

Bias and discrimination: what do we know?

Marina Della Giusta* and Steven Bosworth**

Abstract: The paper presents the economic literature on gender bias, illustrating the underpinnings in the psychology of bias and stereotyping; the incorporation of these insights into current theoretical and empirical research in economics; and the literature on methods to contrast bias, presenting evidence (where it exists) of their effectiveness. The second part of the paper presents results of an experiment in revealing unconscious bias.

Keywords: implicit bias, stereotyping, moral licensing, interventions

JEL classification: B54, C90, D91, J16, J70

I. Discrimination in economics and in psychology

There is a vast empirical literature in economics documenting discrimination in many settings, from education and labour markets to access to credit, housing, and health services, offers for products and services, politics, and law enforcement (Knowles *et al.*, 2001; Bertrand and Mullainathan, 2004; List, 2004; Nelson, 2009; Rodgers, 2009; Wood *et al.*, 2009; Ewens *et al.*, 2014; Alesina *et al.*, 2018). Economists view discrimination as a situation in which individuals with identical productive characteristics are treated differently from each other on the basis of observable personal characteristics (such as age, ethnicity, gender, BMI, etc.) that are unrelated to their productivity.

The tools economists have used to analyse discriminations have been mostly those of competitive markets in which discrimination arises from the behaviour of utility-maximizing individuals, and the efforts have been directed both at understanding why it may arise, why it may persist, as well as the consequences it has on individuals and society at large. The motivation for economists in studying the problem of discrimination have been thus both on grounds of equity and of efficiency: in regards to the latter, standard

*University of Reading; e-mail: m.dellagiusta@reading.ac.uk

**University of Reading; e-mail: s.j.bosworth@reading.ac.uk

The authors acknowledge financial support from the University of Reading. We wish to thank Almudena Sevilla, Margaret Stevens, and our second referee for thorough and insightful comments, and all participants in the editorial seminar for this issue of the *Oxford Review of Economic Policy*, the seminar series of the Department of Economics at the University of Reading, and the 3rd Reading Experimental and Behavioural Economics Workshop for very useful suggestions. We also wish to thank Sahira van de Wouw, who assisted with the Gorilla programming.

doi:10.1093/oxrep/gra045

© The Author(s) 2021. Published by Oxford University Press.

For permissions please e-mail: journals.permissions@oup.com

theory would predict that a well-functioning market should gradually eliminate any selection based on non-economically relevant characteristics as it would be inefficient to discriminate, so theory has also sought to explain the persistence of discrimination in terms of imperfect markets (Stiglitz, 1973).

The seminal contribution in Gary Becker's PhD dissertation (Becker, 1957) arose in a context in which discrimination against black and women workers in the US was legal (the Equal Pay Act dates to 1963 and the Civil Rights Act to 1964), and it was mostly the subject of sociological rather than economic research, although it was the latter and especially the applied studies that were conducted on wage discrimination (Oaxaca, 1973) that was to have the strongest impact on actual litigation cases from then on (Ashenfelter and Oaxaca, 1987). Becker explained discrimination as taste-based, that is based on the dislike of a group on the part of discriminating individuals, who may be employees, co-workers, or customers. In this case, discriminating employers would be effectively willing to pay the price of their choices by allowing a restriction in the pool of talent from which they select hires and promotions. In taste-based discrimination models those who dislike other groups are willing to pay a price to avoid interaction with them (Becker, 1957) and discrimination may thus persist even when it is inefficient. Over time, it may also persist because it becomes a social norm that is costly to break, as discussed by Akerlof in the context of explaining the causes of unemployment (Akerlof, 1980). The model of course assumes that individuals are rational and know their preferences, so that they deliberately discriminate in full knowledge of the costs this will have for themselves (if not those faced by those against whom they discriminate). The solution in this case is for discriminated groups to become more skilled (thus costlier to exclude), so that discrimination becomes so expensive that even die-hard racists and sexists will be willing to forgo their exogenously given and perfectly known preferences in exchange for financial compensation.

The other classical explanation for discrimination in economics is that attributed to Arrow (1973) and Phelps (1972), who proposed statistical discrimination, suggesting that when individuals have insufficient information about a particular person they will attribute to them average group perceived or real characteristics. These attributions amount to what psychologists describe as stereotypes, which are cognitive shortcuts used by the brain to generate expectations of others' behaviour and that can drive decisions based on unverified information about the group an individual belongs to rather than specific information about them, something Kahnemann and Tversky called the heuristics of representativeness (Kahnemann and Tversky, 1973). How these stereotypes affect belief formation has been the subject of work by both Coffman (2014) and Bordalo *et al.* (2016). These models assume that significant biases in beliefs can arise from stereotypes or other conjunction fallacies (Tversky and Kahneman, 1983) which consist in exaggerating small differences in some parts of the distribution of attributes of one group relative to another (for example, believing that because men are over-represented in the top tail of the mathematics GPA scores distribution, this holds true across the whole distribution of scores). Economists have also modelled theoretically how the vicious cumulative effects of these decisions play out in aggregate outcomes. Oxoby (2014) has shown that the process of forming beliefs about one's own ability, incorporating irrelevant information on observable types, can bias downward one's perception of one's own ability (or upward if the type-based biases are positive), and lead to inefficient allocations of agents across more and less skilled sectors in the labour market and a growing segregation over time through the feedback to agents

from increased type-based biases in their beliefs. This type of discrimination can be corrected when information is supplied about the individual in question as these models expect rational individuals to correct their beliefs when they find them to be distorted. For economists, therefore, stereotyping based on group membership that results from imperfect information can be corrected through the provision of more information (Guryan and Charles, 2013).

Psychologists, though, have demonstrated that this is not so straightforward and, in particular, stereotypes related to personal characteristics that are important for identity (as ethnicity and gender typically are) are not so easily corrected. In fact, there is emerging neuroscientific research suggesting that cerebral networks used to process self-identity are different to those used to process more general knowledge and much harder to change, and thus correcting stereotypes with direct experience is difficult, particularly as other types of biases innate in the way we think come into play (Rippon, 2019). Confirmation bias makes us pay much greater attention to the information that confirms what we already believe than to new information, and belief bias makes us forget information that is contrary to our beliefs (Kahneman, 2011).

Evidence of the lengths we are willing to go to protect our biases is provided in the recent paper by Bohren *et al.* (2019) who have conducted an experiment aiming to identify separate causes of discrimination by observing discrimination developing in a dynamic setting (a large online Q&A forum used by students and researchers in STEM (science, technology, engineering, and mathematics)). They have formally tested three hypotheses on the sources of discrimination: preference-based (*à la* Becker), belief based with correct beliefs (*à la* Arrow–Phelps), and belief-based with incorrect beliefs (*à la* Coffman). They find that decreasing subjectivity in judgement (through the provision of more information on the quality of the answers provided by female and male profiles) mitigates belief-based discrimination but does not affect taste-based discrimination, which persists even when quality is perfectly observable. In particular, when discrimination is belief based (beliefs about a group's average ability), observing prior evaluation will reduce discrimination. Conversely, if discrimination is taste-based then, even receiving similar evaluation to men, the women will continue to face discrimination. The mechanism they identify to explain this result runs through the different interpretations given to the signals received observing prior evaluations: in the case of biased beliefs, evaluators become aware that the woman had to produce work of a higher standard to be positively evaluated; in the second, instead, the evaluator may believe she has been 'helped' and that the evaluations do not reflect her true quality (in other words they think the process is rigged). The paper shows both theoretically and empirically that reversals of beliefs can occur when discrimination is based on biased beliefs (stereotypes), which provides an important explanation for results observed in STEM where accomplished female academics are favoured over males (Williams and Ceci, 2015) and discrimination is instead found among female students (Reuben *et al.*, 2015), and in labour markets where discrimination occurs at hiring but reverses at promotion (Lewis, 1986; Groot and Van den Brink, 1996; Booth, 2009), a result corroborated in the leadership literature that finds a gender premium at the top of organizations and discrimination at the bottom. The paper also provides evidence corroborating both Becker and contemporary psychology literature: when people are strongly prejudiced they would rather believe the system is rigged than change their beliefs.

The psychology of stereotyping (Schneider, 2004; Jussim *et al.*, 2015) is fundamentally based on the fact that we are ‘wired to be social’ (Lieberman, 2013): our brains continuously work out extremely fast predictions on who to interact with and how, as described in the ‘thinking fast and slow’ model of decision-making (Kahnemann, 2011). The fast mode unconsciously produces the statistics of repetition to recognize which groups of people we like (in-groups) and which we do not like (out-groups) on the basis of stereotypes: shortcuts that incorporate a range of expectations of how someone else will behave. This is true also of self-expectations, that incorporate social rules of what is expected by someone like us. The slow mode conversely refers to deliberate thinking and works through problems more systematically, comparing and contrasting ideas and thus exercising more effort to arrive at more reliable decisions. Given the vast number of decisions we make in everyday life, a lot of our interactions are driven by fast rather than slow thinking and incorporate the biases that it generates. A wide body of experimental evidence from both psychology and economics shows both that stereotyping and self-stereotyping are present, and that they can be artificially engineered in a wide variety of settings (Anderson *et al.*, 2006). While economists would like to think that system two is always able to correct biases through belief updating, the evidence in psychology is cumulatively demonstrating that biases persist and, even when revealed, strong mean reversion occurs in order to protect self-identity and reduce uncertainty. For a comparative summary of the economics and psychology literatures on bias, refer to Table 1.

II. Gender differences and gender bias

Gender is one important dimension along which stereotypes are formed (and interacts, of course, with other dimensions intersectionally), affecting behaviours and the directions of research itself (Rippon, 2019). Gender stereotyping includes both descriptive stereotypes, beliefs about what men and women typically do—which derive from

Table 1: Concepts of discrimination

Economics			Incorrect beliefs-based (Coffman, Bordalo, Gennaioli)
We know what our tastes are and we know what our beliefs are. We can fix biased beliefs through information provision, we cannot fix tastes.	Taste-based (Becker) <i>You can show me women who are good at top maths, but I still won't hire them.</i>	Correct beliefs-based or Statistical (Arrow-Phelps) <i>I think women are bad at top maths, and if you show me a woman who is good at maths I change my rating of her accordingly but keep my belief about women in general.</i>	<i>I think women are bad at maths in general, but would change my mind if I encounter many typically feminine women who are good at maths.</i>
Psychology	Conscious bias (prejudiced)	Unconscious bias	
We don't know our beliefs or our tastes very well. We can somewhat mitigate biased beliefs through information.	Revealing unconscious bias activates immediate confirmation and belief bias that turns it into a conscious belief.	Revealing unconscious bias activates short-term re-evaluation and even attempts to over-compensate. Long-run effects unknown (reversion to mean possible).	

contact with each other (Fiske and Stevens, 1993), and prescriptive ones, beliefs about what men and women should do (Cialdini and Trost, 1998) which include both prescriptions and proscriptions (Koenig, 2018; Prentice and Carranza, 2002). For example, women are supposed to be communal (warm, sensitive, cooperative; a prescription for women) and avoid dominance (e.g. aggressive, intimidating, arrogant; a proscription for women), and men are supposed to be agentic (assertive, competitive, independent; PPS for men) and avoid weakness (e.g. weak, insecure, emotional; NPS for men). The psychological literature, moreover, tends to find that generally backlash is stronger for men and for boys transgressing the norms (Brown and Stone, 2016; Sullivan *et al.*, 2018).

The evidence from the literature on psychological traits suggests that women on average are expected (both by men and by other women) to be more conscientious and compliant (Carter, 2014; Eswaran, 2014), and the self-reports that generate the data on personality traits (Big Five Inventory) show women reporting on average higher levels of neuroticism, extraversion, agreeableness, and conscientiousness than men across most nations (Costa *et al.*, 2001; Schmitt *et al.*, 2008). When it comes to evaluations of own ability, men on average perceive their general intellect as higher and they tend to overestimate it, while women on average tend to do the opposite (Karwowski *et al.*, 2013).¹ Goals reporting differs, too: women on average declare that social objectives are more important than the goals connected with achievements, while men do the opposite (Piiro, 1991; Kuhn and Villeval, 2015). Recent findings from the Global Preference Survey (Falk *et al.*, 2015) also suggest that women tend to exhibit a stronger social predisposition than men, and to be more responsive to social cues (Zetland and Della Giusta, 2011; Eckel and Fullbrunn, 2015).

Gender stereotypes emerge in early childhood and have important consequences (Bian *et al.*, 2017). La Ferrara (2019) shows, making use of PISA data across OECD countries, that gender is strongly and robustly correlated with the probability of finishing university, with an effect that amounts in the most conservative specifications to a 15 per cent increase over the mean. Worldwide, the achievements and choices in maths by girls have been found to be strongly connected with the wider gender norms of societies (more gender-equal societies have lower gaps (Guiso *et al.*, 2008); in gender segregated schools girls choose STEM more and boys choose humanities more (Favara, 2012)); with the gender of professors (Carrell *et al.*, 2010); teachers' gender views (teachers who have positive expectations increase the performance of pupils (Figlio, 2005; Sprietsma, 2009; Hanna and Linden, 2012; Campbell, 2015); more gender-egalitarian teachers increase the performance and uptake of STEM by girls (Alan *et al.*, 2018; Carlana, 2018); parental beliefs (more gender-egalitarian parents have daughters that do better at maths—Eccles *et al.*, 1990; Fryer and Levitt, 2010; Cornwell *et al.*, 2013; Krings *et al.*, 2014; de San Román and De La Rica, 2016); and self-stereotyping (Coffman, 2014). Beyond the realm of human capital accumulation and its very real welfare consequences, Criado Perez (2019) has documented in her book *Invisible Women* the wide range of significant and persistent effects in the world of technology, medicine, and politics arising from a design bias which takes men as the norm, and women as a deviation from it.

¹ Parents also perceive their sons' intelligence to be higher than their daughters', while children perceive the intelligence of their fathers to be higher than that of their mothers (Karwowski *et al.*, 2013).

III. The effect of the exposure to stereotypes

Exposure to stereotypes effectively acts to hijack the brain by diverting resources to defending one's identity, activating a process of monitoring for failure, and attempting to suppress negative thoughts both of which load working memory that would instead be needed to perform the task at hand (Schmader, 2010). This generates cognitive overload that diminishes performance. Poverty, for example, has also been shown to impede economic decision-making (Shah *et al.*, 2012; Carvalho *et al.*, 2016; Adamkovič and Martončík, 2017) and performance on cognitively demanding tasks (Evans and Schamberg, 2009; Mani *et al.*, 2013; Lichand and Mani, 2020).

Stereotype threat thus affects performance, and evidence suggests it can be manipulated—for example, to convince pupils to believe they belong to a group that has a natural advantage in a particular subject ahead of a test. More often than not it acts in negative ways: if pupils are reminded of their gender they do worse in the subjects in which they are expected to do badly, mathematics for girls and English for boys (Johns *et al.*, 2005; Jussim *et al.*, 2015). The process starts very early: girls aged 4 have a worse performance in a spatial skills test if they colour in a girl playing with a doll before taking the test (Shenouda and Danovitch, 2014).

The literature has shown that exposure to bias toward one's group affects effort, self-confidence, and productivity (Bordalo *et al.*, 2016; Glover *et al.*, 2017; Carlana, 2018). Recent contributions in the development literature reviewed in La Ferrara (2019) explicitly link bias and stereotyping to the process of aspirations formation and aspirations as key contributing factors to poverty traps. Aspirations are strongly correlated to expectations (Carlana *et al.*, 2018), and expectations have been shown to affect performance, for example in the case of education, independently of previous attainment and parental and other characteristics (Jacob and Wilder, 2010).

Exposure to counter-stereotypical role models has been shown to help improve maths performance by women when tests are administered by women whose competence is highlighted (Marx and Roman, 2002). Breda *et al.* (2018) have shown that exposure to female role models in schools can increase the proportion of female students who choose STEM subjects, and recently Porter and Serra (2020) have conducted an RCT (randomized controlled trial) showing that exposure to female role models majoring in economics has a positive effect on the choice of economics majors by women.

IV. Combating gender discrimination

While much of the literature in economics has traditionally concentrated on assessing the effectiveness of specific policy interventions aiming to redress gender inequality (from equal pay legislation to parental leave or quotas), and more recently also on educational interventions aiming to increase the confidence, grit, self-efficacy, leadership, or maths ability of girls (Alan *et al.*, 2019, 2020), there remains a large understudied area in establishing both the detection and the mitigation of discrimination arising from the psychological processes that feed into subjective elements of evaluation.

Discrimination often occurs along dimensions that are hard to quantify, such as the language used when engaging with and evaluating members of a targeted group.

Combining evaluators from different backgrounds and making evaluations anonymous continues to be an important pillar of any strategy aiming to combat discrimination (Bohnet, 2016), but detection is anyway key to effective policy.

There is evidence that appraisal reports written for women are different from those written for men: Dutt *et al.* (2016) have investigated recommendation letters for post-doctoral fellowships in geoscience, focusing on letter length and letter tone, and found that women are significantly less likely to receive excellent recommendation letters than their male counterparts at a critical juncture in their career. Wu (2018) has measured gendered language in postings on the Economics Job Market Rumors (EJMR) forum and found that the words most predictive of a post about a woman (*female* words) are generally about physical appearance or personal information, whereas those most predictive of a post about a man (*male* words) tend to focus on academic or professional characteristics. Bohren *et al.* (2019) found in their study a significant difference in the sentiment of answers to questions from male versus female questions. Answers to questions posted by female accounts score significantly higher on both negative and positive sentiment; responses to female users contain more opinion words, both positive and negative, than responses to male users. The detection of bias in all these cases is likely to involve multiple dimensions, and research is still at the early stages, with algorithms offering promise (Kleinberg *et al.*, 2019) though the latter can incorporate bias in their design or in fact might inadvertently exacerbate it (Schlesinger *et al.*, 2018).

Unconscious bias training (UBT) has been extensively adopted in many organizations, under the premise that by bringing our unconscious biases into conscious awareness, we can determine whether they are still appropriate behaviours and we can find ways to mitigate their impact on our behaviour and decisions (Equality Challenge Unit, 2013). Most UBT interventions include an implicit association test accompanied by a debrief, education on the psychological principles underpinning unconscious bias and its natural presence, information on its impacts, and suggested techniques to mitigate it, ranging from exposure to non-stereotypical situations to aid the removal of group characterization in favour of individuation (Burns *et al.* (2017) present evidence from South Africa), to the use of strategies to reduce the impact by requiring justification of choices, reducing the importance of subjective elements in the evaluation providing very clear performance indicators. There is evidence that appropriately designed UBT works for doctors, managers, and school pupils (Teal *et al.*, 2012; Devine *et al.*, 2012; Campbell, 2015; Gilliam *et al.*, 2016; Morris and Perry, 2017; Atewologun *et al.*, 2018). The principles of well-designed training have been identified as those that include an understanding of what stereotyping and unconscious bias are, that they have real effects and that we all have them and must accept this, that we must mitigate their effects creating a culture of recognition and promoting a culture of motivation of decision based as much as possible on a combination of clearly measurable indicators and bringing diverse perspectives to neutralize each other's biases. A recent meta-review of Implicit Association Test (IAT)-based interventions (FitzGerald *et al.*, 2019) found that while perspective taking was useful but only in the short run, exposure to counter-stereotypical examples was a quite effective tool.

The actual detection of unconscious bias in individuals has been widely done making use of Implicit Association Tests (<https://implicit.harvard.edu/implicit/>). The test typically requires assigning words to categories following both stereotypical (congruent) and

non-stereotypical (incongruent) associations and measuring the speed with which such associations are made. For example, to test for a gender and leadership implicit association one would have to assign female and male first names to words associated with being a leader or a follower in both congruent and incongruent scenarios.

Results of the test have been linked to negative expectations of employees (Reuben *et al.*, 2015), and, at the receiving end of the bias, to the performance of employees of biased managers and pupils' performance and self-confidence in maths ability (Glover *et al.*, 2017; Carlana, 2018). Revealing bias on its own can elicit subsequent moderating behaviours (Alesina *et al.*, 2018), but is criticized by psychologists because of the possible negative reactions it can elicit (Howell *et al.*, 2015), and, most importantly, because of likely subsequent effects that occur through moral licensing (Mazar and Zhong, 2010; Merritt *et al.*, 2012; Cascio and Plant, 2015). This is the process of behaving in a moral way but later being more likely to display behaviours that are immoral, which in the case of bias might lead, subsequent to revelation, to reversion to previous biased beliefs and even backlash in subsequent longer-term behaviours.²

It remains an open question in the social psychology literature whether exposing people to information about prior (im)moral behaviour leads to compensating or reinforcing subsequent behaviour. While the moral licensing effect has been replicated in many studies (see Blanken *et al.*, 2015; Simbrunner and Schlegelmilch, 2017), Conway and Peetz (2012) point out that the moral licensing effect is at odds with results showing that people strive for consistency with past behaviour. Mullen and Monin (2016) argue that information is likely to elicit consistency when it is more abstract and temporally removed, when it is more relevant to the individual's identity, and when the individual's motives underlying their initial behaviour were ambiguous.³

The extant research only, however, speaks to licensing effects arising from interventions to explicit pro- or anti-social behaviour. An evolving public discourse around unconscious bias has yielded a consensus that people may be morally responsible for their biased attitudes (Kelly and Roedder, 2008; Holroyd, 2015; Brownstein, 2016; Holroyd *et al.*, 2017), especially to the extent that they are aware of them (Holroyd, 2015; Madva, 2017). Counteracting biased attitudes is seen as socially desirable, at least among certain population segments (Whiteout, 2018).

A crucial question then, for practitioners seeking to counteract unconscious bias by making people aware of theirs, is whether giving them an opportunity to improve (e.g. Carlana, 2018) is likely to generate licensing behaviour which counteracts the intervention, or identity-affirming behaviour which reinforces it. In order to test whether bias revelation that gives rise to subsequent corrective behaviours also leads to further licensing, we have designed an experiment in bias revelation, the results of which are presented in section V.

V. Revealing gender bias: an experiment

The aim of our study was to reveal experimentally to a selected group of participants their degree of bias, using a common measure from the implicit association test (IAT)

² Bagues and Esteve-Volart (2010) and Schwarz (2011) show, for example, that including more women on hiring panels can bias the selection process *more* towards male candidates.

³ Clot *et al.* (2016) show an interaction between intended motive and social identity in the moral licensing effect: mandatory acts are more likely to elicit compensatory behaviour when the act is identity-relevant.

administered on a computer. These participants were then asked to repeat the IAT exercise on paper. This opportunity for redemption was intended to provide a moral licence, which we then tested by asking participants to hand in their paper forms to two confederate researchers, one female and one male, whom they were told would tally the responses. This final behaviour (handing in the form) is our measure of moral licensing.

A total of 106 students completed participation in an undergraduate introductory microeconomics module at a large university in the south-east of England. The experimenters (also employees at this university) were introduced by the module convener during lecture time. Though participation was voluntary, most in attendance completed at least part of the study. The study consisted of two parts. In the first part, a link to an online IAT was sent to students' email addresses, and those with laptops present were asked to click on the link and complete it in class. The computer IAT was programmed in and hosted on the Gorilla.sc platform. The IAT first asked students to categorize a list of nouns as either 'Male' or 'Female'⁴ and then a separate list of nouns as relating to 'Science' or 'Liberal Arts'.⁵ Incorrect categorizations were communicated as such and participants were not permitted to move on to the next word until they assigned the word to the correct category. The two lists were then mixed in a random order and subjects had to allocate them to a combined category of 'Male or Science' or 'Female or Liberal Arts'. This condition is called the *Congruent* task since it conforms with gender stereotypes about STEM subjects being more male-oriented. A final task asked students to assign a re-mixing of these words to a 'Male or Liberal Arts' or 'Female or Science' category. This is the *Incongruent* task since it goes against the gender stereotype above. At the end of the computer IAT, half the students, randomly allocated by the server, were shown the difference in milliseconds of their average reaction times between the Incongruent and Congruent tasks. These subjects were also told that the difference in reaction times is used as a measure of their bias for associating science-based subjects with men. The other half of the participants were not shown their reaction times, but were similarly informed that their reaction times were recorded, and that this is a measure of psychological bias towards associating men with science-based subjects.

Following completion of the online IAT, a paper booklet was distributed to the students. This booklet was the same for all students and comprised a cover sheet on which students were asked for their student ID number and an instruction to wait for the experimenter before opening the booklet. Once all booklets had been handed out an experimenter instructed them to turn to the first page, which contained a vertical list of the above nouns, again in random order, with bubbles to fill in indicating whether this word was associated with a 'Science/Male' category or a 'Liberal-arts/Female' category. Participants were asked not to begin the task until the experimenter started a 30-second timer, and told that they should stop when the 30 seconds allocated were announced as over. The students were then told to turn to the next page, which again had the same list of words in a new randomized order, but the categories were reversed to be 'Liberal-arts/Male' and 'Science/Female'. Participants had 30 seconds again to complete the categorization task.

⁴ These words were Grandpa, Aunt, Son, Father, Sister, Wife, Girl, Woman, Mother, Man, Boy, Grandma, Husband, Uncle, Brother, and Daughter.

⁵ These words were Literature, History, Geology, Music, Physics, Biology, English, Engineering, Astronomy, Humanities, Chemistry, and Philosophy.

Once all participants had finished, they were informed that we would use their student ID number to match their paper responses to their computer responses for comparison (individualized links had been sent by the platform and email addresses were associated with responses). We told students that once the matching had been done all identifiers would be stripped from the data for analysis, storage, and reporting purposes. Finally, students were asked to hand in their paper booklets to one of two class tutors, one male and one female, who were positioned at the exit of the lecture theatre. Students were also informed that each tutor would tally the booklets handed into them.

(i) Data

We recorded 123 complete responses to the computer IAT and 171 complete responses to the paper IAT. Of these, we were able to match 101 paper and computer responses.⁶ This is because some students did not complete the computer IAT and some students either did not complete the paper IAT or did not record their student ID number on their paper booklet. We report results only from the 101 matched cases. Among the matched sample, 48 participants received feedback on their computer IAT and 53 subjects did not receive feedback. Of the participants 74 were male whereas only 27 were female; this is unbalanced but is reflective of the gender balance in the economics degree programmes at this university, which are approximately two-thirds male. Fifty participants handed in their paper IAT to the female tutor and 51 to the male tutor.

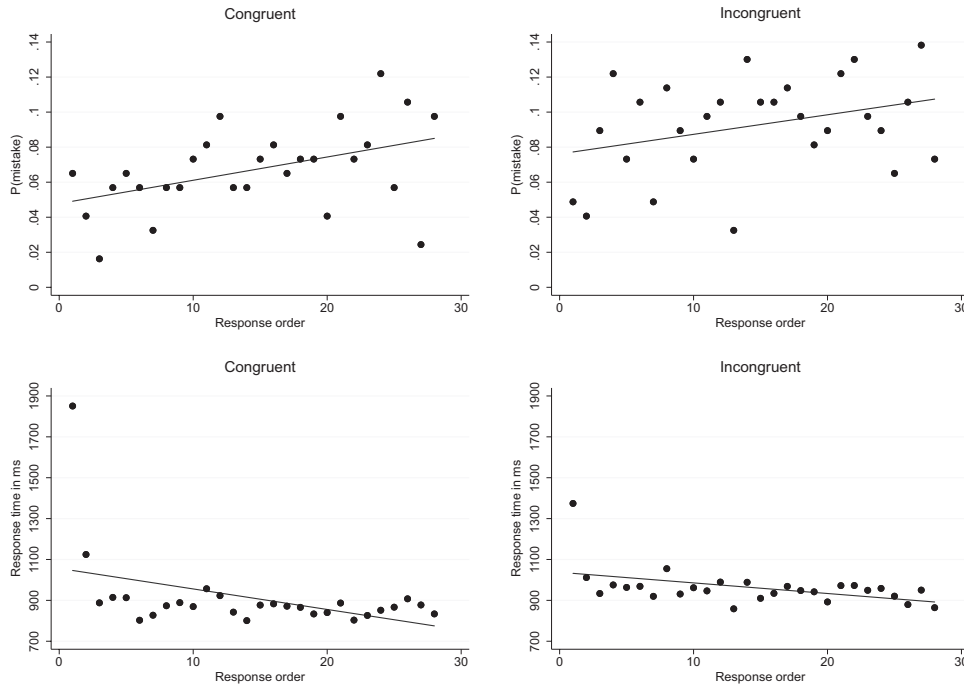
While the computerized version of the IAT enforces that subjects eventually choose the correct response to the categorization question, subjects were able to make any number of mistakes before selecting the correct of the two responses. The average number of mistakes per response and average total reaction time per response are shown in Figure 1. Note that under both the congruent and incongruent conditions, subjects were given 28 words to categorize. The horizontal axis represents the order they progressed through these words, which could have been different for every subject. We see that as subjects progress through the computerized task they gradually make more mistakes as they answer the questions slightly faster. The first word presented under each condition takes markedly longer for subjects to categorize than all the others, possibly due to attention lags. Unconscious bias on computerized IATs is most frequently measured through reaction times, incorporating the time taken for mistakes, so we use these going forward.

Participants' total reaction times to both the Congruent and Incongruent computer tasks were recorded in milliseconds. From these numbers we compute a bias score in the computerized task, $bias_{Ci}$ for subject i by

$$bias_{Ci} = \frac{\overline{\text{incongruent task reaction time}}_i - \overline{\text{congruent task reaction time}}_i}{\sqrt{\sigma_{\text{incongruent}, i}^2 + \sigma_{\text{congruent}, i}^2}}$$

⁶ A further eight subjects have paper IATs that can be matched to computerized IATs, but these subjects completed the computer IAT *after* the classroom experiment had finished and therefore their paper IAT responses are not commensurate with the others. These subjects are not included in the analysis.

Figure 1: Average probability of mistakes (top) and total reaction times for items (bottom) across congruent (left) and incongruent (right) conditions. Horizontal axis is the order in which each item is presented



dividing the subject's mean difference in reaction times between the congruent and incongruent conditions by their standard deviation of reaction time to control for heterogeneity among participants in ability and task engagement.⁷

From the paper IATs, recall that subjects were under time pressure and reaction times aren't observed, so our primary measure of cognitive association here is correct answers.⁸ Subjects' total number of correct categorizations in both the Congruent and Incongruent tasks were recorded and from these numbers we compute a bias score in the paper task, $bias_P$ by

$$bias_P = (\text{congruent correct categorisations} - \text{incongruent correct categorisations}) \\ \times \left(\frac{\text{congruent correct categorisations}}{\text{incongruent correct categorisations}} \right)^2,$$

again to control for any participant heterogeneity unrelated to gender bias. The bias scores were then made comparable by ranking them from lowest to highest among participants and computing percentile scores $\%bias_C$ and $\%bias_P$. Relative improvement between the computer and paper IAT tasks was taken as the difference in percentiles:

$$\text{improvement} = (\%bias_C - \%bias_P) / 2,$$

⁷ Greenwald *et al.* (2003) suggest this method for scoring computerized IATs is preferable.

⁸ This is suggested by Lemm *et al.* (2008) as the appropriate analogue to reaction time in paper IATs. We also adopt in calculating $bias_P$ the scoring method they suggest as having the highest correlation with Greenwald *et al.*'s computerized IAT score.

Figure 2: Bias in response to feedback

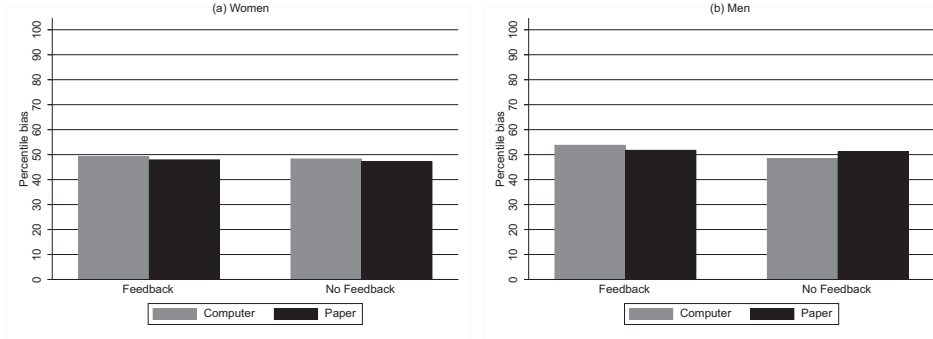
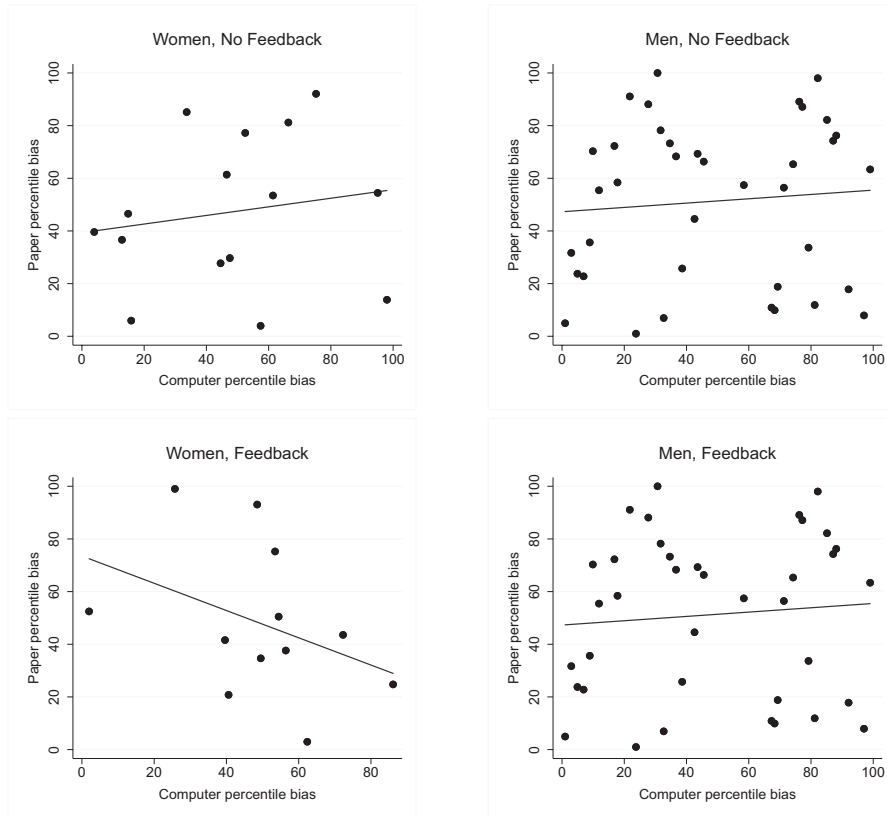


Figure 3: Relationship between computerized and paper IAT percentiles for women and men by feedback condition



with *improvement* > 0 corresponding to those who moved down in the bias ranking between computer and paper IATs and those with *improvement* < 0 moving up in the bias rankings. This was rather variable, with the minimum improvement being negative 37 percentiles and the maximum being positive 45 percentiles, with a standard deviation of 19.5. By construction mean improvement was zero.

VI. Results

Figure 2 shows the percentile bias scores across the Feedback and No Feedback treatments among women (left panel) as well as men (right panel). First, let us remark that bias is similar among both men and women. There are no significant differences in bias rank across genders for either the computer or paper tests. Second, participants exposed to feedback about their gender bias are very slightly more likely to improve in the bias ranking when completing the paper test than those who do not receive feedback, and this is not statistically significant (ranksum $p = .45$).

Figure 3 plots the relationship between the percentiles of both computerized and paper IATs by feedback condition for both women and men. There is surprisingly low correlation between the two IAT measures, likely due to the information presented to subjects that this was a measure of gender bias. The correlations are positive for both women and men under no feedback ($\rho = .17$ and $.08$ respectively, neither significant) as well as for men who receive feedback about their gender bias ($\rho = .18$, also insignificant) but negative for women who receive feedback ($\rho = -.37$), though this is also

Figure 4: Likelihood to return the paper IAT to female tutor

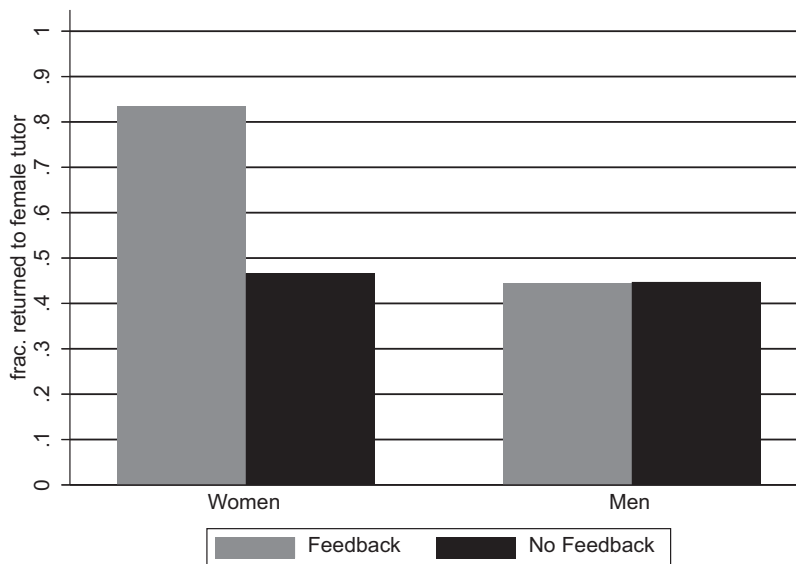


Table 2: Linear probability regressions of returning IAT to female tutor on feedback

	Women	Men
feedback	.373** (.171)	.008 (.117)
improvement	.005 (.008)	-.004 (.004)
feedback × improvement	-.013 (.009)	-.002 (.006)
constant	0.467 (.137)	0.442 (.082)

only suggestive ($p = .23$). We would cautiously interpret this result as meaning female subjects exposed to feedback try (but on average fail) to improve their gender bias between tests.

When we examine participants' likelihood to return their booklet to the female rather than the male tutor, however, we see a stark pattern. This is indicated in [Figure 4](#). Women receiving feedback were much more likely to return their IAT booklet to the female rather than the male tutor, compared with women who did not receive feedback on the computer IAT (ranksum $p = .05$). Men on the other hand were just as likely to return their booklets to the female tutor under both conditions ($p = .98$).

It seems from these aggregate patterns that, among women at least, participants receiving feedback about their implicit gender bias made efforts to improve this when measured again; but ended up giving more work to the female tutor once they had demonstrated to themselves their improvement. The data in fact paint a more nuanced picture than this. [Table 2](#), which reports the coefficients from linear probability regressions of returning the IAT booklet to the female tutor on feedback treatment, participant improvement score, and their interaction, broken out by participant gender, illustrates. We see among women that the main effect of feedback for those whose bias percentile score did not improve is positive, as per [Figure 2](#). Note, however, the interaction between feedback and improvement. It is precisely those women whose bias percentile did *not* improve—or indeed worsened—who were most likely to hand their paper in to the female tutor. Male subjects, on the other hand, display no noticeable difference in behaviour between those whose bias scores improved or did not, though these effects are very imprecisely estimated. It could be that since men on average did not improve with feedback, no moral licensing resulting from a 'second chance' at the paper IAT would have been sought.

VII. Concluding remarks

As ways of detecting the role of bias in everyday decisions and the cumulative effects it exerts on unequal outcomes continue to develop, we are beginning to understand better how to identify crucial points for policy intervention to contrast gender bias and stereotyping. The next stage for research is to conduct thorough assessments of the policies that work to counter unconscious bias, and their appropriateness for different policy areas. Our experiment in revealing bias, licensing bias (by eliciting it again), and exercising biased choices shows that revealing gender bias does not lead to corrective behaviour by men in our student sample, but it does on average lead to correction and thereafter to a larger gender-biased choice by women in the sample. The experiment also suggests that the effects are different depending on the initial level of bias of the subjects as well as their gender. This kind of work needs to be replicated to much larger scales to be meaningful, but it does illustrate the need for caution when advocating bias revelation, as it might inadvertently cause people to act as if they have overcome their bias when they have not. Together with the gains that can be achieved by contrasting unconscious bias, work will anyway need to continue in combating structural barriers and institutional 'gendered bottlenecks' in order to realize progress towards realizing gender equality.

References

- Adamkovič, M., and Martončík, M. (2017), 'A Review of Consequences of Poverty on Economic Decision-making: A Hypothesized Model of a Cognitive Mechanism', *Frontiers in Psychology*, **8**, 1784.
- Akerlof, G. A. (1980), 'A Theory of Social Custom, of which Unemployment may be one Consequence', *The Quarterly Journal of Economics*, **94**(4), 749–75.
- Alan, S., Boneva, T., and Ertac, S. (2019), 'Ever Failed, Try Again, Succeed Better: Results from a Randomized Educational Intervention on Grit', *The Quarterly Journal of Economics*, **134**(3), 1121–62.
- Ertac, S., and Mumcu, I. (2018), 'Gender Stereotypes in the Classroom and Effects on Achievement', *Review of Economics and Statistics*, **100**(5), 876–90.
- — Kubilay, E., and Loranth, G. (2020), 'Understanding Gender Differences in Leadership', *The Economic Journal*, **130**(626), 263–89.
- Alesina, A., Carlana, M., Ferrara, E. L., and Pinotti, P. (2018), 'Revealing Stereotypes: Evidence from Immigrants in Schools', No. w25333, National Bureau of Economic Research.
- Anderson, L., Fryer, R.M and Holt, C. (2006), 'Discrimination: Experimental Evidence from Psychology and Economics', *Handbook on the Economics of Discrimination*, 97–118.
- Arrow, K. J. (1973), 'The Theory of Discrimination', in O. Ashenfelter and A. Rees (eds), *Discrimination in Labor Markets*, Princeton, NJ, Princeton University Press.
- Ashenfelter, O., and Oaxaca, R. (1987), 'The Economics of Discrimination: Economists Enter the Courtroom', *The American Economic Review*, **77**(2), 321–5.
- Atewologun, D., Cornish, T., and Tresh, F. (2018), 'Unconscious Bias Training: An Assessment of the Evidence for Effectiveness', Equality and Human Rights Commission Research Report Series.
- Bagues, M., and Esteve-Volart, B. (2010), 'Can Gender Parity Break the Glass Ceiling? Evidence from a Repeated Randomized Experiment', *Review of Economic Studies*, **77**(4), 1301–28.
- Becker, G. S. (1957), *The Economics of Discrimination*, Chicago, IL, University of Chicago Press.
- Bertrand, M., and Mullainathan, S. (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.
- Bian, L., Leslie, S. J., and Cimpian, A. (2017), 'Gender Stereotypes about Intellectual Ability Emerge Early and Influence Children's Interests', *Science*, **355**(6323), 389–91.
- Blanken, I., van de Ven, N., and Zeelenberg, M. (2015), 'A Meta-analytic Review of Moral Licensing', *Personality and Social Psychology Bulletin*, **41**(4), 540–58.
- Bohnet, I. (2016), *What Works: Gender Equality by Design*, Cambridge, MA, Belknap Press of Harvard University Press.
- Bohren, J. A., Imas, A., and Rosenberg, M. (2019), 'The Dynamics of Discrimination: Theory and Evidence', *American Economic Review*, **109**(10), 3395–436.
- Booth, A. L. (2009), 'Gender and Competition', *Labour Economics*, **16**(6), 599–606.
- Bordalo, P., Coffman, K., Gennaioli, N., and Shleifer, A. (2016), 'Stereotypes', *The Quarterly Journal of Economics*, **131**(4), 1753–94.
- Breda, T., Grenet, J., Monnet, M., and Van Effenterre, C. (2018), 'Can Female Role Models Reduce the Gender Gap in Science? Evidence from Classroom Interventions in French High Schools', PSE Working Papers, halshs-01713068, HAL.
- Brown, C. S., and Stone, E. A. (2016), 'Gender Stereotypes and Discrimination: How Sexism Impacts Development', *Advances in Child Development and Behavior*, **50**, 105–33.
- Brownstein, M. (2016), 'Attributionism and Moral Responsibility for Implicit Bias', *Review of Philosophy and Psychology*, **7**, 765–86.
- Burns, M. D., Monteith, M. J., and Parker, L. R. (2017), 'Training Away Bias: The Differential Effects of Counterstereotype Training and Self-regulation on Stereotype Activation and Application', *Journal of Experimental Social Psychology*, **73**, 97–110.
- Campbell, T. (2015), 'Stereotyped at Seven? Biases in Teacher Judgement of Pupils' Ability and Attainment', *Journal of Social Policy*, **44**(3), 517–47.

- Carlana, M. (2018), 'Implicit Stereotypes: Evidence from Teachers' Gender Bias', HKS Working Paper No. RWP18-034.
- La Ferrara, E., and Pinotti, P. (2018), 'Goals and Gaps: Educational Careers of Immigrant Children', Working Paper Series rwp18-036, Harvard University, John F. Kennedy School of Government.
- Carrell, S. E., Page, M. E., and West, J. E. (2010), 'Sex and Science: How Professor Gender Perpetuates the Gender Gap', *Quarterly Journal of Economics*, **125**(3), 1101–44.
- Carter, M. J. (2014), 'Gender Socialization and Identity Theory', *Social Sciences*, **3**(2), 242–63.
- Carvalho, L. S., Meier, S., and Wang, S. W. (2016), 'Poverty and Economic Decision-making: Evidence from Changes in Financial Resources at Payday', *American Economic Review*, **106**(2), 260–84.
- Cascio, J., and Plant, A. (2015), 'Prospective Moral Licensing: Does Anticipating Doing Good Later Allow You to be Bad Now?', *Journal of Experimental Social Psychology*, **56**, 110–16.
- Cialdini, R. B. and Trost, M. R. (1998), 'Social Influence: Social Norms, Conformity and Compliance', in D. T. Gilbert, S. T. Fiske, and G. Lindzey (eds), *The Handbook of Social Psychology*, McGraw-Hill, 151–92.
- Clot, S., Grolleau, G., and Ibanez, L. (2016), 'Do Good Deeds Make Bad People?', *European Journal of Law and Economics*, **42**(3), 491–513.
- Coffman, K. B. (2014), 'Evidence on Self-stereotyping and the Contribution of Ideas', *The Quarterly Journal of Economics*, **129**(4), 1625–60.
- Conway, P., and Peetz, J. (2012), 'When Does Feeling Moral Actually Make You a Better Person? Conceptual Abstraction Moderates Whether Past Moral Deeds Motivate Consistency or Compensatory Behaviour', *Personality and Social Psychology Bulletin*, **38**(7), 907–19.
- Cornwell, C., Mustard, D. B., and Van Parys, J. (2013), 'Noncognitive Skills and the Gender Disparities in Test Scores and Teacher Assessments: Evidence from Primary School', *Journal of Human Resources*, **48**(1), 236–64.
- Costa Jr, P., Terracciano, A., and McCrae, R. R. (2001), 'Gender Differences in Personality Traits across Cultures: Robust and Surprising Findings', *Journal of Personality and Social Psychology*, **81**(2), 322–33.
- Criado Perez, C. (2019), *Invisible Women: Exposing Data Bias in a World Designed for Men*, Random House.
- de San Román, A. G., and de La Rica, S. (2016), 'Gender Gaps in PISA Test Scores: The Impact of Social Norms and the Mother's Transmission of Role Attitudes', *Estudios de Economía Aplicada*, **34**(1), 79–108.
- Devine, P. G., Forscher, P. S., Austin, A. J., and Cox, W. T. (2012), 'Long-term Reduction in Implicit Race Bias: A Prejudice Habit-breaking Intervention', *Journal of Experimental Social Psychology*, **48**(6), 1267–78.
- Dutt, K., Pfaff, D. L., Bernstein, A. F., Dillard, J. S., and Block, C. J. (2016), 'Gender Differences in Recommendation Letters for Postdoctoral Fellowships in Geoscience', *Nature Geoscience*, **9**(11), 805.
- Eccles, J. S., Jacobs, J. E., and Harold, R. D. (1990), 'Gender-role Stereotypes, Expectancy Effects, and Parents' Role in the Socialization of Gender Differences in Self Perceptions and Skill Acquisition', *Journal of Social Issues*, **46**, 182–201.
- Eckel, C. C., and Fullbrunn, S. C. (2015), 'Thar SHE Blows? Gender, Competition, and Bubbles in Experimental Asset Markets', *American Economic Review*, **105**(2), 906–20.
- Equality Challenge Unit (2013), *Unconscious Bias in Higher Education*, London, Equality Challenge Unit.
- Eswaran, M. (2014), *Why Gender Matters in Economics*, Princeton, NJ, Princeton University Press.
- Evans, G. W., and Schamberg, M. A. (2009), 'Childhood Poverty, Chronic Stress, and Adult Working Memory', *Proceedings of the National Academy of Sciences of the United States of America*, **106**(16), 6545–9.
- Ewens, M., Tomlin, B., and Wang, L. C. (2014), 'Statistical Discrimination or Prejudice? A Large Sample Field Experiment', *Review of Economics and Statistics*, **96**(1), 119–34.
- Falk, A., Becker, A., Dohmen, T. J., Enke, B., and Huffman, D. (2015), 'The Nature and Predictive Power of Preferences: Global Evidence', IZA DP No. 9504.

- Favara, M. (2012), 'The Cost of Acting "Girly": Gender Stereotypes and Educational Choices', IZA DP No. 7037.
- Figlio, D. (2005), 'Names, Expectation and the Black-White Test Score Gap', NBER Working Paper 11195.
- Fiske, S. T., and Stevens, L. E. (1993), *What's So Special About Sex? Gender Stereotyping and Discrimination*, Sage Publications.
- FitzGerald, C., Martin, A., Berner, D., and Hurst, S. (2019), 'Interventions Designed to Reduce Implicit Prejudices and Implicit Stereotypes in Real World Contexts: A Systematic Review', *BMC Psychology*, **7**(1), 29.
- Fryer Jr, R. G., and Levitt, S. D. (2010), 'An Empirical Analysis of the Gender Gap in Mathematics', *American Economic Journal: Applied Economics*, **2**(2), 210–40.
- Gilliam, W. S., Maupin, A. N., Reyes, C. R., Accavitti, M., and Shic, F. (2016), 'Do Early Educators' Implicit Biases Regarding Sex and Race Relate to Behavior Expectations and Recommendations of Preschool Expulsions and Suspensions', Research Study Brief, Yale Child Study Center, New Haven, CT, Yale University.
- Glover, D., Pallais, A., and Pariente, W. (2017), 'Discrimination as a Self-fulfilling Prophecy: Evidence from French Grocery Stores', *The Quarterly Journal of Economics*, **132**(3), 1219–60.
- Greenwald, A. G., Nosek, B. A., and Banaji, M. R. (2003), 'Understanding and Using the Implicit Association Test: I. An Improved Scoring Algorithm', *Journal of Personality and Social Psychology*, **85**(2), 197–216.
- Groot, W., and van den Brink, H. M. (1996), 'Glass Ceilings or Dead Ends: Job Promotion of Men and Women Compared', *Economics Letters*, **53**(2), 221–6.
- Guiso, L., Monte, F., Sapienza, P., and Zingales, L. (2008), 'Culture, Gender, and Math', *Science*, **320**(5880), 1164–5.
- Guryan, J., and Charles, K. K. (2013), 'Taste-based or Statistical Discrimination: The Economics of Discrimination Returns to its Roots', *The Economic Journal*, **123**(572), F417–F432.
- Hanna, R. N., and Linden, L. L. (2012), 'Discrimination in Grading', *American Economic Journal: Economic Policy*, **4**(4), 146–68.
- Holroyd, J. (2015), 'Implicit Bias, Awareness and Imperfect Cognitions', *Consciousness and Cognition*, **33**, 511–23.
- Scaife, R., and Stafford, T. (2017), 'Responsibility for Implicit Bias', *Philosophy Compass*, **12**(3), e12410.
- Howell, J. L., Gaither, S. E., and Ratliff, K. A. (2015), 'Caught in the Middle: Defensive Responses to IAT Feedback Among Whites, Blacks, and Biracial Black/Whites', *Social Psychological and Personality Science*, **6**(4), 373–81.
- Jacob, B. A., and Wilder, T. (2010), 'Educational Expectations and Attainment', No. w15683, National Bureau of Economic Research.
- Johns, M., Schmader, T., and Martens, A. (2005), 'Knowing is Half the Battle: Teaching Stereotype Threat as a Means of Improving Women's Math Performance', *Psychological Science*, **16**(3), 175–9.
- Jussim, L., Crawford, J., Anglin, S., Chambers, J., Stevens, S., and Cohen, F. (2015), 'Stereotype Accuracy: One of the Largest and Most Replicable Effects in All of Social Psychology', in T. Nelson (ed.), *The Handbook of Prejudice, Stereotyping, and Discrimination*, Mahwah, NJ, Lawrence Erlbaum.
- Kahneman, D. (2011), *Thinking, Fast and Slow*, New York, Allen Lane.
- Tversky, A., (1973), 'On the Psychology of Prediction', *Psychological Review*, **80**, 237–51.
- Karwowski, M., Lebuda, I., Wisniewska, E., and Gralewski, J. (2013), 'Big Five Personality Traits as the Predictors of Creative Self-efficacy and Creative Personal Identity: Does Gender Matter?', *The Journal of Creative Behavior*, **47**(3), 215–32.
- Kelly, D., and Roedder, E. (2008), 'Racial Cognition and the Ethics of Implicit Bias', *Philosophy Compass*, **3**(3), 522–40.
- Kleinberg, J., Ludwig, J., Mullainathan, S., and Sunstein, C. R. (2019), 'Discrimination in the Age of Algorithms', No. w25548, National Bureau of Economic Research.

- Knowles, J., Persico, N., and Todd, P. (2001), 'Racial Bias in Motor Vehicle Searches: Theory and Evidence', *Journal of Political Economy*, **109**(1), 203–29.
- Koenig, A. M. (2018), 'Comparing Prescriptive and Descriptive Gender Stereotypes about Children, Adults, and the Elderly', *Frontiers in Psychology*, **9**, 1086.
- Krings, F., Johnston, C., Binggeli, S., and Maggiori, C. (2014), 'Selective Incivility: Immigrant Groups Experience Subtle Workplace Discrimination at Different Rates', *Cultural Diversity and Ethnic Minority Psychology*, **20**(4), 491.
- Kuhn, P., and Villeval, M. C. (2015), 'Are Women More Attracted to Cooperation than Men?', *The Economic Journal*, **125**, 115–40.
- La Ferrara, E. (2019), 'Aspirations, Social Norms, and Development', *Journal of the European Economic Association*, **17**(6), 1687–722.
- Lemm, K. M., Lane, K. A., Sattler, D. N., Khan, S. R., and Nosek, B. A. (2008), 'Assessing Implicit Cognitions with a Paper-format Implicit Association Test', in M. A. Morrison and T. G. Morrison (eds), *The Psychology of Modern Prejudice*, Hauppauge, NY, Nova Science Publishers, 123–46.
- Lewis, G. B. (1986), 'Gender and Promotions: Promotion Chances of White Men and Women in Federal White-collar Employment', *Journal of Human Resources*, 406–19.
- Lichand, G., and Mani, A. (2020), 'Cognitive Droughts', University of Zürich Working Paper No. 341.
- Lieberman, M. D. (2013), *Social: Why our Brains are Wired to Connect*, Oxford, Oxford University Press.
- List, J. A. (2004), 'The Nature and Extent of Discrimination in the Marketplace: Evidence from the Field', *The Quarterly Journal of Economics*, **119**(1), 49–89.
- Madva, A. (2017), 'Implicit Bias, Moods, and Moral Responsibility', *Pacific Philosophical Quarterly*, **99**(S1), 53–78.
- Mani, A., Mullainathan, S., Shafir, E., and Zhao, J. (2013), 'Poverty Impedes Cognitive Function', *Science*, **341**(6149), 976–80.
- Marx, D. M., and Roman, J. S. (2002), 'Female Role Models: Protecting Women's Math Test Performance', *Personality and Social Psychology Bulletin*, **28**(9), 1183–93.
- Mazar, N., and Zhong, C. B. (2010), 'Do Green Products Make Us Better People?', *Psychological Science*, **21**(4), 494–8.
- Merritt, A., Effron, D., Fein, S., Savitsky, K., Tuller, D., and Monin, B. (2012), 'The Strategic Pursuit of Moral Credentials', *Journal of Experimental Social Psychology*, **48**(3), 774–7.
- Morris, E. W., and Perry, B. L. (2017), 'Girls Behaving Badly? Race, Gender, and Subjective Evaluation in the Discipline of African American Girls', *Sociology of Education*, **90**(2), 127–48.
- Mullen, E., and Monin, B. (2016), 'Consistency versus Licensing Effects of Past Moral Behaviour', *Annual Review of Psychology*, **67**, 363–85.
- Nelson, T. D. (2009), *Handbook of Prejudice, Stereotyping, and Discrimination*, Psychology Press.
- Oaxaca, R. (1973), 'Male–Female Wage Differentials in Urban Labor Markets', *International Economic Review*, 693–709.
- Oxoby, R. J. (2014), 'Social Inference and Occupational Choice: Type-based Beliefs in a Bayesian Model of Class Formation', *Journal of Behavioral and Experimental Economics*, **51**, 30–7.
- Phelps, E. (1972), 'The Statistical Theory of Racism and Sexism', *The American Economic Review*, **62**(4), 659–61.
- Piirto, J. (1991), 'Why are There so Few? (Creative Women: Visual Artists, Mathematicians, Musicians)', *Roeper Review*, **13**(3), 142–7.
- Porter, C., and Serra, D. (2020), 'Gender Differences in the Choice of Major: The Importance of Female Role Models', *American Economic Journal: Applied Economics*, **12**(6), 226–54.
- Prentice, D. A., and Carranza, E. (2002), 'What Women and Men Should Be, Shouldn't Be, are Allowed to Be, and Don't Have to Be: The Contents of Prescriptive Gender Stereotypes', *Psychology of Women Quarterly*, **26**(4), 269–81.
- Reuben, E., Sapienza, P., and Zingales, L. (2015), 'Taste for Competition and the Gender Gap Among Young Business Professionals', National Bureau of Economic Research Working Paper 21695.
- Rippon, G. (2019), *The Gendered Brain: The New Neuroscience that Shatters the Myth of the Female Brain*, Random House.

- Rodgers, W. M. (2009), *Handbook on the Economics of Discrimination*, Cheltenham, Edward Elgar.
- Schlesinger, A., O'Hara, K. P., and Taylor, A. S. (2018), 'Let's Talk About Race: Identity, Chatbots, and AI', in *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*, ACM, 315.
- Schmader, T. (2010), 'Stereotype Threat Deconstructed', *Current Directions in Psychological Science*, **19**(1), 14–18.
- Schmitt, D. P., Realo, A., Voracek, M., and Allik, J. (2008), 'Why Can't a Man be More Like a Woman? Sex Differences in Big Five Personality Traits across 55 Cultures', *Journal of Personality and Social Psychology*, **94**(1), 168.
- Schneider, D. (2004), *The Psychology of Stereotyping*, New York, Guilford Press.
- Schwarz, J. A. (2011), 'Essays on the Role of Gender and Groups in Economic Decision Making', unpublished doctoral dissertation, University of Pittsburgh.
- Shah, A. K., Mullainathan, S., and Shafir, E. (2012), 'Some Consequences of Having Too Little', *Science*, **338**(6107), 682–5.
- Shenouda, C. K., and Danovitch, J. H. (2014), 'Effects Of Gender Stereotypes and Stereotype Threat on Children's Performance on a Spatial Task', *Revue Internationale de Psychologie Sociale*, **27**(3), 53–77.
- Simbrunner, P., and Schlegelmilch, B. B. (2017), 'Moral Licensing: A Culture-moderated Meta-analysis', *Management Review Quarterly*, **67**, 201–25.
- Sprietsma, M. (2009), 'Discrimination in Grading? Experimental Evidence from Primary School', ZEW Discussion Paper 09–074.
- Sullivan, O., Gershuny, J., and Robinson, J. P. (2018), 'Stalled or Uneven Gender Revolution? A Long-term Processual Framework for Understanding Why Change Is Slow', *Journal of Family Theory & Review*, **10**(1), 263–79.
- Stiglitz, J. E. (1973), 'Approaches to the Economics of Discrimination', *The American Economic Review*, **63**(2), 287–95.
- Teal, C. R., Gill, A. C., Green, A. R., and Crandall, S. (2012), 'Helping Medical Learners Recognise and Manage Unconscious Bias Toward Certain Patient Groups', *Medical Education*, **46**(1), 80–8.
- Tversky, A., and Kahneman, D. (1983), 'Extensional versus Intuitive Reasoning: The Conjunction Fallacy in Probability Judgment', *Psychological Review*, **90**, 293–315.
- Whiteout, S. (2018), 'Popularizing Wokeness', *Harvard Journal of African American Public Policy*, **2017–18**, 63–70.
- Williams, W. M., and Ceci, S. J. (2015), 'National Hiring Experiments Reveal 2:1 Faculty Preference for Women on STEM Tenure Track', *Proceedings of the National Academy of Sciences*, **112**(17), 5360–5.
- Wood, M., Hales, J., Purdon, S., Sejersen, T., and Hayllar, O. (2009), 'A Test for Racial Discrimination in Recruitment Practice in British Cities', Department for Work and Pensions Research Report, 607.
- Wu, A. H. (2018), 'Gendered Language on the Economics Job Market Rumors Forum', *AEA Papers and Proceedings*, **108**, 175–9.
- Zetland, D., and Della Giusta, M. (2011), 'Focal Points, Gender Norms and Reciprocation in Public Good Games', No. em-dp2011-01, Henley Business School, Reading University.