. Because inferences from statistical hypotheses to research ones are often not easily formalized, it is important to identify which are mathematically valid.

#### 3.1 Scales and Scale Types

Length can be quantified in inches and centimeters. Mass is quantified in kilograms and tonnes. In general, anything that is quantifiable can be quantified in many ways. Roughly, a *scale* is a way of quantifying a property. Scales are rarely unique.

But scales are often related. For example, an inch is 2.54 centimeters; a yard is three feet, and generally, one can convert any unit of length into another by multiplying by a constant. When all scales for an attribute are multiples of one another, the attribute is called *ratio-scaled*.

Not all scales are ratio-scales. Consider calendar date. In all calendar systems, there is an arbitrarily chosen “zeroth” year, and calendar date is determined by counting from that zero. Different zeroes can be chosen, e.g., in Islamic calendars, Muhammad’s pilgrimage fixes the zeroth year. And instead of counting years, one could count days, weeks, or units of time determined by lunar rather than solar events. Thus, in converting calendar date in one system to another, one must first multiply (e.g., to convert years to days) and then add

---

<sup>8</sup>As I am not a psychologist, I will not assess whether there is evidence that ‘memory strength’ admits an interval scale.

another number (e.g., to correct for choices of “zeroth” year). When all scales for a property are related in this way, one says the property is *interval-scaled*.

The reader might ask, “what determines which scales are ‘permissible’ ways of quantifying an attribute?” For the purposes of this paper, my answer is, “consult *Foundations of Measurement*” [Krantz and Tversky, 2006, Krantz et al., 1971, 2006]. There, the reader will find a body of mathematical theorems that show that, if certain relations hold among objects (or events) with a given attribute, then then set of permissible scales must always be of one of the few types identified by Stevens [1946]. Importantly, as Michell [1997] observes, the theorems in *Foundations of Measurement* specify the quantitative structure of an attribute even if the attribute cannot be measured, in any realistic sense. This is important because Krantz et al. [1971] are often said to endorse “positivist” assumptions, e.g., that “empirical relations be directly observable, or ‘identifiable’ ” [Mari et al., 2023, p. 94]. Those philosophical assumptions, however, play no role in the mathematical results about scale types.

What is now important for us is to understand how scale classifications can clarify questions about meaningfulness.

### 3.2 Meaningfulness

Contrast two claims: “Ada is more than twice as tall as Boris” and “Ada’s height in inches is more than twice that of Boris.” Notice that the first sentence is true if and only if the second is true. That may be surprising because the first is a *scale-free* assertion – it contains no mention of *units* of length – whereas the latter is *scale-specific*. But according to an influential definition of “meaningfulness”, one should not be surprised at all: the first sentence has a truth-value if and only if its truth-value matches that of the second.<sup>9</sup>

To understand the proposed theory of meaningfulness, consider the scale-free assertion “Ada’s is three taller than Boris.” That claim is nonsense. If Ada is three inches taller than Boris, then she is not three feet taller than Boris. Units matter. These examples motivate the following proposal: a scale-free sentence about an attribute is *meaningful* (i.e., it has a truth-value) if and only if all the scale-specific instances of the statement have the same truth-value. In other words, a scale-free sentence is meaningful if the units do not matter.

Scale-free hypotheses are ubiquitous in science. Consider Galileo’s law of free fall, which asserts that the distance traveled by an object in free fall is proportional to the square of the time of the descent. Galileo’s law does not require that distance be measured in a specific unit like meters, nor that time be measured in a unit like seconds. Similarly, Boyle’s law about pressure and volume is scale free: neither units of neither pressure nor volume are mentioned. Scale-free hypotheses also occur in the social sciences. For example, economists do not mention a specific currency when they claim that profits are maximized when marginal revenue equals marginal costs. These examples show it is important to understand when scale free hypotheses are meaningful.

---

<sup>9</sup>See references in footnote 3.

To see how the theory of meaning works, consider the scale-free hypothesis “memory strength decreases more rapidly in condition  $A$  than in condition  $B$  between times  $t_1$  to  $t_0$ .” Loftus argued that hypothesis could not be inferred from the observed recall effects.

Now, that scale-free hypothesis is meaningful, according to the above theory of meaning, if for any two scales for memory  $M_1$  and  $M_2$  the following biconditional holds:

$$\begin{aligned} M_1(t_1, A) - M_1(t_0, A) &> M_1(t_1, B) - M_2(t_0, B) \text{ if and only if} \\ M_2(t_1, A) - M_2(t_0, A) &> M_2(t_1, B) - M_2(t_0, B) \end{aligned} \quad (1)$$

where  $M_j(t, x)$  represents the memory strength along scale  $j \in \{1, 2\}$  at recall time  $t \in \{1, 2\}$  in condition  $x \in \{A, B\}$ . Some quick algebra shows that Equation 1 holds if there is a positive number  $c > 0$  and some number  $d$  (possibly negative) such that  $M_2(t, x) = c \cdot M_1(t, x) + d$  for all times  $t$  and all conditions  $x$ . That is, the assertion is meaningful if memory is an interval-scaled attribute.

This example suggests that there is some relationship between (1) meaningfulness and (2) the validity of inferences that have scale-free conclusions. Understanding that relationship is important because whereas *statistical* hypotheses are almost always scale-specific (as they describe the data of a particular experiment, which must be measured in specific units), scientists’ *research* hypotheses are often scale-free.

### 3.3 Mathematical Validity and Research Hypotheses

The theory of meaningfulness allows us to immediately identify a set of mathematically valid inferences that have scale-free conclusions. Let  $M$  be a scale; let  $\varphi_M$  be some scale-specific proposition about the attribute  $A$ , and let  $\varphi_A$  be the corresponding scale-free proposition. For instance, if  $M$  is inches and  $\varphi_M$  is the assertion “Ada’s height in inches is twice that of Boris”, then  $\varphi_A$  is the assertion “Ada’s height is twice that of Boris.” Here’s a theorem (stated imprecisely).

**Fact:** The inference from  $\varphi_M$  to  $\varphi_A$  is valid if (1)  $\varphi_A$  is meaningful and (2)  $M$  is a permissible scale for the attribute  $A$ .

The fact follows immediately from definitions. Suppose 1 and 2 hold. Since  $\varphi_A$  is meaningful (by 1),  $\varphi_A$  is true if and only if  $\varphi_S$  is true for any scale  $S$ . Because  $M$  is a scale for  $A$  (by 2), it follows that if  $\varphi_M$  is true, then  $\varphi_A$  must be true (and so the inference is valid).

So what? Recall, statistical hypotheses are about data in a given experimental context *on a fixed scale* (e.g., the CEO’s data is on a 1-7 scale for satisfaction; the memory experiment’s data is probability of recall in a specific context). In contrast, research hypotheses are often scale-free precisely because researchers desire replicable results that do not depend on choice of measurement units. Thus, the inference from a statistical hypothesis (e.g. that men’s and women’s responses differ on average) to the corresponding scale-free research hypothesis



is mathematically valid if (1) the research hypothesis is meaningful and (2) the measurement scale is a permissible way of quantifying the attribute.

The above simple fact is a generalization of Loftus' positive suggestion. It entails that if memory strength admits an interval scale (and so the hypothesis that memory decreases faster in one condition than another is meaningful), then one can validly infer the research hypothesis from the measured results about recall rate if the recall rate is a permissible scale for memory, i.e., it is a linear function of memory strength. However, the above fact is a generalization of Loftus' claim because it applies to *all* scale types, not just interval ones.

The epistemological importance of the fact above, however, should not be overstated. Notice that in Loftus' critique – as in the above fact – (1) the focus is on *validity* of arguments rather than inductive strengths and (2) the conclusion of the inference  $\varphi_A$  is the scale-free hypothesis *corresponding* to the premise  $\varphi_M$ . That is a very restricted form of inference.

First, many strong arguments are not mathematically valid. Second, many valid inferences are not of the above form and yet have scale-free conclusions. For instance, let  $\psi$  and  $\varphi$  be scale-specific and scale-free hypotheses respectively. Suppose  $\varphi$  is meaningful. Then  $\psi \rightarrow \varphi$  and  $\psi$  together entail  $\varphi$ .

In fact, it is possible to describe such a case of modus ponens when the scale of  $\psi$  is not even a permissible scale for the relevant attribute. Suppose two cross country teams race one another. Let  $\varphi$  be the (scale-free) hypothesis that asserts “The average time of Team 1 is faster than that of Team 2.” Notice that hypothesis is meaningful because its truth does not depend on whether times are recorded in seconds, milliseconds, etc. However, suppose the finishing times of runners are not recorded, only the ranks, with 1 being assigned to the first-place runner, 2 to the second-place and so on. The assignment of ordinal ranks is not a permissible scale for time. But let  $\psi$  be the scale-specific proposition “All runners on Team 1 have a lower rank than all runners on Team 2.” Then  $\psi \rightarrow \varphi$  is a mathematical truth, and so if one has evidence for the scale-specific claim  $\psi$ , then one obtains evidence for the scale-free hypothesis  $\varphi$ .<sup>10</sup>

Despite these limitations, the theory of meaningfulness provide a first step in (i) understanding the debate between statistical libertarians and measurement bureaucrats and (ii) identifying a partial resolution.

## 4 Conclusion

As others have convincingly argued, measurement theory helps scientists identify if an attribute is quantifiable at all.<sup>11</sup> I have further argued that if, conditions for quantifiability are met, measurement theory characterizes auxiliary assumptions that are sufficient to facilitate mathematically valid inferences from statistical hypotheses about measured outcomes to research hypotheses about latent attributes. Namely, by the simple fact established in [Section 3.3](#),

<sup>10</sup>More generally, this argument would work if  $\psi$  were replaced with the claim that the ranks of Team 2 stochastically dominate those of Team 1.

<sup>11</sup>In addition to [\[Michell, 1986\]](#), see [\[Heilmann, 2015\]](#) and [\[Wolff, 2020\]](#).

it suffices to show that the measured outcomes are values along a permissible scale for the attribute.

If the measurement scale is not a permissible (or if there is no latent attribute with the relevant mathematical structure), then often, further data and statistical analyses will be necessary to facilitate inference to research hypotheses. This is the most plausible way of describing the CEO case at the outset of this paper. “Satisfaction with sexual harassment policies” is likely not a latent attribute admitting an interval scale, and what researchers are likely interested in is making inferences from the survey to other behaviors. But such inferences would require data that would allow one to explore statistical associations between survey responses and the relevant behaviors.

This paper has only begun to address the question of when scale-specific propositions provide *evidence* for scale-free ones. I have studied only a very narrow set of mathematically valid inferences, and a general theory of *inductive* inference for scale-free hypotheses is still in its infancy [Larroulet Philippi, 2021, 2022].

## 5 Acknowledgments

Thanks to Christian Larroulet-Philippi, audience members at the PSA, and especially to Paul Pedersen and David Kellen for earlier conversations about measurement and statistics.

## 6 Funding Statement

None to declare.

## 7 Declarations

None to declare.

## References

- Ernest W. Adams, Robert F. Fagot, and Richard E. Robinson. A theory of appropriate statistics. *Psychometrika*, 30(2):99–127, June 1965. ISSN 0033-3123, 1860-0980. doi: 10.1007/BF02289443. URL <http://link.springer.com/10.1007/BF02289443>.
- Leslie Atkinson. The measurement-statistics controversy: Factor analysis and subinterval data. *Bulletin of the Psychonomic Society*, 26(4):361–364, October 1988. ISSN 0090-5054. doi: 10.3758/BF03337683. URL <https://doi.org/10.3758/BF03337683>.

- Hubertt Blalock. *Social statistics*. McGraw Hill Book Company, New York, Toronto, and London, 1960.
- John Gaito. Measurement scales and statistics: Resurgence of an old misconception. *Psychological Bulletin*, 87(3):564–567, 1980. doi: 10.1037/0033-2909.87.3.564. Place: US Publisher: American Psychological Association.
- Conrad Heilmann. A new interpretation of the representational theory of measurement. *Philosophy of Science*, 82(5):787–797, 2015. URL <https://www.cambridge.org/core/journals/philosophy-of-science/article/new-interpretation-of-the-representational-theory-of-measurement/C936F03C617926BA7980795829448997>. Publisher: Cambridge University Press.
- David H. Krantz and Amos Tversky. *Foundations of Measurement Volume III: Representation, Axiomatization, and Invariance*. Dover Publications, Mineola, N.Y, December 2006.
- David H. Krantz, R. Duncan Luce, Patrick Suppes, and Amos Tversky. *Foundations of Measurement Volume I: Additive and Polynomial Representations*. Academic Press Inc., New York, NY, 1971.
- David H. Krantz, R. Duncan Luce, Patrick Suppes, and Amos Tversky. *Foundations of Measurement Volume II: Geometrical, Threshold, and Probabilistic Representations*. Dover Publications, Mineola, N.Y, December 2006.
- Insa Lawler and Georg Zimmermann. Misalignment Between Research Hypotheses and Statistical Hypotheses: A Threat to Evidence-Based Medicine? *Topoi*, 40(2):307–318, April 2021. ISSN 1572-8749. doi: 10.1007/s11245-019-09667-0. URL <https://doi.org/10.1007/s11245-019-09667-0>.
- Geoffrey R. Loftus. On interpretation of interactions. *Memory & Cognition*, 6(3):312–319, May 1978. ISSN 0090-502X, 1532-5946. doi: 10.3758/BF03197461. URL <http://link.springer.com/10.3758/BF03197461>.
- Frederic M. Lord. On the Statistical Treatment of Football Numbers. *American Psychologist*, 8(12):750–751, 1953. ISSN 1935-990X. doi: 10.1037/h0063675. Place: US Publisher: American Psychological Association.
- Frederic M. Lord. Further Comment on "Football Numbers". *American Psychologist*, 9(6):264–265, 1954. ISSN 1935-990X. doi: 10.1037/h0059284. Place: US Publisher: American Psychological Association.
- Helen M. Marcus-Roberts and Fred S. Roberts. Meaningless Statistics. *Journal of Educational Statistics*, 12(4):383–394, December 1987. ISSN 0362-9791. doi: 10.3102/10769986012004383. URL <http://journals.sagepub.com/doi/10.3102/10769986012004383>.

- Luca Mari, Mark Wilson, and Andrew Maul. *Measurement across the sciences: Developing a shared concept system for measurement*. Springer Nature, 2023. URL <https://library.oapen.org/handle/20.500.12657/61879>.
- Joel Michell. Measurement scales and statistics: A clash of paradigms. *Psychological Bulletin*, 100(3):398–407, 1986. ISSN 1939-1455. doi: 10.1037/0033-2909.100.3.398. Place: US Publisher: American Psychological Association.
- Joel Michell. Quantitative science and the definition of *measurement* in psychology. *British Journal of Psychology*, 88(3):355–383, August 1997. ISSN 0007-1269, 2044-8295. doi: 10.1111/j.2044-8295.1997.tb02641.x. URL <https://bpspsychub.onlinelibrary.wiley.com/doi/10.1111/j.2044-8295.1997.tb02641.x>.
- Louis Narens. *Theories of meaningfulness*. Psychology Press, 2012. URL <https://www.taylorfrancis.com/books/mono/10.4324/9781410604156/theories-meaningfulness-louis-narens>.
- Cristian Larroulet Philippi. On Measurement Scales: Neither Ordinal nor Interval? *Philosophy of science*, 88(5):929–939, 2021. URL <https://www.cambridge.org/core/journals/philosophy-of-science/article/on-measurement-scales-neither-ordinal-nor-interval/BE9D7DBF7558EA12221776D5E995516A>. Publisher: Cambridge University Press.
- Cristian Larroulet Philippi. Against Prohibition (Or, When Using Ordinal Scales to Compare Groups Is OK). *The British Journal for the Philosophy of Science*, page 721759, July 2022. ISSN 0007-0882, 1464-3537. doi: 10.1086/721759. URL <https://www.journals.uchicago.edu/doi/10.1086/721759>.
- Fred S. Roberts. *Measurement Theory: Volume 7: With Applications to Decisionmaking, Utility, and the Social Sciences*. Cambridge University Press, Cambridge, reissue edition edition, March 2009. ISBN 978-0-521-10243-8.
- Virginia L. Senders. *Measurement and statistics:: A basic text emphasizing behavioral science applications*. Oxford University Press, first edition edition, January 1958.
- Sidney Siegel and N. John Castellan Jr. *Nonparametric Statistics for The Behavioral Sciences*. McGraw-Hill Humanities/Social Sciences/Languages, Boston, Mass., 2nd edition edition, January 1988.
- S. S. Stevens. On the Theory of Scales of Measurement. *Science*, 103(2684):677–680, June 1946. ISSN 0036-8075, 1095-9203. doi: 10.1126/science.103.2684.677. URL <https://www.science.org/doi/10.1126/science.103.2684.677>.

- Patrick Suppes and Joseph L. Zinnes. *Basic measurement theory*. Univ., 1962. URL [https://web.stanford.edu/group/csli-suppes/techreports/IMSSS\\_45.pdf](https://web.stanford.edu/group/csli-suppes/techreports/IMSSS_45.pdf).
- Eran Tal. Measurement in Science. June 2015. URL <https://plato.stanford.edu/Entries/measurement-science/>. Last Modified: 2020-08-07.
- Eran Tal. Two myths of representational measurement. *Perspectives on Science*, 29(6):701–741, 2021. URL <https://direct.mit.edu/posc/article-abstract/29/6/701/107112>. Publisher: MIT Press One Rogers Street, Cambridge, MA 02142-1209, USA journals-info . . .
- Hoben Thomas. Measurement Structures and Statistics. In *Encyclopedia of Statistical Sciences*. John Wiley & Sons, Ltd, 2006. ISBN 978-0-471-66719-3. doi: 10.1002/0471667196.ess1591.pub2. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/0471667196.ess1591.pub2>. \_eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/0471667196.ess1591.pub2>.
- Paul F. Velleman and Leland Wilkinson. Nominal, Ordinal, Interval, and Ratio Typologies are Misleading. *The American Statistician*, 47(1):65–72, February 1993. ISSN 0003-1305, 1537-2731. doi: 10.1080/00031305.1993.10475938. URL <http://www.tandfonline.com/doi/abs/10.1080/00031305.1993.10475938>.
- Eric-Jan Wagenmakers, Angelos-Miltiadis Krypotos, Amy H. Criss, and Geoff Iverson. On the interpretation of removable interactions: A survey of the field 33 years after Loftus. *Memory & Cognition*, 40(2):145–160, February 2012. ISSN 0090-502X, 1532-5946. doi: 10.3758/s13421-011-0158-0. URL <http://link.springer.com/10.3758/s13421-011-0158-0>.
- Thomas P. Wilson. Critique of Ordinal Variables. *Social Forces*, 49(3):432–444, March 1971. ISSN 0037-7732. doi: 10.1093/sf/49.3.432. URL <https://doi.org/10.1093/sf/49.3.432>.
- Jo E. Wolff. *The metaphysics of quantities*. Oxford University Press, 2020. URL <https://books.google.com/books?hl=en&lr=&id=UzjpDwAAQBAJ&oi=fnd&pg=PP1&dq=wolff+metaphysics+of+quantities&ots=W1hfpRagFr&sig=hPrUe76Troygwjrd0tPVpSZFG8g>.