

DO FIELD EXPERIMENTS ON LABOR AND HOUSING MARKETS OVERSTATE DISCRIMINATION? A RE-EXAMINATION OF THE EVIDENCE

DAVID NEUMARK AND JUDITH RICH*

Since 2000, more than 80 field experiments across 23 countries consider the traditional dimensions of discrimination in labor and housing markets—such as discrimination based on race. These studies nearly always find evidence of discrimination against minorities. The estimates of discrimination in these studies can be biased, however, if there is differential variation in the unobservable determinants of productivity or in the quality of majority and minority groups. It is possible that this experimental literature as a whole overstates the evidence of discrimination. The authors re-assess the evidence from the 10 existing studies of discrimination that have sufficient information to correct for this bias. For the housing market studies, the estimated effect of discrimination is robust to this correction. For the labor market studies, by contrast, the evidence is less robust, as just over half of the estimates of discrimination fall to near zero, become statistically insignificant, or change sign.

Field experiments—specifically, audit or correspondence studies—have been used extensively to test for discrimination in markets. In audit studies of labor market discrimination, fake job candidates (“testers”) of different races, ethnicities, and so forth, who are sometimes actors, are sent to interview for jobs (or in some early studies, apply by telephone). The candidates have similar résumés and are often trained to act, speak, and dress similarly. Correspondence studies use fictitious job applicants who exist on

*DAVID NEUMARK is the Chancellor’s Professor of Economics at the University of California, Irvine, a Research Associate at the National Bureau of Economic Research (NBER), and a Research Fellow at the Institute of Labor Economics (IZA). JUDITH RICH is a Reader in Economics at the University of Portsmouth and a Research Fellow at IZA.

We thank the following authors of studies who generously provided their raw data: Ali Ahmed, Lina Andersson, and Mats Hammarstedt; Stijn Baert, Bart Cockx, Niels Gheyle, and Cora Vandamme; Marianne Bertrand and Sendhil Mullainathan; Mariano Bosch, M. Angeles Carnero, and Lidia Farré; Magnus Carlsson and Stefan Eriksson; Dan-Olof Rooth (also with Magnus Carlsson); Nick Drydakis; Michael Ewens, Bryan Tomlin, and Liang Choon Wang; Hwok-Aun Lee and Muhammad Abdul Khalid; and Phil Oreopoulos. We thank Nick Drydakis and Philip Oreopoulos, and three anonymous referees for helpful comments. For information regarding the data and/or computer programs utilized for this study, please address correspondence to the authors at dneumark@uci.edu or judy.rich@port.ac.uk.

KEYWORDS: discrimination, field experiments, bias, minorities, hiring

paper only (or now, electronically) and who differ systematically only on group membership. The response captured in correspondence studies is a “call-back” for an interview or a closely related positive response. By contrast, the final outcome in audit studies is actual job offers. Differences in outcomes between groups are likely attributable to discrimination, although, naturally, some subtle issues of interpretation occur—including that such differences can be attributable to either taste discrimination or statistical discrimination.

Audit and correspondence (AC) studies have also been used to study discrimination in housing markets. In audit studies, the testers of different races, ethnicities, and so forth are sent to inquire about properties for rent or sale. In correspondence studies, the fictitious inquiry is submitted electronically, applying online to advertised properties for rent or sale.

The large literature using AC studies to test for discrimination in labor markets and housing markets leads to remarkably consistent findings. Nearly every study focusing on race or ethnicity finds evidence of race or ethnic discrimination in the labor market or in the housing market, and the conclusions of the smaller number of studies of sexual orientation discrimination are equally consistent.

The question we ask in this article is whether this near-uniform evidence of discrimination from field experiments is an accurate reflection of discriminatory behavior, supporting a conclusion that discrimination really is this consistent and pervasive. The question might seem misplaced, as AC studies are regarded as providing the most compelling evidence on discrimination. A particularly challenging critique of AC studies (the “Heckman-Siegelman critique”), however, claims that the resulting estimate of discrimination can be biased in either direction; or equivalently, discrimination can be unidentified. This problem arises when the variances of the unobservables differ across the groups studied. Moreover, such a difference in variances—and the bias it creates—cannot be ruled out or easily controlled in AC studies, and most of the past literature using AC studies has simply ignored the problem.

A method to correct AC studies for bias from differences in the variance of unobservables requires more and different kinds of data than AC studies typically collect. We have identified 10 studies of discrimination against minorities in labor and housing markets that do include the requisite data. We re-examine the data from these studies to test whether this evidence is robust to confronting the data with the Heckman-Siegelman critique. Specifically, implementing the correction for bias from differences in the variances of unobservables across groups, do these studies still uniformly point to discrimination?

To summarize the results briefly, for the housing market studies the estimated effects of discrimination are robust to this correction. For the labor market studies, the evidence is less robust; in about half of the cases covered in these studies, the estimated effect of discrimination either falls to near zero or becomes statistically insignificant, and in some the sign changes.

The results for the labor market, in particular, suggest that researchers should build into future AC studies the data and experimental design needed to address the Heckman-Siegelman critique, and that further work on different ways to eliminate bias from AC studies' estimates of discrimination is warranted. More substantively, our re-examination of the evidence suggests that the overall body of experimental evidence on labor market discrimination provides a less clear signal of discrimination than one would draw from the results reported in the existing studies.

Key Background Literature

AC studies are widely regarded as providing more rigorous evidence on discrimination than can be obtained from non-experimental evidence in which group membership may be correlated with unobservables.¹ Heckman and Siegelman (1993) and Heckman (1998), however, showed that in the standard implementation, estimates of discrimination from AC studies can be biased in either direction; or equivalently, discrimination can be unidentified. This problem arises not under some unusual or unlikely theoretical conditions but rather, under an assumption that is at the core of early models of statistical discrimination (Aigner and Cain 1977): the variances of the unobservables differ across the groups studied. This criticism of evidence from AC studies, which we refer to as the “Heckman-Siegelman critique,” holds even under quite ideal conditions (detailed later) in which other potential research design flaws that Heckman and Siegelman discuss are absent.

A statistical method that can lead to unbiased estimates of discrimination using data from AC studies, relying on an identifying assumption, was proposed in Neumark (2012). As explained below, most past AC studies do not have the requisite data, which are applicant or other characteristics aside from the group identifier that shift the probability of call-backs or hires.

The 10 studies, conducted over the past couple of decades, that do include the requisite data,² just like nearly all of the far greater number of AC studies that do *not* have the requisite data, find evidence of discrimination against ethnic or racial minorities, immigrants, or gays and lesbians.³ We have obtained the original data from the authors of these studies, and our goal in this article is to test whether this evidence is robust to confronting the data with the Heckman-Siegelman critique. Specifically, after implementing the correction for bias from differences in the variances of unobservables across groups, do these studies still uniformly point to discrimination?

¹The methods and empirical findings from these studies have been reviewed by Riach and Rich (2002), Pager (2007), Rich (2014), and Neumark (forthcoming). Additionally, similar studies of discrimination have been conducted in consumer markets (e.g., Doleac and Stein 2013).

²The studies are Bertrand and Mullainathan (2004)—the data used in Neumark (2012); Carlsson and Rooth (2007); Ahmed, Andersson, and Hammarstedt (2010); Bosch, Carnero, and Farré (2010); Oreopoulos (2011); Carlsson and Eriksson (2014); Drydakis (2014); Ewens, Tomlin, and Wang (2014); Baert, Cockx, Gheyle, and Vandamme (2015); and Lee and Khalid (2016).

³For the most recent review of a large number of AC studies, see Neumark (forthcoming).

Some very recent AC studies have implemented this bias correction.⁴ Our article revisits past studies that do not address the Heckman-Siegelman critique, to assess whether the near-uniform findings of discrimination from the large body of past research are robust to addressing this critique. We cannot re-examine all such studies; but we do, we believe, re-examine the complete set of such studies that focus on traditional demographic dimensions (such as race, sex, or sexual orientation) of discrimination and have the data required to address this critique.

Field Experiments Covered in This Article

The field experiments re-analyzed in this article are one of three broad types: studies of ethnic/immigrant or race discrimination in labor markets; studies of sexual orientation discrimination in labor markets; and studies of ethnic/immigrant or race discrimination in rental housing markets. Many of the details and results of these studies are discussed in Rich (2014) and Neumark (forthcoming). Here we focus only on what is essential to understand the analysis of bias from differences in unobservables that we implement in this article. Readers interested in more details on these specific studies, and the techniques used more generally, should see our surveys (or, of course, the original papers). We do not go into more detail because our goal here is not to compare or critique other dimensions of these studies, but rather just to consider the robustness of the conclusions to addressing the Heckman-Siegelman critique.⁵

What distinguishes these 10 studies from the others in the literature is that they use applicants distinguished not only by race, ethnicity (including immigrant origin), or sexual orientation but also by different levels of qualifications. In these studies, the additional information was used to ask, in a general way, whether the evidence of discrimination by ethnicity, race, or sexual orientation differed for applicants with different levels of qualifications.⁶ The availability of data with variation in applicant qualifications is

⁴See Carlsson, Fumarco, and Rooth (2013); Baert (2014, 2015, 2016); Nunley, Pugh, Romero, and Seals (2015); and Neumark et al. (2016, forthcoming). Baert and Verhofstadt (2015) also do this, although in relation to criminal background (juvenile delinquency), which is outside the scope of discrimination studies covered in the present article.

⁵There are also field experiments investigating differences in hiring outcomes based on other characteristics, such as criminal background, mental or physical illness, facial attractiveness, veteran status, or socioeconomic background or class. Although these kinds of differences are not the focus of our article (even though some could be interpreted as discrimination), the experimental designs in these articles do not generate the data needed to implement this empirical method, with the exception of Baert and Balcaen (2013), who implemented this method in relation to differential treatment based on military service and found no evidence of bias from differences in the variances of unobservables.

⁶The first study of this type (Jowell and Prescott-Clarke 1970) considered this issue. The study compared job offer outcomes for immigrant versus white British applicants, and they gave half the applications in each group higher qualifications with regard to education. (Additional variation among the immigrants included whether they were English-speaking and whether secondary education occurred in Britain, although this kind of variation that does not apply equally to majority and minority groups is not as useful.) The more recent studies with such data that we re-examine in the present article are those for which we could recover the data from authors.

exactly what is needed to implement the empirical method that addresses the Heckman-Siegelman critique.

Bertrand and Mullainathan (2004), Carlsson and Rooth (2007), Drydakis (2014), Baert et al. (2015), and Lee and Khalid (2016) all used matched pairs (sets) of applicants, with two (or more) applications sent to each job vacancy. Oreopoulos (2011) considered differences for many different ethnic groups (relative to native Canadians), in some cases also signaling immigrant status, and sent multiple résumés for each job vacancy. Across these studies, the authors used either real résumés they had found or résumés they generated randomly. Names used on the résumés signaled race or ethnicity, and education or work experience in a foreign country signaled immigrant status (Oreopoulos 2011). Participation in an organization active on behalf of the gay community or a gay organization signaled sexual orientation.

There have been fewer studies of discrimination in housing markets in the broader literature. Of the ones we re-examine, only Bosch et al. (2010) used matched pairs, whereas the other three (Ahmed et al. 2010; Carlsson and Eriksson 2014; Ewens et al. 2014) sent a single rental enquiry. An accompanying message provided details on the applicant, in which the researchers manipulated the information provided—ethnicity and race, as well as other qualifications or the applicant's job, which indicated ability to pay. In these studies, signaling was done by name, although Bosch et al. (2010) interpreted their results for Moroccan versus Spanish names as measuring discrimination against immigrants.

Other qualifications also varied across the résumés or applications, and this variation in qualifications is essential for implementing the correction for bias from differences in variances of unobservables. We describe the variables used in each study in Tables 2A, 2B, and 3, which report our results from re-analyzing the data from these studies (discussed in detail below). For example, Bertrand and Mullainathan (2004) generally sent four applications to each job. They created two matched pairs of applicants, one pair with low-quality backgrounds and another pair with high-quality backgrounds. The quality of the applicant varied based on labor market experience, career profiles, employment history, and skills such as employment experience gained over the summer or while at school, volunteering, extra computer skills, certification degrees, foreign language skills, honors, or some military experience. Carlsson and Rooth (2007) signaled similar additional information on applicants including different spells of unemployment, work experience over the summer, overqualified or not, personality traits, and cultural and sporting activities listed as hobbies and interests. Oreopoulos (2011) varied the information provided on the extent of foreign education and foreign experience as well as language skills and certification and master's degrees. Drydakis (2014) used an accompanying cover letter to provide more favorable information about applicants in some cases, including a mention of grades, previous job responsibilities and tasks, and personality characteristics associated with work commitment; these

same applicants also included letters of references that more strongly signaled positive work traits such as teamwork and loyalty to the firm. Lee and Khalid (2016) varied factors such as private versus public university, grades, and English proficiency.

In the housing market tests, researchers manipulated the information on the applicant, using an accompanying message, to explore the impact of basic, negative, or positive information, such as habits (smoking, exercise, and nightclub attendance, in Carlsson and Eriksson 2014), variation in smoking and credit rating (in Ewens et al. 2014), and information on positive characteristics such as work history, education, lack of payment complaints, and so on (Ahmed et al. 2010) or stable occupations and contracts (Bosch et al. 2010).

The richness and number of qualifications that researchers chose to vary across the applicants differed quite a bit across these studies. For the labor market studies, these qualifications generally pertained to education, experience, and skills but sometimes extended to attempts to convey something about the applicant's personality or hobbies, the order of the application, and other things. One of the housing studies (Carlsson and Eriksson 2014) tried to provide information on the applicant's lifestyle, which could be relevant to a potential landlord. We do not discuss the different qualifications used in each study in detail, but list them for each study in the tables reporting the statistical analysis (Tables 2A and 2B for the labor market studies, and Table 3 for the housing market studies). Note that we also list other features of the ads that could affect the probability of a call-back, such as characteristics of the job or the apartment. We include these because, as explained in the next section, the statistical method is informed by differences in the coefficients between the two groups studied in *any* of the factors that can affect call-backs.

Findings from the Field Experiments Covered in This Article

Table 1 summarizes the conventional results from the 10 studies we re-examine and provides basic information about them, including the years covered, the groups covered, and the outcomes. The original studies reported results in different ways, varying between chi-square/Fisher exact tests, binomial tests, or tests of the null hypothesis of no difference in the call-back rate between the groups, typically controlling for other aspects of the résumés. Here, we report results on a consistent basis for all studies—marginal effects from probit models using the full set of résumé characteristics included in the data—which we have estimated from data provided by the authors of these studies.⁷

⁷Details on the control variables, the standard errors, and so on are provided in tables discussed below. Not surprisingly, the results in Table 1 closely parallel the conclusions of the original papers—however they report their results—although they are not always identical.

Table 1. Experimental Studies of Discrimination in Labor and Housing Markets: Summary and Key Results

Study (1)	Country (2)	Years (3)	Minority (4)	Outcome (5)	Majority call-back rate (6)	Estimated differential for minority (7)
A. Labor market field experiments						
Baert et al. (2015)	Belgium	2011–12	Turkish	Call-back	0.329	-0.082 (0.034)
				Immediate interview	0.190	-0.056 (0.026)
Carlsson and Rooth (2007)	Sweden	2005–06	Middle Eastern	Call-back	0.269	-0.095 (0.009)
Drydakis (2014)	Cyprus	2010–11	Gay	Call-back	0.554	-0.410 (0.010)
			Lesbian	Call-back	0.523	-0.411 (0.011)
Lee and Khalid (2016)	Malaysia	2011	Malay (vs. Chinese)	Call-back	0.222	-0.152 (0.018)
Oreopoulos (2011)	Canada	2008	Chinese	Call-back	0.142	-0.053 (0.007)
			Indian	Call-back	0.142	-0.056 (0.007)
			Chinese-Canadian	Call-back	0.142	-0.063 (0.008)
			Pakistani	Call-back	0.142	-0.073 (0.009)
			Greek	Call-back	0.142	-0.035 (0.017)
			British	Call-back	0.142	-0.031 (0.011)
Bertrand and Mullainathan (2004)	United States	2001–02	Black-sounding names	Call-back	0.097	-0.030 (0.006)

(continued)

Table 1. Continued

<i>Study</i> (1)	<i>Country</i> (2)	<i>Years</i> (3)	<i>Minority</i> (4)	<i>Outcome</i> (5)	<i>Majority call-back rate</i> (6)	<i>Estimated differential for minority</i> (7)
B. Housing market field experiments						
Ahmed et al. (2010)	Sweden	2008	Arab/Muslim	Positive response	0.514	-0.171 (0.033)
Bosch et al. (2010)	Spain	2009	Moroccan immigrants	Immediate showing	0.254	-0.091 (0.024)
Carlsson and Eriksson (2014)	Sweden	2010–11	Arab	Positive response	0.590	-0.133 (0.014)
Ewens et al. (2014)	United States	2009	Black	Immediate showing	0.541	-0.135 (0.014)
				Positive response	0.387	-0.130 (0.012)
				Immediate showing	0.271	-0.110 (0.011)
				Positive response	0.503	-0.090 (0.019)

Notes: All studies are correspondence studies. Column (7) reports marginal effect from probit models, our estimates, from following tables, with standard errors in parentheses. In the Orcopoulos study, “Chinese-Canadian” means there was an English first name.

As reported in Table 1, the six labor market experiments covered in panel A all found statistically significant evidence of discrimination against ethnic minorities, blacks, or gays and lesbians. The estimated differentials by racial and ethnic groups were in the same range—an approximately 0.03 to 0.15 lower probability of a call-back. These values were on somewhat different baseline rates of call-backs, but the call-back rates did not vary that much across these studies.⁸ The two estimates from Drydakis (2014), for discrimination against gays and lesbians in Cyprus, were much larger (although the baseline call-back rates were much higher, too).

The four housing market studies similarly found consistent evidence of discrimination against minorities. The range of estimates was fairly tight (a 0.09 to 0.17 lower call-back rate). Thus, every one of these studies pointed to evidence of discrimination against the minority group.

Conclusions from these studies strongly echo the broader literature, in which nearly every study found evidence of discrimination in the labor market or the housing market on the basis of race or ethnicity (Rich 2014; Neumark, forthcoming; Zschirnt and Ruedin 2016; Quillian, Pager, Hexel, and Midtboen 2017), as do the smaller number of studies of discrimination based on sexual orientation (Neumark, forthcoming). The question our article addresses is whether this near-uniform evidence of discrimination from field experiments is an accurate reflection of discriminatory behavior, supporting a conclusion that discrimination really is this consistent and pervasive, or whether the evidence in at least some of these studies might reflect biases stemming from differences in the variance of unobservables across groups—the problem highlighted by the Heckman-Siegelman critique.

Some of the studies also include female and male applicants, or more broadly test for discrimination along multiple dimensions, including sex and age (Carlsson and Eriksson 2014). We do not focus, in this article, on evidence of discrimination based on sex or age. The broader literature focuses far more on race and ethnicity (and more recently on sexual orientation), and, as we have noted, delivers a near-uniform finding of discrimination against minorities. The evidence of sex discrimination is less robust, and tends to point less to discrimination against women and more to the importance of sex norms for jobs in whether male or female applicants received more call-backs (Neumark, forthcoming). And recent evidence from a large-scale correspondence study of age discrimination yields ambiguous results for men, but not for women (Neumark et al., forthcoming).

Addressing the Heckman-Siegelman Critique

Quite a few critiques of AC studies can be found aside from the one we focus on here. Most of them are described in Heckman and Siegelman

⁸One might wonder about apparent evidence of discrimination against British immigrants in Canada; indeed, we will see in implementing the correction for the Heckman-Siegelman critique that this evidence appears to be spurious.

(1993) and are discussed further in Neumark (2012) in the context of the framework laid out in this section. Some of the more important critiques—such as the possibility of “experimenter effects,” and the small differences between applicants that can matter a lot when applicants are matched on so many characteristics—can be addressed by using correspondence studies instead of audit studies, and indeed most recent research uses the correspondence study technique. The Heckman-Siegelman critique is of particular importance because it applies equally well to correspondence studies, even under otherwise ideal conditions such as no *mean* differences in unobservables between groups, but only differences in the *variances* of unobservables. This critique is salient because nothing in the research design rules out differences in the variances of unobservables, and indeed, as noted earlier, these differences are foundational in models of statistical discrimination. We first lay out a basic framework for the analysis of data from an audit or correspondence study, and then explain the bias and the correction.⁹

Non-experimental regression-based approaches testing for and measuring discrimination use data on the groups in question in a population, introducing regression controls to try to remove the influence of group differences in the population that can affect outcomes (Altonji and Blank 1999). Correspondence (and audit) studies, by contrast, create an artificial pool of labor market participants among whom there are supposed to be no average differences by group. This is clearly a potentially powerful strategy, because if we have, for example, a sample of blacks and whites who are identical *on average*, because race is randomly assigned to a subset of similar résumés, then in a regression of the form

$$(1) \quad Y = \alpha + \beta B + \varepsilon,$$

where Y is the outcome and B is a dummy variable for blacks, ε is uncorrelated with B , so that the OLS estimate $\hat{\beta}$ (or simply the mean difference in Y) provides an estimate of the effect of race discrimination on Y .¹⁰

Of course, most of the earlier regression studies focus on wages, whereas AC studies focus on hiring. If an employer is free to pay a lower wage to blacks, for example, then in the context of the Becker (1971) employer discrimination model, why discriminate in hiring? One common interpretation is that there is an equal wage constraint, perhaps because of a minimum wage, or because anti-discrimination laws are more effective at rooting out wage discrimination than hiring discrimination. Alternatively, in the simple model, employers with stronger discriminatory tastes than the marginal employer will discriminate in hiring. As we make clear below, however, this framework does not only detect taste discrimination à la Becker (1971).

⁹This section draws heavily on Neumark (2012), yet avoids many details that a reader can find in that article.

¹⁰For simplicity, we couch the discussion here solely in terms of blacks and whites.

To provide a more formal framework, suppose that productivity depends on two individual characteristics (standing in for a larger set of relevant characteristics), $X = (X^I, X^{II})$, so that productivity is $P(X)$. X^I is what the firm observes, and X^{II} is unobserved by firms. It is simplest, for now, to think of Y as continuous, such as the wage offered, although in AC studies we should think of it as latent productivity leading to a decision to hire/call-back or not.

Define discrimination as

$$(2) \quad Y(P(X'), B = 1) \neq Y(P(X'), B = 0).$$

Assume that $P(.,.)$ is additive, so

$$(3) \quad P(X') = \beta_I X^I + X^{II},$$

where the coefficient of X^{II} is normalized to 1 as it is unobservable, and

$$(4) \quad Y(P(X'), B) = P + \gamma B.$$

Discrimination against blacks implies that $\gamma < 0$, so that blacks are paid less than or are perceived as less productive than whites, when actually they are equally productive.

In correspondence studies, researchers create résumés that standardize the productivity of applicants at some level. Denote expected productivity for blacks and whites, based on what the firm observes, as P_B^* and P_W^* . Y is observed for each tester, so each test—the outcome of applications to a firm by one black and one white tester/applicant—yields an observation

$$(5) \quad Y(P_B^*, B = 1) - Y(P_W^*, B = 0) = P_B^* + \gamma - P_W^*.$$

Given that the correspondence study design sets $P_B^* = P_W^*$, we should be able to estimate γ easily from these data, by simply running a regression of Y on the dummy variable B and a constant. (Some potential complications are discussed in Neumark 2012.)

A correspondence study can preclude systematic differences between groups in observables and experimenter effects. But there can still be assumed differences in means between groups despite the groups using matched résumés. In Equation (5) above, $P_B^* = E(\beta_I X_B^I + X_B^{II} | X_B^I, B = 1)$, and similarly for P_W^* . Assuming randomization, and with $X_B^I = X_W^I = X^I$, the right side of Equation (5) reduces to $\gamma + E(X_B^{II} | X^I, B = 1) - E(X_W^{II} | X^I, B = 0)$, implying that we only identify γ if $E(X_B^{II} | X^I, B = 1) = E(X_W^{II} | X^I, B = 0)$. Employers may have different expectations about the mean of X^{II} for blacks and whites, conditional on what they observe, which a labor economist would label statistical discrimination. Although economists are interested in distinguishing between statistical and taste discrimination, both are illegal under US law and both also appear to be illegal under European Union

law.¹¹ Moreover, it is challenging to distinguish between the two models. Thus, we put this issue aside, and interpret the discrimination estimates from the studies considered in this article as the sums of taste and statistical discrimination.¹²

That is not to suggest that researchers using AC methods have not tried to distinguish between taste and statistical discrimination. The idea exploited in most studies is that when the applications include a richer set of applicant characteristics, it is less likely that statistical discrimination plays much of a role in group differences in outcomes (e.g., Ewens et al. 2014). Effectively, one tries to eliminate the term $E(X_B^H|X^I, B=1) - E(X_W^H|X^I, B=0)$ from the estimated difference in hiring rates to see how much of the overall difference in hiring rates is accounted for by this difference in expectations, which corresponds to statistical discrimination.¹³

Oreopoulos (2011) and Ewens et al. (2014) presented perhaps the most thorough attempts at discerning between these hypotheses about discrimination in AC studies. Oreopoulos used the approach of adding information (e.g., on country of education, to signal English language skills) to see whether estimated hiring gaps fall, as well as examining differences in hiring gaps for occupations across which the importance of statistical discrimination likely varies. In many cases, he does not find evidence consistent with statistical discrimination, despite evidence from a survey of participating employers that they used name, or country of education or experience, as a signal of potential language problems.

Ewens et al. (2014) specifically allowed for the means and variances of unobservables to differ across groups (as in Aigner and Cain 1977) and examined whether the differential treatment by race is more consistent with statistical discrimination (both first- and second-moment) or taste-based discrimination. Although they did not correct for differences in variances of unobservables, they demonstrated that group differences in outcomes may decrease when more information is provided, and they argued that the evidence is consistent with statistical discrimination. In particular, they

¹¹As discussed in Neumark (forthcoming), the U.S. Code of Federal Regulations (29, § 1604.2) defines illegal discrimination as “the refusal to hire an individual because of the preferences of coworkers, the employer, clients or customers.” But it also states, “The principle of nondiscrimination requires that individuals be considered on the basis of individual capacities and not on the basis of any characteristics generally attributed to the group.” There is not as explicit a prohibition of statistical discrimination in the European Union (EU). Article 2 of the EU’s Directive 2000/43/EC prohibits both “direct” and “indirect” discrimination, but these appear to line up, respectively, with disparate treatment and disparate impact in the US context (accessed at <http://eur-lex.europa.eu/legal-content/en/ALL/?uri=CELEX%3A32000L0043>, December 2, 2015). Other material suggests that statistical discrimination is covered by direct discrimination (OECD 2013: 195).

¹²Indeed, it seems that we could also include implicit discrimination (e.g., Bertrand, Chugh, and Mullainathan 2005). Implicit discrimination posits a different reason for undervaluing the productivity of a group of workers, which can lead to different policy levers to combat it. But if it arises when employers evaluate applicants in AC studies, the empirical implication for the framework developed here would likely be the same as the implication of taste discrimination.

¹³Neumark (forthcoming) provided many examples and also some criticisms of this approach.

demonstrated that the differences in outcomes across groups vary with the differences in racial composition across neighborhoods in a way that is consistent with the hypothesized differences in variances of unobservables across groups.

One could presumably use the method described below for résumés with varying amounts of information to recover unbiased estimates under different information treatments and hence try to gauge the relative importance of taste and statistical discrimination. However, this issue is not the focus of our analysis here. Instead, our focus is re-examining the 10 studies identified earlier and investigating whether the uniform evidence of discrimination from these studies persists once account is taken of the Heckman-Siegelman critique.

The issue raised by the Heckman-Siegelman critique arises from the potential for differences across groups in the variances of the unobservables, which is equally problematic even in the ideal condition of no assumed mean difference. To see how the difference in variances can drive differences in the results of the analysis of data from an AC study, it is most natural to think of Equation (1) as a latent variable model for productivity, with applicants having to exceed some productivity threshold with sufficiently high probability (where α in Equation (1) can also include observables that vary across individuals that affect productivity, which we have denoted X^I).

To isolate the problem, consider the best-case scenario where $E(X_B^{II}|X^I, B = 1) = E(X_W^{II}|X^I, B = 0)$, that is, no statistical discrimination regarding levels. But the standard deviations of the unobservables, denoted σ_B^{II} and σ_W^{II} , need not be equal.¹⁴

Assume the applicant is called back (hired) if there is a sufficiently high probability that their productivity exceeds a given threshold. In this case, the inequality $\sigma_B^{II} \neq \sigma_W^{II}$ combined with the design of AC studies results in a biased estimate of discrimination; worse, we cannot necessarily even sign the bias.

To see the intuition, recall that the key feature of the usual design of AC studies is using similar résumés on the applicants in different groups. This approach requires choosing a particular level of the quality of the résumés. Suppose, for example, that the research design standardizes X^I at a low level, denoted X^{I*} . Employers care about how likely it is that the sum $\beta_1 X^I + X^{II}$ exceeds some threshold. Given the low value X^{I*} , this is more likely for a group with a high variance of X^{II} . Thus, even in the case of no discrimination ($\gamma = 0$), the employer will favor the high-variance group. Conversely, if standardization is at a high level of X^{I*} , the employer will favor the low-variance group. Because researchers do not have information on the population of real applicants to the jobs studied, there is no definitive way to know whether X^{I*} is high or low relative to the actual

¹⁴As in Neumark (2012), we assume homoskedasticity within groups, and thus suppress conditioning on X_B^I and X_W^I .

distribution, and hence no way to sign the bias. As discussed in more detail below, note that the variances of unobservables affect which group receives more call-backs only because of the research design standardizing the résumés at a particular level (when the level of standardization is not at the central tendency of the distribution).

The technique developed in Neumark (2012) to correct for the bias from differences in the variances of unobservable characteristics relies on the experimental study having extra information that explores the impact of different productivity or quality characteristics (creating applicants who have different levels of qualifications, for example). As long as some of these characteristics have the same effects in the latent variable model for the probability of a call-back—the key identifying assumption—this extra information allows the effect of the difference in variances between the groups' unobserved characteristics on the responses to be isolated from the role of discrimination in evaluating applicants. That is, it allows separate identification of the relative variances in the unobservables and the discrimination coefficient, γ .¹⁵

Correspondence studies rarely include variables that shift the call-back probability, because these studies typically create one “type” of applicant for which random variation occurs only in characteristics that are not intended to affect outcomes. However, the 10 studies discussed earlier in this article have this information—as in Bertrand and Mullainathan (2004), whose data Neumark (2012) used to illustrate this method for correcting for the bias in AC studies. Applying this method to the studies re-examined in this article therefore allows us to determine whether the measures of discrimination from conventional analyses of the data in these studies provided unbiased estimates of discrimination, or instead either overstated or understated discrimination.¹⁶

The intuition behind the solution stems from the fact that a higher variance for one group (say, whites) implies a smaller effect of observed characteristics on the probability that a white applicant meets the standard for hiring. Thus, information from a correspondence study on how variation in observable qualifications is related to call-backs can be informative about the relative variance of the unobservables, and this, in turn, can identify the effect of discrimination. Based on this idea, the identification problem highlighted by the Heckman-Siegelman critique is solved by invoking an identifying assumption—specifically, that the effects of applicant characteristics that affect perceived productivity and hence call-backs are equal across groups—along with the testable requirement that some applicant characteristics

¹⁵To reiterate, for the purposes of simplification, we assume $E(X_B^H|X^L, B = 1) = E(X_W^H|X^L, B = 0)$. Without this assumption, references to γ in the remainder of this section should be read as references to $\gamma + E(X_B^H|X^L, B = 1) - E(X_W^H|X^L, B = 0)$, that is, the sum of taste and statistical discrimination.

¹⁶For recent code to implement the estimator, we direct readers to the code used in Neumark et al. (2016) on the website of the *American Economic Review* (click on “Data Set” on the webpage accessed at <https://www.aeaweb.org/articles?id=10.1257/aer.p20161.008>).

affect the call-back probability (since if all the effects are zero we cannot learn about $\sigma_B^{II}/\sigma_W^{II}$ from these coefficient estimates).

In a probit specification, for example, we know that we can identify only the coefficients of the latent variable model for productivity relative to the standard deviation of the unobservable. In this case, we effectively have two probit models, one for blacks and one for whites. If we normalize σ_W^{II} to 1, then for a characteristic (Z) that affects the call-back rate, we identify its coefficient (δ_W) relative to σ_W^{II} , or δ_W/σ_W^{II} . However, if we assume that $\delta_W = \delta_B$, then we do not need to impose the normalization that $\sigma_B^{II} = 1$, but instead can identify $\sigma_B^{II}/\sigma_W^{II}$ from the ratio of the coefficients on Z in the probit for whites versus blacks, which in turn allows us to identify γ . The estimation can be done using a heteroskedastic probit model. Finally, when *multiple* productivity-related characteristics shift the call-back probability Z_k ($k = 1, \dots, K$), there is an overidentification test because the ratio of coefficients on each Z , for whites relative to blacks, should equal $\sigma_B^{II}/\sigma_W^{II}$.¹⁷

The heteroskedastic probit model estimates can be decomposed into the estimated differential due to differences in γ , and the estimated differential due to differences in the variance of the unobservables. In generic notation, let the latent variable depend on a vector of variables S and coefficients ψ , and the variance depend on a vector of variables T , which includes S , with coefficients θ . The elements of S are indexed by k . For a standard probit model, we translate coefficient estimates into estimates of the marginal effects of a continuous variable S using

$$(6) \quad \partial P(\text{call-back})/\partial S_k = \psi_k \phi(S\psi)$$

where S_k is the variable of interest with coefficient ψ_k , $\phi(\cdot)$ is the standard normal density, and the standard deviation of the unobservable is normalized to 1. Typically, this partial derivative is evaluated at the means of S . When S_k is a dummy variable such as race, the difference in the cumulative normal distribution functions is often used instead, although the difference is usually trivial.

The marginal effect is more complicated in the case of the heteroskedastic probit model, because if the variance of the unobservable differs by race, then when race “changes” both the variance and the level of the latent variable that determines hiring can shift. As long as we use the continuous version of the partial derivative to compute marginal effects from the heteroskedastic probit model, a unique decomposition exists of the effect of a change in a variable S_k (which also appears in T) into these two components. In particular, denoting the variance of the unobservable $[\exp(T\theta)]^2$,

¹⁷Indeed, the identifying restriction $\delta_W = \delta_B$ only has to hold for subsets of the characteristics that shift the call-back probability, and one can rely only on this subset if the overidentification test for a larger set of résumé characteristics fails (see Neumark 2012).

with the variables in T arranged such that the k^{th} element of T is S_k , then the overall partial derivative of $P(\text{call-back})$ with respect to S_k is

$$(7) \quad \partial P / \partial S_k = \phi(S\psi / \exp(T\theta)) \{\psi_k / \exp(T\theta)\} + \phi(S\psi / \exp(T\theta)) \cdot \{(-S\psi \cdot \theta_k) / \exp(T\theta)\}.$$
¹⁸

The first part of the sum in Equation (7) is the partial derivative with respect to changes in S_k affecting only the level of the latent variable—corresponding to the counterfactual of S_k changing the valuation of the worker without changing the variance of the unobservable. The second part is the partial derivative with respect to changes through the variance of the unobservable. In the analysis below, we report these two separate effects as well as the overall marginal effect, and we calculate standard errors using the delta method.¹⁹

This discussion raises the issue of what we are trying to measure in audit and correspondence studies. Focusing on γ , the structural effect of race, captures the potential discounting by employers of black workers' productivity à la Becker (and possibly statistical discrimination about the mean of X^H). But, as shown, employers could treat blacks and whites differently in hiring because of different variances of the unobservable. If the latter is accepted as a meaningful measure of discrimination, we might not want to eliminate it.

The coefficient γ is the focus of interest for two reasons. First, to the best of our knowledge, differential treatment based on assumptions (true or not) about variances have not been viewed as discriminatory in the legal literature. Second, and probably more important, the taste discrimination (and possibly “first-moment” statistical discrimination) that correspondence studies capture in γ generalizes from the correspondence study to the real economy. By contrast, the difference in treatment based on differences in the variances of unobservables is an artifact of the design of correspondence (or audit) studies—in particular, the standardization of applicants to particular, and similar, values of the observables, relative to the actual distribution of observables among real applicants. If, instead, a study used applicants that replicated the actual distribution of applicants to the employers in the study, there would be no bias—in the setting described here—from different variances of the unobservables (see detailed discussion in Neumark 2012).

That is not to say, however, that there cannot be discrimination based on second moments with, for example, risk averse firms. In that sense, one can potentially interpret the bias correction and decomposition not as separating out real versus spurious discrimination, but rather first-

¹⁸See Cornelißen (2005).

¹⁹Because the formula for the derivative based on a continuous variable yields this unique decomposition, it is used below. One can decompose the partial derivative from the heteroskedastic probit model based on the partial derivative for discrete variables calculated from differences in the cumulative normal distribution functions, but then the decomposition is not unique.

moment versus second-moment discrimination. We could imagine, for example, that risk-averse firms are less likely to call back (or hire) workers with more uncertain productivity, even when on average they are as productive as another group. The potential difficulty with this interpretation, however, is that we do not uniformly find that the minority group that experiences discrimination according to the conventional analysis generally has a higher variance of the unobservable; indeed, in both the labor market studies and the housing market studies we analyze, this is the case in just about half of the estimates. This finding is a further reason for why, in the remainder of the article, we interpret the evidence as isolating discrimination by adjusting for differences in the variances of unobservables.

Results from Re-examination of Field Experiments with Quality Variation across Résumés

Labor Market Field Experiments

We report the results for the re-analysis of the data sets from the labor market field experiments in Tables 2A and 2B. Turning to the first set of labor market studies covered in Table 2A, we first report the estimated discrimination coefficient (γ , in the equations from above) in the first row of the table (panel A). These match the estimates in the last column of Table 1 and have already been summarized.

Panel B turns to the heteroskedastic probit estimates that correct for biases from differences in the variance of unobservables. The “Controls” entry toward the bottom of the table lists the résumé characteristics, including those likely to shift the call-back rate (such as education, skills, and so forth).²⁰ The first row of panel B reports the overall effect from the heteroskedastic probit estimates, which are similar to the probit estimates. The next two rows of the table report the key results from the decomposition of the heteroskedastic probit estimates. The “level” effect (labeled “Marginal effect through level (unbiased)” in the table) is the unbiased estimate, and the “variance” effect reflects the bias from the correspondence study design, arising because of the interaction between the quality of the résumés sent out (relative to the actual distribution) and differences in the variances of unobservables.

Looking at these estimates, for the first study—the Baert et al. (2015) experiment on discrimination against Turkish job applicants relative to

²⁰Some studies include résumé characteristics that are not independent of minority group status. For example, Oreopoulos (2011) indicated, for some of his ethnic groups, that some education or experience occurred in a foreign country. This detail is useful for asking what might explain variation in the amount of discrimination immigrants face, which was the focus of his study. But it does not fit into the narrower question considered in this article of discrimination against the minority group per se. Hence, we only use résumé characteristics that are constructed to be orthogonal to minority group status.

Table 2A. Estimates for Labor Market Discrimination Studies: Full Specifications

Study	Baert et al. (2015), Belgium		Carlsson and Rooth (2007), Sweden		Drylakis (2014), Cyprus		Lee and Khalid (2016), Malaysia	
	Call-back	Immediate interview	Call-back	Middle Eastern, males	Gay	Lesbian	Call-back	Malay
Outcome	Turkish, males		(3)		(4)	(5)	(6)	
Minority group	(1)	(2)						
A. Estimates from basic probit								
Minority, marginal effect	-0.082 (0.034)	-0.056 (0.026)	-0.095 (0.009)		-0.410 (0.010)	-0.411 (0.011)		-0.152 (0.018)
B. Heteroskedastic probit model								
Minority, marginal effect	-0.096 (0.034)	-0.072 (0.028)	-0.095 (0.009)		-0.384 (0.040)	-0.304 (0.091)		-0.201 (0.038)
Marginal effect through level (unbiased)	0.044 (0.068)	0.073 (0.087)	-0.102 (0.024)		-0.476 (0.029)	-0.499 (0.016)		0.244 (0.108)
Marginal effect through variance	-0.141 (0.065)	-0.145 (0.093)	0.007 (0.026)		0.093 (0.065)	0.195 (0.104)		-0.445 (0.142)
Standard deviation of unobservables, minority/non-minority	0.49	0.55	1.03		1.59	2.27		0.11
Wald test, overidentification, ratios of coefficients equal (β value)	0.97	0.93	0.87		0.09	0.64		0.94
LR test: standard vs. heteroskedastic probit (β value)	0.01	0.10	0.80		0.06	0.01		0.01
Wald test, ratio of standard deviations = 1 (β value)	0.00	0.03	0.79		0.18	0.16		0.00
Controls (jobs or applicants)	High education, over-educated, distance, vacancy duration, vacancies/unemployed, unemployment, % foreign, % Turkish, city, multiple jobs, average occupation wage, job quality, intensive/moderate customer contact		Unemployment spells, cultural activities, sport, personality, summer experiences, US high school, high education, multiple employers, occupation		Enhanced cover letters, enhanced reference letters, first applicant, résumé type, reference type, occupation tester,		Occupation, cover letter with good English, extracurricular skills, BA from private university, grades, language and writing skills stated (Malay, Chinese), MS Office, software/accounting skills, high-quality CV, degree or degree project on CV, pre-university institution,	
Clustered (within-pair design)	Yes	Yes	Yes		Yes	Yes		Yes
N	736	736	5,636		4,846	4,194		3,009

Notes: In panel A, the marginal effect is based on the standard formula for a discrete variable, with other variables set at sample means. In panel B, the continuous approximation for marginal effects is used, with the decomposition in Equation (7) immediately below. Standard errors for the two components of the marginal effects, reported in parentheses, are computed using the delta method. The only individual controls for which interactions are not introduced are for other demographic groups. LR, likelihood ratio. Values set in bold are key estimates to focus on.

Table 2B. Estimates for Labor Market Discrimination Studies: Full Specifications

Study	Oropoulos (2011), Canada				Bertrand and Mullainathan (2004), US	
	Chinese	Indian	Chinese-Canadian	Pakistani	Greek	British
Outcome	Call-back					
Minority group	(1)	(2)	(3)	(4)	(5)	(6)
A. Estimates from basic probit						
Minority, marginal effect	-0.053 (0.007)	-0.056 (0.007)	-0.063 (0.008)	-0.073 (0.009)	-0.035 (0.017)	-0.031 (0.011)
B. Heteroskedastic probit model						
Minority, marginal effect	-0.046 (0.009)	-0.050 (0.008)	-0.068 (0.009)	-0.083 (0.014)	-0.066 (0.073)	-0.038 (0.013)
Marginal effect through level (unbiased)						
Marginal effect through variance	-0.131 (0.046)	-0.101 (0.041)	-0.029 (0.054)	-0.076 (0.078)	-0.169 (0.208)	-0.031 (0.045)
	0.086 (0.052)	0.052 (0.046)	-0.040 (0.054)	-0.007 (0.070)	0.102 (0.139)	-0.068 (0.052)
Standard deviation of unobservables, minority/non-minority	1.46	1.26	0.84	0.97	1.54	0.75
Wald test, overidentification, ratios of coefficients equal (<i>p</i> value)	0.72	0.85	0.78	0.48	0.66	0.20
LR test: standard vs. heteroskedastic probit (<i>p</i> value)	0.07	0.22	0.46	0.92	0.33	0.21
Wald test, ratio of standard deviations = 1 (<i>p</i> value)	0.19	0.32	0.42	0.92	0.55	0.13
Controls (job or applicants)						
			Extracurricular activities, top-ranked bachelor's, master's, occupation, speaking/social/writing skills required, female			
					Bachelor's, experience and square, volunteer, military service, email address, gaps in work history, work during school, academic honors, computer and other skills, female; in zip code (% high school dropout, college graduate, black, and white, log median household income)	
Clustered (within-pair design)	Yes	Yes	Yes	Yes	Yes	Yes
<i>N</i>	5,866	6,373	4,468	3,978	3,388	3,934
						4,784

Notes: In panel A, the marginal effect is based on the standard formula for a discrete variable, with other variables set at sample means. In panel B, the continuous approximation for marginal effects is used, with the decomposition in Equation (7) immediately below. Standard errors for the two components of the marginal effects, reported in parentheses, are computed using the delta method. The only individual controls for which interactions are not introduced are for other demographic groups. Some skills are specific to immigrant groups and used to distinguish among immigrants (such as specific language fluencies or where experience obtained) and are not included. LR, likelihood ratio. Values set in bold are key estimates to focus on.

natives in Belgium—the evidence of discrimination completely disappears in the heteroskedastic probit estimates. In both columns (1) and (2), the first for a call-back and the second for an immediate interview, the negative and significant coefficient estimate on the indicator for Turkish applicants becomes positive and statistically insignificant.

By contrast, the estimated effect through the variance is negative and significant, implying that the study design generates bias toward finding evidence of discrimination. The next row of the table reports that the ratio of the estimated standard deviations of the unobservables for minority versus non-minority candidates is approximately 0.5, indicating a lower variance of unobservables for the Turkish applicants. In terms of the model, the reduction in estimated discrimination coupled with a lower variance of unobservables for minorities implies that on average the résumés in this study were of relatively low quality compared to what employers see; thus, the low-variance group is less likely to be of sufficiently high quality on the unobservables to merit a call-back, and the difference in variance creates a bias toward finding discrimination against Turkish applicants.

Below the decomposition estimates, the table reports some additional diagnostic test results. On the one hand, it reports the p value from the overidentification test that the ratios of the skill coefficients between (in this case) Turkish and native applicants are equal across all of the skills/résumé characteristics. The p value is 0.97 in column (1) and 0.93 in column (2), indicating that we do not reject the overidentifying restrictions. On the other hand, in this case, as reported in the next row, the data tend to reject the restriction to the homoskedastic specification; the p value from a likelihood ratio test is 0.01 in column (1) and 0.10 in column (2). The final test result reported is whether the ratio of variances of the unobservables equals 1, which is rejected strongly in both columns (a result we expect would parallel to some extent the likelihood ratio test).

Thus, for the Baert et al. (2015) study, application of this method of correcting for bias from differences in the variances of unobservables very much overturns the evidence of ethnic discrimination. We have one additional point to make with reference to the more general earlier discussion about interpreting the effect through the variance. One might refer to the negative (and significant) estimates on “Marginal effect through variance” as suggesting that the evidence of discrimination has not gone away, but simply been “displaced” to show up in the variance. We have already explained why, in the context of the method and the underlying model used in this article, the estimated effect through the variance is an artifact of the study and would not be expected to be replicated in the real world. Similarly, it would not be replicated if the study had used high-quality résumés, or a distribution of résumés that matched the distribution employers actually see. An alternative hypothesis though is that the effect of variance is real, and it reflects employer risk aversion rather than how the employer evaluates the likelihood that an applicant exceeds a call-back/

hiring threshold, given the résumé. If there is risk aversion, however, then high-variance groups would be penalized. That pattern is inconsistent with the evidence from the Baert et al. (2015) data, since the minority applicants are estimated to have lower variance.²¹

Having gone through the results for the first study in detail, the results for the other labor market studies can be covered more succinctly. The Carlsson and Rooth (2007) study of discrimination against Middle Easterners in Sweden asks a question very similar to the one asked in Baert et al. (2015). In this case, however, the conclusions are scarcely affected by addressing the Heckman-Siegelman critique. The estimated marginal effect through the level (-0.102) is very similar to the simple heteroskedastic probit estimate (-0.095), and the estimated marginal effect through the variance is close to zero (0.007) and estimated precisely. In this case the ratio of the estimated standard deviations of the unobservable for minorities relative to non-minorities is very close to one (1.03), which implies—in terms of the Heckman-Siegelman critique—that bias is unlikely regardless of the quality of the artificial résumés relative to the population of résumés that the employer sees, which is consistent with the robustness of the evidence for this study. Note also that the data do not reject the overidentifying restrictions, nor do they reject the restriction to the homoskedastic model or that the ratio of standard deviations equals 1, which is not surprising given the estimates.

The Drydakis (2014) study looked at discrimination against gays and lesbians. In this study, also, correcting for potential bias from differences in the variances of the unobservables does not alter the conclusion by much. Indeed, the estimated effect of being gay or lesbian is a larger negative (-0.476 or -0.499) after correcting for this bias, relative to the overall effect of -0.384 for gays and -0.304 for lesbians. For both groups, the estimated variance of the unobservable is quite a bit larger than for straight men or women, with a ratio of standard deviations of 1.59 for gay versus straight men, and 2.27 for lesbian versus straight women. The combination of a higher variance for gays or lesbians with a larger estimate of discrimination would imply that the résumés were of low quality relative to the distribution, which would lead employers to favor the high-variance group and generate a bias toward zero in the estimate of discrimination.

Note that for the Drydakis analyses there is strong evidence against the homoskedastic probit model and marginally significant evidence against equal standard deviations. Also, for the analysis of gay versus straight men the overidentifying restrictions are rejected at the 10% level. This last result prompted us to estimate a less constrained model that did not restrict the

²¹This may be too strong a statement, since if employers actually evaluate applicants based on their assumed variance of the unobservable, the statistical model might be different. We are not aware of any field experiments that have tried to incorporate risk aversion, although this might be fruitful. Dickinson and Oaxaca (2009) provided a laboratory experiment study of this type of discrimination in labor markets.

effects of two of the résumé characteristics to be the same across gay and straight men—chosen based on the estimates indicating that these interactions did not fit the expected pattern if the coefficients in the latent variable model were equal and only the variances of the unobservables varied.²² In this case, the overidentification restrictions were no longer rejected (the p value was 0.751), yet the estimates were very similar to those reported in column (4) of Table 2A.

Lee and Khalid (2016) studied discrimination against Malays (compared to Chinese) in the private sector in Malaysia.²³ In this case, the conclusions are dramatically affected by addressing the Heckman-Siegelman critique, as the estimated marginal effect through the level changes sign and becomes significant and positive, consistent with discrimination in favor of Malays.²⁴ By contrast, the estimated marginal effect through the variance is large, negative, and significant (-0.445). In this case, the ratio of the estimated standard deviations of the unobservable for Malays relative to Chinese is very low (0.11). The combination of a lower variance for Malays with a smaller (indeed, opposite-signed) estimate of discrimination would imply that the résumés were of low quality relative to the distribution for jobs included in the study, which would lead employers to favor the high-variance group and generate a bias toward discrimination in favor of Chinese applicants. Note also that the data do not reject the overidentifying restrictions.

Turning to the remaining labor market studies, in Table 2B, Oreopoulos (2011) studied outcomes for six immigrant groups relative to native Canadians. It turns out that for two of these groups—Chinese and Indian—the evidence of discrimination remains significant after addressing the Heckman-Siegelman critique and is actually stronger, with estimates changing from approximately -0.05 to -0.10 or greater. For both groups, the estimated variance of the unobservable is larger for immigrants than for natives, which appears to interact with the applicants being low quality so

²²These résumé characteristics were the indicators for a high-quality résumé (more experience) and for résumé type. These were chosen because the estimated signs of the interactions relative to the signs of the main effects were rather strongly inconsistent with what would be predicted based on the higher estimated variance of the unobservable for gays. Note that the model is identified as long as the effects of *some* variables that shift the call-back probability are restricted to be equal across the two groups; this restriction does not have to hold for all of them and can be relaxed by adding interactions between the group indicator and the résumé characteristic to the heteroskedastic probit model.

²³Malays are not the minority group, although we retain that label in the table to be consistent with other studies. Lee and Khalid (2016) discussed issues related to potential discrimination against Malays in the private sector, including affirmative action for Malays in public education that may lead Malay graduates to be less preferred. Their sample size with controls was a bit smaller than ours (see their table 4), because they also included data on the companies in the study; these data were not always available, and the company data were not provided to us.

²⁴Although this change in results is striking, other findings in the Lee and Khalid article do not cleanly fit the expected story of discrimination against Malays. In particular, they found stronger anti-Malay discrimination in hiring for private university graduates, where affirmative action in education is *not* implemented.

that the higher variance biases the estimate of discrimination from the standard probit toward zero. But for the other four groups—Chinese-Canadian,²⁵ Pakistani, Greek, and British—there is no longer significant evidence of discrimination. Note that in two cases—Pakistani and Greek—the point estimate of the marginal effect of minority group membership through the level is still a large negative number but is insignificant. By contrast, for the British, the point estimate is no longer negative.

Turning to the other diagnostics, in every case for the Oreopoulos analysis, the overidentification restrictions are not rejected. Similarly, with the exceptions of the analysis for the Chinese applicants, the data do not reject the restriction to the homoskedastic model. Thus, in this case we are sometimes failing to find evidence of discrimination because we are estimating a more flexible model even when the data do not reject a more restrictive model that provides evidence of discrimination, and the results for the Pakistani and Greek applicants are notable in this regard. Estimating a more flexible model poses the usual trade-off of bias versus precision, although generally speaking labor economists are willing to estimate less restrictive models that eliminate bias at the risk of decreased precision. Regardless, it seems reasonable to conclude that the re-analysis of the Oreopoulos data indicates far less robust evidence of discrimination than did the original study.

Finally, column (7) of Table 2B repeats the re-analysis of the Bertrand and Mullainathan (2004) data from Neumark (2012). In this case, the evidence of discrimination becomes a bit stronger, and we estimate the variance of the unobservable to be larger for blacks. These findings are consistent with low-quality résumés generating a bias against finding discrimination, although the qualitative conclusions are unchanged.

Thus, the conclusion from our re-examination of the labor market experiments is that the findings from the existing studies of discrimination against ethnic, racial, or sexual orientation minorities are not always robust to addressing the Heckman-Siegelman critique. All 13 estimates based on the existing studies, using the conventional approach, point to evidence of discrimination. But only six (or just less than one-half) of the corrected estimates provide evidence of discrimination.²⁶

This conclusion that the analysis of data from field experiments on labor market discrimination is not always robust is echoed in the findings reported in Neumark et al. (forthcoming). They studied age discrimination in hiring and found that the evidence of discrimination against older women is robust to addressing the Heckman-Siegelman critique, but the evidence of discrimination against older men is not robust. Some other recent papers using this technique do not find large differences. Carlsson et al.

²⁵This categorization refers to an English first name and a Chinese last name.

²⁶This finding includes the evidence from Carlsson and Rooth (2007), Drydakis (2014, for both gays and lesbians), Oreopoulos (2011, for Chinese and Indian), and Bertrand and Mullainathan (2004, significant at 10% level).

(2013) re-examined data from four previous studies of the Swedish labor market, each of which included some form of the data required to implement the bias correction. Their re-analysis did not lead to large changes in the estimates of discrimination, although sometimes the estimated discrimination (against those with Arabic names, and in favor of women) becomes smaller. Three recent studies by Baert, all on the Belgian labor market, found no change in the estimates of discrimination in these experimental studies. Baert (2015) implemented this method in a study of sex discrimination in Belgium for jobs entailing a promotion, using information on distance from the worker's residence to the workplace to identify the heteroskedastic probit model, and reported that this correction does not alter the conclusions (although the estimated effect of discrimination does become smaller and statistically insignificant).²⁷ Baert (2014: 551, note 15) applied the bias correction in an investigation of discrimination based on sexual orientation and family responsibilities and found no bias or difference in reported results. Baert (2016: 83–84) found similar results in a study of hiring discrimination against disabled individuals. Nunley et al. (2015) studied racial discrimination in hiring of recent college graduates in the United States. Applying the bias correction to their finding of a significant, lower interview rate to black graduates indicated that the baseline estimate of discrimination was understated, although the resulting estimated marginal effects through the level and variance were not statistically significant (p. 1118). Thus, among these latter studies, there is again sometimes an indication that the results are not robust to addressing the Heckman-Siegelman critique, although less clear is whether ignoring this critique leads to overstating discrimination.

Housing Market Field Experiments

Table 3 presents the results from re-examination of evidence from the housing discrimination studies. Ahmed et al. (2010) studied discrimination against Arab applicants in Sweden, looking at both positive responses and offers of immediate showings—as do three of the four housing studies. In this study, correcting for potential bias from differences in the variances of the unobservables did very little to change the conclusions. The estimates of lower positive responses or offers of immediate showings to Arab applicants became, if anything, more negative—most notably for immediate showing, where the estimate changes from -0.074 to -0.146 —and both estimates are statistically significant. The estimated effects of Arab ethnicity through the variance are positive, and larger for immediate showings, corresponding to the larger negative estimate on the marginal effect through the level. The estimated variance of the unobservable is larger for Arab applicants, so

²⁷Baert, De Pauw, and Deschacht (2016) used these same data, but they did not include the bias correction. Since the data were used in the 2015 paper to do the bias correction, these data are not included in our re-analysis.

combined, the estimates imply that the applications were lower quality than the population of applications to these landlords, biasing toward zero the conventional probit estimate of discrimination in immediate showings. Turning to the other diagnostics, in neither analysis are the overidentification restrictions, the restriction to a homoskedastic probit model, or equality of the standard deviations rejected. Thus, in the Ahmed et al. (2010) study, evidence of discrimination persists.

These same conclusions are echoed in the remaining columns of the table: for the Bosch et al. (2010), Carlsson and Eriksson (2014), and Ewens et al. (2014) studies. In all cases, the bias-corrected estimates still lead to statistically significant evidence of discrimination based on race and ethnicity. And, in most cases, the point estimate for the marginal effect through the level is very close to the overall heteroskedastic probit estimate, whereas the estimates of the effect of race or ethnicity through the variance are very small.²⁸

In one case, Ewens et al. (2014), our analysis rejected the overidentifying restrictions at the 10% level (and the p values for the other tests were fairly low). We therefore carried out an additional analysis, paralleling what we did with the Drydakis (2014) data on gay and straight male applicants. In this case, we estimated a less constrained model that did not restrict the effects of percent black in the area or city to be the same across black and white applicants, based on the estimates indicating that these interactions did not fit the pattern of equal coefficients in the latent variable model with probit coefficients differing because of differences in the variances of unobservables. In this case, the overidentification restrictions were no longer rejected (the p value was 0.877), yet the conclusions were similar to those in column (7) of Table 3. Our overall estimate (standard error) of discrimination from the heteroskedastic probit model was -0.064 (0.023), and the unbiased estimated effect through the level was -0.067 (0.023).

Thus, from our re-examination of the housing market studies, we conclude that the findings from the existing studies of discrimination against ethnic or racial minorities are robust to addressing the Heckman-Siegelman critique. With one minor exception, these past studies found evidence of discrimination, and our corrected estimates are qualitatively and usually quantitatively very similar.

Why might the housing market tests of call-backs for rental enquiries be more robust to addressing the Heckman-Siegelman critique? One possibility is that the information provided in the housing market tests is sufficiently

²⁸One reason for the robustness of the results in Carlsson and Eriksson (2014) could be because they use applications with substantial variation in applicant characteristics. The authors do this because by avoiding standardizing applicants to a very narrow range, the bias identified by the Heckman-Siegelman critique can be reduced, although this cannot ensure that the range of quality of actual applicants is not larger. It is also the case that, especially for the positive response outcome, the variances are nearly equal (the ratio of estimated standard deviations is 1.02), so that using a narrow range of applicant quality would not introduce bias.

complete that there is little scope for a role for unobservables, and hence little impact of any differences in the variances of unobservables across groups. In housing markets, there may be little more that matters to agents than ability to pay, and the information in the applications may convey this quite reliably. By contrast, an employer has an ongoing relationship with a worker, as do the employer's customers, so that many factors that are not conveyed in an online job application could potentially weigh on an employer's decision, and hence, correspondingly, differences in the variances of these unobserved factors across groups could matter much more.

Conclusions

The goal of this article is to re-examine evidence from field experiments on labor market and housing market discrimination (experiments that, in general, identify the combined effect of taste discrimination and statistical discrimination). Specifically, our goal is to see if the near-uniform findings of discrimination against minorities hold up after correcting for an important source of bias originally identified in Heckman and Siegelman (1993), which we refer to as the Heckman-Siegelman critique. This critique emphasizes that even under quite ideal conditions for these studies, the evidence can be biased in either direction—or, equivalently, discrimination can be unidentified—if the variances of the unobservables differ across the groups studied. This concern is plausible given that a difference in the variances of unobservables across groups cannot be ruled out and indeed is at the core of early theoretical models of statistical discrimination (Aigner and Cain 1977). We re-examine evidence from 10 studies that have the requisite data—applicant or other characteristics aside from the identifier for the group in question, which shift the probability of call-backs or hires—implementing a correction for this bias proposed in Neumark (2012).

We find that for the housing market studies, the estimated effect of discrimination is robust to this correction. For the labor market studies, the evidence is less robust; in about half of the cases the estimated effect of discrimination either falls to near zero or becomes statistically insignificant, and in some cases the sign changes.

We, of course, cannot definitively extrapolate from the 10 studies we were able to re-examine to the broader set of field experiments on discrimination by race, ethnicity, and sexual orientation. Nonetheless, given that about half of the estimates of labor market discrimination that we could re-examine no longer provide statistical evidence of discrimination (or discrimination in the same direction) after correcting for bias from differences in the variance of unobservables, it seems reasonable to suggest that the overall (and overwhelming) evidence of labor market discrimination from field experiments is likely less robust than it seems. We have no doubt that in many countries discrimination occurs in labor and housing markets against many groups, and that—like the subset of studies we re-examine in this article—

the evidence of discrimination would frequently be robust to addressing the Heckman-Siegelman critique. But our evidence also indicates that in some cases a research design that enables a researcher to address this critique would not find evidence of labor market discrimination.

If nothing else, this conclusion implies that we need three types of research to draw more definitive conclusions from field experiments on labor and housing market discrimination: 1) more evidence using this kind of research design and methods; 2) more analysis of how best to implement these methods, what kinds of quality shifters provide the most informative estimates, and so forth; and 3) further consideration of whether other methods can address the Heckman-Siegelman critique and whether they would generate similar answers. Moreover, given the non-robustness of the experimental evidence on labor market discrimination, in particular, to addressing the Heckman-Siegelman critique, one could reasonably argue that future experimental studies of labor market discrimination (and perhaps of discrimination in any market) must take account of this critique to be regarded as credible.

References

- Ahmed, Ali, Lina Andersson, and Mats Hammarstedt. 2010. Can discrimination in the housing market be reduced by increasing the information about the applicants? *Land Economics* 86(1): 79–90.
- Aigner, Dennis, and Glen C. Cain. 1977. Statistical theories of discrimination in labor markets. *Industrial and Labor Relations Review* 30(2): 175–87.
- Altonji, Joseph G., and Rebecca M. Blank. 1999. Race and gender in the labor market. In Orley C. Ashenfelter and David Card (Eds.), *Handbook of Labor Economics*, Volume 3, pp. 3143–259. Amsterdam: Elsevier.
- Baert, Stijn. 2014. Career lesbians. Getting hired for not having kids? *Industrial Relations Journal* 45(6): 543–61.
- . 2015. Field experimental evidence on gender discrimination in hiring: Biased as Heckman and Siegelman predicted? *Economics* 9(25). Accessed at <http://dx.doi.org/10.5018/economics-ejournal.ja.2015-25>.
- . 2016. Wage subsidies and hiring chances for the disabled: Some causal evidence. *European Journal of Health Economics* 17(1): 71–86.
- Baert, Stijn, and Pieter Balcaen. 2013. The impact of military work experience on later hiring chances in the civilian labour market. Evidence from a field experiment. *Economics* 7(37). Accessed at <http://www.economics-ejournal.org/economics/journalarticles/2013-37>.
- Baert, Stijn, and Elsy Verhofstadt. 2015. Labour market discrimination against former juvenile delinquents: Evidence from a field experiment. *Applied Economics* 47(11): 1061–72.
- Baert, Stijn, Bart Cockx, Niels Gheyle, and Cora Vandamme. 2015. Is there less discrimination in occupations where recruitment is difficult? *ILR Review* 68(3): 467–500.
- Baert, Stijn, Ann-Sophie De Pauw, and Nick Deschacht. 2016. Do employer preferences contribute to sticky floors? *ILR Review* 69(3): 714–36.
- Becker, Gary S. 1971. *The Economics of Discrimination*, 2nd edition. Chicago: University of Chicago Press.
- Bertrand, Marianne, and Sendhil Mullainathan. 2004. Are Emily and Greg more employable than Lakisha and Jamal? A field experiment on labor market discrimination. *American Economic Review* 94(4): 991–1013.
- Bertrand, Marianne, Dolly Chugh, and Sendhil Mullainathan. 2005. Implicit discrimination. *American Economic Review Papers and Proceedings* 95(2): 94–98.

- Bosch, Mariano, M. Angeles Carnero, and Lidia Farre. 2010. Information and discrimination in the rental housing market: Evidence from a field experiment. *Regional Science and Urban Economics* 40(1): 11–19.
- Carlsson, Magnus, and Stefan Eriksson. 2014. Discrimination in the rental market for apartments. *Journal of Housing Economics* 23: 41–54.
- Carlsson, Magnus, and Dan-Olof Rooth. 2007. Evidence of ethnic discrimination in the Swedish labor market using experimental data. *Labor Economics* 14(4): 716–29.
- Carlsson, Magnus, Luca Fumarco, and Dan-Olof Rooth. 2013. Artifactual evidence of discrimination in correspondence studies? A replication of the Neumark method. IZA Discussion Paper No. 7619. Bonn, Germany: Institute of Labor Economics.
- Cornelißen, Thomas. 2005. Standard errors of marginal effects in the heteroskedastic probit model. Institute of Quantitative Economic Research, Discussion Paper No. 230. Hanover, Germany: University of Hanover.
- Dickinson, David L., and Ronald L. Oaxaca. 2009. Statistical discrimination in labor markets: An experimental analysis. *Southern Economic Journal* 71(1): 16–31.
- Doleac, Jennifer L., and Luke C. D. Stein. 2013. The visible hand: Race and online market outcomes. *Economic Journal* 123(572): F469–92.
- Drydakis, Nick. 2014. Sexual orientation discrimination in the Cypriot labor market: Distastes or uncertainty? *International Journal of Manpower* 35(5): 720–44.
- Ewens, Michael, Bryan Tomlin, and Liang Choon Wang. 2014. Statistical discrimination or prejudice? A large sample field experiment. *Review of Economics and Statistics* 96(1): 119–34.
- Heckman, James J. 1998. Detecting discrimination. *Journal of Economic Perspectives* 12(2): 101–16.
- Heckman, James J., and Peter Siegelman. 1993. The Urban Institute audit studies: Their methods and findings. In Michael Fix and Raymond J. Struyk (Eds.), *Clear and Convincing Evidence: Measurement of Discrimination in America*, pp. 187–258. Washington, DC: Urban Institute Press.
- Jowell, Roger, and Patricia Prescott-Clarke. 1970. Racial discrimination and white-collar workers in Britain. *Race* 11(4): 397–417.
- Lee, Hwok-Aun, and Muhammed Abdul Khalid. 2016. Discrimination of high degrees: Race and graduate hiring in Malaysia. *Journal of the Asia Pacific Economy* 21(1): 53–76.
- Neumark, David. 2012. Detecting discrimination in audit and correspondence studies. *Journal of Human Resources* 47(4): 1128–57.
- . Forthcoming. Experimental research on labor market discrimination. *Journal of Economic Literature*.
- Neumark, David, Ian Burn, and Patrick Button. 2016. Experimental age discrimination evidence and the Heckman critique. *American Economic Review Papers and Proceedings* 106(5): 303–8.
- . Forthcoming. Is it harder for older workers to find jobs? New and improved evidence from a field experiment. *Journal of Political Economy*.
- Nunley, John M., Adam Pugh, Nicholas Romero, and R. Alan Seals. 2015. Racial discrimination in the labor market for recent college graduates: Evidence from a field experiment. *B.E. Journal of Economic Analysis and Policy* 15: 1093–125.
- OECD. 2013. *International Migration Outlook 2013*. Paris: Organisation for Economic Co-operation and Development.
- Oreopoulos, Philip. 2011. Why do skilled immigrants struggle in the labor market? A field experiment with thirteen thousand resumes. *American Economic Journal: Economic Policy* 3(4): 148–71.
- Pager, Devah. 2007. The use of field experiments for studies of employment discrimination: Contributions, critiques, and directions for the future. *Annals of the American Academy of Political and Social Science* 609(1): 104–33.
- Quillian, Lincoln, Devah Pager, Ole Hexel, and Arnfinn H. Midtboen. 2017. Meta-analysis of field experiments shows no change in racial discrimination in hiring over time. *Proceedings of the National Academy of Sciences of the United States of America* 114(41): 10870–75.

- Riach, Peter A., and Judith Rich. 2002. Field experiments of discrimination in the market place. *Economic Journal* 112(483): F480–518.
- Rich, Judith. 2014. What do field experiments of discrimination in markets tell us? A meta-analysis of studies conducted since 2000. IZA Discussion Paper No. 8584. Bonn, Germany: Institute of Labor Economics.
- Zschirnt, Eva, and Didier Ruedin. 2016. Ethnic discrimination in hiring decisions: A meta-analysis of correspondence tests 1990–2015. *Journal of Ethnic and Migration Studies* 42(7): 1115–34.