

# EVALUATING STUDENT EVALUATIONS OF TEACHING

## MEASUREMENT AND EQUITY BIAS IN SETS AND RECOMMENDATIONS FOR REFORM

Rebecca J. Kreitzer (UNC-Chapel Hill) and Jennie Sweet-Cushman (Chatham University)

2/15/2021

Standard Evaluations of Teaching (SETs) are commonly used in critical personnel decisions in academia. Despite their ubiquity, SETs are problematic because of their known measurement and equity biases. This brief summarizes the literature on SETs and makes recommendations on how to make SETs less problematic.

### MEASUREMENT BIAS

Measurement bias occurs where variables unrelated to teaching effectiveness systematically influence the results in SETs. Common examples of [course characteristics that cause measurement bias](#) include class time, class size, if the class is a required or elective class, course difficulty, and discipline. Common examples of [individual characteristics of students that cause measurement bias](#) include the student's interest in the course material and their previous coursework.

Some of the ways in which evaluations vary:

- Classes with lighter workloads or higher grading distributions have higher scores.<sup>i</sup>
- Evaluations are lower for non-elective and quantitative courses.<sup>ii</sup>
- Evaluations are higher for upper-level, discussion-based classes compared to larger, introductory courses.<sup>iii</sup>
- Evaluations vary across disciplines; natural science courses receive the lowest scores and humanities the highest.<sup>iv</sup>
- Bringing cookies or chocolate to class increases evaluation scores.<sup>v</sup>

“[T]he best evidence – the meta-analyses of SET/learning correlations when prior learning/ability are taken into account – indicates that the SET/learning correlation is zero.”<sup>vi</sup>

Computational simulation models demonstrate that “even under ideal conditions, under ideal circumstances, even careful and judicious use of SETs to assess faculty can produce an unacceptably high error rate.”<sup>vii</sup>

In summary, evaluations are shaped by course and individual characteristics unrelated to actual instructor quality. On this basis alone, universities and colleges should reconsider the use of these evaluations in high stakes employment decisions, such as hiring and for promotion.

### EQUITY BIAS

Equity bias occurs when variables outside the instructor's control systematically influence the results. Common examples of [instructor characteristics that cause equity bias](#) include instructor's gender, race, ethnicity, accent, sexual orientation or disability. The evidence of equity bias is strongest in the qualitative comments about the course or instructor.

[Evidence of gender equity bias includes:](#)

- Male instructors are perceived as more accurate in their teaching, have more education, are less sexist, more enthusiastic, competent, organized, professional, effective, easier to understand, prompt in providing feedback, and are less penalized for being tough graders.<sup>viii</sup>
- Experimental designs that manipulate the gender of online instructors find that instructors receive lower evaluations when students believe their instructor is a woman, despite identical course delivery.<sup>ix</sup>
- Women receive lower scores in the social sciences and higher scores in the humanities.<sup>x</sup> There is no discipline where women receive higher evaluative scores compared to men in online evaluations.<sup>xi</sup>
- There is a gender affinity effect, whereby students prefer instructors of the same gender.<sup>xii</sup> This affinity is also likely the case with race, though there is not currently any research on this.

[Conforming to prescribed gender roles has a more significant effect than gender itself.](#)<sup>xiii</sup> Women faculty are rated highly for exhibiting traditionally feminine traits, like warmth and sensitivity, by male and female students.<sup>xiv</sup> Men are evaluated

positively on traditionally male traits, like perceptions of intelligence.<sup>xv</sup> These gender stereotypes harm women instructors because students prefer professors with masculine traits, yet penalize women for not conforming to feminine stereotypes.<sup>xvi</sup> Because of the conditional nature of equity bias, there are not many estimates of the size of its effect. One study found that, controlling for other factors, **female instructors received average ratings .5 standard deviation lower** than male instructors' ratings.<sup>xvii</sup>

There is far less research on equity bias in teaching evaluations for faculty of color, in part because of their severe underrepresentation in academia. **Faculty of color are evaluated worse than their male colleagues**, especially Black and Asian professors, with Black male professors fairly particularly poorly.<sup>xviii</sup> Faculty with accents and Asian last names fare worse than their native English-speaking counterparts.<sup>xix</sup> People of color are also punished for not conforming to intersectional stereotypes.<sup>xx</sup>

There is some evidence to suggest that **LGBTQ faculty fare worse** than their straight colleagues.<sup>xxi</sup> Some research indicates seniority decreases bias,<sup>xxii</sup> while other research finds that younger professors are more popular and receive higher evaluations.<sup>xxiii</sup> We know almost nothing about biases against other relevant intersectional identities, such as pregnancy or disability.<sup>xxiv</sup>

## RECOMMENDATIONS FOR BETTER USE OF STUDENT EVALUATIONS

1. **Contextualize evaluations of students' experiences, not a measure of teaching.**<sup>xxv</sup> Students cannot, and arguably should not, evaluate teaching. When contextualized as surveys of students' perceptions and experiences, they can provide useful feedback for faculty and administrators.
2. **Be proactive about increasing the validity of these assessments by improving the response rate.** A lower response rate is more likely to be unrepresentative.<sup>xxvi</sup>
3. **Administrators should interpret the results of student ratings with caution.** Student evaluations were not designed to be used as a comparative metric across faculty.<sup>xxvii</sup> Evaluations should be used to compare a faculty member's trajectory of teaching over time, and ideally, within a single course.<sup>xxviii</sup> Because the distribution of most faculty members' reviews have a negative skew,<sup>xxix</sup> administrators should look at the median or modal response, rather than the mean. Mean the distribution of evaluations is not normally distributed, means may be biased.
4. **Restrict or eliminate the use of qualitative comments, which have the strongest evidence of equity bias.** Women faculty and faculty of color are more likely to receive negative comments about personality traits, appearance, mannerisms, competence, and professionalism.<sup>xxx</sup> Instead of asking for general "comments," assessments should direct students to provide feedback in response to specific prompts.

**Qualitative comments are problematic for many reasons**, including being difficult to aggregate because of small sample sizes<sup>xxxi</sup>; they are often contradictory and not reliable;<sup>xxxii</sup> suffer from novelty bias (people are more likely to remember unexpected or uncommon comments) and negativity bias (people are more likely to remember negative information).

5. **Administrators should not rely on SETs as the sole method of assessing teaching.** There are several alternatives or supplements to SETs, including: peer observation<sup>xxxiii</sup>, comprehensive evaluations of teaching portfolios,<sup>xxxiv</sup> and reviews of course materials.<sup>xxxv</sup> While these alternatives may also be susceptible to biases, they are not systematically biased in the same way.<sup>xxxvi</sup> Several imperfect measures are better than using just one.
6. **There should be more research on interventions to reduce bias.** There are only a few articles on testing interventions to reduce equity bias. Reducing the size of the scale can mitigate gender bias.<sup>xxxvii</sup> One random control trial (RCT) finds that making students aware of biases can mitigate the gender gap in SETs,<sup>xxxviii</sup> while another RCT finds the opposite.<sup>xxxix</sup> Anecdotal evidence suggests there may be a backlash effect when under-represented groups discuss equity bias in SETs with students.

---

This policy brief provides a summary of a peer-reviewed article.

Suggested citation: Kreitzer, R. J., & Sweet-Cushman, J. (2021). Evaluating Student Evaluations of Teaching: a Review of Measurement and Equity Bias in SETs and Recommendations for Ethical Reform. *Journal of Academic Ethics*, 1-12.

- 
- <sup>i</sup> Greenwald & Gillmore, 1997; Miles & House, 2015; Rosen, 2018; Sinclair & Kunda, 2000
- <sup>ii</sup> Benton & Cashin, 2012; Boring et al, 2016; Chamberlin & Hickey, 2001; Elmore & LaPointe, 1975; Greenwald & Gillmore, 1997; Mengel et al 2018; Uttel et al, 2017, 2013
- <sup>iii</sup> Hamermesh & Parker 2005; Miles & House, 2015; Sidanius & Crane, 1989; Sporen et al, 2013; Centra & Gaubatz, 1998
- <sup>iv</sup> Basow & Montgomery, 2005; Basow & Silberg, 1987
- <sup>v</sup> Hessler et al, 2018; Youmans & Jee, 2007
- <sup>vi</sup> Uttel et al, 2017, 19
- <sup>vii</sup> Esarey & Valdes, 2020, 1
- <sup>viii</sup> Abel & Meltzer, 2007; Arbuckle & Williams, 2004; Basow, 1995; Basow & Silberg, 1987; Boring et al, 2016; Elmore & LaPointe, 1975; MacNell et al, 2015; Rivera & Tilcsik, 2019; Miller & Chamberlin, 2000; Sidanius & Crane, 1989; Sinclair & Kunda, 2000; Sprague & Massoni, 2005; Storage et al, 2016
- <sup>ix</sup> Boring et al, 2016; MacNell et al 2015
- <sup>x</sup> Basow & Montgomery, 2005
- <sup>xi</sup> Rosen, 2018
- <sup>xii</sup> Bachen et al 1999, Martin, 1984, Young et al, 2009; Basow & Silberg, 1987; Basow, 1995; Burns-Glover & Veith, 1995; Fan et al, 2019; Kaschak, 1981, 1978; Mengel et al, 2018; Bray & Howard, 1980; Centra, 2000, Rowden & Carlson, 1996
- <sup>xiii</sup> Basow & Silberg, 1987; Freeman, 1994
- <sup>xiv</sup> Wallisch & Karau, 2002
- <sup>xv</sup> Bian et al, 2017; Basow, 2000; Boring, 2017; Leslie et al, 2015
- <sup>xvi</sup> Burns-Glover & Veith, 1995; Bennett, 1982; Boring, 2017; Kierstead et al, 1988
- <sup>xvii</sup> Hamermesh & Parker, 2005
- <sup>xviii</sup> Reid 2010; Chavez & Mitchell, 2020
- <sup>xix</sup> Fan et al, 2019; Subtirelu, 2015
- <sup>xx</sup> Anderson, 2010, Anderson & Smith, 2005
- <sup>xxi</sup> Anderson & Kanner 2011; Ewing et al, 2013
- <sup>xxii</sup> Mengel et al, 2018; Wiginton et al, 1989
- <sup>xxiii</sup> Arbuckle & Williams, 2003; McPherson et al, 2009
- <sup>xxiv</sup> But see, Baker & Copp, 1997
- <sup>xxv</sup> Linse, 2017; Abrami, 2001; Arreola, 2004
- <sup>xxvi</sup> Chapman & Joines, 2017; Adams & Umbach 2012
- <sup>xxvii</sup> Franklin 2001
- <sup>xxviii</sup> Linse, 2017
- <sup>xxix</sup> Linse, 2017; Arreola, 2004; Hativa 2013a, b
- <sup>xxx</sup> Wallace et al, 2019
- <sup>xxxi</sup> Himelein, 2018
- <sup>xxxii</sup> Linse, 2017
- <sup>xxxiii</sup> Miller & Seldin, 2014
- <sup>xxxiv</sup> Centra, 2000; Seldin et al 2010
- <sup>xxxv</sup> Chism, 2007
- <sup>xxxvi</sup> Esarey & Valdes, 2020
- <sup>xxxvii</sup> Rivera & Tilcsik, 2019
- <sup>xxxviii</sup> Peterson et al, 2019
- <sup>xxxix</sup> Key and Ardoin, 2019