

Reply to Rodgers and Rowe (2002)

Seth Roberts
University of California, Berkeley

Harold Pashler
University of California, San Diego

That a theory fits data is meaningful only if it was plausible that the theory would not fit. S. Roberts and H. Pashler (2000) knew of no enduring theories initially supported by good fits alone (good fits, i.e., where it was not clear that the theory could have plausibly failed to fit). J. L. Rodgers and D. C. Rowe (2002) claimed to provide 6 examples. Their 3 nonpsychological examples (Kepler, etc.) are instances of good practice: How the theory constrained outcomes was evident, so it was clear that the theory could have plausibly failed to fit. Their 3 psychological examples are flawed in various ways. It remains possible that no examples exist.

Were we to sit down at a table with Rodgers and Rowe, we like to think that the four of us could eventually agree on many things—for instance, that there are important differences between Mendel's support for his theory and goodness-of-fit evidence for mathematical learning theories. We might not agree, however, on a point that Roberts and Pashler (2000) did not make clear enough: When used to evaluate complex psychological theories, goodness-of-fit tests have been too easy to pass. The theories are too flexible, the data too variable, and contradictory outcomes too implausible. We fear we did not state this bluntly enough.

In many cities, taxi drivers must pass a written test. Current practice in major journals, such as *Psychological Review*—which often publish theories supported by no more than good fits—resembles allowing a driver to pass the test if he can write his name. "Theory Development Should Begin (But Not End) With Good Empirical Fits" is the title of Rodgers and Rowe's (2002) article. We agree. Taxi-driver tests should begin (but not end) with the test taker writing her name on the test. "Unless a researcher can demonstrate that a theory matches the data to begin with, there is little point in considering how it constrains the possible outcomes, how variable the data are, or whether results are surprising" (Rodgers & Rowe, 2002, p. 602). We agree. Unless a test taker can write his name, there is little point in considering whether he has more difficult-to-acquire skills. Rodgers and Rowe (2002, p. 600) noted "the many uses of goodness of fit." We agree—it does have many legitimate uses, just as it is often helpful to learn if a person can write his or her name.

The proof of a test is the competence of those who pass it. Many people who can write their name would not be good taxi drivers. What about the psychological theories that have passed goodness-of-fit tests and nothing else (tests without the additional informa-

tion that we argue is necessary)? Many such theories have been published. How have they fared? Roberts and Pashler (2000) emphasized that "there seem to be no examples of a theory supported mainly by good fits that has led to demonstrable progress" (p. 358)—progress, that is, measured by a better metric than goodness of fit. This is an important difference, we hope everyone can agree, between our critique and criticisms of null hypothesis significance testing. Not even the severest critics of null hypothesis significance testing have claimed that it has never been associated with work of lasting importance (a claim which is obviously false). Even good tools can be misused. In Roberts and Pashler (2000), we were not saying that the use of good fits as evidence for complex theories had been misused. We were saying—perhaps too politely, perhaps not—that the practice is rotten to the core, in the sense that a taxi-driver test that can be passed by writing your name is rotten to the core.

Rodgers and Rowe (2002) disagreed, of course. This is why their examples are the most important part of their argument. If current practice is reasonable—if passing a goodness-of-fit test by itself implies that a theory should be taken seriously by other scientists—then many theories that initially passed such tests (without the additional information we argue is necessary) should have turned out to have lasting importance. It is telling that Rodgers and Rowe's first three examples (Kepler, Mendel, and Watson and Crick) are outside psychology—an indication, we suppose, that examples within psychology are hard to find.

In Roberts and Pashler (2000), we tried to contrast what we considered bad practice (testing goodness of fit) with what we considered good practice (testing predictions). From Rodgers and Rowe (2002), we see that we might have been clearer. We presented the two as disjoint (nonoverlapping), but the real relation is that testing predictions (i.e., meaningful tests) is a subset of testing goodness of fit. Here is what we should have said: Consistency of theory and data is meaningful only if inconsistency was plausible. With simple theories—for example, theories that predict straight lines—that the theory might not have fit may be obvious, but with more complex theories it is usually not obvious and therefore needs to be made explicit. However, authors of complex theories routinely omit the necessary information. To show that a theory could have plausibly failed to fit requires three steps:

Seth Roberts, Department of Psychology, University of California, Berkeley; Harold Pashler, Department of Psychology, University of California, San Diego.

We thank Joseph Rodgers and Saul Sternberg for helpful comments.

Correspondence concerning this article should be addressed to Seth Roberts, Department of Psychology, University of California, Berkeley, California 94720-1650. E-mail: roberts@socrates.berkeley.edu

1. Determine a way in which the theory constrains possible outcomes. Any prediction can be described in terms of a *forbidden zone* (Roberts & Sternberg, 1993, p. 638) consisting of outcomes the theory cannot generate no matter what parameter values are used. For example, if the prediction is that a certain value computed from the results will be more than zero, then the forbidden zone is zero and less.

2. Shrink the forbidden zone to allow for the variability of the data. If the initial forbidden zone—not taking variability into account—is zero and less and the variability in the data (e.g., between-subject variation) gives a 95% confidence interval of length 10 for the value, then (assuming symmetric confidence intervals) the forbidden region shrinks to -5 and less, because an observed value of (say) -3 would be compatible with the theory.

3. Determine if any plausible outcomes lie in the shrunken forbidden zone. For example, do plausible alternative theories generate outcomes in that zone? Only if the forbidden zone contained plausible outcomes should consistency of theory and data be impressive.

We might have used Rodgers and Rowe's (2002) nonpsychological examples as instances of good practice. Kepler, Mendel, and Watson and Crick made explicit how their theories constrained the data. Kepler derived Kepler's Laws, which are constraints. Mendel pointed out that his theory predicted certain ratios of phenotypes. Watson and Crick noted that their theory predicted Chargaff's Rules. Because these constraints were made explicit (Step 1 of the three-step process described above), it was possible to take variability (Step 2) and plausibility (Step 3) into account. In each case—Kepler, Mendel, and Watson and Crick—readers of the initial work could see that the theory narrowly constrained outcomes in cases where a much wider range of outcomes was plausible. Readers could see, in other words, that the theory could have easily failed the test. This is what many modelers do not make clear.

The first of Rodgers and Rowe's (2002) three psychological examples is "the Weber-Fechner psychophysical function" (p. 601). We are unsure what theory this refers to; functions are not theories. Fechner (1860/1912) proposed that sensation grows as the logarithm of stimulus strength, but this was apparently a definition rather than a theory because sensation strength was not independently defined. Rodgers and Rowe's reference here (Falmagne, 1974) contained several theoretical ideas, but as far as we know none of them has been influential.

Rodgers and Rowe's (2002) second psychological example is Shepard's (1987) theory of stimulus generalization, built to explain an empirical "law" (Shepard, 1987, p. 1317) that Shepard noticed. We are less sure than Rodgers and Rowe that this theory has lasting value, but the more important point is that Shepard made explicit how his theory constrains the data (it implies the empirical law). The data were given in enough detail to take variability into account, and a careful reader can see that plausible outcomes would have contradicted the theory.

Rodgers and Rowe's (2002) final psychological example is Tversky and Kahneman's (1992) prospect theory. The curve fitting to which Rodgers and Rowe referred (Tversky & Kahneman, 1992, Figures 1 and 2, pp. 310 and 311) involved an equation not derived from the theory. The theory was supported not by the good fit of this equation but by more general features of the data.

Equation fitting has never been used to support prospect theory, as far as we know (Kahneman & Tversky, 2000).

Rodgers and Rowe (2002) did readers a valuable service by providing the six examples. None of their examples, we believe, supported their case. If any examples exist of what we could not find—a theory initially supported only by goodness of fit (without the necessary additional information) that has turned out to be important—they are plainly few and far between.

Rodgers and Rowe (2002) also defended their own work (Rodgers & Rowe, 1993), which Roberts and Pashler (2000) mentioned briefly. Because the practice we criticized is so common, it is a little unfair to spend much time on any one instance. Rodgers and Rowe (1993) were more sophisticated than most modelers. They understood the problem, or came close to understanding it: "We were concerned that the model we built to explain prevalence curves for adolescent sexual development might ... explain any developmental process" (Rodgers & Rowe, 2002 p. 601). However, they did not correctly solve the problem, at least in our judgment, and why they failed is worth pointing out so that others can avoid their mistakes. "We fit a two-gender sexual development model to data for smoking onset and to several other artificial data sets... In each case the model was rejected" (Rodgers & Rowe, 2002, p. 601). One mistake is that they failed to consider variability. It must be shown that the theory could not have fit outcomes of the procedure that produced the actual results (including sample sizes). The smoking onset data came from the same survey as the sexual development data, but the sample sizes may have been different (due to nonresponses). Maybe smoking questions were answered more often than sexual questions; a larger sample size makes it easier to reject the theory. The sample sizes of the artificial data sets were not stated. Second, and perhaps more importantly, they did not consider plausibility. Whether the smoking data and the artificial data were plausible outcomes of the sexual part of the survey is unclear. Rodgers and Rowe did not present the smoking data and the artificial data, so the reader cannot decide. Their third mistake was more subtle. To show that the theory cannot fit one particular outcome (or a small number of outcomes) is not enough. The probability of any one outcome is usually very low. To show that one particular outcome would falsify the theory is only to show that one very unlikely event would do so. The data used by Rodgers and Rowe (1993) can be described as a set of 96 integers, the number of students who responded "yes" in each of the 96 cells of the data table. The sample size for each cell was about 50, so there were about 51^{96} possible outcomes. (A sample size of n allows $n + 1$ possible outcomes. For instance, with a sample size of 2, there are 3 possible outcomes: 0, 1, and 2 "yes" answers.) Any one was very unlikely, of course. To show the theory might plausibly fail, it must be shown that the theory cannot fit a range of plausible outcomes, not just 1 or 10 or 1,000.

Rodgers and Rowe's (2002) discussions of aspects of variance (p. 600) and plausibility (p. 600) are beside the point. To show that the good fits of Rodgers and Rowe (1993) support the theory requires showing it was plausible the theory might not have fit—nothing more, nothing less. The discussion of irrelevant issues is probably our fault. We failed to make clear that the three issues we raised in Roberts and Pashler (2000)—the need to determine what the theory cannot fit, to consider variability, and to consider

plausibility—were three parts of one task: showing that the theory might plausibly not have fit.

Roberts and Pashler (2000) pointed out that typical uses of good fits to support theories have lacked crucial information. This meant that this type of evidence (good fit) could be abused—used to claim more support for a theory than it deserves. We also pointed out that this type of evidence had not been the initial evidence for any enduring theory, as far as we knew. This suggested that this type of evidence *had* been abused. In spite of ample opportunity, Rodgers and Rowe (2002) have failed to come up with even one clear counterexample. This suggests that the abuse has been extensive.

References

- Falmagne, J. C. (1974). Foundations of Fechnerian psychophysics. In D. H. Krantz, R. C. Atkinson, R. D. Luce, & P. Suppes (Eds.), *Contemporary developments in mathematical psychology: Vol. II. Measurement, psychophysics, and neural information processing* (pp. 127-159). San Francisco: Freeman.
- Fechner, G. T. (1912). Elements of psychophysics: Sections VII and XIV. In B. Rand (Ed.), *The classical psychologists* (pp. 562-572). Boston: Houghton Mifflin. (Original work published 1860)
- Kahneman, D., & Tversky, A. (2000). *Choices, values, and frames*. Cambridge, England: Cambridge University Press.
- Roberts, S., & Pashler, H. (2000). How persuasive is a good fit? A comment on theory testing. *Psychological Review*, *107*, 358-367.
- Roberts, S., & Sternberg, S. (1993). The meaning of additive reaction-time effects: Tests of three alternatives. In D. E. Meyer & S. Kornblum (Eds.), *Attention and performance XIV: Synergies in experimental psychology, artificial intelligence, and cognitive neuroscience—A silver jubilee* (pp. 611-653). Cambridge, MA: MIT Press.
- Rodgers, J. L., & Rowe, D. C. (1993). Social contagion and adolescent sexual behavior: A developmental EMOSA model. *Psychological Review*, *100*, 479-510.
- Rodgers, J. L., & Rowe, D. C. (2002). Theory development should begin (but not end) with good empirical fits: A comment on Roberts and Pashler (2000). *Psychological Review*, *109*, 599-603.
- Shepard, R. N. (1987, September 11). Toward a universal law of generalization for psychological science. *Science*, *237*, 1317-1323.
- Tversky, A., & Kahneman, D. (1992). Advances in prospect theory: Cumulative representation of uncertainty. *Journal of Risk and Uncertainty*, *5*, 297-323.

Received September 24, 2001
Accepted December 17, 2001