



Constructs, Tape Measures, and Mercury

Journal:	<i>Perspectives on Psychological Science</i>
Manuscript ID	PPS-21-412.R1
Manuscript Type:	Commentary
Date Submitted by the Author:	28-Mar-2022
Complete List of Authors:	Boykin, Malik; Brown University, Department of Cognitive, Linguistic, and Psychological Sciences
Keywords:	Individual Differences, Intergroup Relations, Application: Education, Culture / Diversity, Social Cognition
User Defined Keywords:	graduate admissions, bias, discrimination, higher education

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Constructs, Tape Measures, and Mercury

C. Malik Boykin

Department of Cognitive, Linguistic, and Psychological Sciences, Brown University

For Review Only

Abstract

This is a Lewinian field theory approach to understanding the GRE in the context of racism to contribute to the debate about whether graduate schools should remove GRE scores from admissions processes. Woo and colleagues (this issue) review the empirical literature on bias from a psychometric perspective. In this commentary, I challenge the definition of the underlying construct measured by the GRE and offer alternative definitions of what is measured. Next, drawing on an analogy from gene-wide association studies, I discuss how genomic models that predict height that are trained on data from European ancestral populations systematically underpredict the height of West Africans (Martin et al, 2017). Our access to data from tape measures, and their correlation with height, provide objective opportunities to audit our prediction. I discuss the implications of this when the criterion variable for validating the GRE is 1st-year grades. I then probe an analogy used by Woo and colleagues, where they assert that blaming the GRE for disparities in scores across groups is akin to blaming the thermometer for global warming. I describe racism as context for a field theory approach to thinking about the limitations of this misguided analogy. Lastly, I suggest pathways forward.

Constructs, Tape Measures, and Mercury

Woo, LeBreton, Keith, and Tay (2022) published an article in this journal intending to clarify the meanings of *bias* and *fairness* from a psychometric perspective. The article was framed in the context of psychology departments considering the removal of the GRE in their admissions decision process. Their expressed goal was to define *bias* and *fairness* psychometrically and to create common ground for discussion. Using results from a previously published meta-analysis, Woo and colleagues determine that the GRE is a *fair* and *unbiased* estimator of ability and state that criticizing standardized tests as discriminatory is akin to “blaming a thermometer for global warming.” Additionally, they critically examine the sometimes-limited empirical evidence about *bias* and *fairness* in six types of information used to evaluate students during the graduate admissions process including: grade point average, personal statements, resumes/CVs, letters of recommendation, interviews, and the GRE. I believe that within a world circumscribed in a commonly accepted toolkit of psychometric approaches to thinking about *bias* and *fairness* – a world that has perpetuated the standardized testing industry - this examination was well executed by Woo et al. (2022). Their review of *bias* and *fairness* in prediction regarding the GRE and other materials evaluated for graduate student admissions demonstrates a wealth of information worth reviewing. However, I take issue with the narrowness of their definitions of *bias* and *fairness*, and their choice to define the GRE and its’ subsections as tests of ability. Therefore, I find their conclusion that the GRE is a *fair* and *unbiased* estimator of *ability* to be used in graduate admissions decisions to be flawed. I argue that *choosing* these *definitions* for these terms ignores important historical context, neglects important features of the GRE data, overvalues what the GRE actually predicts, and in doing so, serves to perpetuate structural racism.

Construct Definition

We create the constructs we theorize about from our imperfect minds and then design measurements to bring these unseen attributes into view. No matter how you slice it, psychometrics and latent variable modeling approaches still cause us, as psychological scientists, to subjectively define what we intend to measure and how we interpret the obtained. Cronbach and Meehl (1955) wrote that, “A construct is some postulated attribute of people, assumed to be reflected in test performance.” Thus, the act of *defining* the GRE as a measurement of *ability* is a human choice that is susceptible to human error and human *biases*, as it is a postulation about a latent construct *assumed* to exist within people and *assumed* to be brought into view by our questions. There is the potential for bias in the choice of how to interpret the data, regardless of whether there is bias in the coefficient of any particular model. Choosing to define the GRE as a test of ability can never cross over from the realm of assumptive reasoning to the realm of factual reasoning until we no longer need constructs to attempt to elucidate portions of ability through measurement. Historically, defining what the test measures from an array of alternative choices is a site of serious differences in perspective, and as Woo and colleagues (2022) define the GRE as a test of ability, this is a major point where my perspective diverges from that of Woo et al. and many other psychometricians. I repeat the words “*choice*” and “*define*” purposefully, for they are the most important words to focus on to understand my arguments. Regardless of whether there is bias in any given statistical estimator analyzing the data we have, there is *bias* toward structural racism inherent to choosing to define the GRE as a *fair* and *unbiased* estimator of *ability*.

I would like to start by inviting readers, as I invited Woo and colleagues during my review of their manuscript, to consider what James Jones (2003) described as the Universal

CONSTRUCTS, TAPE MEASURES, & MERCURY

5

Context of Racism. Racism is inextricable from our culture, and Kurt Lewin (1941) described culture in terms of a force field that bends and warps our respective realities and impacts our behavior. In this view, racism can serve as an ever-present social reality that introduces a 3rd variable problem that compromises whether the test is measuring the same thing across groups who live diverging racial realities. This complicates whether the underlying meaning of what is measured by the test generalizes well across groups and contexts. I revisit this concept later in this commentary.

The GRE was conceived in 1936 in a legally segregated United States, an era where most colleges, let alone graduate school programs in psychology, had yet to admit their first minoritized student. Less than 100 Black people obtained a Ph.D. in psychology prior to 1966, compared to nearly 10,000 White people (Williams, 1970). This was a GRE designed by White psychometricians for White selection committees to select assumed-to-be-White (male) students. The latent constructs that underlie the GRE were first defined in this context. This was an era of both legalized racism and psychometric eugenics. The test was not designed to be devoid of cultural bias or racism; rather it was designed to be monocultural and used in the context of legally institutionalized racism. Attempts have been made to retrofit the test for fairness, but what are we missing when we try to retrofit such an instrument for fairness across groups as opposed to developing new measures to be validated, in their own terms, for the now diverse populations we're evaluating for selection?

The Power to Define

Fifty years ago, Robert Williams (1972) asked many of these same questions, among others, when he developed the Black Intelligence Test of Cultural Homogeneity (BITCH). The BITCH was designed to test aptitude couched within the language of African-American

CONSTRUCTS, TAPE MEASURES, & MERCURY

6

Vernacular English, with references relevant to Black culture and shared experiences. In an article published in the *Journal of Applied Psychology*, Mattarzo and Weins (1977) presented results from a study of applicants to a police officer selection program in Portland, Oregon wherein Black respondents unsurprisingly scored considerably higher than White respondents on the BITCH. Among several interesting findings in this study, White respondents who had more experience with Black people, Black language, and Black culture scored higher than their White counterparts who had less experiences with Black people, language, and culture. The irony here is that GRE scores in a distribution of Black people is also likely predicted by experiences with White people, White language, and White culture, which should not have any impact on measures of the latent ability in Black people. It would be racist to think that experiences with White people improved the latent ability within Black people, but I do concede that experiences with White people and White culture may help Black students navigate graduate school in predominantly White universities. However, the latent constructs defined by the GRE and discussed by Woo et al. (2022) are not defined with this kind of flexibility of meaning and I, among others, argue that they should be (Sireci, 2021). It is a choice, an irresponsible one, to define these constructs as measures of ability with meaning that generalizes across groups.

Curiously, Mattarazo and Wiens (1977) determined that since the BITCH did not correlate with Wechsler's Adult Intelligence Test (WAIS), that it was likely an inappropriate measure for selection of police officers. Here, the WAIS, by way of its historical use and representation in literature refereed by mostly White editorial boards (Roberts et al., 2020), was defined as an acceptable criterion to invalidate the use of the BITCH, as opposed to some sort of job-relevant performance outcome which was not measured. However, in an unpublished dissertation, Hammond (1987) demonstrated that the BITCH predicted education placement

CONSTRUCTS, TAPE MEASURES, & MERCURY

7

1
2
3 outcomes for low-income Black students. Here, Hammond demonstrates a similar-in-concept
4
5 predictive validity criterion for low-income Black people that many accept, logically, as criterion
6
7 for using the GRE. I am not arguing for the validity or the use of the BITCH at all, and to
8
9 conclude that I am would demonstrate a willful commitment to missing the point. I am, however,
10
11 using the BITCH to highlight the power asymmetry demonstrated by White people at the
12
13 systemic and societal level in the process of designing, defining, and validating a self-referential
14
15 construct and then imposing it on aspiring graduate students from across the world as a standard
16
17 unbiased measure of ability that generalizes to all of humanity.
18
19
20

Retrofitted for Fairness?

21
22
23
24 In 1969, members of the Association of Black Psychologists stormed the stage while
25
26 George Miller was delivering his Presidential Address at the annual conference of the American
27
28 Psychological Association to express grievances about racism inherent to standardized testing
29
30 and to the field of psychology broadly (Williams, 2008). Five years later, **for the first time in its**
31
32 **history**, the 3rd edition of the Published Professional Standards for Educational and
33
34 Psychological Testing (produced as a joint committee of the APA, NCME, and the AERA)
35
36 included an acknowledgement of the idea that standardized tests should consider the potential for
37
38 the discrimination against minoritized people and work to ensure fairness across groups (Sireci
39
40 & Randall, 2021). To reiterate, the 1954 and 1966 versions of the Standards document made no
41
42 mention of fairness or discrimination, and one could argue that they likely would never have
43
44 considered this in the absence of pressure and critique from a literal insurrection of Black
45
46 Psychologists in the 1960's and 1970's.
47
48
49
50

51
52 In a reality that is not ahistorical, it has been argued that the GRE could be defined as a
53
54 better predictor of the race of the test-taker than the 1st year grades that supposedly validate the
55
56
57
58
59
60

CONSTRUCTS, TAPE MEASURES, & MERCURY

8

1
2
3 test (Miller & Stassun, 2014). This interpretation is actually a better fit to the data. In a world
4
5 where standardized testing is still a correlational method of knowledge production, built on
6
7 subjective definitions of constructs that are used to validate each other, one could reasonably ask,
8
9 is the GRE just an elaborate way of predicting the test-taker's race? In this *definition* of the latent
10
11 variable, after controlling for race which takes up most of the variance, the validating criterion
12
13 variable is reduced to a residual prediction. Willfully ignoring this question or this reality is a
14
15 *choice*. The power of this question is amplified in historical context, where psychometrics and
16
17 the statistics that underlie latent class analyses of subpopulations performance on ability
18
19 measures were birthed from the expressed motivation to demonstrate the intellectual superiority
20
21 of White people above other groups (Saini, 2019; Thomas & Sillen, 1972). So then, the fact that
22
23 the racial stratification in scores that the GRE produces being similar to its overtly racist
24
25 predecessors, becomes a suspicious indicator that it is potentially the same wolf in sheep's
26
27 clothing. This clothing is tailored from a Professional Standards for Educational and
28
29 Psychological Testing document outlining a psychometric approach to *bias* and *fairness* written
30
31 to convince society that a testing tradition originally designed to uncover a latent variable that
32
33 demonstrated White intellectual supremacy has been sufficiently retrofitted (starting in 1974) for
34
35 fairness even though it produces a similar racial stratification (Gould, 1996; Sireci & Randall,
36
37 2021). Seriously, what are we doing?

No Tape Measure for Ability

44
45
46
47 In a potentially analogous example from genome-wide association studies (GWAS),
48
49 Martin and colleagues (2017) demonstrated models that predicting height from the human
50
51 genome do not generalize with the same accuracy across groups. Models derived from
52
53 individuals with predominately European ancestry still predicted height in West African
54
55
56
57
58
59
60

CONSTRUCTS, TAPE MEASURES, & MERCURY

9

1
2
3 populations, which is an amazing feat of science. We could stop there, the model predicts height
4
5 in both populations, we're done. However, the point of their paper is that the accuracy of the
6
7 models derived from predominantly European ancestral samples diminished in their prediction
8
9 the further away a population was from being genetically European. Here, the less ancestrally
10
11 European the population was, the shorter they were predicted to be, and West Africans were
12
13 predicted to be shorter on average than Europeans, an outcome that is completely incongruent
14
15 with reality (Gustafson & Lindenfors, 2004; N'Diaye et al., 2011). Martin et al. (2017) further
16
17 show that models trained on West African populations predicted the height of West African
18
19 people with improved accuracy over those trained on European populations. This is a profound
20
21 finding and the implications of it are still reverberating through the genomics literature.
22
23
24 Obviously, with height, we have the ability to independently validate how well our statistical
25
26 models perform, aggregated and disaggregated, by way of the tape measure — a luxury we do
27
28 not have to validate our models of ability. So then, we can demonstrate the underprediction of
29
30 the height of West African people inherent to aggregate models of height prediction derived from
31
32 heavily biased European ancestry data and avoid erroneously defining the European derived
33
34 model as a generalizable predictor for height across populations.
35
36
37
38
39

40 Luckily for descendants of West Africans, teams in the National Basketball Association
41
42 (and other basketball leagues around the world) get to consider selecting players based, in part,
43
44 on information about height obtained from tape measures as opposed to prediction models from
45
46 genome-wide association studies derived from European ancestral populations. With the GRE, it
47
48 seems our best offering for criterion validation is the comparatively weak correlations between
49
50 test scores and first year grades, and not much else, in the absence of a tape-measure equivalent
51
52 (Sedlacek, 2017). Squaring the coefficients - since we don't know the causal direction (only the
53
54
55
56
57
58
59
60

CONSTRUCTS, TAPE MEASURES, & MERCURY

10

1
2
3 temporal one since both GRE scores and grades may be indicators of privilege taken at two time
4
5 points) - yields the knowledge of a 3% to 15% known variance space about what is supposed to
6
7 be a measure of *ability* because it predicts 1st year grades but is a better predictor of race.
8
9
10 Therefore, most of what is predicting 1st year grades is unknown and we have no idea whether
11
12 removing bias from predicting 3% to 15% of the variance space generalizes to the other 85% to
13
14 97% (Cohen et al., 2003; Sternberg & Williams, 1997; Woo et. al., 2022). The correlation
15
16 between an individual's height and their tape measured height is 1 to 1 (100% for a coefficient of
17
18 1.00), thus we can confidently audit these models for biased prediction across groups.
19
20

21
22 A logical extension of Martin and colleagues' work (2017) highlights points made
23
24 recently by Jennifer Randall (2021), and historically by Robert Williams (1972) among countless
25
26 others. Essentially, from a bottom-up approach, the algorithms are finding unique components of
27
28 the genomic architecture of West Africans that better predict height in these ancestral
29
30 populations than bottom-up models derived from European ancestry skewed samples. The
31
32 European ancestry derived models are missing key components and areas of the genome for
33
34 height prediction in people descended from West Africans. So then, I again ask a question that
35
36 many before me have asked: What competencies would we measure and what items would we
37
38 create if we sought to create an appropriate test of latent ability for graduate school selection for
39
40 minoritized populations? Would this look or read anything like the GRE?
41
42
43

44
45 If we took this project seriously, would we cull sets of reading passages from the same
46
47 sources or different sources? Would we test for the same esoteric GRE vocabulary words that I
48
49 never used in graduate school or would we use a different set of words? What linguistic styles
50
51 would the passages be written in? Would we assume the same baseline cultural knowledge or
52
53 different cultural knowledge to define and assess ability? Would we test cultural flexibility?
54
55
56
57
58
59
60

CONSTRUCTS, TAPE MEASURES, & MERCURY

11

1
2
3 Would we repeat the same questions to be considered in different contexts to see how people's
4 answers shifted with new information? Would the knowledge of boating and nautical terms (e.g.,
5 flotsam and jetsam) an aspiring test-taker learns in Manhattan Prep GRE's Verbal Strategies
6 workbook be replaced with the vocabulary of critical race theory to help students prepare for
7 these new tests? Intersectionality anyone? Would these tests be timed or untimed? For the new
8 GRE Psychology Subject Tests, would we ask questions about E.L. Thorndike, Francis Galton,
9 Carl Jung, and Wilhelm Wundt (each who published explicitly racist ideas about Black people
10 and other minorities) or questions about Frantz Fanon, Frances Cecil Sumner, Albert Beckham,
11 Mamie Phipps Clark, Vonnie McLoyd, Derald Wing Sue, Rodolfo Mendoza-Denton, Joseph
12 Gone, Diane Sanchez, Stephanie Fryberg, Kevin Nadal, Mahzarin Banaji, Ramiswammi
13 Mahalingham, Belinda Campos, Sapna Cheryan, Inez Prosser, Jennifer Richeson, Bill Cross, or
14 A.Wade Boykin? Let's throw in a section on foundational papers about the psychology of
15 racism, or testing bias, or eugenics, or White privilege, and then let's validate it with criterion
16 prediction for the success of minoritized people in their 1st year in graduate school. Then the
17 follow-up question would be how well this assessment predicts *ability* or 1st year graduate school
18 grades in psychology for White students who were aspiring to pursue Ph.D. studies. Would there
19 be a gap in scores between White test-takers and the minoritized populations it was validated
20 for? Would the scores have the same meaning across groups? Would tests of differential item
21 functioning and the psychometric elimination of statistical bias in prediction of 1st year grades
22 ever make this test fair to the White students who had to show up to take it? Should we force
23 them to take it regardless of the answer to these questions? Why or why not? Think about this,
24 because it is important to consider the context of sitting at the computer to take the GRE as a
25 minoritized person and to subsequently think about the meaning of the score that is obtained for
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 consideration. It's an important consideration because presently, a Black person aspiring to be a
4
5 psychologist has to cognitively and emotionally suffer through the legitimization of eugenic racists
6
7 like Raymond Cattell on their road to a competitive GRE score.
8
9

10 Questions remain as to whether *ability* can be measured at all, or if these psychometric
11
12 approaches to measure it tell us anything worth knowing. They predict first year grades, but they
13
14 do not predict career success, innovation, or productivity. Hence, people are *able* to succeed,
15
16 innovate, and produce in science at relatively random levels of GRE measured *ability*.
17
18 Borrowing an argument from Navarro (2019), I am unconvinced that the chosen psychometric
19
20 methods of measuring *ability* or statistically *debiasing* the models are getting us any closer to
21
22 what we'd want to know to understand *ability* or how to select graduate students.
23
24
25

26 **On GRE Fairness and Race Neutrality**

27

28 The history of psychology and the psychological testing of ability is a problematic one. I
29
30 will not lay out that entire history here, but I will refer readers to Randall's (2021) arguments
31
32 about how attempts to make the GRE race neutral actually just make the test even more
33
34 culturally White and, by their broad acceptance and institutionalization for selection, structurally
35
36 racist. In short, the illusion that a race -or culture- neutral test could even exist is a privilege that
37
38 only White people experience. In general, White people and White psychologists; live in this
39
40 illusion by themselves and undertake the project of actively explaining how and why their
41
42 positionality doesn't matter to their science (Dupree & Kraus, 2021). This is done while mostly
43
44 White editorial boards of psychology journals have a track record of explaining why everyone
45
46 else's positionality relative to their research makes them and their work biased (Roberts, 2020).
47
48
49 The collective power of the White majority throughout the field and the broader power structure
50
51
52
53
54
55
56
57
58
59
60

CONSTRUCTS, TAPE MEASURES, & MERCURY

13

1
2
3 of American education imposes this illusion on everyone's reality, though most non-White
4
5 people do not subscribe to such an illusion.
6

7
8 Many Black people, and members of other minoritized groups, know of this problematic
9
10 history of psychology and testing, which means we have a healthy and history-informed
11
12 skepticism of positivism in general and of psychometrics specifically. So then, preparation for
13
14 the GRE for an aspiring graduate student who is Black, or who is from another historically
15
16 marginalized group, is a preparation for an exercise that has a very different meaning for a Black
17
18 or minoritized test taker than a White one. Would a differential item functioning examination
19
20 help us to control for such a context altering reality or the demoralization that accompanies it? A
21
22 question remains about whether these psychometric definitions of *bias* and *fairness* are a class of
23
24 mathematically artful explanations designed to maintain racial hierarchy and justify systemic
25
26 racism? Just because predominately White psychometricians tend to agree on a set of
27
28 explanations doesn't mean that they're correct, especially when the explanations are in service of
29
30 making their advantages in society sound fair and their tests seem legitimate. Once upon a time,
31
32 a completely White male Supreme Court of the United States came to an overall consensus in
33
34 deciding that separate but equal was fair for trains, schools, bathrooms, and water fountains. Our
35
36 definitions of *fairness* evolve and have shifted with time in standards for standardized testing as
37
38 well (Sireci & Randall, 2021).
39
40
41
42
43

44
45 Students who attend high schools that offer calculus to most attending students are more
46
47 likely to be affluent, are more likely to score higher on the GRE quantitative section, and are
48
49 more likely to be White – who cares? Does this score predict whether they complete their Ph.D.
50
51 program? No. Does it predict that they will publish more? No. Does it predict whether they'll ask
52
53 novel scientific questions? No. Does it predict whether they will produce more sound science?
54
55
56
57
58
59
60

CONSTRUCTS, TAPE MEASURES, & MERCURY

14

1
2
3 No. Does it predict whether they will be good teaching assistants? No. Does it predict whether
4 they'll be good psychologists? No. Does it predict whether they'll be positive contributors to
5 their lab community? No. Does it predict the applicability of their work to the real world? No.
6
7 Well, what does it predict? First year grades? Who cares? Wait, you care? Why? Because it's
8 fair? Fair to whom? Oh, fair to students who are more likely to score higher on the GRE who are
9 more likely to be affluent and/or White. And this matters because? Oh, right, because they have
10 the power to say it does. I *choose* not to participate.
11
12
13
14
15
16
17
18

On Mercury

19
20
21 Race, a social construct, is made real in the lives of minoritized people (Black, Filipino,
22 Mexican, etc.) through the political levers of institutional violence (Sidanius & Pratto, 1999),
23 segregation (Rothstein, 2017), discrimination experiences, (Lee et al., 2019), economic
24 oppression (Massey & Denton, 1993), devaluation of cultural expressions (Boykin, 1986), the
25 overestimation of progress (Kraus et al., 2019), differential opportunities in higher education,
26 (Boykin & Dupree, 2021), as well as racial biological realism in science history (Gould, 1996;
27 Guthrie, 2004; Saini, 2019; Thomas & Sillen, 1972). Intelligence testing and standardized tests
28 derive from stated attempts to show that race was biologically real and that some races were
29 smarter than others (Guthrie, 2004; Saini, 2019; Thomas & Sillen, 1972), and the legacy of this
30 exists in the minds of lay people in society (Sidanius & Pratto, 1999; Guthrie, 2004; Saini, 2019;
31 Thomas & Sillen, 1972; Zou & Cheryan, 2017). Through the process that Fields and Fields
32 (2014) describe as Racecraft, wherein something the oppressive system does (e.g. develop ability
33 tests using eugenic methodologies in a segregated society), by virtue of a linguistic sleight of
34 hand, becomes something that oppressed people are (groups with lower ability). Standardized
35 testing is an inextricable part of the cultural forcefield that helps to cause and maintain inequality
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

CONSTRUCTS, TAPE MEASURES, & MERCURY

15

1
2
3 in society. If you *choose to communicate* that a standardized test was sufficiently retrofitted for
4
5 *fairness* in the testing of *ability* while it still reflects the inequality that standardized tests were
6
7 initially designed to demonstrate, I'm simply going to request that you make better *choices*.
8
9

10
11 The points made throughout this commentary help to explain why I find Woo et al.'s
12
13 thermometer analogy to be especially misguided. In most natural, contexts a mercury
14
15 thermometer objectively measures temperature as advertised. Much like Martin et al.'s (2017)
16
17 height example, latent construct validity is not needed to measure temperature and track global
18
19 warming with a thermometer, as the results generalize across contexts. We have a tape measure
20
21 for temperature, the thermometer, like we have a tape measure for height, the actual tape
22
23 measurer. However, it is known that mercury becomes highly magnetic at super cool
24
25 temperatures (e.g., below 4 degrees Kelvin) that would render mercury thermometer readings
26
27 invalid in the presence of magnetic fields (van Delft, 2012). In an imagined world, where super
28
29 cool temperatures were more prevalent, mercury thermometers would become differentially
30
31 valid depending on the presence, strength, and pull of magnetic fields in a given context. As
32
33 mentioned earlier, Kurt Lewin explained that the social environment and cultural context can
34
35 operate as a magnetic field that differentially alters people's behavior and their interaction with
36
37 society. The lasting histories of housing segregation, school segregation, economic inequality,
38
39 and other social disparities which have been the direct result of deliberate policy choices that
40
41 have impacted marginalized Black people and other minorities inform the differential pull the
42
43 force field has on our differentiated racial realities (Bonam et al. 2015; Clark & Clark, 1939;
44
45 Lewin, 1941; Rothstein, 2017; McCall et al., 2017). As in the case with exposure to pollution
46
47 (Bonam et al., 2015), if minoritized populations were caused by the attitudes and actions of the
48
49
50
51
52
53
54
55
56
57
58
59
60

CONSTRUCTS, TAPE MEASURES, & MERCURY

16

1
2
3 majority group to live in zones where magnetic fields are present, then mercury thermometers
4
5 would be systematically invalid across these social contexts.
6

7
8 Minoritized people form competencies and even pursue doctoral studies in psychology
9
10 partially in adaptive response to the force field of the Universal Context of Racism (Jones, 2003).
11
12 From this perspective, a world where the force fields of race and stigma do not operate to
13
14 confound the measurements of the GRE and other tests born out of a tradition intended to prove
15
16 Black inferiority is also an imagined one. If this is taken into careful consideration then neither
17
18 the thermometer nor the test is being blamed for global warming or disparity. The target of
19
20 blame would be the *choice* of *defining* the test as an *unbiased* generalizable measure of *ability*
21
22 across various confounding magnetic fields in our culture that differentially pull on minoritized
23
24 people. Endorsing these definitions and then using the GRE as a selection tool is also a *choice*.
25
26 These choices serve to assert power and to deny the reality of the force field's impact and call it
27
28 fair and unbiased. These choices tacitly serve to locate the cause of group disparities in scores
29
30 within the groups and their members, or define differences in temperature at a given location as
31
32 the cause of the reading on the thermometer, when in each case the reading would be partly
33
34 attributable to the differential influence of magnetic force fields. Among the questions that
35
36 Cronbach and Meehl (1955) listed that might be of interest to psychometricians was, "To what
37
38 extent is this test of intelligence culture free?" I would modify this to ask, to what extent is this
39
40 test of ability racism free? I would then ask to what extent this test of ability is useful for
41
42 predicting things we care about? I believe these are also questions for selection committees in
43
44 psychology departments as they consider whether to use, or how to interpret, GRE scores when
45
46 evaluating prospective students from minoritized populations for opportunities to learn how to
47
48 advance our collective knowledge of psychological science.
49
50
51
52
53
54
55
56
57
58
59
60

So Now What?

Our path forward needs to include radically changing our understanding of what standardized tests are (predictors of race and grades) and potentially aren't (predictors of things that matter for selecting students). We need to be honest about why statistically controlling for race in using GRE scores as predictors makes many people so uncomfortable. Is it because statistically controlling for race is unfair, or because it relinquishes hierarchy, power, and privilege? Then we need to ask whether what remains still matters. Once we do this, the critical thing we need to do next is figure out what actually does matter. Is it research productivity, career success, taking feedback with humility, completing the doctoral program, innovation, orientation toward growth, practicality, ability to be an effective psychologist with training, or problem solving in the context of science (as opposed to figuring the area of Rhombus)? Answering the question of what matters will allow us to better capture what we truly care about and better predict what we are hopeful for when selecting students into our programs and research labs. Once we figure this out, then we can begin the process of formulating assessments that help us make these predictions more effectively and fairly. On this point, I am in full agreement with Woo and colleagues; we'll need to collect as much data as we can to best predict the outcomes we care about, outcomes that actually matter, with as little psychometric *bias* as possible in our estimators.

It is here where William Sedlacek (2017) is likely decades ahead of most of the field in his research of non-cognitive factors in student selection that predict meaningful outcomes. His focus has been on assessments of predictors such as perseverance, leadership experience, and long-range goal setting among others, which have shown to predict outcomes that are associated with long term student success (e.g. degree attainment). Illustrating the value of centering non-

CONSTRUCTS, TAPE MEASURES, & MERCURY

18

1
2
3 cognitive factors for selection criteria, Casey Miller and Keivan Stassun (2014) have used
4
5 Sedlacek's recommendations in selecting students into a master's degree program that prepares
6
7 students, predominately minoritized students who did not perform well on the GRE, for doctoral
8
9 studies in the mathematical and physical sciences. Upwards of 80% of students who have entered
10
11 their program have attained doctoral degrees and several of these students have published
12
13 innovative science in high impact journals and have made it to tenure in the professoriate. Their
14
15 program uses non-cognitive factors to select students into opportunities that GRE-based selection
16
17 criteria would have missed and demonstrates the possibility to better predict key outcomes
18
19 expected from doctoral student success.
20
21
22

23
24 Similarly, in applied mathematics, Carlos Castillo-Chavez and his colleagues ran a
25
26 summer program for 24 years that provided minoritized students the opportunities to collaborate
27
28 and apply scientific problem solving to real-world problems of their choosing – with scaffolding
29
30 (Castillio-Chavez et al., 2017; Castillo-Garsaw et al, 2013). Their program sent the majority of
31
32 their students on the trajectory toward graduate studies and produced many scientific innovations
33
34 and publications. Castillo-Chavez sought “diamonds in the rough” in his selections criteria.
35
36 Castillo-Chavez and his colleagues program used non-cognitive factors, such as indicators of
37
38 sustained interest, reviewing grades in key math courses, and multiple readers per application in
39
40 the selection process (Christopher Kribs, March 13th, 2022 - personal communication). Stassun
41
42 and Castillo-Chavez, a physicist and an applied mathematician, are demonstrating that we can
43
44 assess ability to succeed in scientific careers without considering the GRE.
45
46
47
48

49 Similar models already exist in psychology graduate programs if we *choose* to learn from
50
51 and emulate them, which I argue that we should. In his 41 years as a professor at Howard
52
53 University, my father, A. Wade Boykin, mentored an endless stream of Black (and other
54
55
56
57
58
59
60

CONSTRUCTS, TAPE MEASURES, & MERCURY

19

1
2
3 minoritized) doctoral students in psychology. This includes selecting and then mentoring
4
5 students who demonstrated *ability* through scholarly productivity and success in rising to the
6
7 ranks of tenured associate and full professors at universities such as the University of Michigan,
8
9 Bowie State University, Pomona College, Smith College, the University of Kentucky, North
10
11 Carolina A&T, and North Carolina Central University, among others, without ever ceding a
12
13 shred of merit to GRE scores. Much like Carlos Castillo-Chavez, my father sought diamonds in
14
15 the rough. Additionally, for 5 years, and with great success, Michael Kraus has run an internship
16
17 program at Yale University for aspiring doctoral students in psychology and organizational
18
19 behavior that has provided an on-ramp for minoritized students to pursue doctoral studies. He
20
21 focused the selection criteria on the quality and kinds of scientific questions students wanted to
22
23 ask and answer, their passion for wanting to pursue the questions, and indicators that they would
24
25 persevere through the humbling challenges of research. Kraus has used similar criteria to select
26
27 his several productive graduate students who have published numerous articles in high impact
28
29 journals while wholly ignoring GRE scores. Regarding GRE use in selection decisions, Kraus
30
31 simply stated “I don’t care about triangles or know what memorizing the Pythagorean theorem
32
33 has to do with psychology (personal communication, March, 16th, 2022).” I could not agree more
34
35 and I hope others will join me in asking similar questions about the value of these standardized
36
37 tests.
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

References

- Bonam, C. M., Bergsieker, H. B., & Eberhardt, J. L. (2016). Polluting black space. *Journal of Experimental Psychology: General*, *145*(11), 1561.
- Boykin, A. W. (1986). The triple quandary and the schooling of Afro-American children. In U. Neisser (Ed.), *The school achievement of minority children* (pp. 57-93). Lawrence Erlbaum Associates.
- Castillo-Garsow, C., Castillo-Chavez, C., & Woodley, S. (2013). A preliminary theoretical analysis of a research experience for undergraduates community model. *PRIMUS*, *23*(9), 860-880.
- Castillo-Chavez, C., Kribs, C., & Morin, B. (2017). Student-driven research at the mathematical and theoretical biology institute. *The American Mathematical Monthly*, *124*(9), 876-892.
- Clark, K. B., & Clark, M. K. (1939). Segregation as a factor in the racial identification of Negro pre-school children: A preliminary report. *The Journal of Experimental Education*, *8*(2), 161-163.
- Cohen, J.C., Cohen, P., West, S.G., & Aiken, L.S. (2003) *Applied multiple regression/correlation analysis for the behavioral sciences*. Routledge.
- Cronbach, L. J., & Meehl, P. E. (1955). Construct validity in psychological tests. *Psychological Bulletin*, *52*(4), 281.
- Dupree, C. H., & Boykin, C. M. (2021). Racial inequality in academia: Systemic origins, modern challenges, and policy recommendations. *Policy Insights from the Behavioral and Brain Sciences*, *8*(1), 11-18.
- Dupree, C. H., & Kraus, M. W. (2020). Psychological science is not race neutral. *Perspectives on Psychological Science*, 1745691620979820.

CONSTRUCTS, TAPE MEASURES, & MERCURY

21

Fields, K. E., & Fields, B. J. (2014). *Racecraft: The soul of inequality in American life*. Verso Books.

Gould, S. J., & Gold, S. J. (1996). *The mismeasure of man*. WW Norton & company.

Gustafsson, A., & Lindenfors, P. (2004). Human size evolution: no evolutionary allometric relationship between male and female stature. *Journal of human evolution*, 47(4), 253-266.

Guthrie, R. V. (2004). *Even the rat was white: A historical view of psychology* (2nd ed.). Pearson Education.

Hammond, D. P. (1987). *Accuracy of prediction of the Culture Fair Intelligence Test, Short Form Test of Academic Aptitude, and Black Intelligence Test of Cultural Homogeneity of the current educational functioning of low-income black students* (Doctoral dissertation, Southern Illinois University at Carbondale).

Jones, J. M. (2003). TRIOS: A psychological theory of the African legacy in American culture. *Journal of Social Issues*, 59(1), 217-242.

Kraus, M. W., Onyeador, I. N., Daumeyer, N. M., Rucker, J. M., & Richeson, J. A. (2019). The misperception of racial economic inequality. *Perspectives on Psychological Science*, 14(6), 899-921.

Lee, R. T., Perez, A. D., Boykin, C. M., & Mendoza-Denton, R. (2019). On the prevalence of racial discrimination in the United States. *PloS one*, 14(1), e0210698.

Lewin, K. (1941). Self-hatred among Jews. In K. Lewin (Ed.), *Resolving social conflicts: Selected papers on group dynamics* (pp. 186-200). New York, NY: Harper & Brothers.

CONSTRUCTS, TAPE MEASURES, & MERCURY

22

1
2
3 Lewin, K. (1946). Action research and minority problems. *Journal of Social Issues*, 2, 34–46.

4
5 Reprinted in G. W. Lewin (1997) (Ed.), *Resolving social conflicts: Selected papers on*
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

group dynamics (pp. 143–152). Washington, DC: American Psychological Association.

Martin, A. R., Gignoux, C. R., Walters, R. K., Wojcik, G. L., Neale, B. M., Gravel, S., ... &
Kenny, E. E. (2017). Human demographic history impacts genetic risk prediction across
diverse populations. *The American Journal of Human Genetics*, 100(4), 635-649.

Massey, D., & Denton, N. A. (1993). *American apartheid: Segregation and the making of the*
underclass. Harvard university press.

McCall, L., Burk, D., Laperrière, M., & Richeson, J. A. (2017). Exposure to rising inequality
shapes Americans' opportunity beliefs and policy support. *Proceedings of the National*
Academy of Sciences, 114(36), 9593-9598.

Matarazzo, J. D., & Wiens, A. N. (1977). Black Intelligence Test of Cultural Homogeneity and
Wechsler Adult Intelligence Scale scores of Black and White police applicants. *Journal*
of Applied Psychology, 62(1), 57.

Miller, C. & Stassun, K. (2014). A test that fails. *Nature*, 510(7504), 303-304.

Navarro, D. J. (2019). Between the Devil and the Deep Blue Sea: Tensions Between Scientific
Judgment and Statistical Model Selection. *Computational Brain & Behavior*, 2, 28-34.

N'Diaye, A., Chen, G. K., Palmer, C. D., Ge, B., Tayo, B., Mathias, R. A., ... & Haiman, C. A.
(2011). Identification, replication, and fine-mapping of Loci associated with adult height
in individuals of african ancestry. *PLoS genetics*, 7(10), e1002298.

Randall, J. (2021). “Color-neutral” is not a thing: Redefining construct definition and
representation through a justice-oriented critical antiracist lens. *Educational*
Measurement: Issues and Practice.

CONSTRUCTS, TAPE MEASURES, & MERCURY

23

- 1
2
3 Roberts, S. O., Bareket-Shavit, C., Dollins, F. A., Goldie, P. D., & Mortenson, E. (2020). Racial
4
5 inequality in psychological research: Trends of the past and recommendations for the
6
7 future. *Perspectives on Psychological Science*, *15*(6), 1295-1309.
8
9
- 10 Rothstein, R. (2017). *The color of law: A forgotten history of how our government segregated*
11
12 *America*. Liveright Publishing.
13
- 14 Saini, A. (2019). *Superior: the return of race science*. Beacon Press.
15
- 16 Sedlacek, W. (2017). *Measuring noncognitive variables: Improving admissions, success and*
17
18 *retention for underrepresented students*. Stylus Publishing, LLC.
19
- 20 Sidanius, J., & Pratto, F. (1999). *Social dominance: An intergroup theory of social hierarchy and*
21
22 *oppression*. New York, NY: Cambridge University Press.
23
- 24 Sireci, S. G. (2021). NCME Presidential address 2020: Valuing educational measurement.
25
26 *Educational Measurement: Issues and Practice*, *40*(1), 7-16. DOI: 10.1111/emip.12415.
27
28
- 29 Sireci, S. G., & Randall, J. (2021). Evolving notions of fairness in testing in the United States.
30
31 *In The History of Educational Measurement* (pp. 111-135). Routledge.
32
33
- 34 Sternberg, R. J., & Williams, W. M. (1997). Does the Graduate Record Examination predict
35
36 meaningful success in the graduate training of psychology? A case study. *American*
37
38 *Psychologist*, *52*(6), 630–641. <https://doi.org/10.1037/0003-066X.52.6.630>
39
40
- 41 Thomas, A. & Sillen, S. (1972) *Racism and psychiatry*. New York, NY: Brunner/Mazel.
42
- 43 van Delft, D. (2012). History and significance of the discovery of superconductivity by
44
45 Kamerlingh Onnes in 1911. *Physica C: Superconductivity*, *479*, 30-35.
46
47
- 48 Williams, R. L. (1970). Black pride, academic relevance, and individual achievement.
49
50 *Counseling Psychologist*, *2*(1), 18-22.
51
52
53
54
55
56
57
58
59
60

CONSTRUCTS, TAPE MEASURES, & MERCURY

24

1
2
3 Williams, R. L. (1972). The BITCH test (Black Intelligence Test of Cultural Homogeneity). St.

4
5 Louis: Williams & Associates.

6
7 Williams, R. L. (2008). *History of the Association of Black Psychologists: Profiles of*

8
9
10 *outstanding Black psychologists*. AuthorHouse.

11
12 Woo, S. E., LeBreton, J., Keith, M., Tay, L. (2022) Bias, fairness, and validity in graduate

13
14 admissions: A psychometric perspective. *Perspectives on Psychological Science*.

15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Review Only