

ALGORITHMIC ACCOUNTABILITY:
A LEGAL AND ECONOMIC FRAMEWORK

Robert P. Bartlett, III^{*}
Adair Morse[‡]
Nancy Wallace[†]
Richard Stanton[°]

Abstract

Despite the potential for machine learning and artificial intelligence to reduce face-to-face bias in decision-making, a growing chorus of scholars and policymakers have recently voiced concerns that if left unchecked, algorithmic decision-making can also lead to unintentional discrimination against members of historically marginalized groups. These concerns are being expressed through Congressional subpoenas, regulatory investigations, and an increasing number of algorithmic accountability bills pending in both state legislatures and Congress. To date, however, prominent efforts to define policies whereby an algorithm can be considered accountable have tended to focus on output-oriented policies and interventions that either may facilitate illegitimate discrimination or involve fairness corrections unlikely to be legally valid.

We provide a workable definition of algorithmic accountability that is rooted in the caselaw addressing statistical discrimination in the context of Title VII of the Civil Rights Act of 1964. Using instruction from the burden-shifting framework, codified to implement Title VII, we formulate a simple statistical test to apply to the design and review of the inputs used in any algorithmic decision-making processes. Application of the test, which we label the *input accountability test*, constitutes a legally viable, deployable tool that can prevent an algorithmic model from systematically penalizing members of protected groups who are otherwise qualified in a target characteristic of interest.

^{*} I. Michael Heyman Professor of Law & Faculty Co-Director of the Berkeley Center for Law and Business - UC Berkeley School of Law.

[‡] Associate Professor & Solomon P. Lee Chair in Business Ethics – UC Berkeley Haas School of Business.

[†] Professor & Lisle and Roslyn Payne Chair in Real Estate Capital Markets, Co-Chair, Fisher Center for Real Estate and Urban Economics – UC Berkeley Haas School of Business.

[°] Professor & Kingsford Capital Management Chair in Business Economics – UC Berkeley Haas School of Business.

ALGORITHMIC ACCOUNTABILITY:
A LEGAL AND ECONOMIC FRAMEWORK

TABLE OF CONTENTS

	<u>Page</u>
I. INTRODUCTION	2
II. ACCOUNTABILITY UNDER U.S. ANTIDISCRIMINATION LAW	9
A. <i>Accountability and the Burden-Shifting Framework of Title VII</i>	9
B. <i>The Input Accountability Test (IAT) Versus Outcome-Oriented Approaches</i>	14
C. <i>The Input Accountability Test Versus the ‘Least Discriminatory Alternative’ Test</i>	19
D. <i>The Input Accountability Test Versus HUD’s Mere Predictive Test</i> . 24	24
III. THE INPUT ACCOUNTABILITY TEST	25
A. <i>The Test</i>	26
B. <i>The Test in Regression Form</i>	30
C. <i>Challenges in Implementing the Test</i>	33
i. <i>Unobservability of the Target Variable</i>	33
ii. <i>Measurement Error in the Target</i>	34
iii. <i>Testing for ‘Not Statistically Correlated’</i>	38
iv. <i>Nonlinearities or Interactions Among Proxies</i>	40
D. <i>Simulation</i>	40
i. <i>Set Up</i>	40
ii. <i>Applying the Input Accountability Test</i>	42
IV. APPLICATIONS BEYOND EMPLOYMENT	46
A. <i>Domains with Court-Defined Business Necessity Targets</i>	47
B. <i>Domains Without Court-Defined Business Necessity Targets</i>	49
C. <i>Self-Determining Business Necessity</i>	50
V. CONCLUSION.....	52
APPENDIX	54

I. INTRODUCTION

In August 2019, Apple Inc. debuted its much-anticipated Apple Card, a no fee, cash-rewards credit card “designed to help customers lead a healthier financial life.”¹ Within weeks of its release, Twitter was abuzz with headlines that the card’s credit approval algorithm was systematically biased against

¹ Press Release, Apple Inc., *Introducing Apple Card, A New Kind of Credit Card Created by Apple* (March 25, 2019), <https://www.apple.com/newsroom/2019/03/introducing-apple-card-a-new-kind-of-credit-card-created-by-apple/>.

women.² Even Apple co-founder Steve Wozniak weighed in, tweeting that the card gave him a credit limit that was ten times higher than what it gave his wife, despite the couple sharing all their assets.³ In the days that followed, Goldman Sachs—Apple’s partner in designing the Apple Card—steadfastly defended the algorithm, insisting that “we have not and will not make decisions based on factors on gender.”⁴ Yet doubts persisted. By November, the New York State Department of Financial Services had announced an investigation into the card’s credit approval practices.⁵

Around that same time, buzz spread across the media about another algorithm, that of health insurer UnitedHealth.⁶ The algorithm was used to inform hospitals about patients’ level of sickness so that hospitals could more effectively allocate resources to the sickest patients. However, an article appearing in *Science* showed that because the company used cost of care as the metric for gauging sickness and because African-American patients historically incurred lower costs for the same illnesses and level of illness, the algorithm caused them to receive substandard care as compared to white patients.⁷

Despite the potential for algorithmic decision-making to eliminate face-to-face biases, these episodes provide vivid illustrations of the widespread concern that algorithms may nevertheless engage in objectionable discrimination.⁸ Indeed, a host of regulatory reforms have emerged to contend with this challenge. For example, New York City has enacted an algorithm accountability law, which creates a task force to recommend procedures for determining whether automated decisions by city agencies disproportionately impact protected groups.⁹ Likewise, the Washington State House of Representatives introduced an algorithm accountability bill, which would require the state’s chief information officer assess whether any automated decision system used by a state agency “has a known bias, or is untested for

² See Sridhar Natarajan & Shahien Nasiripour, *Viral Tweet About Apple Card Leads to Goldman Sachs Probe*, BLOOMBERG (Nov. 19, 2019), <https://www.bloomberg.com/news/articles/2019-11-09/viral-tweet-about-apple-card-leads-to-probe-into-goldman-sachs>.

³ See Isobel Asher Hamilton, *Apple Cofounder Steve Wozniak Says Apple Card Offered His Wife a Lower Credit Limit*, BUSINESSINSIDER (Nov. 11, 2019), <https://www.businessinsider.com/apple-card-sexism-steve-wozniak-2019-11>.

⁴ *Id.*

⁵ See Neil Vigdor, *Apple Card Investigated After Gender Discrimination Complaints*, NY TIMES (Nov. 10, 2019), <https://www.nytimes.com/2019/11/10/business/apple-credit-card-investigation.html>.

⁶ Melanie Evans & Anna Wilde Mathews, *New York Regulator Probes UnitedHealth Algorithm for Racial Bias*, WSJ (Oct. 26, 2019), <https://www.wsj.com/articles/new-york-regulator-probes-unitedhealth-algorithm-for-racial-bias-11572087601>

⁷ Ziad Obermeyer, Brian Powers, Christine Vogeli, and Sendhil Mullainathan, *Dissecting Racial Bias in an Algorithm Used to Manage the Health of Populations*, 366 SCIENCE 447 (2019).

⁸ See, e.g., Salon Barocas & Andrew D. Selbst, *Big Data’s Disparate Impact*, 104 CAL L. REV. 671, 673 (2016) (“If data miners are not careful, the process can result in disproportionately adverse outcomes concentrated within historically disadvantaged groups in ways that look a lot like discrimination.”).

⁹ See Zoë Bernard, *The First Bill to Examine ‘Algorithmic Bias’ in Government Agencies Has Just Passed in New York City*, BUSINESSINSIDER (Dec. 19, 2017), <http://www.businessinsider.com/algorithmic-bias-accountability-bill-passes-in-new-york-city-2017-12?IR=T>.

bias.”¹⁰ Federally, the Algorithmic Accountability Act of 2019, which is currently pending in Congress, would require large companies to audit their algorithms for “risks that [they] may result in or contribute to inaccurate, unfair, biased, or discriminatory decisions impacting consumers.”¹¹

Yet, a notable absence in these legislative efforts is a formal standard for courts or regulators to deploy in evaluating algorithmic decision-making, raising the fundamental question: *What exactly does it mean for an algorithm to be accountable?* The urgency of this question follows from the meteoric growth in algorithmic decision-making, spawned by the availability of unprecedented data on individuals and the accompanying rise in techniques in machine learning and artificial intelligence.¹²

In this Article, we provide an answer to the pressing question of what accountability is, and we put forward a workable test that regulators, courts, and data scientists can apply in examining whether an algorithmic decision-making process complies with long-standing antidiscrimination statutes and caselaw. Central to our framework is the recognition that, despite the novelty of artificial intelligence and machine learning, existing U.S. antidiscrimination law has long provided a workable definition of accountability dating back to Title VII of the Civil Rights Act of 1964.¹³

Title VII and the caselaw interpreting it define what it means for any decision-making process—whether human or machine—to be accountable under U.S. antidiscrimination law. At the core of this caselaw is the burden-shifting framework initially articulated by the Supreme Court in *Griggs v. Duke Power Co.*¹⁴ Under this framework, plaintiffs putting forth a claim of unintentional discrimination under Title VII must demonstrate that a particular decision-making practice (e.g., a hiring practice) lands disparately on members of a protected group.¹⁵ If successful, the framework then demands that the burden shift to the defendant to show that the practice is “consistent with business necessity.”¹⁶ If the defendant satisfies this requirement, the burden returns to the plaintiff to show that an equally valid and less discriminatory practice was available that the employer refused to use.¹⁷ The focus of Title VII is on discrimination in the workplace, but the analytical framework that emerged from the Title VII context now spans

¹⁰ House Bill 1655, 66th Leg., Reg. Sess. (Wash. 2019), <http://lawfilesex.leg.wa.gov/biennium/2019-20/Htm/Bills/House%20Bills/1655-S.htm>.

¹¹ H.R. 2231, 116th Cong. (2019).

¹² See C. Scott Hemphill, *Disruptive Incumbents: Platform Competition in an Age of Machine Learning*, 119 *COL. L. REV.* 1973, 1975-1979 (2019) (surveying rapid deployment of machine learning technologies).

¹³ 42 U.S.C. § 2000e (2012).

¹⁴ *Griggs v. Duke Power Co.*, 401 U.S. 424, 432 (1971).

¹⁵ See *Dothard v. Rawlinson*, 433 U.S. 321, 329 (1977).

¹⁶ 42 U.S.C. § 2000e-2(k); see also *Griggs*, 401 U.S. at 431 (in justifying employment practice that produces disparate impact, [t]he touchstone is business necessity”).

¹⁷ See *Albemarle Paper Co. v. Moody*, 422 U.S. 405, 425 (1975).

other domains and applies directly to the type of unintentional, statistical discrimination utilized in algorithmic decision-making.¹⁸

Despite the long tradition of applying this framework to cases of statistical discrimination, it is commonly violated in the context of evaluating the discriminatory impact of algorithmic decision-making. Instead, for many, the legality of any unintentional discrimination resulting from an algorithmic model is presumed to depend on simply the accuracy of the model—that is, the ability of the model to predict a characteristic of interest (e.g., productivity or credit risk) generally referred to as the model’s “target.”¹⁹ An especially prominent example of this approach appears in the Department of Housing and Urban Development’s 2019 proposed rule revising the application of the disparate impact framework under the Fair Housing Act (FHA) for algorithmic credit scoring.²⁰ The proposed rule provides that, after a lender shows that the proxy variables used in an algorithm do not substitute for membership in protected group, the lender may defeat a discrimination claim by showing that the model is “predictive of risk or other valid objective.”²¹ Yet this focus on predictive accuracy ignores how courts have applied the *Griggs* framework in the context of statistical discrimination.

To see why, consider the facts of the Supreme Court’s 1977 decision in *Dothard v. Rawlinson*.²² There, a prison system desired to hire job applicants who possessed a minimum level of strength to perform the job of a prison guard, but the prison could not directly observe which applicants satisfied this requirement.²³ Consequently, the prison imposed a minimum height and weight requirement on the assumption that these observable characteristics were correlated with the requisite strength required for the job.²⁴ In so doing, the prison was thus engaging in statistical discrimination: It was basing its hiring decision on the statistical correlation between observable proxies (an applicant’s height and weight) and the unobservable variable of business necessity (an applicant’s job-required strength).

¹⁸ For example, this general burden-shifting framework has been extended to other domains where federal law acknowledges the possibility for claims of unintentional discrimination under a disparate impact theory. *See, e.g.*, *Texas Department of Housing and Community Affairs v. Inclusive Communities Project, Inc.*, 135 S. Ct. 2507 (2015), (adopting burden-shifting framework for disparate impact claims under the Fair Housing Act); *Ferguson v. City of Charleston*, 186 F.3d 469, 480 (4th Cir. 1999) (discussing cases adopting the Title VII burden-shifting framework in Title VI disparate impact cases), *rev’d on other grounds*, 532 U.S. 67 (2001).

¹⁹ *See infra* Part 2(C).

²⁰ *See* Department of Housing and Urban Development, *HUD’s Implementation of the Fair Housing Act’s Disparate Impact Standard*, 84 FR 42,854 (August 19, 2019) [hereinafter “2019 HUD Proposal”]. The rulemaking was intended to amend HUD’s interpretation of the disparate impact standard “to better reflect” the Supreme Court’s 2015 ruling in *Texas Department of Housing and Community Affairs v. Inclusive Communities Project, Inc.*, 135 S. Ct. 2507 (2015), which upheld the ability of plaintiffs to bring disparate impact cases of discrimination under the FHA.

²¹ *Id.*

²² 433 U.S. 321 (1977).

²³ *Id.* at 331-32.

²⁴ *Id.*

Because this procedure resulted in adverse hiring outcomes for female applicants, a class of female applicants brought suit under Title VII for gender discrimination.²⁵ Deploying the burden-shifting framework, the Supreme Court first concluded that the plaintiffs satisfied the disparate outcome step,²⁶ and it also concluded that the prison had effectively argued that hiring applicants with the requisite strength could constitute a business necessity.²⁷ However, the Court ultimately held that the practice used to discern strength—relying on the proxy variables of height and weight—did not meet the “consistent with business necessity” criteria.²⁸ Rather, absent evidence showing the precise relationship between the height and weight requirements to “the requisite amount of strength thought essential to good job performance,”²⁹ height and weight were noisy estimates of strength that risked penalizing females over-and-above these variables’ relation to the prison’s business necessity goal. In other words, height and weight were likely to be biased estimates of required strength whose use by the prison risked systematically penalizing female applicants who were, in fact, qualified.

The Court thus illustrated that in considering a case of statistical discrimination, the “consistent with business necessity” step requires the assessment of two distinct questions. First, is the unobservable “target” characteristic (e.g., requisite strength) one that can justify disparities in hiring outcomes across members of protected and unprotected groups? Second, even with a legitimate target variable, are the proxy “input” variables used to predict the target noisy estimates that are biased in a fashion that will systematically penalize members of a protected group who are otherwise qualified? In this regard, the Court’s holding echoes the long-standing prohibition against redlining in credit markets. A lender who engages in redlining refuses to lend to residents of a majority-minority neighborhood on the assumption that the average unobservable credit risk of its residents is higher than those of observably-similar but non-minority neighborhoods.³⁰ Yet while differences in creditworthiness can be a legitimate basis for racial or ethnic disparities to exist in lending under the FHA,³¹ courts have consistently held that the mere fact that one’s neighborhood is correlated with predicted credit risk is insufficient to justify red-lining.³² By assuming that all residents

²⁵ *Id.* at 323.

²⁶ *Id.* at 330-31.

²⁷ *Id.* at 332.

²⁸ *Id.*

²⁹ *Id.*

³⁰ The term red-lining derives from the practice of loan officers evaluating home mortgage applications based on a residential map where integrated and minority neighborhoods are marked off in red as poor risk areas. Robert G. Schwemm, *Housing Discrimination* 13–42 (Release # 5, 1995).

³¹ See *infra* note 170.

³² See *Laufman v. Oakley Building & Loan Company*, 408 F. Supp. 489 (S.D. Ohio 1976)(redlining on the basis of race violates the “otherwise make unavailable or deny” provision of § 3604(a) of the FHA); *Wai v. Allstate Ins. Co.*, 75 F. Supp. 2d 1, 7 (D.D.C. 1999)(interpreting identical language in § 3604(f)(2)

of minority neighborhoods have low credit, redlining systematically penalizes minority borrowers who actually have high credit worthiness.

These two insights from *Dothard*—that statistical discrimination must be grounded in the search for a legitimate target variable and that the input proxy variables for the target cannot systematically discriminate against members of a protected group who are qualified in the target—remain as relevant in today’s world of algorithmic decision-making as they were in 1977. The primary task for courts, regulators, and data scientists is to adhere to them in the use of big data implementations of algorithmic decisions (e.g., in employment, performance assessment, credit, sentencing, insurance, medical treatment, college admissions, advertising, etc.).

Fortunately, the caselaw implementing the Title VII burden-shifting framework, viewed through basic principles of statistics, provides a way forward. This is our central contribution: We recast the logic that informs *Dothard* and courts’ attitude towards redlining into a formal statistical test that can be widely deployed in the context of algorithmic decision-making. We label it the *Input Accountability Test (IAT)*.

As we show, the IAT provides a simple and direct diagnostic that a data scientist or regulator can apply to determine whether an algorithm is accountable under U.S. antidiscrimination principles. For instance, a statistician seeking to deploy the IAT could do so by turning to the same training data that she used to calibrate the predictive model of a target. In settings such as employment or lending where courts have explicitly articulated a legitimate business target (e.g., a job required skill or creditworthiness),³³ the first step would be determining that the target is, in fact, a business necessity variable. Second, taking a proxy variable (e.g., height) that her predictive model utilizes, she would next decompose the proxy’s variation across individuals into that which correlates with the target variable and an error component. Finally, she would test whether that error component remains correlated with the protected category (e.g., gender). If a proxy used to predict a legitimate target variable is unbiased with respect to a protected group, it will pass the IAT, even if the use of the proxy disparately impacts members of protected groups. In this fashion, the test provides a concrete method to harness the benefits of statistical discrimination with regard to predictive accuracy while avoiding the risk that it systematically

of the FHA as prohibiting insurance redlining); Laufman, 408 F. Supp. at 496–497 (mortgage redlining); Nationwide Mut. Ins. Co. v. Cisneros, 52 F.3d 1351 (6th Cir. 1995)(insurance redlining); American Family Mut. Ins., 978 F.2d at 297 (insurance redlining); Lindsey v. Allstate Ins. Co., 34 F. Supp. 2d 636, 641–643 (W.D. Tenn. 1999)(same); Strange v. Nationwide Mut. Ins. Co., 867 F. Supp. 1209, 1213–1214 (E.D. Pa. 1994)(same). The regulatory agencies charged with interpreting and enforcing the lending provisions of the FHA have defined redlining to include “the illegal practice of refusing to make residential loans or imposing more onerous terms on any loans made because of the predominant race, national origin, etc. of the residents of the neighborhood in which the property is located. Redlining violates both the FHA and ECOA.” Joint Policy Statement on Discrimination in Lending, 59 Fed. Reg. 18266 (1994).

³³ See Part 4(A).

penalizes members of a protected group who are, in fact, qualified in the target characteristic of interest.

We provide an illustration of the IAT in the *Dothard* setting, not only to provide a clear depiction of the power of the test, but also to introduce several challenges in implementing it and suggested solutions. These challenges include multiple incarnations of measurement error in the target, as exemplified by the UnitedHealth use of cost as a target, rather than the degree of illness, mentioned previously. These challenges also include understanding what “significantly correlated” means in our era of massive datasets. We offer an approach that may serve as a way forward. Beyond the illustration, we also provide a simulation of the test using a randomly constructed training dataset of 800 prison employees.

Finally, we illustrate how the IAT can be deployed by courts, regulators, and data scientists. In addition to employment, we list a number of other sectors – including credit, parole determination, home insurance, school and scholarship selection, and tenant selection – where either caselaw or statutes have provided explicit instructions regarding what can constitute a legitimate business necessity target.³⁴ We also discuss other domains such as automobile insurance and health care where claims of algorithmic discrimination have recently surfaced, but where existing discrimination laws are less clear whether liability can arise for unintentional discrimination. Businesses in these domains are thus left to self-regulating and have generally professed to adhering to non-discriminatory business necessity targets.³⁵ For firms with an express target delineation (whether court-formalized or self-imposed), our IAT provides a tool to pre-test their models.

We highlight, however, that firm profit margins and legitimate business necessity targets can easily be confounded in the design of machine learning algorithms, especially in the form of exploiting consumer demand elasticities (e.g., profiling consumer shopping behavior).³⁶ In lending, for instance, courts have repeatedly held that creditworthiness is the sole business necessity target that can justify outcomes that differ across protected and unprotected groups.³⁷ Yet, newly-advanced machine learning techniques make it possible to use alternative targets, such as a borrower’s proclivity for comparing loan products, that focus on a lender’s profit margins in addition to credit risk. In other work, we provide empirical evidence consistent with FinTech algorithms’ engaging in such profiling, with the result that minority borrowers face higher priced loans, holding constant the price impact of borrowers’ credit risk.³⁸ As such, these findings illustrate how the incentive

³⁴ *See Id.*

³⁵ *See* Part 4(B).

³⁶ *See* Part 4(C).

³⁷ *See infra* note 170.

³⁸ Robert Bartlett, Adair Morse, Richard Stanton, and Nancy Wallace, Consumer Lending Discrimination in the FinTech Era, NBER Working Paper No. 25943, available at <https://www.nber.org/papers/w25943>.

of firms to use shopping behavior as a target can lead to discrimination in lending—a practice that could be detected by application of the IAT.³⁹ Profiling for shopping behavior is a subject applicable to many settings beyond the lending context and a leading topic for future research and discourse.

Our approach differs from other approaches to “algorithmic fairness” that focus solely on ensuring fair outcomes across protected and unprotected groups.⁴⁰ As we show, by failing to distinguish disparities that arise from a biased proxy from those disparities that arise from the distribution of a legitimate target variable, these approaches can themselves run afoul of U.S. antidiscrimination law. In particular, following the Supreme Court’s 2009 decision in *Ricci v. DeStefano*,⁴¹ efforts to calibrate a decision-making process to equalize outcomes across members of protected and unprotected groups—regardless of whether individuals are qualified in a legitimate target of interest—are likely to be deemed impermissible intentional discrimination.⁴²

This Article proceeds as follows. In Part 2, we begin by articulating a definition for algorithmic accountability that is at the core of our input accountability test. As we demonstrate there, our definition of algorithmic accountability is effectively a test for “unbiasedness,” which differs from various proposals for “algorithmic fairness” that are commonly found in the statistics and computer science literatures. Building on this definition of algorithmic accountability, Part 3 formally presents the IAT. The test is designed to provide a workable tool for data scientists and regulators to use to distinguish between legitimate and illegitimate discrimination. The test is directly responsive to the recent regulatory and legislative interest in understanding algorithmic accountability, while being consistent with long-standing U.S. antidiscrimination principles. Part 4 follows by exploring how the IAT can likewise be applied in other settings where algorithmic decision-making has come to play an increasingly important role. Part 5 concludes.

II. ACCOUNTABILITY UNDER U.S. ANTIDISCRIMINATION LAW

A. *Accountability and the Burden-Shifting Framework of Title VII*

We ground our definition of accountability in the antidiscrimination principles of Title VII of the Civil Rights Act of 1964.⁴³ Title VII, which focuses on the labor market, makes it “an unlawful employment practice for an employer (1) to ... discriminate against any individual with respect to his

³⁹ See Part 4(C).

⁴⁰ See Part 2(B).

⁴¹ 557 U.S. 557 (2009).

⁴² We discuss this challenge in more detail in Part 2(B).

⁴³ 42 U.S.C. § 2000e (2012).

compensation, terms, conditions, or privileges of employment, because of such individual's race, color, sex, or national origin; or (2) to limit, segregate, or classify his employees or applicants for employment in any way which would deprive or tend to deprive any individual of employment opportunities ... because of such individual's race, color, religion, sex, or national origin."⁴⁴ Similar conceptualizations of antidiscrimination law were later written to apply to other settings, such as the prohibition of discrimination in mortgage lending under the FHA.⁴⁵

In practice, Title VII has been interpreted as covering two forms of impermissible discrimination. The first and "the most easily understood type of discrimination"⁴⁶ falls under the *disparate-treatment* theory of discrimination and requires that a plaintiff alleging discrimination prove "that an employer had a discriminatory motive for taking a job-related action."⁴⁷ Additionally, Title VII also covers practices which "in some cases, ... are not intended to discriminate but in fact have a disproportionately adverse effect on minorities."⁴⁸ These cases are usually brought forth under the *disparate-impact* theory of discrimination and allow for an employer to be liable for "facially neutral practices that, in fact, are 'discriminatory in operation,'" even if unintentional.⁴⁹

Critically, in cases where discrimination lacks an intentional motive, an employer can be liable only for disparate outcomes that are unjustified. The process of how disparities across members of protected and unprotected groups might be justified is articulated in the *burden-shifting framework* initially formulated by the Supreme Court in *Griggs v. Duke Power Co.*⁵⁰ and subsequently codified by Congress in 1991.⁵¹ This delineation is central to the definition of accountability in today's era of algorithms.

Under the burden-shifting framework, a plaintiff alleging discrimination under a claim without intentional motive bears the first burden. The plaintiff must identify a specific employment practice that causes "observed statistical disparities"⁵² across members of protected and unprotected groups.⁵³ If the plaintiff succeeds in establishing this evidence, the burden shifts to the

⁴⁴ 42 U.S.C. § 2000e-2(a) (2012).

⁴⁵ 42 U.S.C. § 3605 (2012) ("It shall be unlawful for any person or other entity whose business includes engaging in residential real estate-related transactions to discriminate against any person in making available such a transaction, or in the terms or conditions of such a transaction, because of race, color, religion, sex, handicap, familial status, or national origin.")

⁴⁶ *Int'l Bhd. of Teamsters v. United States*, 431 U.S. 324, 335 n.15 (1977).

⁴⁷ *Ernst v. City of Chi.*, 837 F.3d 788, 794 (7th Cir. 2016).

⁴⁸ *Ricci v. DeStefano*, 557 U.S. 557, 577 (2009).

⁴⁹ *Id.* at 577-78 (quoting *Griggs*, 401 U.S. at 431).

⁵⁰ *Griggs*, 401 U.S. at 432.

⁵¹ Civil Rights Act of 1991, Pub. L. No. 102-66, 105 Stat. 1071 (1991).

⁵² *Watson v. Fort Worth Bank & Trust*, 487 U.S. 977, 979 (1988).

⁵³ *See also Albemarle Paper Co. v. Moody*, 422 U.S. 405, 425 (1975) (holding that the plaintiff has the burden of making out a prima facie case of discrimination).

defendant.⁵⁴ The defendant must then “demonstrate that the challenged practice is job related for the position in question and consistent with business necessity.”⁵⁵ If the defendant satisfies this requirement, then “the burden shifts back to the plaintiff to show that an equally valid and less discriminatory practice was available that the employer refused to use.”⁵⁶

This overview highlights two core ideas that inform what it means for a decision-making process to be accountable under U.S. antidiscrimination law. First, in the case of unintentional discrimination, disparate outcomes must be justified by reference to a legitimate “business necessity.”⁵⁷ In the context of employment hiring, for instance, this is typically understood to be a job-related skill that is required for the position.⁵⁸ Imagine, for instance, an employer who made all hiring decisions based on applicant’s level of a direct measure of the job-related skill necessary for the job. Even if the outcome of these decision-making processes results in disparate outcomes across minority and non-minority applicants, these disparities would be justified as nondiscriminatory with respect to a protected characteristic.

Second, in invalidating a decision-making process, U.S. antidiscrimination law does so because of invalid “inputs” rather than invalid “outputs” or results. This feature of U.S. antidiscrimination law is most evident in the case of disparate treatment claims involving the use by a decision-maker of a protected category in making a job-related decision. For instance, Section (m) of the 1991 Civil Rights Act states that “an unlawful employment practice is established when the complaining party demonstrates that race, color, religion, sex, or national origin was a motivating factor for any employment practice, even though other factors also motivated the

⁵⁴ See *Albermarle*, 422 U.S. at 425 (noting that the burden of defendant to justify an employment practice “arises, of course, only after the complaining party or class has made out a prima facie case of discrimination.”)

⁵⁵ 42 U.S.C. § 2000e-2(k)(1)(A)(i); see also *Griggs*, 401 U.S. at 432 (“Congress has placed on the employer the burden of showing that any given requirement must have a manifest relationship to the employment in question.”)

⁵⁶ *Puffer v. Allstate Ins. Co.*, 675 F.3d 709, 717 (7th Cir. 2012); see also 42 U.S.C. § 2000e-2(k)(1)(A)(ii), (C).

⁵⁷ 42 U.S.C. § 2000e-2(k)(1)(A)(i). Likewise, even in the case of claims alleging disparate treatment, an employer may have an opportunity to justify the employment decision. In particular, absent direct evidence of discrimination, Title VII claims of intentional discrimination are subject to the burden-shifting framework established in *McDonnell Douglas Corp. v. Green*, 411 U.S. 792 (1973). Under the *McDonnell Douglas* framework, a plaintiff must first “show, by a preponderance of the evidence, that she is a member of a protected class, she suffered an adverse employment action, and the challenged action occurred under circumstances giving rise to an inference of discrimination.” *Bennett v. Windstream Communications, Inc.*, 792 F.3d 1261, 1266 (10th Cir. 2015). If the plaintiff succeeds in establishing a prima facie case, the burden of production shifts to the defendant to rebut the presumption of discrimination by producing some evidence that it had legitimate, nondiscriminatory reasons for the decision. *Id.* at 1266.

⁵⁸ See, e.g., *Griggs*, 401 U.S. at 432 (holding that the employer’s practice or policy in question must have a “manifest relationship” to the employee’s job duties); see also *Albermarle*, 422 U.S. at 425 (“If an employer does then meet the burden of proving that its tests are ‘job related,’ it remains open to the complaining party to show that other tests or selection devices, without a similarly undesirable racial effect, would also serve the employer’s legitimate interest in ‘efficient and trustworthy workmanship.’”)

practice.”⁵⁹ However, this focus on inputs is also evident in cases alleging disparate impact, notwithstanding the doctrine’s initial requirement that a plaintiff allege disparate outcomes across members of protected and unprotected groups. Recall that even with evidence of disparate outcomes, an employer that seeks to defend against a claim of disparate impact discrimination must demonstrate why these outcomes were the result of a decision-making process based on legitimate business necessity factors (i.e., the disparate outcomes were the result of legitimate decision-making inputs).⁶⁰ This focus on “inputs” underscores the broader policy objective of ensuring a decision-making process that is not discriminatory.

The practical challenge in implementing this antidiscrimination regime is that the critical decision-making input—an individual’s possession of a job-related skill—cannot be perfectly observed at the moment of a decision, inducing the decision-maker to turn to proxies for it. However, the foregoing discussion highlights that the objective in evaluating these proxy variables should be the same: ensuring that qualified minority applicants are not being systematically passed over for the job or promotion. As summarized by the Supreme Court in *Ricci v. DeStefano*, “[t]he purpose of Title VII ‘is to promote hiring on the basis of job qualifications, rather than on the basis of race or color.’”⁶¹

This objective, of course, is the basis for prohibiting the direct form of statistical discrimination famously examined by economists Kenneth Arrow⁶² and Edmund Phelps.⁶³ In their models, an employer uses a job applicant’s race as a proxy for the applicant’s expected productivity because the employer assumes that the applicant possesses the average productivity of his or her race. If the employer also assumes the average productivity of minority applicants is lower than non-minorities (e.g., because of long-standing social and racial inequalities), this proxy will ensure that above-average productive minorities will systematically be passed over for the job despite being qualified for it. Because this practice creates a direct and obvious bias against minorities, this practice is typically policed under the disparate treatment theory of discrimination.⁶⁴

Beyond this clearly unlawful form of statistical discrimination, a decision-maker can use statistical discrimination to incorporate not just the protected-class variable but also other proxy variables for the business-necessity unobservable attributes. For instance, an employer might seek to

⁵⁹ 42 U.S.C. § 2000e-2(m).

⁶⁰ *See, e.g.*, *Dothard*, 433 U.S. at 331 (holding that, to satisfy the business necessity defense, an employer must show that a pre-employment test measured a characteristic “essential to effective job performance” given that the test produced gender disparities in hiring).

⁶¹ *Ricci*, 557 U.S. at 582 (citing *Griggs*, 401 U.S. at 424).

⁶² Kenneth J. Arrow, *The Theory of Discrimination*, in *DISCRIMINATION AND LABOR MARKETS* 3 (Orley Ashenfelter & Albert Rees eds., 1973).

⁶³ Edmund S. Phelps, *The Statistical Theory of Racism and Sexism*, 62 *AM. ECON. REV.* 659 (1972).

⁶⁴ *See* text accompanying note 59.

predict a job applicant's productivity based on other observable characteristics that the employer believes are correlated with future productivity, such as an applicant's level of education or an applicant's performance on a personality or cognitive ability test.⁶⁵ Indeed, it is the possibility of using data mining to discern new and unintuitive correlations between an individual's observable characteristics and a target variable of interest (e.g., productivity or creditworthiness) that has contributed to the dramatic growth in algorithmic decision-making.⁶⁶ The advent of data mining has meant that thousands of such proxy variables are sometimes used.⁶⁷

As the UnitedHealth algorithm revealed, however, the use of these proxy variables can result in members of a protected class experiencing disparate outcomes. The problem arises from what researchers call "redundant encodings"—the fact that a proxy variable can be predictive of a legitimate target variable *and* membership in a protected group.⁶⁸ Moreover, there are social and economic factors that make one's group membership correlated with virtually any observable proxy variable. As one proponent of predictive policy declared, "If you wanted to remove everything correlated with race, you couldn't use anything. That's the reality of life in America."⁶⁹ At the same time certain proxy variables may predict membership in a protected group over-and-above their ability to predict a legitimate target variable; relying on these proxy variables therefore risks penalizing members of the protected group who are otherwise qualified in the legitimate target variable.⁷⁰ In short, algorithmic accountability requires a method to limit the use of redundantly encoded proxy variables to those that are consistent with the anti-discrimination principles of Title VII of the Civil Rights Act and to prohibit the use of those that are not.⁷¹

⁶⁵ See, e.g., Neal Schmitt, *Personality and Cognitive Ability as Predictors of Effective Performance at Work*, 1 ANNUAL REVIEW OF ORGANIZATIONAL PSYCHOLOGY AND ORGANIZATIONAL BEHAVIOR 45, 56 (2014) (describing web-based tests pre-employment tests of personality and cognitive ability).

⁶⁶ See Barocas & Selbst, *supra* note 8, at 677 ("By definition, data mining is always a form of statistical (and therefore seemingly rational) discrimination.")

⁶⁷ See, e.g., Mikella Hurley & Julius Adebayo, *Credit Scoring in the Era of Big Data* 18 YALE J. L. TECH. 148, 164 (2020) (describing how ZestFinance uses an "all data is credit data" approach to predict an individual's creditworthiness based on "thousands of data points collected from consumers' offline and online activities").

⁶⁸ See Barocas & Selbst, *supra* note 8, at 691 (citing Cynthia Dwork et al., *Fairness Through Awareness*, 3 PROC. INNOVATIONS THEORETICAL COMPUTER SCI. CONF. 214 app. at 226 (2012)).

⁶⁹ Nadya Labi, *Misfortune Teller*, ATLANTIC (Jan.–Feb. 2012), <http://www.theatlantic.com/magazine/archive/2012/01/misfortune-teller/308846> (quoting Ellen Kurtz, Director of Research for Philadelphia's Adult Probation and Parole Department).

⁷⁰ As noted in the Introduction, redlining represents a classic example: An individual's zip code may be somewhat predictive of one's creditworthiness, but given racialized housing patterns, it is almost certainly far more accurate in predicting one's race. Assuming that all residents in a minority-majority zip code have low creditworthiness will therefore result in systematically underestimating the creditworthiness of minorities whose actual creditworthiness is higher than the zip code average.

⁷¹ In theory, there are statistical methods that would estimate the precise degree to which a redundantly encoded proxy variable predicts a legitimate target variable that is independent of the degree to which it predicts membership in a protected classification. We discuss these methods and their shortcomings *infra* at notes 144 to 146 and in the Appendix.

Our central contribution is in developing accountability input criteria that speak directly to the process demanded by Title VII. Specifically, we use these accountability input criteria to develop a statistical test for whether a proxy variable (or each proxy variable in a set of proxy variables) is being used in a way that causes illegitimate statistical discrimination and should therefore not be used in an algorithmic model. Fundamentally, it is a test for “unbiasedness” designed to ensure that the use of a proxy input variable does not systematically penalize members of a protected group who are otherwise qualified with respect to a legitimate-business-necessity objective. We refer to this test as the *input-accountability test*. We illustrate the test and its application with a simple pre-employment screening exam designed to infer whether a job applicant possesses sufficient strength to perform a particular job. Before doing so, however, we differentiate the input-accountability test from other approaches to algorithmic accountability.

B. The Input Accountability Test Versus Outcome-Oriented Approaches

Our input-based approach differs significantly from that of other scholars who have advanced outcome-oriented approaches to algorithmic accountability. For instance, Talia Gillis and Jann Spiess have argued that the conventional focus in fair lending on restricting invalid inputs (such as a borrower’s race or ethnicity) is infeasible in the machine-learning context.⁷² The reason, according to Gillis and Spiess, is because a predictive model of default that excludes a borrower’s race or ethnicity can still penalize minority borrowers if one of the included variables (e.g., borrower education) is correlated with both default and race.⁷³ Gillis and Spiess acknowledge the possibility that one could seek to exclude from the model some of these correlated variables on this basis, but they find this approach infeasible given that “a major challenge of this approach is the required articulation of the conditions under which exclusion of data inputs is necessary.”⁷⁴ They

⁷² See Talia B. Gillis and Jann L. Spiess, *Big Data and Discrimination*, 86 U. CHICAGO L. REV 459 (2019).

⁷³ *Id.* at 468-69.

⁷⁴ *Id.* at 469. Elsewhere in their article, Gillis and Spiess also suggest that input-based analysis may be infeasible because “in the context of machine-learning prediction algorithms, the contribution of individual variables is often hard to assess.” *Id.* at 475. They illustrate this point by showing how in a simulation exercise, the variables selected by a logistic lasso regression in a predictive model of default differed each time the regression was run on a different randomly-drawn subsample of their data. However, this evidence does not speak to how an input-based approach to regulating algorithms would be deployed in practice. A lasso regression—like other models that seek to reduce model complexity and avoid overfitting—seeks to reduce the number of predictors based on the underlying correlations among the full set of predictor variables. Thus, it can be used in training a model on a set of data with many proxy variables, and running a lasso regression multiple times on different subsamples of the data should be expected to select different variables with each run. However, once a model has been trained and the model’s features are selected, the model must be deployed, allowing the features used in the final model to be evaluated and tested for bias. That is, regardless of the type of model fitting technique one uses in the training procedure (e.g., lasso regression, ridge regression, random forests, etc.), the model that is ultimately deployed will utilize a set of features that can be examined.

therefore follow the burgeoning literature within computer science on “algorithmic fairness”⁷⁵ and advocate evaluating the outcomes from an algorithm against some baseline criteria to determine whether the outcomes are fair.⁷⁶ As examples, they suggest a regulator might simply examine whether loan prices differ across members of protected or unprotected groups, or a regulator might look at whether “similarly situated” borrowers from the protected and unprotected groups were treated differently.⁷⁷

Gillis and Spiess are, of course, correct that simply prohibiting an algorithm from considering a borrower’s race or ethnicity will not eliminate the risk that the algorithm will be biased against minority borrowers in a way that is unrelated to their creditworthiness (which is a legitimate-business-necessity variable).⁷⁸ Indeed, we share this concern about redundant encodings, and it motivates our empirical test. However, we part ways with these authors in that we do not view as insurmountable the challenge of articulating the conditions for excluding variables that are correlated with a protected classification, as we illustrate in Part 3.

Equally important, it is with an outcome-based approach rather than with an input-based approach where one encounters the greatest conceptual and practical challenges for algorithmic accountability. As Richard Berk and others have noted, efforts to make algorithmic outcomes “fair” pose the challenge that there are multiple definitions of fairness, and many of these

⁷⁵ For a summary, see Sam Corbett-Davies and Sharad Goel, *The Measure and Mismeasure of Fairness: A Critical Review of Fair Machine Learning* (arXiv.org, August 2018), available at <https://arxiv.org/pdf/1808.00023.pdf>. In particular, a common approach to algorithmic fairness within computer science is to evaluate the fairness of a predictive algorithm by use of a “confusion matrix.” *Id.* at 4. A confusion matrix is a cross-tabulation of actual outcomes by the predicted outcome. For instance, the confusion matrix for an algorithm that classified individuals as likely to default on a loan would appear as follows:

	Default Predicted	No Default Predicted
Default Occurs	# Correctly Classified as Defaulting = N_{TP} (True Positives)	# Incorrectly Classified as Non-Defaulting = N_{FN} (False Negatives)
Default Does Not Occur	# Incorrectly Classified as Defaulting = N_{FP} (False Positives)	# Correctly Classified as Non-Defaulting = N_{TN} (True Negatives)

Using this table, one could then evaluate the fairness of the classifier by inquiring whether classification error is equal across members of protected and unprotected groups. *Id.* at 5. For example, one could use as a baseline fairness criterion a requirement that the classifier have the same false positive rate (i.e., $N_{FP} / (N_{FP} + N_{TN})$) for minority borrowers as for non-minority borrowers. Alternatively, one could use as a baseline a requirement of treatment equality (e.g., the ratio of False Positives to False Negatives) across members of protected and unprotected groups.

⁷⁶ See Gillis & Spiess, *supra* note 72, at 480 (“In the case of machine learning, we argue that outcome analysis becomes central to the application of antidiscrimination law.”)

⁷⁷ *Id.* at 484-85.

⁷⁸ See also Jon Kleinberg, et al, *Algorithmic Fairness*, 108 AEA PAPERS AND PROCEEDINGS 22, 22–23 (2018) (“Our central argument is that across a wide range of estimation approaches, objective functions, and definitions of fairness, the strategy of blinding the algorithm to race inadvertently detracts from fairness.”)

definitions are incompatible with one another.⁷⁹ The central challenge is that an outcome test will often result in *some* form of residual discrimination, raising the inevitable question: *how much* discrimination should be permissible in the outcomes?⁸⁰

In a concrete illustration of this challenge, Richard Berk and a team of researchers at the University of Pennsylvania describe a criminal-risk-assessment tool they designed for a jurisdiction that was concerned about racial bias in the pre-trial release rates among criminal defendants who were awaiting trial.⁸¹ In general, when a defendant was arraigned in this jurisdiction, a magistrate judge was required to decide whether the defendant should be released until the trial date, considering (among other things) the defendant's threat to public safety.⁸² Berk and his team developed a forecasting algorithm of a defendant's risk, using a subsequent arrest for a violent crime within 21 months of release as a proxy for the defendant's threat to public safety.⁸³

To reduce racial disparities, Berk and his team tuned the algorithm so that it was equally accurate at predicting release across racial categories; that is, the rate of re-arrest for a violent crime among minority and non-minority defendants was the same.⁸⁴ However, the base rate of re-arrest among minority defendants was higher than non-minority defendants, meaning that the chosen fairness objective could be accomplished only by making the algorithm biased. In particular, the algorithm had to classify more "violent" non-minority defendants as "nonviolent" (thus resulting in their release), and it had to classify more "nonviolent" minority defendants as "violent" (thus resulting in their detention).⁸⁵ The need to bias the algorithm in this fashion

⁷⁹ See Richard Berk et al., *Fairness in Criminal Justice Risk Assessments: The State of the Art* 33 (arXiv.org, May 30, 2017), available at <https://arxiv.org/pdf/1703.09207.pdf> (arguing that "[t]here are different kinds of fairness that in practice are incompatible").

⁸⁰ See, e.g., Gillis & Spiess, *supra* note 72, at 486 (advocating an outcome test in which a regulator evaluates whether lending outcomes differ by race among "similarly situated" borrowers "should include a degree of tolerance set by the regulator").

⁸¹ See Berk et al., *supra* note 79, at 31-33.

⁸² *Id.* at 31.

⁸³ *Id.* at 31-33.

⁸⁴ *Id.*

⁸⁵ As Sam Corbett-Davies and Sharad Goel note, it is the different underlying distribution of risk (or other unobservable characteristic of interest) among minority and non-minority populations that gives rise to the alternative and incompatible definitions of fairness based on classification errors. See Corbett-Davies & Goel, *supra* note 75, at 2 ("When the true underlying distribution of risk varies across groups, differences in group-level error rates are an expected consequence of algorithms that accurately capture each individual's risk."). Within the antidiscrimination literature, this statistical challenge is known as the problem of infra-marginality and has long plagued outcome tests of discrimination in human decisions. See Ian Ayres, *Outcome Tests of Racial Disparities in Police Practices*, 4 JUSTICE RESEARCH AND POLICY 131 (2002). The central problem is that an inquiry into whether a decision is unbiased is concerned with what happens at the margin (i.e., is the same standard being applied to everyone?). Error rates, however, are evaluated away from the margin as they rely on evaluating outcomes following the application of a cut-off standard to all individuals (both those who might be near the cut-off and those who might be far from it). If the risk distributions differ across minority and non-minority individuals, lumping together both marginal and infra-marginal individuals will produce error rates that differ by race.

arose from the fact that minority defendants had a higher baseline re-arrest rate.⁸⁶ As a result, the algorithm could achieve its particular definition of fairness—equality of accuracy, conditional on release—only by releasing more non-minority defendants that were more likely be arrested (to compensate for the overall lower rate of arrest for non-minority defendants) and by not releasing some minority defendants that were unlikely to be arrested (to compensate for their overall higher rate of arrest).⁸⁷ As they note, in sacrificing one form of fairness for another, the resulting “differences [in error rates] can support claims of racial injustice.”⁸⁸

To be sure, applying our test for “unbiasedness” does not solve the challenge of addressing concerns about fairness. A decision that passes our test might still be objectionable for other distributional reasons. In the case of lending, for instance, creditworthiness is a well-recognized target variable, but the determinants of creditworthiness (e.g., income, income growth, wealth) also reflect long-standing racial and economic inequalities, ensuring that creditworthiness will likewise reflect these racial and economic inequalities. Thus, even an unbiased lending rule would result in lending outcomes that reflect these structural inequalities, and rectifying them would require an additional intervention, such as through subsidized loan programs and other policies designed to encourage lending to low and moderate-income families. Indeed, this approach is reflected in existing U.S. housing programs such as the Federal Housing Administration mortgage program (which seeks to provide mortgages to low and moderate-income borrowers)⁸⁹ and the Community Reinvestment Act (which seeks to encourage lenders to provide loans to residents of low and moderate-income neighborhoods).⁹⁰

This two-step approach—ensuring that decision-making processes are unbiased and then subsequently addressing distributional concerns directly through transfers and subsidies—is consistent with democratic principles. As we show, it is conceivable to design a decision-making process that is unbiased with respect to a legitimate business necessity. This is the objective of the IAT. But as Berk’s study illustrates, it is not possible to design a decision-making process that satisfies every possible definition of “fairness.” Evaluating an algorithm for whether it is “fair” rather than “unbiased” thus risks enforcing a particular vision of fairness and doing so in a way that lacks transparency. Indeed, Berk et al. themselves provide no explanation for why

⁸⁶ *Id.* at 32.

⁸⁷ *Id.*

⁸⁸ *Id.*

⁸⁹ See James H. Carr, Katrin B. Anacker, *The Complex History of the Federal Housing Administration: Building Wealth, Promoting Segregation, and Rescuing the U.S. Housing Market and The Economy*, 34 BANKING & FIN. SERVICES POL’Y REP. 10 (2015) (describing program).

⁹⁰ See Keith N. Hylton, *Banks and Inner Cities: Market and Regulatory Obstacles to Development Lending*, 17 YALE J. ON REG. 197 (2000) (describing Community Reinvestment Act).

they opted to implement their particular definition of fairness.⁹¹ Likewise, an algorithm that seeks to “fix” disparate outcomes that arise from an unbiased decision-making process can risk diminishing the ability to identify the source of the underlying structural inequalities and/or measurement error in the decision-making process. In Berk’s setting, for instance, a risk-assessment algorithm that results in equality of release rates across minority and non-minority defendants could hide the possibility that minority defendants have a higher re-arrest rate because of prejudice among police, which in turn would raise the question of whether re-arrests are truly a decent proxy for a defendant’s level of violence. For all of these reasons, determination of distributional equity is accordingly best left to institutions that can evaluate the relevant trade-offs in a transparent fashion.

Regardless of these conceptual and practical challenges, outcome-based approaches to algorithmic fairness would almost certainly be deemed legally problematic following the Supreme Court’s 2009 decision in *Ricci v. DeStefano*.⁹² The facts giving rise to *Ricci* involved a decision by the city of New Haven to discard the results of an “objective examination” that sought to identify city firefighters who were the most qualified for promotion.⁹³ The city justified its decision to discard the results on the basis that there was a statistical racial disparity in the results, raising the risk of disparate impact liability under Title VII.⁹⁴ A group of white and Hispanic firefighters sued, alleging that the city’s discarding of the test results constituted race-based disparate-treatment.⁹⁵ In upholding their claim, the Court emphasized the extensive efforts that the city took to ensure the test was job-related⁹⁶ and that there was “no genuine dispute that the examinations were job-related and consistent with business necessity.”⁹⁷ Nor did the city offer “a strong basis in evidence of an equally valid, less-discriminatory testing alternative.”⁹⁸ Prohibiting the city from discarding the test results was therefore required to prevent the city from discriminating against “qualified candidates on the basis of their race.”⁹⁹

The Court’s assumption that the promotion test identified the most qualified firefighters makes it difficult to see a legal path forward for explicit race-based adjustments of algorithmic outcomes. Assuming the algorithm

⁹¹ Berk et al, *supra* note 79, at 32 (describing their choice of error metric as “conditional use accuracy equality, which some assert is a necessary feature of fairness.”)

⁹² 557 U.S. 557 (2009)

⁹³ *Id.* at 562.

⁹⁴ *Id.*

⁹⁵ *Id.* at 562-63.

⁹⁶ *Id.* 586-588.

⁹⁷ *Id.* at 587; see also *id.* at 589 (“The City, moreover, turned a blind eye to evidence that supported the exams’ validity.”)

⁹⁸ *Id.* at 589.

⁹⁹ *Id.* at 584 (“Restricting an employer’s ability to discard test results (and thereby discriminate against qualified candidates on the basis of their race) also is in keeping with Title VII’s express protection of bona fide promotional examinations.”)

properly functions to identify individuals who are qualified in a specified target, such race-based adjustments would appear to be no different than what the city of New Haven attempted to do with the promotion test results. Rather, *Ricci* underscores the fundamental importance of ensuring that decision-making processes do not systematically discriminate against qualified individuals because of their race. And as noted previously, this is the objective of the IAT.

C. The Input Accountability Test Versus the “Least Discriminatory Alternative” Test

We differ also from scholars and practitioners who focus only on the final step in the disparate-impact burden-shifting framework. Recall that according to this burden-shifting framework, an employer who establishes that a business practice can be justified by a legitimate business necessity shifts the burden back to the plaintiff to show that an equally valid and less discriminatory practice was available that the employer refused to use.¹⁰⁰ Some commentators have mistakenly assumed that this test implies that the critical question for an algorithm that produces a disparate impact is whether the algorithm uses the least discriminatory predictive model for a given level of predictive accuracy. Of course, in using machine learning over thousands of variables, it is easy to run many models and decide which creates the least disparate impact for a given level of accuracy in prediction. But this approach will not address whether any of the variables used in the model are systematically penalizing members of a protected group that are otherwise qualified in the skill or characteristic the model is seeking to predict.

Nonetheless, a number of commentators have mistakenly argued that the central test for whether an algorithm poses any risk of illegitimate discrimination should be whether there are alternative models that can achieve the same level of predictive accuracy with lower levels of discrimination.¹⁰¹ For instance, in an oft-cited discussion paper regarding fair lending risk of credit cards, David Skanderson and Dubravka Ritter advocate that lenders should focus on this step of the disparate-impact framework when evaluating the fair-lending risk of algorithmic credit-card models.¹⁰² Specifically, Skanderson and Ritter note that “a model or a model’s predictive

¹⁰⁰ See text accompanying note 56.

¹⁰¹ See, e.g., Nicholas Schmidt and Bryce Stephens, *An Introduction to Artificial Intelligence and Solutions to the Problems of Algorithmic Discrimination*, (arXiv.org, Nov. 2019), available at <https://arxiv.org/pdf/1911.05755.pdf> (advocating for using “a ‘baseline model’ that has been built without consideration of protected class status, but which shows disparate impact, and then search[ing] for alternative models that are less discriminatory than that baseline model, yet similarly predictive.”)

¹⁰² See, e.g., David Skanderson & Dubravka Ritter, Discussion Paper, Fair Lending Analysis of Credit Cards, Federal Reserve Bank of Philadelphia (August 2014), available at <https://www.philadelphiafed.org/-/media/consumer-credit-and-payments/payment-cards-center/publications/discussion-papers/2014/d-2014-fair-lending.pdf?la=en> (last visited February 2, 2020).

variable with a disproportionate adverse impact on a prohibited basis may still be legally permissible if it has a demonstrable business justification and there are no alternative variables that are equally predictive and have less of an adverse impact.”¹⁰³ For Skanderson and Ritter, the business necessity defense for an algorithmic decision-making process therefore boils down to whether it is the most accurate possible test in predicting a legitimate target variable of interest. As they summarize in the context of lending, “If a scoring system is, in fact, designed to use the most predictive combination of available credit factors, then it should be unlikely that someone could demonstrate that there is an equally effective alternative available, which the lender has failed to adopt.”¹⁰⁴

To see how validating an algorithm based entirely on the fact that it is the most predictive model available would validate algorithms that are clearly biased against members of a protected group, we offer an example. Consider an employer who needs employees that can regularly lift 40 pounds as part of their everyday jobs. Imagine this employer designs a one-time test of whether applicants can lift 70 pounds as a proxy for whether the applicant can repetitively lift 40 pounds. The employer can show that this test has 90% prediction accuracy. However, those applicants that fail the test who in fact could regularly lift 40 pounds are disproportionately female. Thus, the test, because it is not a perfect proxy, causes a disparate impact on female applicants. Now assume that it can be shown that a one-time test of whether applicants can lift 50 pounds produces no disparate impact on females but has an accuracy rate of just 85%. Under Skanderson and Ritter’s approach, the employer would have no obligation to consider the latter test, despite the fact that a 70-pound test will systematically penalize female applicants that can in fact satisfy the job requirement.

Not surprisingly, this approach to pre-screening employment tests has been expressly rejected by courts. In *Lanning v. Southeastern Pennsylvania Transportation Authority*,¹⁰⁵ for instance, the Third Circuit considered a physical fitness test for applicants applying to be transit police officers. The fitness test involved a 1.5 mile run that an applicant was required to complete

¹⁰³ *Id.* at 38.

¹⁰⁴ *Id.* at 43. This line of reasoning also informs Barocas and Selbst’s conclusion that Title VII provides a largely ineffective means to police unintentional discrimination arising from algorithms. *See* Barocas & Selbst, *supra* note 66, at 701-714. According to Barocas and Selbst, the business necessity defense requires that an algorithm is “predictive of future employment outcomes.” *Id.* If this is correct, it would logically follow that an employer will have no disparate impact liability from using the most predictive algorithmic model for a legitimate job-related quality since an equally predictive, less discriminatory alternative would not be available. However, this conclusion relies on an assumption that predictive accuracy is a necessary and sufficient condition to justify a decision-making process that produces a disparate impact. As we show, this is an incorrect assumption as courts have been careful not to conflate the business necessity defense with predictive accuracy. A predictive model may be accurate in predicting whether an individual is likely to have a legitimate target characteristic but nevertheless be biased against members of a protected group who are otherwise qualified in the target characteristic.

¹⁰⁵ 181 F.3d 478 (3rd Cir. 1999), cert. denied, 528 U.S. 1131 (2000).

within 12 minutes; however, the 12 minute cut-off was shown to have a disparate impact on female applicants.¹⁰⁶ The transit authority acknowledged that officers would not actually be required to run 1.5 miles within 12 minutes in the course of their duties, but it nevertheless adopted the 12 minute cut-off because the transit authority's expert believed it would be a more "accurate measure of the aerobic capacity necessary to perform the job of a transit police officer."¹⁰⁷

In considering the transit authority's business-necessity defense, the court agreed that aerobic capacity was related to the job of a transit officer.¹⁰⁸ It also agreed that by imposing a 12 minute cut-off for the run, the transit authority would be increasing the probability that a job applicant would possess high aerobic capacity.¹⁰⁹ Nonetheless, the court rejected this "more is better" approach to setting the cutoff time:

Under the District Court's understanding of business necessity, which requires only that a cutoff score be "readily justifiable," [the transit authority], as well as any other employer whose jobs entail any level of physical capability, could employ an unnecessarily high cutoff score on its physical abilities entrance exam in an effort to exclude virtually all women by justifying this facially neutral yet discriminatory practice on the theory that more is better.¹¹⁰

Accordingly, the court required "that a discriminatory cutoff score be shown to measure the minimum qualifications necessary for successful performance of the job in question in order to survive a disparate impact challenge."¹¹¹ In other words, in determining whether disparate outcomes are justified, the question to ask in evaluating a predictive model of a legitimate target variable is not simply whether the model is accurate in predicting the target variable. Rather, the inquiry should be both whether the model is accurate and whether the cutoff score used to classify individuals serves the employer's legitimate business goals.¹¹²

¹⁰⁶ *Id.* at 482

¹⁰⁷ *Id.*

¹⁰⁸ *Id.* at 492.

¹⁰⁹ *Id.* ("The general import of these studies is that the higher an officer's aerobic capacity, the better the officer is able to perform the job.")

¹¹⁰ *Id.* at 493.

¹¹¹ *Id.*

¹¹² The Third Circuit was even more explicit that setting the cut-off was effectively about calibrating the predictive accuracy of the employment test. *See* Lanning, 308 F.3d at 292 ("It would clearly be unreasonable to require SEPTA applicants to score so highly on the run test that their predicted rate of success be 100%. It is perfectly reasonable, however, to demand a chance of success that is better than 5% to 20%."); *see also* E.E.O.C. vs. Simpson Timber Co., 1992 WL 420897 (finding that a pre-employment step test accurately measured strength and endurance, which were legitimate business goals, but an equally

An even stronger rejection of the “more is better” approach to predictive accuracy appeared in *Murphy v. Derwinski*.¹¹³ There, the plaintiff, Mary Murphy, applied to become a Roman Catholic chaplain at hospitals operated by the United States Veterans Administration (VA).¹¹⁴ The VA rejected Murphy’s application on the ground that VA guidelines required that all applicants be ordained in the relevant religion and receive an ecclesiastical endorsement from their churches.¹¹⁵ However, within the Roman Catholic religion, only men can be ordained as priests, making it impossible for Murphy to satisfy these requirements.¹¹⁶ In her subsequent Title VII lawsuit, the district court determined that Murphy made out a prima facie case of discrimination based on this policy and that the defendant articulated a legitimate business justification for it.¹¹⁷ In particular, the court agreed with the VA that the agency’s interest in providing a full range of ritual services to its Catholic patients creates a legitimate purpose for requiring ordination for VA chaplains.¹¹⁸ The court nevertheless rejected the VA’s argument that if the ordination requirement were eliminated, the VA would be unable to accommodate the needs of its patients by providing the full range of religious services.¹¹⁹ Rather, the court held that by removing the ordination requirement and requiring only ecclesiastical endorsement, the VA could still ensure that its patients received the religious services that the Catholic Church deemed sufficient.¹²⁰

On appeal, the Tenth Circuit affirmed and elaborated on why removing the ordination requirement would not impair the VA’s legitimate interests.¹²¹ Citing the VA’s own administrative materials, the court noted that the chaplain service’s primary objective was to “provide for the spiritual welfare”¹²² of patients such as through establishing relationships with patients and providing patients and family members with ministry in crisis situations, and “[p]riests are not needed to administer these functions.”¹²³ The court acknowledged that only priests could administer sacraments to patients subscribing to the Roman Catholic faith,¹²⁴ but it concluded that the VA would still be able to accommodate the religious needs of its Roman Catholic patients:

effective, less discriminatory alternative existed in the form of using a lower cut-off score to determine if an applicant passed the test).

¹¹³ 990 F.2d 540 (10th Cir. 1993).

¹¹⁴ *Id.* at 542.

¹¹⁵ *Id.*

¹¹⁶ *Id.* at 542 n. 5.

¹¹⁷ *Murphy v. Derwinski*, 776 F. Supp. 1466, 1470 (D. Colo. 1991).

¹¹⁸ *Id.*

¹¹⁹ *Id.* at 1472-73

¹²⁰ *Id.*

¹²¹ *Murphy* 990 F.2d at 545-547

¹²² *Id.* at 546

¹²³ *Id.*

¹²⁴ *Id.* at 545.

The experience of the VA hospital in Denver where Murphy sought to work suggests that removal of the ordination requirement will not disrupt services only priests may perform. Of the hospital's six chaplains at the time of this lawsuit, two were Catholic priests. Thus, four of the chaplains could not provide Roman Catholics with services unique to that religion. Similarly, none of the six could administer unique religious services to members of nonrepresented faiths. When a priest is needed but, for whatever reason, is unavailable, the VA Manual calls for supplementing its full-time chaplain services through contract help or other arrangements.¹²⁵

Thus, the court held that requiring only the ecclesiastical endorsement was an alternative, nondiscriminatory requirement that could serve the VA's legitimate interest in providing the full range of religious services to its patients.¹²⁶ In so doing, note the inconsistency with the approach to the "less discriminatory alternative" inquiry as interpreted by Skanderson and Ritter. Like the transit authority in *Lanning*, the VA in *Murphy* was concerned about identifying job applicants who, at any given moment during their job performance, were likely to serve the VA's legitimate interest.¹²⁷ That is, the VA's hiring guidelines were designed to provide an answer to the question: When a Roman Catholic patient requires religious services, will this applicant be able to provide them? The two requirements—ordination and ecclesiastical endorsement—were clearly accurate in predicting whether an applicant could provide these services. And the requirement that applicants have both characteristics made it virtually certain that a VA chaplain could, in fact, provide any and all of these religious services, any time of the day (morning, noon or night). But like the court in *Lanning*, the *Murphy* court also concluded that setting the probability threshold so high—in this case, imposing an application requirement that made it close to 100% certain that a chaplain would be available to provide any and all Catholic religious services—was simply too high. As the court emphasized, most of the services required of chaplains did not require ordination. Thus, eliminating the ordination requirement might lessen the probability that a VA chaplain would actually be available to administer the sacraments if a patient happened to require them, but the probability would nonetheless remain high enough to

¹²⁵ *Id.* at 546.

¹²⁶ *Id.* at 545-546.

¹²⁷ *See Murphy*, 776 F. Supp. at 1472 ("The VA asserts that if ordination were not required, it would not be able to accommodate the needs of its patients by providing the full range of religious services. VA chaplains must be able to administer the various sacraments, and only ordained priests are qualified for these duties.")

satisfy the VA's legitimate interest in accommodating the religious needs of its Roman Catholic patients.

In short, in the era of algorithmic decision-making, we view the need to inquire into whether there exists a "less discriminatory alternative" to be fundamentally about the cut-off threshold applied to an algorithm that otherwise passes our test. Whether an algorithm is screening for acceptable job applicants or acceptable borrowers, the end result is both to estimate the probability that an individual has a legitimate target characteristic and then to apply a probability cut-off to make the ultimate accept/reject classification. In setting this cut-off, *Lanning* and *Murphy* are reminders of the need to consider whether the cut-off has been set at the minimum level required to advance a legitimate business interest, such as successful performance of the job in question.¹²⁸ As we show below, doing so can help ensure that a decision-making process that passes our test is not inappropriately biased against members of a protected group simply because of the unequal distribution of a legitimate target variable (e.g., strength or speed) across protected and unprotected groups.¹²⁹

D. The Input Accountability Test Versus HUD's Mere Predictive Test

Finally, we consider the IAT against HUD's 2019 proposed rulemaking regarding the application of the disparate impact framework under the FHA.¹³⁰ Given the increasing role of algorithmic credit scoring, the proposed rule-making expressly provides for a new defense for disparate impact claims where "a plaintiff alleges that the cause of a discriminatory effect is a model used by the defendant, such as a risk-assessment algorithm..."¹³¹ In particular, the proposed rule provides that in these cases, a lender may defeat the claim by "identifying the inputs used in the model and showing that these inputs are not substitutes for a protected characteristic and that the model is predictive of risk or other valid objective."¹³² In other words, so long as a

¹²⁸ See *Lanning* F.3d. 481 ("[U]nder the Civil Rights Act of 1991, a discriminatory cutoff score on an entry level employment examination must be shown to measure the minimum qualifications necessary for successful performance of the job in question in order to survive a disparate impact challenge."); see also *Association of Mexican-American Educators v. State of California*, Nos. 96-17131 and 97-15422, 1999 WL 976720 (9th Cir. Oct. 28, 1999) (upholding, against a disparate-impact challenge under Title VII, a requirement that public-school teachers "demonstrate basic reading, writing and mathematics skills in the English language as measured by a basic skills proficiency test" and holding as not clearly erroneous the district court's finding that the cutoff scores "reflect[ed] reasonable judgments about the minimum levels of basic skills competence that should be required of teachers.").

¹²⁹ This interpretation of the third prong of the Title VII burden-shifting framework is also consistent with the common view that it is effectively a test for whether an articulated business necessity defense is a pretext for discrimination; that is, as noted in *Lanning*, one could purposefully set a threshold at a sufficiently high level to ensure that members of protected groups will fail the test. See, e.g., *Murphy* 990 F.2d at 545 ("The focus on appeal is whether the VA's business justification for requiring an ordained clergy person constitutes a pretext for gender discrimination.")

¹³⁰ See 2019 HUD Proposal, *supra* note 20.

¹³¹ *Id.* at 42,862.

¹³² *Id.*

variable is not an undefined “substitute” for a protected characteristic, any model that predicts creditworthiness is sufficient to defeat a claim of disparate impact discrimination.

This approach to algorithmic accountability, however, suffers from the same defect noted previously with regard to those who have misapplied the “least discriminatory alternative” test. Specifically, by focusing solely on whether a model is “predictive of risk or other valid objective,” HUD’s test leaves open the possibility that a lender can adopt a model that systematically discriminates against borrowers who are, in fact, creditworthy. Recall that in our hypothetical strength test, the ability to lift 70 pounds was, in fact, predictive of whether an applicant could regularly lift 40 pounds; however, it systematically discriminated against women who were qualified for the job. Worse still, by not even requiring that a model have any particular level of accuracy, HUD’s test would seemingly permit the use of any proxy so long as it has *some* correlation with credit risk. Indeed, this approach would even appear to permit the use of explicit redlining in a predictive model so long as a lender could show that the average credit risk of a majority-minority neighborhood is marginally higher than that of non-majority-minority neighborhoods.

In contrast, a central goal of the IAT is to ensure that in evaluating whether a model is consistent with a decision-maker’s legitimate business necessity, it incorporates only those proxy variables that are not correlated with a protected characteristic beyond the proxy variables’ correlation with a legitimate target variable.

III. THE INPUT ACCOUNTABILITY TEST

In this section, we formally present our *input accountability test* (IAT) for unintentional discrimination. We begin with some nomenclature. The design of a decision-making algorithm rests fundamentally on the relationships between a set of variables, referred to as “features,” and an underlying latent skill or attribute of interest (creditworthiness, productivity, etc.), referred to as a “target.” Today, the relationships between targets and features are increasingly analyzed and developed within artificial-intelligence and machine-learning processes, but it is just as likely that an organization uses an algorithmic decision process based on human-selected data or even on personal intuition. The IAT applies to both machine learning and human learning.

Our core contribution is a test that informs when a feature’s (a proxy variable’s) use has correlations with a target that produce statistical discrimination against a protected class that is unjustified according to the criteria developed in Part 2. That is, the IAT detects if the use of a feature results in systematically penalizing members of a protected group who are otherwise qualified in the target variable of interest. After illustrating the

IAT, we extend our analysis to consider the mis-assertion of a target cutoff that does not reflect the true level of the target that is required, reflecting the prior example we gave of requiring job applicants to lift 70 pounds as a mis-asserted target.

A. *The Test*

We illustrate our test throughout with the facts giving rise to the 1977 Supreme Court decision in *Dothard v. Rawlinson*.¹³³ As noted previously, in *Dothard*, female applicants for prison guard positions challenged a prison's minimum height and weight requirements as inconsistent with Title VII.¹³⁴ Because the average height and weight of females was less than that for males, the female applicants argued that the requirement created an impermissible disparate impact for females under Title VII.¹³⁵ In response, the prison argued that a height and weight requirement was a justified job requirement given that an individual's height and weight are predictive of strength, and strength was required for prison guards to perform their jobs safely.¹³⁶ In short, the prison took the position that the general correlation between one's height/weight and strength was sufficient to justify the disparate outcomes this requirement caused for women. The Supreme Court, however, rejected this defense.¹³⁷ Rather, to justify gender differences in hiring outcomes, the prison would need to show that it had tested for the *specific type of strength* required for effective job performance;¹³⁸ in other words the prison would have to be concerned with the aspects of strength that the proxy variables were and were not picking up that related to a prison guard's need for strength.

We use this setup and some hypothetical applicants to lay out the IAT. Imagine for example that twelve individuals apply for an open prison guard position, of which six applicants are male and six are female. In evaluating the applicants, the prison seeks to select those applicants who possess the actual strength required for successful job performance. For simplicity, assume that an individual's strength can be measured on a scale of zero to one hundred, and that a strength score of at least sixty is a true target for job effectiveness (in the Court's language a strength of sixty is a legitimate-business-necessity criterion). The challenge the prison faces in evaluating job applicants is that each applicant's actual strength is unobservable at the time of hiring, thus inducing the prison to rely on height as a proxy.

¹³³ 433 U.S. 321 (1977).

¹³⁴ *Id.* at 323-24.

¹³⁵ *Id.*

¹³⁶ *Id.* 331.

¹³⁷ *Id.* at 332.

¹³⁸ *Id.* at 332 ("If the job-related quality that the appellants identify is bona fide, their purpose could be achieved by adopting and validating a test for applicants that measures strength directly.")

Assume that the use of a minimum height requirement results in the following distribution of applicants according to their actual but unobservable strength (Figure 1).

Figure 1

Actual Strength	Results with Height Test		
	Meets Height Requirement	Fails Height Requirement	
100	x		
90	x		
80	x		
70	x	x	
60		x	<i>Minimum Required Strength</i> ↓
50	x		
40	x	x	
30		x	
20		x	
10		x	
0			

x = applicant

Consistent with the prison’s argument, there is a clear correlation between an applicant’s height and actual strength. However, when we examine the gender of the applicants, we discover that only the six male applicants satisfy the minimum height requirement (Figure 2).

Figure 2

Applicant’s Strength	Results with Height Test		
	Meets Height Requirement	Fails Height Requirement	
100	□		
90	□		
80	□		
70	□	•	
60		•	<i>Minimum Required Strength</i> ↓
50	□		
40	□	•	
30		•	
20		•	
10		•	
0			

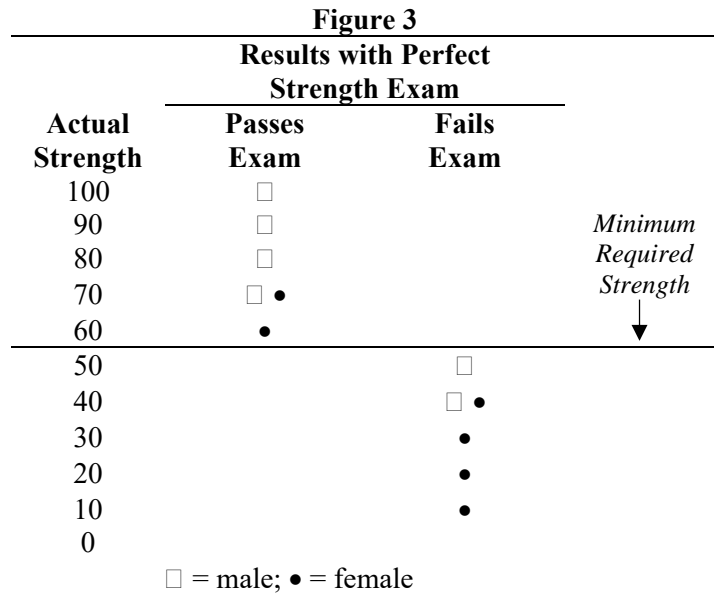
□ = male; • = female

In this situation, a basic correlation test between height and strength has produced exactly the same injury noted in Part 2: The imperfect relationship between height and strength results in penalizing otherwise qualified female applicants and benefiting unqualified male applicants. This can be seen from the fact that the only applicants who possessed sufficient strength but failed the height test were female. Likewise, the only applicants who met the height test but lacked sufficient strength were male. The screening test is thus systematically biased against female applicants for reasons unrelated to a legitimate business necessity.

This example points to the crux of the IAT. In general, the objective of the test is to ensure that a proxy variable is excluded from use if the imperfect relationship between the proxy variable and the target of interest results in systematically penalizing members of a protected group that are otherwise qualified in the target of interest. In other words, since the proxy variable (height) is not a perfect predictor of having the target strength, there is some residual or unexplained variation in height across applicants that is unrelated to whether one has the required strength. The question is whether that residual is correlated with gender. In Figure 2, it is.

To avoid this result in *Dothard*, the Supreme Court therefore required a better proxy for required strength. In particular, the prison would be required to “adopt[] and valida[te] a test for applicants that measures strength directly” in order to justify disparities in hiring outcomes.¹³⁹ For example, assume that the prison implemented as part of the job application a physical examination that accurately assessed required strength, which produced the following results (Figure 3).

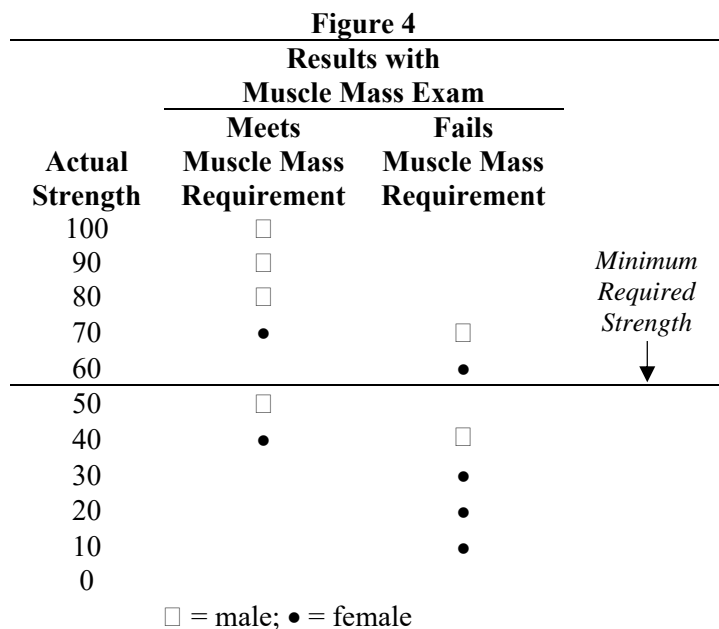
¹³⁹ *Id.* at 332.



The examination was perfect in classifying all individuals – male and female – as qualified if they in fact were so. Note that, even under this perfect exam, more males than females would be deemed eligible for the position. This disparity, however, arises solely through differences in actual strength (a legitimate business necessity).

Figure 3 is an ideal outcome in the sense that the prison was perfect in measuring each applicant’s actual strength, but perfect proxy variables are rarely available. Imagine instead that the prison asks the applicants to complete a simple muscle-mass index assessment (Figure 4).¹⁴⁰

¹⁴⁰ For instance, imagine the prison assesses each applicant’s mid-arm muscle circumference (MAMC) and requires a minimum measure which the prison believes is associated with having a minimum strength of 60. The MAMC is one of several techniques to measure muscle mass. See Julie Mareschal et al., *Clinical Value of Muscle Mass Assessment in Clinical Conditions Associated with Malnutrition*, 8 J. CLIN. MED. 1040 (2019).



As can be seen, muscle mass proxies for required strength with a positive, significant correlation, but it does so with error. In particular, there are applicants who are sufficiently strong but fail the muscle mass requirement, and there are applicants who meet the muscle mass requirement but are not sufficiently strong. The difference from Figure 2, however, is that the proxy is unbiased: Neither male applicants nor female applicants are favored by the fact that the proxy does not perfectly measure required strength. This is illustrated by the fact that one male and one female fail the muscle mass requirement but possess sufficient strength for the job, and one male and one female meet the muscle mass requirement but lack sufficient strength. Because the proxy is unbiased with respect to gender, an employer should therefore be permitted to use muscle mass as a proxy for required strength.

B. The Test in Regression Form

Moving from concepts to practice, standard regression techniques provide a straightforward means to implement the IAT. In keeping with the foregoing example, we return to the modified facts of *Dothard*, in which a prison uses a job applicant's height as a proxy for whether they have the required strength to perform the job of a prison officer.¹⁴¹ The prison does so based on the assumption that required strength is manifested in an

¹⁴¹ Of course, there might be multiple proxies. For instance, imagine the job requirements were strength and IQ, in some combination. Such a specification could be handled by more complex formations on the right-hand side of the regression framework that we discuss in this subsection.

individual's height. However, height is also determined by a host of other causes that are unrelated to strength. If we represent this group of non-strength determinants of height for a particular individual i as ε_i , we can summarize the relationship between the height and strength as follows:

$$\text{Height}_i = \alpha \cdot \text{Strength}_i + \varepsilon_i,$$

where α is a transformation variable mapping the relationship of strength to expected height. If ε_i was zero for each individual i , the equation becomes $\text{Height}_i = \alpha \cdot \text{Strength}_i$. In such a setting, an individual's height would be precisely equal to the individual's strength, multiplied by the scalar α . Therefore, one could compare with perfect accuracy the relative strength of two individuals simply by comparing their heights.

Where ε is non-zero, using height as a proxy for strength will naturally be less accurate; however, using height in this fashion will pose no discrimination concerns if ε (the unexplained variation in height that is unrelated to strength) is uncorrelated with a protected classification. This was precisely the case in Figure 4: Strength was *somewhat* manifested through the muscle mass index. Thus, it would be a useful variable for predicting which job applicants had the required strength for the job. Moreover, while it was error-prone in measuring actual strength (i.e., $\varepsilon_i \neq 0$), using one's muscle mass index to infer strength would pass the IAT:

$$\varepsilon_i \perp \text{gender};$$

the errors were not statistically correlated with gender, the protected category in our example. To implement this test empirically, the prison would use the historical data it holds concerning its existing employees' measured height and strength and regress employee height on employee strength to estimate α , which can be used to estimate ε_i for each employee.¹⁴² Using these ε_i estimates, the prison would then examine whether they are correlated with employee gender.

How would the IAT be used in a setting where the proxy is not a continuous measure (such as one's height or muscle mass) but rather a binary outcome of whether an individual possesses a specified level of the measure? Recall that this was the case in our hiring example where the prison first assessed an applicant's height and then applied a cut-off score to eliminate from consideration those applicants who fell below it. As reflected in

¹⁴² The regression will also estimate a constant term that is utilized in calculating the relationship between strength and height.

Dothard and *Lanning*, applying a minimum cut-off score to a proxy variable is a common decision-making practice, including within machine learning.¹⁴³

The application of the IAT would use the same framework, but using as the left-hand-side variable an indicator variable for whether an individual i was above or below the cutoff—for our example, $Height_i = 1$ for applicants above the cutoff and $Height_i = 0$ for applicants below it. To estimate a discrete 0-1 variable ($Height$) as a function a target (e.g., $Strength$), the preferred model is a logistic regression (or a comparable model for use with a dichotomous outcome variable). Logistic regression is a transformation that takes a set of zeros and ones representing an indicator variable and specifies them in terms of the logarithm of the odds ratio of an outcome (in our example, the odds ratio is the probability of $Height_i$ being above the cut-off divided by the probability that it is below the cut-off). This formulation is then regressed on the target measure ($Strength$). To generate the residuals (ε_i) for the IAT test, one predicts the probability of a positive outcome and then generates the error as the true outcome minus the predicted probability. As above, to pass the test, the residuals should not be significantly correlated with gender.

Finally, we conclude this overview with a discussion of what happens when a proxy variable fails the input accountability test: exclusion from the model. If the residuals (ε_i) are correlated with a protected classification (e.g., gender), it may be possible to “de-bias” a model that predicts strength from height, most notably by adding an individual’s membership (or lack of membership) in a protected class as an input in the predictive model.¹⁴⁴

However, as shown in the Appendix, the fact that de-biasing requires us to include *Gender* in the predictive model impairs the utility of this approach. A predictive model that explicitly scores individuals differently according to gender constitutes disparate treatment, making it a legally impermissible means to evaluate individuals. To avoid this challenge, proponents of this approach have therefore advocated that, in making predictions, the model should assign all individuals to the mean of the protected classification,¹⁴⁵ in our example, one would do so by treating all individuals as if $Gender = 0.5$ (i.e., $(1 + 0) / 2$) when estimating the effect of *Gender* on predicted strength. Doing so introduces prediction error, however, and as demonstrated by Kristen Altenburger and Daniel Ho, this error can be especially problematic

¹⁴³ See, e.g., Elizabeth A. Freeman & Gretchen G. Moisen, *A Comparison of the Performance of Threshold Criteria for Binary Classification in Terms of Predicted Prevalence and Kappa*, 217 *ECOLOGICAL MODELING* 48 (2008) (reviewing criteria for establishing cutoffs in ecological forecasting).

¹⁴⁴ This approach to de-biasing proxy variables has been advanced by several scholars. See Devin G. Pope & Justin R. Sydnor, *Implementing Anti-Discrimination Policies in Statistical Profiling Models*, 3 *AM. ECON. J.* 206, 206 (2011); Crystal Yang & Will Dobbie, *Equal Protection Under Algorithms: A New Statistical and Legal Framework*, John M. Olin Center For Law, Economics, and Business Discussion Paper No. 1019 (October 2019). We provide an example of this approach, as well as its limitations, in the Appendix.

¹⁴⁵ See, e.g., Pope & Sydnor, *supra* note 144, at 212.

when the approach is deployed in common machine-learning models.¹⁴⁶ More troublesome, these prediction errors can themselves be systematically biased against members of a protected group who are otherwise qualified in the target. We illustrate this challenge in the Appendix, which presents a simple example showing that this “de-biasing” procedure may actually have almost no effect on the extent of bias in the final outcome.

These considerations reinforce our conclusion that any variable that fails our test should be excluded from a decision-making model. While this approach risks sacrificing some degree of predictive accuracy in favor of ensuring an unbiased decision-making process, our discussion in Part 2(C) illustrates that U.S. antidiscrimination law has long made this trade-off. Additionally, a rule of exclusion also creates obvious incentives to seek out observable variables that can more accurately capture the target variable of interest, consistent with the holding of *Dothard* that the prison should adopt a test that more directly measured applicant’s strength.¹⁴⁷ Indeed, in the machine learning context, this history of U.S. antidiscrimination law provides an independent reason to adhere to a rule of exclusion given the capacity of machine-learning processes to analyze an ever-increasing volume of data to identify proxy variables that enhance accuracy while remaining unbiased with respect to a protected classification.

C. Challenges in Implementing the Test

Implementing the IAT faces several challenges, which we list below and then discuss in the context of the hiring test ($Height_i = \alpha \cdot Strength_i + \varepsilon_i$), where the target variable is *Strength*.

i. Unobservability of the Target Variable

The problem of an unobservable target variable of interest is always the starting point for constructing an algorithm to screen an applicant (or make some other decision), since the motivation for using statistical inference in the first place is the challenge of measuring unobservable characteristics such as creditworthiness, productivity, longevity, or threat to public safety.¹⁴⁸ In designing a machine-learning algorithm, the need to solve this problem arises in the training procedure, where data on a target variable are required to determine which features predict the target. In practice, the solution is to turn to historical data, which can be used to train the predictive model.¹⁴⁹ In the

¹⁴⁶ See Kristen M. Altenburger and Daniel Ho, *When Algorithms Import Private Bias into Public Enforcement: The Promise and Limitations of Statistical Debiasing Solutions*, 175 JOURNAL OF INSTITUTIONAL AND THEORETICAL ECONOMICS 98 109-118 (2018).

¹⁴⁷ See *supra* note 138.

¹⁴⁸ See Jon Kleinberg, Jens Ludwig, Sendhil Mullainathan and Cass R. Sunstein, *Discrimination in the Age of Algorithms*, 10 J. LEG. ANALYSIS 113, 132 (2019) (“One way to think about the goal of prediction is to overcome a missing information problem.”).

¹⁴⁹ *Id.*

employment setting, for instance, an employer seeking to predict the future productivity of job applicants could train a model with data concerning the productivity of existing employees along with data concerning the characteristics of these employees at the time of application. The data may suffer from selection bias given that the employer will not observe applicants who were not hired, which is why in both training a model and in running the IAT, one must be attendant to measurement error—a point we discuss in subsection (ii).

Nonetheless, the threshold challenge for the IAT—that the target is unobservable—is in many ways one of transparency. That is, data concerning the target variable exist (after all, these data were required to train the model), but they may not necessarily be available. As Jon Kleinberg, Jens Ludwig, Sendhil Mullainathan, and Cass Sunstein emphasize, transparency in the training data is therefore an important step in ensuring the ability to evaluate whether algorithmic decision-making facilitates discrimination.¹⁵⁰ We agree. The ability to examine the training data used in designing a model would allow a regulator, litigant or data scientist to conduct the empirical test we describe in this section. In the UnitedHealth example discussed in the Introduction, one could apply the IAT using actual morbidity data to assess whether the substitute measure of the target—the cost of healthcare—has a discriminatory effect. Indeed, the availability of actual morbidity data was what enabled researchers to quantify the racial bias in *Science*.¹⁵¹

Even with data on the target variable of interest, however, this last example highlights the problem of measurement error: Do the data on the target measure what they purport to measure with error? We address this problem in the following subsection.

ii. Measurement Error in the Target

In addressing the unobservability problem of the target, one can inadvertently mis-measure it. This challenge of measurement error—or what is alternatively referred to as “label bias”¹⁵²—has been studied in the computer science and economics literatures, providing useful guidance for addressing it when applying our test.¹⁵³

Consider, for instance, judicial bail decisions where data scientists have used past judicial bail decisions to train algorithms to decide whether a

¹⁵⁰ *Id.* at 114 (arguing that harnessing the benefits of algorithmic decision-making while avoiding the risk of discrimination “will only be realized if policy changes are adopted, such as the requirement that all the components of an algorithm (including the training data) must be stored and made available for examination and experimentation”).

¹⁵¹ Obermeyer, et al., *supra* note 7, at 447 (“Because we have the inputs, outputs, and eventual outcomes, our data allow us a rare opportunity to quantify racial disparities in algorithms and isolate the mechanisms by which they arise.”).

¹⁵² Corbett-Davies & Goel, *supra* note 75, at 3.

¹⁵³ See *id.* at 17-18.

defendant should be released on bail pending trial.¹⁵⁴ In many states, judges are required to consider the risk that a defendant poses for public safety, and in training the model, the target variable is often defined to be whether a defendant who was released was later arrested prior to the trial.¹⁵⁵ However, heavier policing in minority neighborhoods might lead to minority defendants being arrested more often than non-minorities who commit the same offense.¹⁵⁶ Consequently, Sam Corbett-Davies and Sharad Goel have warned that this form of label bias risks causing a model to estimate a positive relationship between a defendant's race (and correlates of race) and whether the defendant poses a risk to public safety, simply due to the correlation of race with measurement error.¹⁵⁷

Likewise, as Jon Kleinberg and others have noted, an employer who seeks to measure employee productivity through the number of hours that an employee spends at work will likely be using a biased measure of productivity if there are gender differences in how efficiently an employee works at the office (for example, to attend to childcare obligations before or after work).¹⁵⁸ Similar to the bail example, this form of label bias is problematic because the measurement error may be correlated with a protected characteristic, in this case, gender.

These examples illustrate a more general point, which is that measurement error in a target variable will create discriminatory bias when the measurement error is correlated with membership in a protected group. This result occurs because a statistical model that seeks to estimate the predictors of a true target y that is mis-measured as $y + \mu$ will inevitably discover that the protected classification (and any correlate of it) predicts the level of the mis-measured target.

For similar reasons, when measurement error in a target variable is correlated with a protected classification, application of our test may fail to detect this bias. Returning to the *Dothard* example, imagine that we applied the IAT to *Height* as before, but we use a measure for strength, *Strength**, that has measurement error μ that is correlated with gender. Formally, the test would be:

$$Height_i = \alpha \cdot Strength_i^* + \varepsilon_i$$

which is equivalent to:

$$Height_i = \alpha \cdot (Strength_i + \mu_i) + \varepsilon_i$$

¹⁵⁴ See, e.g., Berk et al., *supra* note 79, at 31-33.

¹⁵⁵ *Id.* at 31.

¹⁵⁶ Corbett-Davies & Goel, *supra* note 75, at 18.

¹⁵⁷ *Id.*

¹⁵⁸ See Kleinberg et al., *supra* note 148, at 139.

In such a setting, the IAT may fail to reveal that the errors (ε_i) are correlated with the protected classification of gender. The reason is because the unexplained variation between “true” *Strength* and *Height* is $(\mu_i + \varepsilon_i)$, but the IAT will not be able to detect how gender is correlated with μ_i because it is part of *Strength**, the mis-measured target. In short, measurement error in a target variable is a critical issue to consider regardless of whether one is calibrating a model or running our test.

Recognition of this latter point is implicit in Kleinberg, Ludwig, Mullainathan, and Sunstein’s argument for making training datasets transparent. Often, the data for a target will reveal fairly obvious risks that the measurement error is biased with respect to a protected classification (such as the example cited earlier when an employer uses hours-worked as a measure for productivity). At the same time, other instances when this problem arises may be less obvious. In the example we provided in the Introduction, that of UnitedHealth, it may not have been immediately obvious that patient costs—the substitute measure of the target of interest—had measurement error for the true target (severity of illness) that was correlated with race. Yet this correlated measurement error was nevertheless revealed when researchers used an alternative estimate for severity of illness.

This last example thus underscores the need to run the IAT with alternative measures of the target which may reveal problematic measurement error in the primary target data. Moreover, opening up the possibility of running the IAT with alternative measures of the target variable should also encourage the use of theory-based models of target characteristics. Theory-based estimates of target variables may be especially valuable in addressing the measurement error that arises from estimating targets based on binary outcomes. Common approaches to estimating target variables often rely on estimating a predictive model based on a binary outcome variable, such as whether a borrower defaults on a loan or whether a defendant who was released on bail later commits a crime prior to trial. Yet estimating unobservable characteristics such as “creditworthiness” or “risk” based on these binary behaviors necessarily implicates the risk of measurement error in the true target of interest.

Consider, for instance, a model that seeks to predict creditworthiness based solely on whether a borrower defaults in the training data. By construction, the training dataset consists only of those borrowers who received a loan; borrowers who do not get a loan provide no information. Thus, it is infeasible to estimate actual creditworthiness within the broader group of all applicants. This is the “selective labels” problem that has been

studied in the computer science and economics literatures.¹⁵⁹ The literature on selective labels in training a model has suggested a process of interventions to correct the misestimations.¹⁶⁰ Another approach would be to implement the IAT through a structural estimation of theoretic representations of the target business necessity.¹⁶¹

Another version of the problem of measurement error comes in the context of threshold analysis. In our example, the prison asserted that it needed a minimum required level of strength. As a result, the target was not the continuous variable of strength, but the applicant possessing a strength level of at least 60, which we assumed was a legitimate business necessity threshold for a prison guard job. But what if the level of strength needed is not obvious? What if the prison erroneously thought the true level of required strength was 80? We previously referred to this setting as a mis-asserted target threshold. Cases such as *Lanning v. Southeastern Pennsylvania Transportation Authority* underscore the potential for these target thresholds to be mis-asserted in a way that results in intentional discrimination, such as when they are purposefully set at a level that will adversely affect members of a protected group.

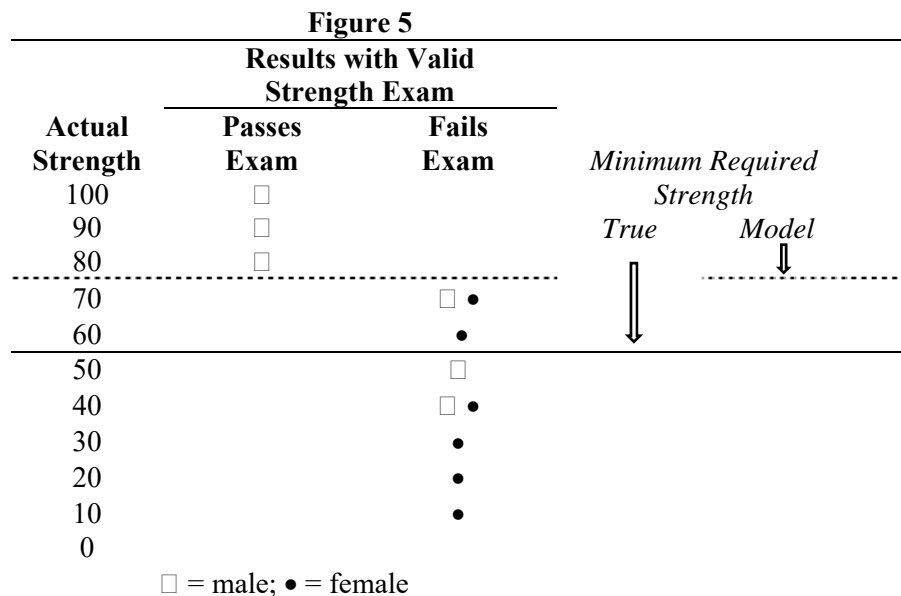
In Figure 5, we assume that, as in Figure 3, the prison implements a physical exam that perfectly measures actual strength. If the prison mistakenly sets the minimum required strength threshold at 80 (the dashed line), the resulting problem is that more women cluster in the just-failed space (between the dashed and straight line), which is the region of between the mis-asserted target threshold relative to the true required strength level. In fact, if an employer did not want to hire women, it could intentionally

¹⁵⁹See Himabindu Lakkaraju, et al., *The Selective Labels Problem: Evaluating Algorithmic Predictions in the Presence of Unobservables*, in KDD Conference Proceedings, 2017; Kleinberg et al., *Human Decisions and Machine Predictions*, 133 Q. J. ECON. 237, 256- (2018).

¹⁶⁰ See, e.g., Maria De-Arteaga, et al., *Learning Under Selective Labels in the Presence of Expert Consistency*, (July 4, 2018), <https://arxiv.org/abs/1807.00905v1> (proposing a data augmentation approach that can be used to leverage expert consistency to mitigate the partial blindness that results from selective labels).

¹⁶¹For instance, consider a credit scoring algorithm that predicts credit risk based on default rates for loans that were previously extended to a group of borrowers. A model built using these target data (i.e., whether or not a borrower defaults) suffers from bias insofar as it only includes default data for loans that were approved by a lender. This selective labels problem can result in bias if the human decision-maker who approved the loans based the approval decision on borrower characteristics that were observable to the loan officer but are unobservable to the data scientist because they do not appear in the dataset. Imagine, for instance, that a loan officer records data on a loan applicant's occupation and, for low-paying occupations, the loan officer also evaluates informally an applicant's attire, which the officer believes is associated with creditworthiness. Assume the loan officer approves loans to well-dressed applicants in occupations that would otherwise make them ineligible for a loan and that these applicants are, in fact, more creditworthy than their occupation would suggest. Training a predictive model using only default data and occupation at the time of application would therefore suggest to the model that "high risk" occupations are actually more creditworthy than they are because they default infrequently. Moreover, given racial, ethnic and gender differences in the composition of certain occupations, this model would likely be biased in addition to being inaccurate. However, evidence of this bias would become apparent in applying the IAT if one were to run the test using an estimate for creditworthiness that was based on borrowers' cash flow data as opposed to default data.

implement a mis-asserted target, knowing that more women would be excluded.



In this setting, the exam would pass the IAT insofar that it was unbiased with respect to gender in predicting whether an applicant had strength of at least 80. However, the employer’s use of the exam would nevertheless fail our definition of accountability set forth in Part 2 because the employer has set the cut-off at a level where qualified females are systematically excluded from the position. As emphasized in *Lanning*, this example underscores the importance of supplementing the IAT with the ability to scrutinize whether a classification threshold has been set at a level that is justified by actual business necessity.

iii. Testing for “Not Statistically Correlated”

The third challenge concerns how to reject the null hypothesis that no correlation exists between a set of proxy variable residuals and a protected category. In our *Dothard* illustration, the use of *Height* as a proxy for *Strength* would pass the IAT if the unexplained variation between *Strength* and *Height* (ϵ_i) is uncorrelated with *Gender*, as given by the test:

$$\begin{aligned} \text{Regression:} & \quad \epsilon_i = \beta_0 + \beta_1 \text{Gender}_i \\ \text{Null Hypothesis:} & \quad \beta_1 = 0. \end{aligned}$$

The tradition in courts and elsewhere is to use a statistical significance level

of 0.05;¹⁶² i.e., we are willing to allow for a 5% probability of making the “Type I” error of rejecting the null hypothesis ($\beta_1 = 0$) by chance, when it is actually true. A related concept is the p-value of an estimate: the probability of obtaining an estimate for β_1 at least as far from zero as the value estimated, assuming the null hypothesis is true. If the p-value is smaller than the statistical significance level, one rejects the null hypothesis.

However, a problem with focusing on p-values is that as the sample size grows increasingly large, realized p-values converge to zero if the sample estimate for β_1 is even trivially different from the null. This is because as the sample size grows larger, the uncertainty of our estimates (usually measured by their “standard error”) gets closer and closer to zero, causing any coefficient (even magnitude-irrelevant ones) to look different from an exact null of $\beta_1 = 0$ in a p-value test. In particular, a company that brings a large dataset to bear on an IAT test might be disadvantaged relative to firms with less data.

The source of the problem is the fact that in any statistical test we are actually trading off the probabilities of making two different errors: Type I errors (when we wrongly reject the null when it is, in fact, true) and Type II errors (when we wrongly fail to reject the null when it is, in fact, false). The “significance level” of a test is the probability of making a Type I error. Keeping this fixed (e.g., at 5%) as the sample size increases means that we are keeping the probability of a Type I error fixed. But at the same time, again because the standard error of our estimates is going to zero as the sample size gets large, the probability of a Type II error is actually converging to zero. If we care about both types of error, it makes sense to reduce the probability of *both* as the sample size increases, rather than fixing the probability of Type I errors and letting that of Type II errors go to zero. This point has been made forcefully by many authors, especially Edward Leamer, and a number of solutions have been proposed for adjusting the significance level as the sample size increases.¹⁶³ A full consideration of these different approaches is

¹⁶² See, e.g., Karen A. Gottlieb, *What Are Statistical Significance and Probability Values?* 1 TOXIC TORTS PRAC. GUIDE § 4:10 (2019) (“Through a half century of custom, the value of 0.05 or 1 in 20 has come to be accepted as the de facto boundary between those situations for which chance is a reasonable explanation (probabilities > 0.05) and those situations for which some alternative is a reasonable explanation (probabilities < 0.05).”); see also *Eastland v. Tennessee Valley Authority*, 704 F.2d 613, 622 n. 12 (1983) (in employment discrimination lawsuit, noting that “a probability level of .05 is accepted as statistically significant” in determining whether racial disparities in pay were statistically significant).

¹⁶³ See, e.g., Edward Leamer, SPECIFICATION SEARCHES: AD HOC INFERENCE WITH NONEXPERIMENTAL DATA (1978) (proposing p-value adjustment to minimize error losses associated with Type I and Type II error); I.J. Good, *Standardized Tail-Area Probabilities*, 16 JOURNAL OF STATISTICAL COMPUTATION AND SIMULATION 65 (1982) (proposing p-value adjustment based on a “Bayes/non-Bayes compromise”); Mingfeng Lin, Henry C. Lucas, Jr., and Galit Shmueli, *Too Big to Fail: Large Samples and the p-Value Problem*, 24 INFORMATION SYSTEMS RESEARCH 906, 908-915 (2013) (surveying approaches to adjusting p-values in large samples and recommending the reporting of effect sizes and confidence intervals and using coefficient/p-value/sample-size plots for interpreting the data along with Monte Carlo simulations); Eugene Demidenko, *The p-value You Can’t Buy*, 70 THE AMERICAN STATISTICIAN 33, 34-37 (2016) (proposing use of d-values for assessing statistical inference in large datasets).

beyond the scope of this Article; however, we provide below an example of one such approach to illustrate how it can be utilized to discern when a seemingly significant result when applying the IAT is actually a function of the large sample size and not evidence of a discriminatory proxy variable.

iv. Nonlinearities or Interactions Among Proxies

Machine learning models are often focused on forming predictions based on nonlinear functions of multiple variables. In introducing the IAT, our specification focused on linear settings, but the IAT could in principle be amended to handle nonlinear models as well. For example, rather than just running the test regression once, we could run it repeatedly, with each of a set of basis functions of the explanatory variables on the left-hand side. Full consideration of this topic is beyond the scope of this Article, but in general, implementation of the IAT could be made part of the type of feature selection and feature analysis protocols that are used in practice with both linear and non-linear machine-learning processes.¹⁶⁴

D. Simulation

To illustrate how the concerns of discrimination enter through proxy variables, we simulate the setting in *Dothard* of hiring a prison worker.

i. Set Up

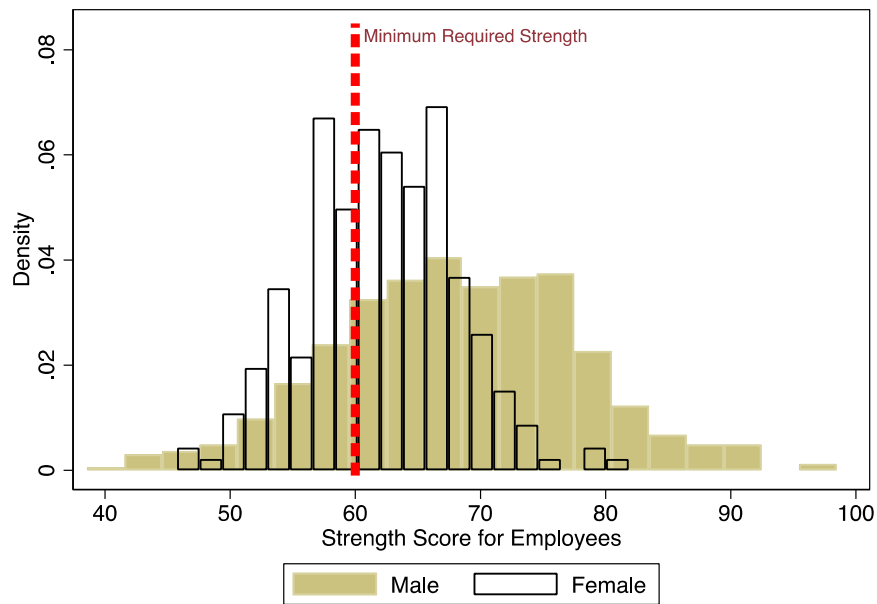
The simulation assumes that the prison has historical records for 800 employees, of which roughly one-third are female (n=256) and two-thirds are males (n=544). We further assume that the prison uses these historical records to develop a sorting algorithm for considering a pool of 1,200 applicants. The 800 employees are endowed with an *unobservable* strength level, which we model as a random variable distributed normally with (i) a mean of 68 and a standard deviation of 10 for male employees and (ii) a mean of 62 and a standard deviation of 6 for female employees. With these modeling assumptions, females have lower mean strength but a smaller standard

¹⁶⁴In particular, a related literature in computer science focuses on feature selection to enhance model interpretability. See Datta et al. *Algorithmic Transparency via Quantitative Input Influence: Theory and Experiments with Learning Systems*, *Proceedings of IEEE Symposium on Security & Privacy 2016*, 598–617, 2016 (proposing a quantitative-input-influence (QII) protocol based upon Shapley values to determine the importance of features and clustering metrics to summarize feature influence); see also Phillippe Bracke et al., *Machine learning explainability in default risk analysis*, Bank of England Staff Working Paper No. 816 (June 5, 2019) (implementing QII method in predicting mortgage defaults). More formally, Lundberg, et al., *Consistent Individualized Feature Attribution for Tree Ensembles*, [arXiv:1802.03888v3](https://arxiv.org/abs/1802.03888v3) [cs.LG], March 7, 2019 and Merrill et al., *Generalized Integrated Gradients: A practical method for explaining diverse ensembles*,” ArXiv 2019, build upon game-theoretic SHAP (Shapley Additive explanation) values and propose new feature credit-assignment algorithms that can handle a broad class of predictive functions with both piecewise-constant (tree-based), continuous (neural-network or radial-basis-function based), and mixed models.

deviation, as plotted below in Figure 6. To be an effective prison guard requires a strength of 60, the business necessity. Hiring is not perfectly effective at sorting which guards will meet this threshold; therefore, even among the employees, there are guards who fall below the required strength for the job. For now, we assume that the prison can implement a costly physical exam to measure true strength for these employees. (We abstract from other aspects of effectiveness such as psychological and managerial skills needed for prison-guard work.)

We assume the strength of applicants is likewise distributed randomly. However, for obvious reasons, the applicant pool has not been previously selected for strength as employees have. Therefore, we model strength across applicants as a random variable distributed normally with a mean of 50 and a standard deviation of 10 for male employees and a mean of 44 and a standard deviation of 6 for female employees.

Figure 6



The prison managers cannot directly observe applicants’ strength, and, as noted, implementing a full physical exam across applicants is costly. Therefore, the prison decides to use height as a proxy variable for an applicant’s strength, since it is easily measured on applications. We model height as a sum of a baseline 50 inches (with a normally-distributed error of 4 inches) plus a concave (quadratic) function increasing in strength. Female height has the same relation to strength but a ten percent lower baseline. The

resulting mean height in the employee training dataset is 5'10" with a standard deviation of 5".

Finally, as in *Dothard*, the prison seeks to filter applicants by imposing a minimum height requirement. To determine the height cut-off, the prison runs a classification analysis. In doing so, the prison determines that they want to ascertain that an individual will be above the strength threshold with an 80% certainty, i.e., they want only a 20% risk of incorrectly classifying an applicant as eligible for hiring (above the strength threshold of 60) when the person in fact has a strength of less than 60. Based on the height and strength of the prison employees, this results in a 5'10" cut-off. The prison applies this cut-off to all 1,200 applicants.

Among the 370 female applicants, 344 (93%) fail the height test. In contrast, among the 830 male applicants, 504 (61%) fail the height test. These disparities suggest that the height cut-off may discriminate against females applicants, but we cannot definitively conclude this from the high rejection rates because, as we saw in Figure 6, females in our samples have lower strength than males on average.

ii. Applying the Input Accountability Test

Assume that in advance of deploying the height test, the prison instead decides to conduct the IAT to ensure that any disparities in hiring would be based on differences in predicted applicant strength. Table 1 presents the results from the test. To run the IAT, the prison would return to the training data it possesses regarding its employees' actual strength and height that it used to determine the 5'10" cut-off. In panel A, we present the first step of regressing the proxy variable on employee strength, the target of interest. Because the prison is focused on using a cutoff for height, we estimate a logistic regression of whether an employee passes the height cut-off as a function of the employee's strength. (To do so, we use as our dependent variable an indicator variable that equals 1 for employees that are at least 5'10" and 0 for all others.) Note that this indicator variable is on the left-hand side of the regression (and not strength) because we want to decompose whether an employee meets the height cut-off into two components – the part that can be predicted from an employee's strength and the part that cannot be predicted from an employee's strength (the "residual"). Stated differently, logistic regression effectively estimates the probability that an employee is 5'10" based on employee strength. Therefore, the residual, which is equal to one minus this predicted probability for each employee, can be viewed as the variation in whether an employee meets the height threshold of 5'10" that is unrelated to an employee's strength. In panel B, we present the results from regressing the residual from panel A onto the indicator variable for female.

Table 1

	(1)	(2)	(3)	(4)	(5)
Panel A: First Step of IAT (DV=Column Heading)					
	Cut-Off Height	Cut-Off Muscle Mass	Muscle Mass	Job Performance	Cut-Off Muscle Mass
<i>Strength</i>	0.0206*** [0.00155]	0.0377*** [0.000747]	0.9965*** [0.0191]		0.0387*** [0.0000138]
<i>Performance Score</i>				0.675*** [0.0307]	
Observations	800	800	800	800	2,000,000
[Pseudo] R-squared	0.111	0.466	0.772	0.376	0.496
Panel B: Second Step of IAT (DV=Residuals from Step 1)					
<i>Female</i>	-0.354*** [0.0327]	-0.013265 [0.02625]	-0.3552 [0.379]	-8.858*** [0.542]	-0.0013*** [0.000505]
Observations	800	800	800	800	2,000,000
R-squared	0.128	0	0	0.25	0
d-value					50%

Standard errors in brackets

*** p<0.01, ** p<0.05, * p<0.1

Panel A of Column (1) reports that strength only accounts for a small part of the variation (R-squared = 0.111) for whether an employee is (or is not) taller than 5'10". In Panel B, our column (1) results show that the residual of the first step regression has a negative, significant correlation with gender, thus failing the IAT. Females incur a penalty because the proxy variable for the business necessity of required strength has residual correlation with gender.

Imagine that the prison realizes this flaw in using a height cut-off and decides instead to consider incurring an extra cost for doing a muscle-mass index evaluation of applicants. Because the evaluation is imperfect in assessing true strength, we assume that the results of a muscle-mass index evaluation is equal to an individual's strength plus random noise.¹⁶⁵ To implement this screening procedure, the prison first applies the muscle-mass index evaluation to existing employees so that it can estimate the minimum muscle mass an individual should have to be above the minimum strength

¹⁶⁵ We model the random noise as a randomly distributed variable with a mean of zero and a standard deviation of 5.

threshold with an 80% certainty. The classification analysis produces a muscle-mass cut-off score of 64. As above, the prison then conducts the IAT.

In column (2) of panel A we present the results of the IAT for the muscle-mass index evaluation based on the employee training data. To implement the IAT, we run the same regressions that we used for testing the height cut-off, but we substitute an indicator variable for whether an employee has a muscle mass of at least 64 for the indicator variable for whether an employee is at least 5'10". In panel A, column (2) shows that the probability that an employee has a muscle mass of at least 64 is (unsurprisingly) related to an employee's strength, resulting in a much larger R-squared. Importantly, the residual should not fail the IAT, because it has no bias against females. In column (2) of panel B, we see that this is indeed the case; the coefficient on female is statistically insignificant and small in magnitude.

In column 3, we instead consider a continuous variable version of muscle mass as a scoring variable rather than a cut-off version of the indicator variable. Perhaps the underlying job-required strength is not a threshold but a strength score that will feed into wage-setting or other profiling of individuals that focus on continuous rather than discrete measures. To implement the IAT in this context, we use the same training data that was used for column (2) of Table 1; however, the regression specification for the first step takes the form of a linear regression of employees' muscle mass scores on their measured strength. As in column (2), column (3) shows that muscle mass is a legitimate business necessity variable. In panel A, we find that muscle mass and strength are very correlated, with strength accounting for almost 80% of the variation in muscle mass. Column (3) of panel B shows that muscle mass again passes the IAT: the residual is uncorrelated with the female indicator variable.

In the final two columns of Table 1, we demonstrate the importance of the challenges we introduced in Part 3(C).

First, we use column (4) to illustrate the concern about measurement error in the target (strength). Thus far, we have been working under the assumption that the prison can take an accurate measurement via a physical exam of the training dataset employees. However, what if instead the prison cannot measure actual strength but uses a job performance assessment made by a manager. (We label this job performance measure an employee's "Performance Score"). As noted above, a central challenge in real world settings is that target variables used to train predictive models are typically estimated in this fashion and may contain measurement error that is correlated with a protected characteristic. We therefore simulate an employee's Performance Score as biased against females.¹⁶⁶ In this regard, the simulation

¹⁶⁶ In particular, for males, we model the job performance measure as strength plus random noise; however, for females, we model job performance as concave in strength (like the height variable)—a quadratic

replicates the same problem illustrated with the UnitedHealth example (where the illness severity measure was inadvertently biased against African Americans).

In addition to employees' Performance Scores, assume that the prison also has at its disposal data from the muscle measure index evaluation used in columns (2) and (3). Even without perfect data regarding employee strength, the prison can still use these data with the IAT to evaluate whether its preferred estimate of the target (an employee's Performance Score) suffers from bias. To implement this test, we treat muscle mass as an alternative measure of the target of interest (strength), and we treat the Performance Score as a proxy for strength. Accordingly, the first step of the IAT is conducted by regressing employees' Performance Scores on the muscle mass evaluation data. The results are shown in column (4) of panel A. Not surprisingly, an employee's muscle mass is closely related to an employee's Performance Score. In column (4) of panel B, we show the results of regressing the residuals from this regression on the gender variable. As shown in the table, Job Performance fails the IAT. In this fashion, the IAT can be used to test whether an estimate for a target suffers from biased measurement error, so long as one has an alternative estimate for the target (even a noisy one) that is believed to be unbiased.

The final column in Table 2 illustrates the concern of large data samples. For this column, we implement the same muscle mass test as in column (2), except that we randomly draw 2 million employees for the training dataset rather than 800 employees. (For all 2 million employees, we model their strength using the same assumptions used for the original 800 employees). For each employee, we likewise calculate muscle mass as employee strength plus a random variable distributed normally with a mean of 0 and a standard deviation of 5. Thus, in our simulated setting, muscle mass is a noisy estimate of employee strength but it has zero bias with respect to gender. Even so, however, the possibility remains that in drawing random measurement error for our sample, very slight differences may exist by chance between the average measurement error of females and males. (This is equivalent to observing that even if a coin is unbiased, it may still return more than 50% heads in a trial of 100 flips). Moreover, as we described in section 3, the p-value may converge to 0 for any small deviation, as sample sizes approach infinity. Thus, even a small (economically non-meaningful) correlation may look significant. This would create a setting of a large-dataset proxy variable failing the IAT, not because of a fundamental problem, but just because of the use of a fixed p-value. This is what we have modeled in column (5). The coefficient on female in column (5) is very small (-0.0013) but statistically

concave function of strength plus random noise. The managers evaluating females do not fairly evaluate them, especially for the stronger females.

significant, notwithstanding the fact that we modeled measurement error from a distribution that had exactly zero gender bias.

As noted in subsection 3(C)(iii), where the IAT is applied to a large dataset, it is therefore critical to check whether a proxy that fails the IAT might have failed the test simply because of the large number of observations in the sample. That the seemingly statistical finding in column (5) may be an artifact of a trivial difference within a large dataset can initially be seen by the fact that the R-squared in column (5) is 0%; if effectively no variation in the residuals can be explained by gender, how can it be that this proxy is penalizing females in a systematic fashion? Additionally, as noted previously, a number of formal solutions also exist to examine this issue more fully. Here, we illustrate one such approach using the concept of the “d-value” proposed by Eugene Demidenko.¹⁶⁷ Rather than focus on a comparison of group means, the d-value is designed to examine how a randomly chosen female fared under this proxy variable relative to a randomly chosen male. Specifically, in the context of the IAT, the d-value answers the question “what is the probability that members of a protected group are being penalized by the proxy?” As shown in the last row of column (5), the d-value is approximately 50%, indicating that the probability that females are penalized by the use of a muscle-mass proxy is effectively a coin-toss; that is, there is no evidence that female applicants are being systematically penalized by the use of this proxy.

This finding, of course, is hardly a surprise given that we designed the simulation to ensure that it was an unbiased proxy. In this fashion, the use of a d-value can highlight when a seemingly significant finding is a function of the large sample size and not evidence of a discriminatory proxy variable.¹⁶⁸

IV. APPLICATIONS BEYOND EMPLOYMENT

The fact that the IAT is rooted in general antidiscrimination principles makes it applicable to any setting where a decision-maker relies on statistical discrimination, regardless of whether conducted by humans or algorithms.

¹⁶⁷ See Demidenko, *supra* note 163.

¹⁶⁸ To the extent one utilizes the d-value in this fashion, a natural question is what level of a d-value would constitute evidence of a discriminatory proxy. Given that the d-value answers the question “what is the probability that members of a protected group are being penalized by the proxy?”, any result that yields a d-value deviating from 50% would presumably be evidence of a discriminatory proxy, allowing for a percentage difference to incorporate a far tail sampling draw. This conclusion follows from the conventional judicial reliance to on p-values, which likewise assumes that any finding with a p-value of less than 0.05 is evidence of discrimination. That said, in adopting such an approach, it would be important to utilize a d-value analysis only upon a finding that a proxy fails the IAT using a conventional statistical test. The reason stems from the fact that in smaller samples, even an unbiased proxy could result in a d-value that is slightly different from 50% due sample variance. For example, the d-value for column (3) is just slightly less than 51%, despite the fact that muscle mass is modeled as an unbiased proxy. However, running the same simulation with 50,000 observations produces a d-value of 50%.

Central to our argument is the idea of using a test to ascertain adherence to business necessity targets when designing a decision-making process. Indeed, even the Equal Employment Opportunity Commission subscribes to a business necessity *test* in its Uniform Guidelines on Employee Selection Procedures, stating that: “[e]vidence of the validity of a test or other selection procedure by a criterion-related validity study should consist of empirical data demonstrating that the selection procedure is predictive of or significantly correlated with important elements of job performance.”¹⁶⁹ Note, however, that even the EEOC’s validity test looks only to the predictive capacity of an employment exam. But as we have emphasized throughout this Article, a simple correlation test leaves open the possibility that a test will penalize members of a protected group who are, in fact, qualified in the job-related skill in question. This, of course, was the lesson of *Dothard* and the cases examined in Part 2. In this regard, a simple method to remedy this defect when conducting a criterion-related validity study would be to incorporate the IAT.

In this section, we discuss additional implementations outside of the employment setting. We first focus on settings where a decision-maker can face liability for claims of unintentional discrimination and where a court or legislature has expressly considered what constitutes a legitimate business necessity target. We then address the application of the IAT in settings where formal liability for claims of disparate impact or other claims of unintentional discrimination are currently less clear, but where firms can use the IAT to self-regulate. Finally, given the latitude firms have to set their own business necessity targets, we conclude with an admonition that firms must be vigilant in monitoring whether a purported target is, in fact, a legitimate one to use.

A. *Domains with Court-Defined Business Necessity Targets*

Consider, for instance, a regulator tasked with evaluating a decision-making algorithm in one of the following domains where claims of unintentional discrimination may be possible, and where courts have expressly defined a legitimate “target” variable that can justify unintended disparities that vary across protected and unprotected groups:

Table 2	
Domain:	Legitimate Target Variable:
Credit Determinations	Creditworthiness ¹⁷⁰

¹⁶⁹ 29 C.F.R. § 1607.5B.

¹⁷⁰ See *A.B. & S. Auto Service, Inc. v. South Shore Bank of Chicago*, 962 F. Supp. 1056 (N.D. Ill. 1997) (“[In a disparate impact claim under the ECOA], once the plaintiff has made the prima facie case, the defendant-lender must demonstrate that any policy, procedure, or practice has a manifest relationship to

Home Insurance Pricing	Risk of Loss ¹⁷¹
Parole Determinations	Threat to Public Safety ¹⁷²
Tenant Selection	Ability to meet lease obligations, ¹⁷³ pay rent, ¹⁷⁴ and resident safety ¹⁷⁵
Post-Secondary School Admission	Predicted academic success ¹⁷⁶
Selection into Special Education	Educational ability ¹⁷⁷
State Merit Scholarship Eligibility	Academic achievement in high school ¹⁷⁸

Just as employers are permitted to make hiring decisions based on the legitimate target variable of a job-required skill, courts in these settings have likewise determined that decision-making outcomes can lawfully vary across protected and unprotected groups only if decisions are based on the target variable noted in Table 2.

the creditworthiness of the applicant..."); *see also* Lewis v. ACB Business Services, Inc., 135 F.3d 389, 406 (6th Cir. 1998) ("The [ECOA] was only intended to prohibit credit determinations based on 'characteristics unrelated to creditworthiness.'"); Miller v. Countrywide Bank, NA, 571 F.Supp.2d 251, 258 (D. Mass 2008) (rejecting argument that discrimination in loan terms among African American and white borrowers was justified as the result of competitive "market forces," noting that prior courts had rejected the "market forces" argument insofar that it would allow the pricing of consumer loans to be "based on subjective criteria beyond creditworthiness.")

¹⁷¹ *See, e.g.,* Owens v. Nationwide Mut. Ins. Co., No. Civ. 3:03-CV-1184-H, 2005 WL 1837959, at *9 (N.D. Tex. Aug. 2, 2005) (minimizing the "risk of loss in homeowner's insurance" was a legitimate business necessity under the Fair Housing Act that justified the use of facially neutral policy of using credit to determine eligibility for homeowner's insurance).

¹⁷² *See, e.g.,* CAL. PENAL. CODE § 3041 (West 2017) (The Board of Prison Term "shall grant parole to an inmate unless it determines that the gravity of the current convicted offense or offenses, or the timing and gravity of current or past convicted offense or offenses, is such that consideration of the public safety requires a more lengthy period of incarceration for this individual."); *see also* Smith v. Sisto, 2009 WL 3294860 at *6 (E.D. Cal. Oct. 13, 2009) (denying claim that denial of parole constituted discrimination and concluding that "[t]he need to ensure public safety provides the rational basis for section 3041").

¹⁷³ *See* 24 C.F.R. § 100.202(c)(1) (permitting under the FHA a landlord's "[i]nquiry into an applicant's ability to meet the requirements of ownership or tenancy").

¹⁷⁴ *See* Ryan v. Ramsey, 936 F.Supp. 417 (S.D. Texas 1996) (noting that under the FHA, "there is no requirement that welfare recipients, or any other individuals, secure apartments without regard to their ability to pay.")

¹⁷⁵ *See* Evans v. UDR, Inc., 644 F.Supp.2d 675, 683 (2009) (permitting landlord to reject tenant based on prior criminal history as the "policy against renting to individuals with criminal histories is thus based concerns for the safety of other residents of the apartment complex and their property").

¹⁷⁶ *See* Kamps v. Baylor University, 592 F. App'x 282 (5th Cir. 2014) (rejecting age discrimination case based on law school admissions criteria that relied on applicant's grade point average (GPA) because GPA is quantitative predictor of academic success in law school and thus a "reasonable factor other than age").

¹⁷⁷ *See* Ga. State Conf. of Branches of NAACP v. Georgia, 775 F.2d 1403, 1420 (11th Cir. 1985) (finding that, in Title VI case alleging that school district achievement grouping caused disparate impact on minority students, school district's effort to classify students based on assessment of ability was justified because it bore "a manifest demonstrable relationship to classroom education").

¹⁷⁸ *See* Sharif by Salahuddin v. New York State Educ. Dept., 709 F. Supp. 345, 362 (SDNY 1989) (finding that state's use of SAT scores did not have a "manifest relationship ... [to] recognition and award of academic achievement in high school" in Title IX claim of disparate impact alleging that state's use of SAT scores to determine student eligibility for merit scholarships had a discriminatory effect on women).

In applying the IAT in these settings, the regulator’s task thus follows the same process noted in Part 3. First, the regulator must evaluate whether the decision-making process does, in fact, seek to produce outcomes based on the legitimate target variable. Second, using historical data for both the target variable and the model’s full set of features, the regulator would then apply the IAT to each feature used in the model. Finally, any feature that failed the test would be required to be excluded from the model.

B. Domains Without Court-Defined Business Necessity Targets

The IAT is equally applicable to domains where antidiscrimination laws do not formally regulate decision-making processes governing disparities across protected and unprotected groups or where the legal risk for unintentional discrimination is presently unclear. We provide an example of each.

The first domain concerns insurance outside the context of home insurance.¹⁷⁹ As Ronen Avraham, Kyle Logue, and Daniel Schwarcz show, a number of jurisdictions do not have any laws restricting providers of automobile or life insurance from discriminating on the basis of race, national origin, or religion.¹⁸⁰ Nor is there a federal antidiscrimination statute applicable to insurance outside of the context of home insurance.¹⁸¹ Consequently, insurers likely have considerable discretion to rely on statistical discrimination to underwrite policies, which may produce unintended disparities across protected and unprotected groups. Yet evidence that racial disparities exist in the pricing of auto loans has routinely been met by the insurance industry with assurances that premiums are based on risk. For instance, following a nationwide study by the Consumer Federation of America in 2015 that found that predominantly African-American neighborhoods pay higher auto premiums,¹⁸² the Property Casualty Insurers Association of America responded with a declaration that “Insurance rates are color-blind and solely based on risk.”¹⁸³ Thus, insurers claim to self-regulate themselves by setting risk as the business necessity target. To the

¹⁷⁹ As noted in Table 2, discrimination in home insurance is governed by the FHA.

¹⁸⁰ See Ronen Avraham, Kyle D. Logue & Daniel Benjamin Schwarcz, *Understanding Insurance Anti-Discrimination Laws*, 87 S. Cal. L. Rev. 195, 239 (2014).

¹⁸¹ *Id.* at 241. Additionally, the few cases alleging discrimination by insurance providers under 42 U.S.C. § 1981—a Reconstruction-era statute that prohibits racial discrimination in private contracting—have required a showing of intentional discrimination. See, e.g., *Amos v. Geico Corp.*, 2008 WL 4425370 (U.S. Minn. 2008) (“To prevail under § 1981, plaintiffs must prove that GEICO intentionally discriminated against them on the basis of race.”).

¹⁸² Consumer Federation of America, *High Price of Mandatory Auto Insurance in Predominantly African American Communities* (2015), available at https://consumerfed.org/wp-content/uploads/2015/11/151118_insuranceinpredominantlyafricanamericancommunities_CFA.pdf.

¹⁸³ Press Release of American Property Casualty Insurers Association of America, *Auto Insurance Rates are Based on Cost Drivers, Not Race*, November 18, 2015, available at <https://www.pciaa.net/pciwebsite/cms/content/viewpage?sitePageId=43349>.

extent insurers are sincere in this claim, the IAT provides them with a ready test to ensure compliance.

An example in the second domain concerns disparities in medical treatment, as motivated by our example in the Introduction concerning UnitedHealth. Discrimination in healthcare provision is covered by Title VI of the Civil Rights Act of 1964, thus making it a more regulated setting than the insurance example. However, in *Alexander v. Sandoval*,¹⁸⁴ the U.S. Supreme Court held that Title VI does not provide for a private right of action to enforce disparate impact claims, greatly diminishing the risk that a provider of healthcare will face a claim of unintentional discrimination. Nonetheless, the UnitedHealth algorithm was designed to determine optimal medical treatment according to an individual's level of illness. Thus, one can presume that "level of illness" is a revealed business necessity target. Here, too, the IAT can provide healthcare providers such as UnitedHealth with a means to test the proxy variables utilized in predicting their target of interest.

C. *Self-Determining Business Necessity*

Regardless of whether an algorithm is based on complex machine-learned insights or on conventional physical exams, the IAT can serve as an important check for consistency with the principles undergirding U.S. antidiscrimination law across a number of decision-making domains. This tool is not simply a utility for courts to evaluate claims of discrimination, but a tool for regulators and self-regulating firms seeking to detect and avoid discrimination in the first place. Before closing, however, we emphasize two considerations. First, the fact that a proxy input variable is predictive of a business necessity target is not sufficient to rule out the possibility that it systematically penalizes members of a protected group who are actually qualified in the target. This is the principle behind the IAT. Second, although we have argued above that often businesses self-regulate themselves to determine business necessity targets (e.g., *risk* for insurers, *illness intensity* for healthcare providers), businesses must be ever vigilant that a purported target is a legitimate one to use. This is especially the case when working in a domain where courts have defined what can (and cannot) constitute a business necessity target.

A case in point comes from the credit markets, whereby lenders may have incentives to deploy predictive algorithms to estimate demand elasticities across different borrowers to engage in price discrimination. Price discrimination is made possible by the fact that certain borrowers are more prone to accept higher priced loans rather than engage in price shopping. These borrowers may not shop around for a host of reasons: They might live in financial desert locations of low competition, lack the knowledge to shop

¹⁸⁴ 532 U.S. 275 (2001).

for the best rate, need to transact in a hurry, have a historical discomfort with financial institutions due to prior discrimination, and/or have a history of being rejected for loans in the past. Empirical studies document that loan officers and mortgage brokers are aware of variation in borrowers' interest rate sensitivity and engage in price discrimination.¹⁸⁵

A loan applicant's "price sensitivity" or "willingness to shop" may therefore be an additional unobserved characteristic that is of interest to a lender. Said another way, a lender's profit margin depends on both creditworthiness (the court-determined legitimate business necessity from Table 2) and shopping profiles. A lender might therefore design an algorithm that seeks to maximize profits by uncovering credit risk and shopping profiles. Furthermore, the lender (if lending were not in a formally-regulated domain) would argue that profits are legitimate business necessity. Yet, as noted in Table 2, lending is a domain where courts have expressly held that if a lending practice creates a disparate impact, "the defendant-lender must demonstrate that any policy, procedure, or practice has a manifest relationship to the creditworthiness of the applicant."¹⁸⁶ That is, while differences in creditworthiness can justify disparate outcomes in lending, differences in shopping behavior cannot.

The concern of algorithmic profiling for shopping behavior is of general concern because empirical evidence, again in lending, finds that profiling on lack-of-shopping almost certainly leads to higher loan prices for minority borrowers. For instance, Susan Woodward and Robert Hall¹⁸⁷ as well as Mark Cohen¹⁸⁸ find that adverse pricing for minority borrowers has generally been the rule when it comes to lenders engaging in price discrimination. In separate work,¹⁸⁹ we likewise find empirical evidence that, even after controlling for borrower credit risk, "FinTech" lenders charge minority homeowners higher interest rates. We interpret these pieces of evidence as consistent with loan originators using a form of algorithmic price discrimination. Were these algorithms subject to an internal or external "accountability audit," it is likely that the proxy variables used would fail the IAT because, no matter how well the algorithm performed in detecting the profitability of a loan, the target for the test would, by law, be creditworthiness—not an outcome that included price sensitivity. In this

¹⁸⁵ See, e.g., Susan E. Woodward, U.S. Dep't of Hous. & Urban Dev., *A Study of Closing Costs for FHA Mortgages* xi (2008), http://www.huduser.org/Publications/pdf/FHA_closing_cost.pdf ("In neighborhoods where borrowers may not be so familiar with prevailing competitive terms, or may be willing to accept worse terms to avoid another application, lenders make higher-priced offers....")

¹⁸⁶ A.B. & S. Auto Service, Inc., 962 F. Supp. at 1056.

¹⁸⁷ Susan Woodward and Robert E. Hall, *Consumer Confusion in the Mortgage Market: Evidence of Less than a Perfectly Transparent and Competitive Market*, 100 AMER. ECON. REV. 511 (2010).

¹⁸⁸ Mark Cohen, *Imperfect Competition in Auto Lending: Subjective Markup, Racial Disparity, and Class Action Litigation*, 8 REV. LAW ECON. 21 (2012)

¹⁸⁹ Bartlett, et al., *supra* note 38.

fashion, simply asking what target variable an algorithm seeks to detect can illuminate illegitimate algorithmic discrimination.

Finally, we want to end this applications section on a positive note. In many discussions with lenders, it has become evident that, at least in the finance realm, firms want to be able to validate what they are doing or what they intend to do before they invest and commit to a predictive algorithm. As we have demonstrated throughout this Article, the standard set by an IAT-accepted environment can provide the valuable consequence of validating the use of proxy variables when their use causes no disparities except through their role in picking up business necessity leveling.

V. CONCLUSION

The era of Big Data places the antidiscrimination mandate at the heart of the Civil Rights Acts of 1964 and 1968 at a critical cross-roads. By relying on data-driven, statistical models, machine learning provides a promising alternative to the type of subjective, face-to-face decision-making that has traditionally been fraught with the risk of bias or outright animus against members of protected groups. Yet left unchecked, algorithmic decision-making can also undermine a central goal of U.S. antidiscrimination law. As we have shown throughout this Article, any decision-making rule that simply maximizes predictive accuracy can result in members of historically marginalized groups being systematically excluded from opportunities for which they are qualified to participate.

Ensuring that algorithmic decision-making promotes rather than inhibits equality thus demands a workable antidiscrimination framework. To date, however, prevailing approaches to this issue have focused on solutions that fail to grapple with the unique challenge of regulating statistical discrimination. Prominent legal approaches (such as reflected in HUD's recent proposed rule-making) have frequently prioritized predictive accuracy despite the fact that such an approach ignores the central risk posed by statistical discrimination demonstrated in our simulation. Conversely, interventions emanating from the field of computer science have largely focused on outcome-based interventions that could themselves lead to claims of intentional discrimination.

Because we derive our input accountability test from caselaw addressing statistical discrimination—in particular, the burden-shifting framework—the IAT advances a vision of algorithmic accountability that is consistent with the careful balance courts have struck in considering the decision-making benefits of statistical discrimination while seeking to minimize their discriminatory risks. By enhancing the predictive accuracy of decision-making, statistical discrimination can greatly enhance the ability of an employer, lender or other decision-maker to identify those individuals who possess a legitimate target characteristic of interest. However, cases such as

Griggs and *Dothard* underscore the danger of simply focusing on predictive accuracy because a proxy that predicts a target variable can nonetheless result in systematically penalizing members of a protected group who are qualified in the target characteristic. That such discriminatory proxies have been consistently declared to be off limits underscores the conclusion that predictive accuracy alone is an insufficient criterion for evaluating statistical discrimination under U.S. antidiscrimination law.

At the same time, our approach is also consistent with the focus in *Griggs* and *Dothard* that differences in a legitimate target can justify disparities that differ across members of protected and unprotected groups. As we show, so long as a proxy used to predict a legitimate target variable is unbiased with respect to a protected group, it will pass the IAT, even if it results in disparate outcomes. The IAT can therefore provide greater transparency into whether disparate outcomes are the result of a biased model or more systemic disparities in the underlying target variable of interest, such as credit risk. In so doing, it can provide vital information about whether the proper way to address observed disparities from an algorithmic model is through de-biasing the model or through addressing disparities in the underlying target variable of interest, such as through targeted subsidies or other transfers. More generally, because the goal of the IAT is to avoid penalizing members of a protected group who are otherwise qualified in a target characteristic of interest, our approach will also be immune to the concern informing cases such as *Ricci v. DeStefano* that our test is biased against qualified individuals.

Finally, our approach provides clear “rules of the road” for how to exploit the power of algorithmic decision-making while also adhering to the antidiscrimination principles at the heart of the Civil Rights Acts of 1964 and 1968. In particular, the IAT offers data scientists a simple test to use in evaluating the risk that an algorithm is producing biased outcomes, mitigating a key source of the regulatory uncertainty surrounding the growing use of algorithmic decision-making. Additionally, our exploration of the early caselaw considering statistical discrimination also reveals that these rules of the road encompass more general concepts to guide both data scientists and regulators when evaluating algorithmic discrimination. These include the notion that, fundamentally, algorithmic decision-making is an effort to assess an unobservable attribute, such as productivity, criminality, longevity, or creditworthiness, through the use of one or more proxy variables. Consequently, evaluating an algorithm must begin with transparency about this target characteristic. And they likewise include the fact that correlation between the unobservable characteristic and the proxy is not, by itself, sufficient to justify the use of the proxy under antidiscrimination principles.

APPENDIX

DE-BIASING PROXY VARIABLES VERSUS DE-BIASING PREDICTIVE MODELS

In this Appendix, we conduct a simulation exercise to illustrate how attempting to de-bias a proxy variable used in a predictive algorithm may do little to de-bias the ultimate predictions. The example we use assumes that a college admissions director wishes to use applicants' standardized test scores (STS) to predict college success (the criterion for allowing an application to continue to the next stage of evaluation.) For this purpose, we assume that a student's performance on the STS is a function of just two factors: *aptitude* and *family wealth*. In our simulation, wealth contributes to test performance because children of wealthier households purchase expensive test preparation classes. To keep the simulation tractable, we assume that wealth does not affect college performance; its only effect is on a student's STS.

Our simulation involves 1,000 college graduates where the admissions director has data on each student's STS at the time of application, student race, and the student's ultimate college performance (e.g., a weighted grade point average or other measure of performance). We divide the race of students, X_i^R , equally so that 500 students are Non-White ($X_i^R = 0$) and 500 are White ($X_i^R = 1$). We assume that wealth and aptitude are distributed as follows:

$$X_i^{Wealth} \sim \begin{cases} N(0,1) & \text{if } X_i^R = 0 \\ N(5,1) & \text{otherwise} \end{cases}$$

$$X_i^{Aptitude} \sim N(0,1)$$

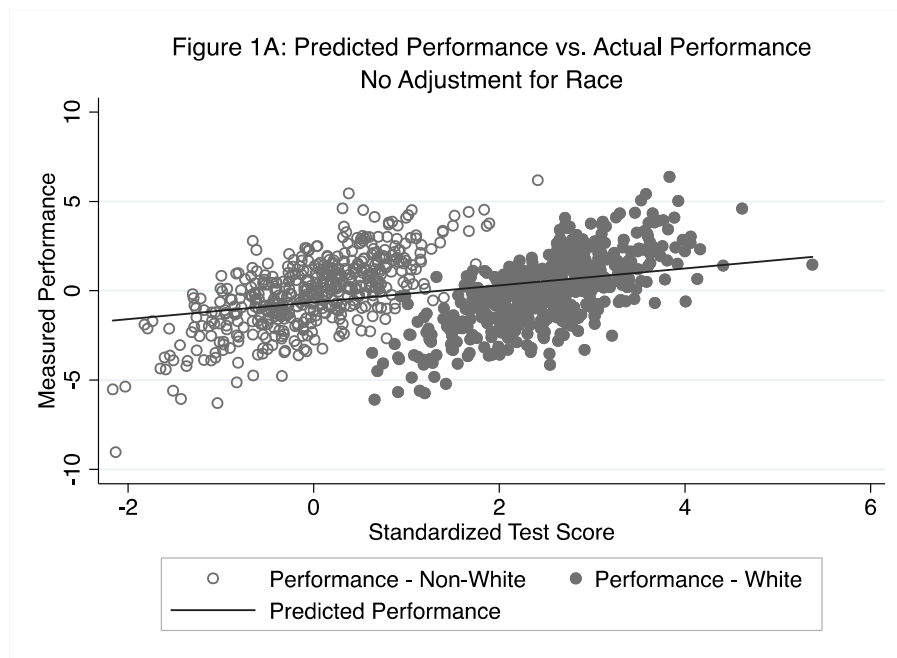
Note that under these distributional assumptions, there is very little common support in wealth across race categories. As noted by Kristen Altenburger and Daniel Ho, it is in these settings where the effort to de-bias proxy variables can produce the largest estimation errors.¹⁹⁰ As noted, a student's STS (X_i^{STS}) is a function of X_i^{Wealth} and $X_i^{Aptitude}$, with each variable given equal weight:

$$X_i^{STS} = 0.5(X_i^{Wealth}) + 0.5(X_i^{Aptitude})$$

Finally, we simulate college performance ($Performance_i$) to be entirely determined by aptitude multiplied by a scalar (which we assume here to be 2).

¹⁹⁰ See Altenburger & Ho, *supra* note 146, at 111. These settings arise "where sharp preexisting demographic differences may exist across groups." *Id.*

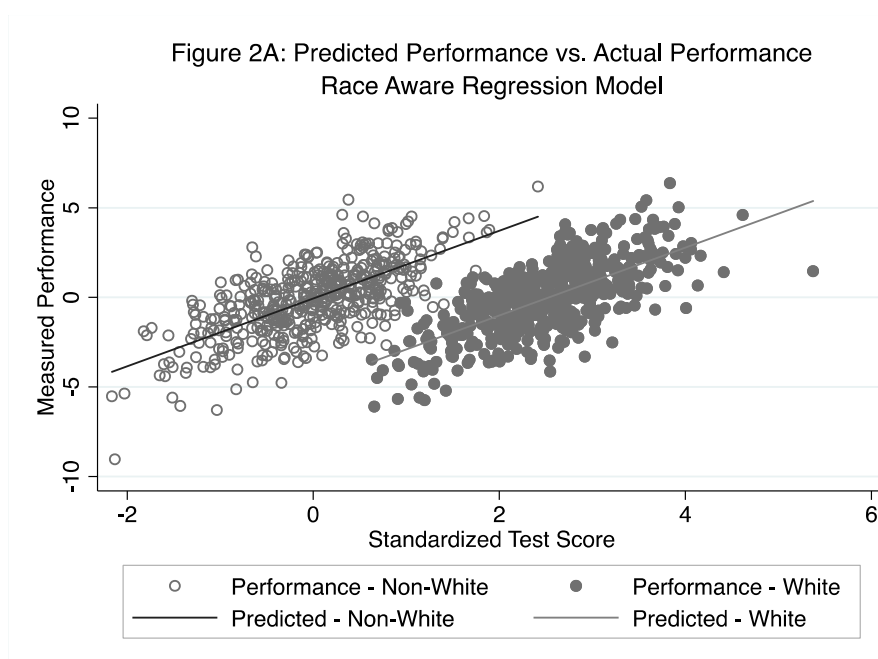
Aptitude is unobservable to the admissions director, inducing her to estimate whether she can use STS to predict college performance. In Figure 1A, we plot separately for White and Non-White graduates the relationship between college performance and STS based on data simulated using the foregoing assumptions. We also include a line that provides the predicted college performance from a simple regression of college performance on STS. As shown in the figure, White graduates had much higher STS scores on average, as would be expected from their much higher family wealth.



The director of admissions would like to admit students that are likely to have a positive measure of college performance (i.e., $Performance > 0$). She therefore runs a simple regression of STS on $Performance$, which produces a regression coefficient ($\hat{\beta}^{STS}$) of 0.47. This estimate indicates that a one-point change in STS is associated with a 0.47 change in $Performance$. Using this regression estimate, the director generates the fitted line shown in Figure 1A, which provides a predicted measure of $Performance$ based solely on STS . The fitted line predicts that $Performance$ is zero at roughly 1.3, suggesting that using a minimum STS of 1.3 would admit students with an expected college performance of at least 0. However, had the admissions director applied this cut-off to these individuals, the bias in STS would result in significant bias against Non-White students owing to their lack of access to test preparation classes:

	Non-White	White
# of Qualified Candidates Predicted by Test Score	13	465

Now assume that the admissions director seeks to control for the greater wealth (and therefore, the greater test preparation bias) among White student applicants. Using the same data, the director expressly adds X_i^R as a control variable in the regression of *STS* on *Performance*. Doing so allows the director to predict *Performance* as a function of both *STS* and *Race*. The results are presented in Figure 2A.



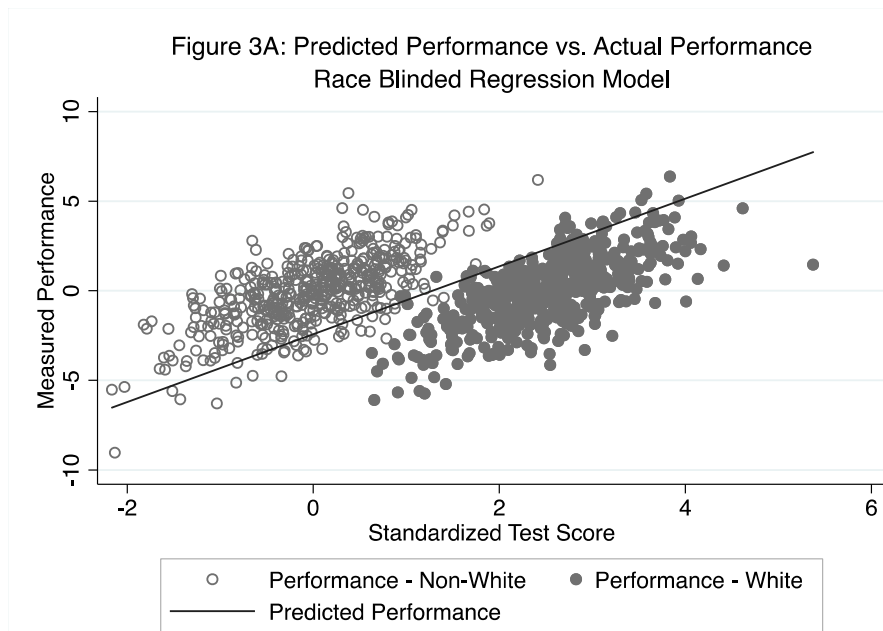
This procedure corrects for the racial bias that arises from using only *STS* to predict *Performance*. This can be seen by the two fitted regression lines, which do a much better job of predicting measured performance across the two racial groups than in Figure 1A. The reason stems from the fact that this regression specification estimates a different y-intercept for each racial category in estimating the relationship between *STS* and *Performance*. Specifically, the regression yields a y-intercept for X_i^R of -4.72, which indicates that in using *STS* to predict *Performance*, it is necessary to deduct 4.72 from the expected performance of White students. (Recall that the difference in average wealth across White and Non-White students is 5.0, so this adjustment eliminates the bias that Wealth creates when using *STS* as a measure of aptitude). With that adjustment, the regression coefficient for *STS* increases from 0.47 to 1.89 because the regression has effectively removed

the confounding effect of wealth on *STS* so that it more cleanly reflects aptitude. As above, the admissions director evaluates each fitted line and determines that the fitted line for Non-White students predicts that *Performance* is zero where *STS* is also zero and that the fitted line for White students predicts that *Performance* is zero at 2.53. Applying a minimum test cut-off of 0 for Non-White students and 2.53 for White students would result in the following students being deemed qualified:

	Non-White	White
# of Qualified Candidates Predicted by Test Score	250	248

This procedure solves the racial bias created by using only *STS* to estimate *Performance*, but it is clearly problematic insofar that it requires a different minimum cut-off for White and Non-White students. This is disparate treatment. To avoid this problem the director therefore turns to the approach advanced by Devin Pope and Justin Sydnor as well as by Crystal Yang and Will Dobbie.¹⁹¹ This procedure involves using the regression estimates generated for Figure 2A but treating all students as if they had the average value of race, or in this example, a race of 0.5. Making this adjustment means that every student receives a deduction of -2.36 (i.e., 0.5×-4.72) after multiplying their exam score by the slope coefficient for *STS* of 1.89, which remains purged of the confounding influence of Wealth. This permits the director to estimate a single fitted regression line as shown in Figure 3A:

¹⁹¹ See *supra* note 71.



The fitted line predicts that *Performance* is zero at approximately 1.28, which the director uses as the minimum cut-off. Had the director applied this cut-off to this group of individuals, the following results would have occurred:

	Non-White	White
# of Qualified Candidates Predicted by Test Score	15	468

In effect, the results are largely identical to those obtained by using only *STS* to predict performance. The reason stems from the lack of common support in wealth across White and Non-White students, resulting in the need for a significant negative adjustment to every White student when estimating performance from *STS*. Applying half of this negative adjustment to *every* student thus works against the de-biasing of the slope coefficient for *STS*. In short, the slope coefficient for *STS* in Figure 3A is unbiased with respect to Non-White students, but the predictive model is not. This problem was significant in this example because there was so little common support in wealth across White and Non-White students—a problem that will exist whenever there are significant demographic differences across protected and unprotected groups.