

# A Behavioral Economics Approach to Discrimination

Fatas, Enrique <sup>1 2</sup>

<sup>1</sup> Behavioral Economics Institute

<sup>2</sup> Universidad Europea de Valencia

## TO CITE

Fatas, E. (2026). A Behavioral Economics Approach to Discrimination. In *Proceedings of the Paris Institute for Advanced Study* (Vol. 21). <https://paris.pias.science/article/a-behavioral-economics-approach-to-discrimination>

## DATE DE PUBLICATION

22/04/2026

## RÉSUMÉ

*How should we measure discrimination? Economics has long relied on the analysis of taste-based discrimination and statistical discrimination, using a well-established empirical toolkit centered on audit studies, correspondence experiments, and regression-based decompositions. These methods have produced landmark findings, yet they face persistent challenges: identifying the precise mechanisms behind differential treatment, capturing the multidimensional nature of real-world discrimination, and accounting for the behavioral complexities that shape both discriminators and those who are discriminated against. This paper reviews the traditional methods used to measure discrimination in economics and the social sciences, and argues that behavioral economics offers a powerful complementary toolkit. Drawing on insights about mental models, implicit biases, adaptive preferences, and social norms, behavioral approaches can illuminate the micro-foundations of discriminatory behavior in ways that standard methods cannot. The paper illustrates these arguments with examples from my own recent fieldwork, five policy-oriented behavioral studies conducted in Ecuador, Spain, Colombia, Pakistan, and Peru, to explain how intersectional discrimination can be studied in the field.*

**Keywords:** discrimination, behavioral economics, field experiments, survey experiments, beliefs-based, taste-based, social exclusion

## Acknowledgments:

I am grateful to the Paris Institute for Advanced Study for its generous fellowship and intellectual hospitality. I thank my co-authors on the various projects discussed here (Lina Restrepo-Plaza, Paulius Yamin, Luis Artavia-Mora, Héctor Solaz, Cristina Escamilla, Thomas Kruijer, Julian Pinazo, Ana M. Rojas, and Lorena Levano) for their contributions to the empirical work reviewed in this paper. Any mistake remains my own responsibility. Special gratitude goes to international organizations (International Finance Corporation, World Bank and the International Labour Organization),

governments (in Spain and Ecuador) and corporations (in all five countries). Without their support, the work described here would have never been possible.

# 1. Foundations

The economic analysis of discrimination begins with Becker's (1957) seminal model, in which some agents may hold a "taste" for discrimination, or a preference-based distaste for interacting with members of certain groups. Employers with such tastes behave as though hiring a worker from the disfavored group imposes an additional cost, leading them to demand a productivity premium or to avoid hiring such workers altogether. In competitive markets, Becker argued, discrimination should erode over time because non-discriminating firms enjoy a cost advantage. Yet this prediction has not been borne out empirically: discrimination persists across labor markets, credit markets, housing, and public services worldwide (Bertrand & Duflo, 2017; Guryan & Charles, 2013).

The resilience of taste-based discrimination has prompted important extensions. Customer discrimination arises when consumers prefer to interact with members of certain groups (many times their own ethnic, religious, or class groups), giving employers an incentive to cater to those preferences. Co-worker discrimination occurs when employees demand a wage premium to work alongside diverse and disfavored colleagues. These extensions explain why market competition alone may not eliminate discrimination: the "tastes" of customers and co-workers can sustain it even if employers themselves are indifferent (Charles & Guryan, 2008).

A fundamentally different account emerged with Phelps (1972) and Arrow (1973). In models of statistical discrimination, decision-makers use group membership as a proxy for unobserved individual characteristics. An employer who cannot perfectly observe a job applicant's productivity may rationally rely on group-level averages or perceived signal precision. This can generate differential treatment even in the absence of prejudices or animus. Critically, it can be self-reinforcing. If members of a group are statistically discriminated against, they may underinvest in human capital, confirming the very beliefs on which discrimination was built.

Recent theoretical work has expanded the statistical discrimination framework in important ways. Bohren et al. (2025) identify a fundamental problem: observed discrimination may reflect either accurate or inaccurate statistical discrimination, and

standard empirical methods cannot distinguish between the two. Campos-Mercade and Mengel (2024) show that non-Bayesian updating can generate systematic patterns of discrimination even when agents lack discriminatory preferences. Onuchic (2025) provides a comprehensive survey of these developments, including models of algorithmic fairness and learning traps.

Neither taste-based nor statistical discrimination models fully capture the behavioral complexity of real-world discrimination. Hoff and Walsh (2018) argue that behavioral economics provides a missing layer of analysis. They identify four key mechanisms through which social exclusion operates below the level of conscious deliberation:

*Mental models and categorical thinking.* People use simplified cognitive frameworks to process social information. These mental models activate automatically and shape expectations, interpretations, and decisions in ways that systematically disadvantage certain groups. In this framework, stereotypes are a canonical example because they reduce cognitive load but introduce systematic bias into judgments about individuals.

*Implicit discrimination.* Much discriminatory behavior operates outside conscious awareness. The large literature on implicit bias, anchored by the Implicit Association Test (Greenwald et al., 1998), demonstrates that individuals often hold negative associations with outgroups that they would explicitly disavow. These implicit biases predict real-world behavior in domains from healthcare to criminal justice (FitzGerald & Hurst, 2017).

*Self-stereotyping and identity threat.* Discrimination does not only affect how individuals are treated by others; it shapes how they see themselves. Members of stigmatized groups may internalize negative stereotypes, leading to reduced aspirations, undermined performance (through stereotype threat), and constrained identity expression. These vicious self-reinforcing dynamics mean that discrimination can persist even when external barriers are removed.

*Adaptive preferences and limited aspirations.* Perhaps the most insidious behavioral mechanism is the adaptation of preferences to constrained circumstances. When individuals face sustained exclusion, they may adjust their desires and expectations downward, coming to "prefer" outcomes that reflect their constrained choice set rather than their unconstrained potential. This process, emphasized in the capabilities approach of Sen (1999), makes discrimination partially invisible: its victims do not report dissatisfaction precisely because they have adapted to their diminished status.

These four mechanisms, and their interactions, suggest that discrimination is not merely a matter of tastes or information but a deeply embedded cognitive and social phenomenon. Measuring it, therefore, requires tools that go beyond revealed preference and self-report.

## 2. Methods

### 2.1. Audit and Correspondence Studies

The gold standard for detecting discrimination in market settings has been the field experiment. Audit studies, pioneered in the 1960s, send matched pairs of real individuals, identical in qualifications but differing in a protected characteristic such as race or gender, to apply for jobs, housing, or services. Differences in outcomes can then be attributed to discrimination because the only difference between participants is the protected characteristic associated with a gender, ethnic, or class identity. Correspondence studies refine this approach by sending fictitious applications (typically résumés or letters of inquiry) that manipulate group membership through names or other signals while holding all other characteristics constant.

Bertrand and Mullainathan's (2004) landmark study, which sent résumés with distinctively White- and African-American sounding names to employers in Boston and Chicago, exemplifies the power of this method. They found that White-sounding names received 50 percent more callbacks for interviews, holding every other relevant socio-economic characteristic of applicants constant, consistent with racial discrimination in hiring. Subsequent correspondence studies have documented discrimination on the basis of ethnicity, gender, age, disability, sexual orientation, and religion across dozens of countries (Bertrand & Duflo, 2017; Rich, 2014; Schaerer et al., 2023).

These methods have clear strengths: they observe actual behavior (not self-reports), they can establish causality through randomization, and they are conducted in real market settings. But they also have important limitations. First, they typically (but not always) test for discrimination along a single dimension, making it difficult to study how multiple identities interact. Second, they measure differential treatment at one point in a process (e.g., the callback stage in hiring) but cannot easily capture discrimination at other stages or across domains. Third, they reveal the existence and magnitude of

discrimination but not its underlying mechanism. The same pattern of differential callbacks could reflect taste-based discrimination, statistical discrimination, or implicit bias.

A complementary approach uses observational data and regression analysis to decompose outcome gaps between groups. The Blinder-Oaxaca decomposition (Blinder, 1973; Oaxaca, 1973), originally developed for wage gaps, partitions the difference in mean outcomes into an "explained" component (due to differences in observable characteristics) and an "unexplained" component (often interpreted as discrimination). Extensions of this method have been applied to employment, earnings, health, and education. While widely used, regression-based methods face a well-known limitation: the "unexplained" component captures not only discrimination but also all omitted variables correlated with group membership. This makes causal interpretation difficult and has led many researchers to prefer experimental methods for establishing the existence of discrimination, while using decomposition methods for descriptive analysis of disparities.

Survey measures of discrimination range from direct questions about experienced discrimination (e.g., the Everyday Discrimination Scale; Williams et al., 1997) to measures of attitudes toward outgroups (e.g., feeling thermometers, social distance scales). These instruments are valuable for capturing the subjective experience of discrimination and for tracking changes in expressed attitudes over time. However, they are subject to well-documented biases: social desirability may lead respondents to underreport discriminatory attitudes, while recall bias and framing effects can distort reports of experienced discrimination (Harnois et al., 2020).<sup>1</sup>

Taken together, these traditional methods have established beyond a reasonable doubt that discrimination exists across markets and domains worldwide. What they have been less successful at is identifying the behavioral mechanisms that produce discrimination, capturing the multidimensional and intersectional nature of discriminatory experiences, and measuring the ways in which discrimination shapes the beliefs, preferences, and behaviors of those who experience it. These are precisely the gaps that behavioral economics can help fill.

### 3. A Behavioral Economics Toolkit

Behavioral economics does not replace traditional methods; it complements them by providing tools that can probe the mechanisms, perceptions, and cognitive processes underlying discriminatory behavior. A distinctive contribution of behavioral economics is the use of incentivized games to classify individuals by their behavioral type—that is, by the underlying motivations that drive their decisions rather than by observed outcomes alone. In the study of discrimination, this approach draws on the literature on conditional cooperation (Fischbacher et al., 2001) and social preferences, which shows that individuals differ systematically in how they respond to others' behavior.

By using the strategy method, in which participants state their decisions for every possible action of their counterpart, researchers can construct a complete behavioral profile and classify individuals as conditional cooperators, free-riders, or other types. When combined with identity manipulation (varying the group membership of the counterpart), this method can reveal whether behavioral types shift across group boundaries and whether discrimination reflects a change in preferences, beliefs, or both.

Restrepo-Plaza and Fatas (2022) illustrate this approach in a lab-in-the-field experiment with victims, non-victims, and ex-combatants of the Colombian conflict. Using a public goods game with the strategy method, they show that behavioral types are remarkably stable across group identities: participants who cooperate conditionally with ingroup members also cooperate conditionally with outgroup members, including former adversaries. This finding challenges the assumption that intergroup discrimination necessarily reflects a shift in social preferences, suggesting instead that beliefs about others' behavior—rather than preferences—mediate group-based differential treatment.

Laboratory experiments offer greater control than the observational paradigms described in the previous section. The four behavioral methods described below go beyond laboratory experiments. Laboratory experiments, while extremely useful to identify mechanisms and test articulated theoretical models, typically do not consider real identities, but minimal groups, artificially created in the lab. <sup>2</sup>

*Survey vignette experiments* present respondents with short scenarios in which key attributes are randomly varied. Unlike traditional survey questions, which ask respondents directly about their attitudes or experiences, vignettes allow researchers to observe how variation in identity characteristics affects judgments, decisions, and attributions under controlled conditions. The factorial structure of these experiments,

where multiple attributes (e.g., gender, ethnicity, socioeconomic status) are independently randomized, enables causal identification of both main effects and interactions.

This approach has several advantages for measuring discrimination. First, it can manipulate multiple identities simultaneously, making it particularly well-suited for studying intersectional discrimination (Hainmueller et al., 2013). Second, it can measure not only treatment (how others respond to an individual) but also perceptions (how individuals expect to be treated), allowing researchers to study both sides of the discrimination equation. Third, it can be deployed at scale through surveys, making it cost-effective compared to audit studies. Fourth, because respondents do not realize which attributes are being tested, it partially mitigates social desirability bias.

*Double randomization* extends the logic of correspondence studies by combining identity randomization (varying the group membership of a target) with treatment randomization (varying the experimental condition or intervention). This two-layer design allows researchers to estimate not only the baseline level of discrimination but also how an intervention differentially affects treatment across groups.

In the context of my fieldwork in b (Fatas et al., 2026), this design was used to evaluate a behavioral training program for civil servants. As will see below, public officials were randomly assigned to treatment and control conditions, and within each condition, the identity of the service user described in survey vignettes was randomized along dimensions of nationality and gender. This double randomization made it possible to measure both the level of discrimination against Venezuelan migrants in access to public services and the causal effect of behavioral drivers on that discrimination.

The behavioral approach can also enhance traditional correspondence experiments by combining field-based measures of discrimination with survey-based measures of behavioral types, beliefs, and norms. The Ecuador study mentioned above included a correspondence experiment, in which fictitious emails from citizens with distinctively Ecuadorian or Venezuelan names were sent to public officials, paired with a survey experiment administered to the same population of officials. This pairing made it possible to correlate revealed discrimination in the field (differential email response rates) with stated attitudes, behavioral types, and reported social norms measured through the survey.

This combined approach addresses a persistent limitation of correspondence studies: the inability to identify mechanisms. By linking field behavior to individual-level survey

measures, researchers can test whether observed discrimination is associated with explicit prejudice, implicit bias, inaccurate beliefs about the target group, or social norms about appropriate treatment.

## 4. Intersectional Discrimination

Real-world discrimination rarely operates along a single dimension. A Latina immigrant may face discrimination based on her ethnicity, her gender, her immigration status, and (many times critically) on the interaction of these identities. Crenshaw's (1989) concept of intersectionality, originally developed in legal theory, captures this insight: the experience of someone who belongs to multiple marginalized groups is not simply the sum of the discrimination associated with each identity but may be qualitatively different (additive if linear, super-additive if the combined effect is above the sum of both individual components). Quantifying this intersectional component is both conceptually and methodologically challenging (Block et al., 2023; Else-Quest & Hyde, 2016).<sup>3</sup>

Standard approaches to intersectionality in quantitative research typically include interaction terms in regression models. While this can test whether the effect of one identity is moderated by another, it treats intersectionality as a statistical interaction and may miss the distinctive experiences of those at the intersection. Moreover, regression-based approaches face the challenge of small cell sizes when multiple identities are crossed, limiting statistical power.

In Fatas et al. (2025), we propose a novel methodology, the Intersectional Discrimination Dual Index (IDDI), that addresses these challenges. The IDDI measures intersectional discrimination from two complementary perspectives: the supply side (discriminatory attitudes and behaviors directed at individuals) and the demand side (the experiences and perceptions of those who are discriminated against). On the supply side, the IDDI uses randomized vignette experiments to measure how decision-makers respond to individuals who vary in multiple identity dimensions. Because identities are independently randomized, the design can estimate both main effects and interactions, providing a causal measure of intersectional discrimination. On the demand side, the IDDI surveys members of potentially discriminated groups about their perceived experiences, expectations, and beliefs regarding discrimination across identity dimensions.

Fatas, E. (2026). A Behavioral Economics Approach to Discrimination. In *Proceedings of the Paris Institute for Advanced Study* (Vol. 21). <https://paris.pias.science/article/a-behavioral-economics-approach-to-discrimination>  
2026/4 - paris-ias-ideas - Article No.2. Disponible <https://paris.pias.science/article/a-behavioral-economics-approach-to-discrimination> - ISSN 2826-2832/© 2026 Fatas E.  
This is an open access article published under the [Creative Commons Attribution-NonCommercial 4.0 International Public License \(CC BY-NC 4.0\)](https://creativecommons.org/licenses/by-nc/4.0/)

The dual nature of the index is crucial for the identification of social exclusion causes and drivers. As the behavioral economics literature on self-stereotyping and adaptive preferences makes clear, discrimination shapes not only how individuals are treated but how they perceive their own treatment and adjust their aspirations accordingly. A measure that captures only one side of this equation will understate the full impact of discrimination. Discrimination is one component of a broader phenomenon that Cuesta et al. (2024) term social exclusion. Their global study estimates that between 2.33 and 2.43 billion people worldwide are at risk of social exclusion, defined as systematic disadvantage across economic, political, and social domains. This framing situates discrimination within a larger set of processes that limit human capabilities and agency, resonating with both the capabilities approach (Sen, 1999) and the behavioral mechanisms identified by Hoff and Walsh (2018).

## 5. Behavioral Methods in the Field: Five Studies

This section illustrates the behavioral economics approach to measuring discrimination through five field studies I have conducted in collaboration with colleagues across five countries (Ecuador, Pakistan, Colombia, Spain, and Peru) and three continents. In each case, I focus on the methodological innovation, the way behavioral tools were deployed to measure discrimination, rather than on specific results, which are reported in the original studies and policy reports.

### Ecuador: Discrimination Against Migrants in Public Services

The Ecuador study (Fatas et al., 2026b) is the most comprehensive of the five, combining three methodological approaches in a single design. The study focused on discrimination against Venezuelan migrants in accessing public services, working with approximately 4,000 civil servants across three government ministries, before scaling up the intervention to all 300,000 civil servants in the country. While access to public services is a constitutional mandate in Ecuador, migrants face serious difficulties in accessing basic public services (like primary healthcare and education). The main goal of the study was to identify barriers and improve migrants' access to an online training course offered to civil servants. The online training was built from scratch, applying

Fatas, E. (2026). A Behavioral Economics Approach to Discrimination. In *Proceedings of the Paris Institute for Advanced Study* (Vol. 21). <https://paris.pias.science/article/a-behavioral-economics-approach-to-discrimination>  
2026/4 - paris-ias-ideas - Article No.2. Disponible <https://paris.pias.science/article/a-behavioral-economics-approach-to-discrimination> - ISSN 2826-2832/© 2026 Fatas E.  
This is an open access article published under the [Creative Commons Attribution-NonCommercial 4.0 International Public License \(CC BY-NC 4.0\)](https://creativecommons.org/licenses/by-nc/4.0/)

lessons from behavioral sciences, both in the contents included in the training course (e.g., on the characterization of biases and individual and group self-efficacy) and in the method used to deliver it (e.g., applying gamification and education tools).

*Design.* The study comprised three components. First, a baseline survey (N = 2,835) used double-randomized vignettes to measure discriminatory attitudes and behavioral types among civil servants. The vignettes varied the nationality and gender of a service user seeking assistance, while the behavioral classification used incentivized cooperation games with the strategy method. Second, a behavioral training intervention was randomly assigned to a subset of officials, and a follow-up survey (N = 1,515) measured changes in attitudes, beliefs, and behavioral types. Third, a correspondence experiment (N = 3,921 emails) sent fictitious service requests to public officials, manipulating the apparent nationality of the sender through name signals.

*Methodological contribution.* The three-component design addresses different layers of the discrimination problem. The double randomization of the survey vignettes (identity × treatment) allows causal identification of both discrimination and intervention effects. The behavioral type classification, based on the conditional cooperation types described above, goes beyond measuring discriminatory outcomes to characterize the motivational structure underlying differential treatment. The correspondence experiment provides a measure of revealed discrimination in actual service delivery, which can be linked back to the survey measures for mechanism identification. This triangulation of methods (including a survey experiment, a behavioral taxonomy, and a field experiment) represented a significant methodological advance over designs that rely on any single approach.

## Colombia: Gender Gaps in Rural Areas

We (Aguilar-Salamanca et al., 2026) investigated the role of discriminatory social norms and institutional trust in shaping the gender digital gap in financial technology adoption among rural women in the Cauca Valley, an area where women face the double challenge of being excluded from accessing basic financial services by the lack of bank branches (the rural area is heavily hit by drug smuggling violence) and by pervasive discriminatory norms supporting the idea that women should not use digital bank applications to improve their income. The goal of the study was to understand how two behavioral elements (discriminatory norms and low institutional trust in banks) prevented

women from installing bank apps on their cell phones, and from using them to improve their household income.

*Design.* Using small teams of trained female facilitators, all of them residing in the area, we run a survey experiment with 30% of all rural women, using a 2×2 factorial vignette design. Participants were presented with a scenario about using a mobile banking application, with two dimensions independently randomized: whether the scenario described discriminatory social norms (most/few residents hold a discriminatory belief discouraging women from using the app), and empirical expectations about trust in banks (most/few women in the area distrusted banks).

*Methodological contribution.* This design isolates the causal role of social norms and institutional trust, two factors that behavioral economics identifies as key drivers of exclusion but that are difficult to manipulate in traditional field experiments. Given that both drivers and discriminatory norms cannot be exogenously manipulated in a randomized control trial, vignette experiments allowed us to learn how residents in this rural area react to different counterfactuals. By experimentally varying the normative environment within a controlled vignette, we could estimate how much of the gender digital gap is attributable to discriminatory social norms (as opposed to, say, differences in digital literacy or access). The design draws on Bicchieri's (2006, 2017, 2021) framework, which distinguishes between empirical expectations (beliefs about what others do) and normative expectations (beliefs about what others think one should do). The vignette manipulation targeted the normative channel, testing whether exposure to discriminatory norms causally reduces women's willingness to adopt financial technology, with a characterization of treatment effects in isolation, or combined, and a taxonomy of heterogeneity effects (including the financial literacy and agency of women).

## Spain: Public Perceptions of Discrimination in Public Services

The Spain study (Fatas et al., 2026c) applied the behavioral toolkit described above by using the IDDI approach to understand both the supply and demand sides of discrimination towards migrants in Spain. The study combined two different methods to study how the general public in the Valencia department perceived discrimination in the delivery of social services (with a survey experiment administered to a large and

representative sample of residents), and how migrants perceived these very same service access (using focus group discussions with regularized and non-regularized migrants, and with members of other vulnerable groups acting as controls).

*Design.* Our survey experiment was administered to a large and representative sample of 1,700 residents, stratified by gender, income, education, and location, combining three methodological innovations. First, relative trust scales measured trust in members of different groups using an average person as a reference point (e.g., "How much do you trust a family member/scientist/politician/civil servant relative to an average Spaniard?"). As the relative trust scale (0-10) was divided into two parts, the "5" value was associated with a level of trust similar to the level of trust participants had in the average Spaniard. Any value above "5" was an indication of a higher level of trust (frequently the case when participants were asked about trust in family members). Any value below "5" could be interpreted as a sign of "relative" mistrust. In other words, we obtained information not only about different levels of trust in absolute terms (as the World Values Survey obtains with their 1-4 scale) but about the prevalence and intensity of (relative) mistrust. Second, the survey experiment randomized the vignette description of an interaction between civil servants (always neutrally presented) and service users with different origins (Latin America, Spain, Africa) in Vignette #1, and different levels of income (high versus low, keeping constant their (Spanish) origin) in Vignette #2. This allowed us to measure the general public's expectations about the prevalence of discriminatory behavior by civil servants. Third, objective knowledge elicitation asked respondents factual questions about inclusion policies and migrants (e.g., the proportion of policy budget allocated to migrants, the requirements to access benefits) to measure both their beliefs' accuracy and their knowledge of social services policies.

Our focus group discussions were open to both regularized migrants (with the right to work in Spain recently acquired) and non-regularized ones, coming from Latin America, Asia, and Africa. Their answers to the questions raised to moderators focused on their perception of being discriminated, the complexity of inclusion policies, their knowledge about how to access specific policy benefits, and their expectations of an equal and fair treatment by civil servants. Their answers were compared with the responses of members of vulnerable groups with the same socioeconomic status.

*Methodological contribution.* The combination of two methodologies, three instruments, and two samples reveals a "metacognitive illusion" that goes in opposite directions in both populations: a systematic disconnect between what people believe about

discrimination, what they know about the groups being discriminated against, who they expect to be discriminated, and how they actually evaluate encounters between civil servants and service users of different identities. Traditional surveys that simply ask about attitudes would miss this gap between stated beliefs and revealed evaluations. The inclusion of objective knowledge measures is particularly innovative, as it allows the study to test whether discriminatory attitudes are associated with inaccurate beliefs about the target group—a prediction of the statistical discrimination model that is rarely tested directly.

## Pakistan: Gender and Workplace Harassment

In collaboration with the International Labour Organization (United Nations), we studied gender-differentiated accountability attributions and workplace harassment exposure in the healthcare sector in Pakistan (see the policy report Fatas et al., 2026a). With a team of local collaborators, we run our study in six hospitals and clinics in Pakistan, aiming to document how three different groups (medical doctors, nurses, technicians) were exposed to three types of harassment (physical, violent, sexual), as exerted by patients and colleagues, and how they perceived the effectiveness of anti-harassment policies in the workplace.

*Design.* The study combined quantitative and qualitative methods, surveying 406 healthcare workers and conducting 8 focus group discussions. The key quantitative innovation was the combination of non-experimental relative scale questions (like the one on relative trust described above), and three randomized vignette experiments using a between subjects factorial design in which participants evaluated scenarios of workplace harassment with escalating severity. Critically, the gender of the victim and the perpetrator was randomized across respondents, while the harassment scenario in each vignette was held constant. This design allowed us to estimate whether identical incidents of harassment elicit different attributions of responsibility, severity, and appropriate institutional response depending on whether the victim is male or female. Moreover, the design also allowed us to estimate the tolerance to discrimination and harassment when the gender of the victim and the perpetrator was randomly manipulated.

*Methodological contribution.* The randomized assignment of victim and perpetrator gender within a constant harassment scenario is a clean application of the behavioral

toolkit to workplace harassment, an area where direct measurement is complicated by strong social desirability pressures and cultural sensitivity. The escalating severity design adds a further dimension: it tests not only whether gender affects overall attributions but whether gender's effect is moderated by the severity of the incident. This interaction is important for policy because it reveals whether gender bias is constant across the spectrum of harassment or concentrated at particular severity levels. The mixed-methods design, combining vignette experiments with focus group data, provides both causal estimates and contextual understanding; a triangulation that is particularly valuable in settings where cultural norms strongly shape both the experience and the reporting of harassment.

## Peru: Financial Inclusion of Migrants

Our Peru experiment (Artavia-Mora et al., 2026) examined discrimination against Venezuelan migrants in accessing financial services, conducted in partnership with a large commercial bank and the International Finance Corporation, part of the World Bank Group. The support of our partners was essential to run not only the qualitative study described below, but several semi-structured interviews and activities (including a behavioral variant of theatre of the oppressed, an interactive, community-based theatre methodology we used to help bank executives to critically examine discriminatory biases and outcomes using performance and improvisation). The Senior Executive Team participated in a two-day in-person training event run in Lima.

*Design.* The study combined the methodology described above and was run in person over three days in Lima, with a large-scale survey experiment with above 80% of all senior bank employees participating in it. In the experiment, we assigned executives randomly to one vignette condition (with a total of three vignettes) to measure discriminatory attitudes toward migrant customers, combined with the behavioral classification of executives' cooperative types also used in previous studies. The vignettes varied the nationality, gender, and socioeconomic cues of a customer seeking to open a bank account and get a loan. The design allowed us to estimate the prevalence and intensity of intersectional discrimination (by comparing their answers with different genders and origin of the vignette protagonist). As we also obtained individual information about the executives' expectations about the capacity of locals versus migrants, and men versus women, of having a large collateral, a sufficient

income, and a good financial literacy, we could also assess whether any discrimination could be explained by statistical discrimination.

*Methodological contribution.* This study applies the behavioral toolkit used in the studies described above in the financial sector, with a well-educated sample of bank executives, demonstrating its generalizability across domains and settings. The partnership with a commercial bank is noteworthy: by measuring discriminatory attitudes among the bank's own employees using behavioral games and randomized vignettes, the study provides the institution with actionable diagnostic information. As migrants in Peru have, on average, higher educational achievements than locals, they are younger and more mobile, the commercial bank understood that closing the door to their financial products to some clients could mean missing a business opportunity. As our study also measured how participants in the in-person event in Lima self-reported the prevalence of intensity of discrimination in the bank, we could document the large mismatch between the data obtained when asking senior employees about a sensitive matter (discrimination) and when their answers came from a randomized study. As predicted by social desirability theory, answers in the survey experiment revealed the extent of discrimination in a very different manner from self-reported responses. Given the small sample of in-person participants, we refrain from making any formal claim about this mismatch in the report.

## 6. Open Questions

The previous sections document that the measurement of discrimination is a fast-moving field. Several recent developments deserve attention, most notably the behavioral kit described in this paper, as the complexity of discrimination requires a careful choice when selecting the methodology used, but the combination of multiple methods to obtain rigorous answers. In this section, I point toward potential future directions in the field, for both theory and methodology.

### Prejudice

In a recent paper, Martin and Marx (2022) develop a clever identification strategy that distinguishes taste-based discrimination from statistical discrimination without requiring knowledge of the decision-maker's information set. Their "robust test of prejudice"

exploits the prediction that taste-based discrimination should persist even when private information would lead a rational statistical discriminator to treat members of different groups equally.<sup>4</sup> This approach provides a tighter test of prejudice than traditional correspondence studies, which can detect differential treatment but cannot rule out statistical motives. The adaptation of this "robust test of prejudice" to the behavioral kit described in this paper is a promising line of research.

## Learning traps

As discussed in the introduction, discrimination may become a self-fulfilling prophecy. Onuchic (2025) surveys a burgeoning literature on how rational learning processes can generate persistent discrimination. The key insight, developed formally by Bardhi et al. (2024), is that decision-makers who interact less frequently with members of certain groups accumulate information about them more slowly. This creates a "learning trap": initial disadvantage reduces the rate at which evidence is gathered, which in turn sustains the initial disadvantage.<sup>5</sup> This mechanism is distinct from both taste-based and standard statistical discrimination and has important implications for the design of anti-discrimination interventions: simply removing barriers may be insufficient if the learning trap has already depressed the information available to decision-makers. As learned in field work run mostly in Latin America (see Restrepo-Plaza & Fatas, 2022, 2023), being exposed or in contact with members of the other group may dramatically alter our expectations and beliefs about them. The design of feasible and actionable interventions based on contact theory may get vulnerable groups out of these learning traps.

## Algorithmic Fairness

The growing use of algorithms in hiring, lending, criminal justice, and other domains raises new questions about discrimination. Algorithms can encode and amplify human biases through biased training data, biased feature selection, or optimization criteria that correlate with protected characteristics. In the paper discussed above, Onuchic (2025) also reviews recent literature on algorithmic fairness, grappling with the fundamental tension between different definitions of fairness (e.g., demographic parity, equalized odds, individual fairness) and the impossibility of satisfying all of them simultaneously. From a behavioral perspective, the interaction between algorithmic recommendations and

human decision-makers is critical: research on "selective adherence" suggests that decision-makers may be more likely to follow algorithmic advice that confirms their prior biases (Rosenthal-von der Pütten & Sach, 2024). Dealing with confirmation bias with information conveyed effectively through education techniques (mixing education with entertainment) is a valuable reality, more than a distant promise. [6](#)

## What Works?

One of the main advantages of behavioral public policy, defined as the design (and evaluation) of public policies applying the lessons learned in behavioral sciences, is that it is deeply evidence based. In her comprehensive synthesis of evidence on discrimination reduction, Bartos (2025) provides a drawing on how both economics and psychology may help us to not only identify but also understand which interventions work better. While the evidence suggests that information interventions (correcting inaccurate beliefs about outgroups), contact interventions (facilitating positive intergroup interaction), and behavioral "nudges" (altering choice architecture to reduce the influence of bias) may all reduce discrimination under certain conditions, effect sizes are typically modest, and there is limited evidence on long-term persistence. The behavioral field studies described in Section 5 contribute to this evidence base by testing specific behavioral interventions in real institutional settings (with large treatment effects).

A particularly promising direction is the combination of diagnostic measurement and tailored intervention. By first using the behavioral toolkit to identify the specific mechanisms driving discrimination in a given context (e.g., inaccurate beliefs, discriminatory norms, implicit bias), policy experts can design interventions that target the relevant mechanism rather than applying one-size-fits-all approaches. The behavioral taxonomy used in the Ecuador, Spain, and Peru studies, based on conditional cooperation theory, for example, allows for the identification of subpopulations that may be particularly responsive to different types of interventions. While conditional cooperators may react to social norms nudges, free riders may react strongly to incentives or deterrence. As learning about the distribution of types in one group is a low-cost component of any intervention, adjusting interventions to specific behavioral types may improve an intervention's effectiveness.

## 7. Conclusion

The measurement of discrimination has come a long way since Becker's (1957) foundational model. Audit and correspondence studies have established the existence and magnitude of discrimination across markets and countries. Regression-based decompositions have documented the scope of disparities. Survey methods have tracked attitudes and subjective experiences. Each of these approaches has contributed essential evidence, and each has intrinsic limitations.

This paper has argued that behavioral economics offers a powerful complementary toolkit, one that can illuminate the micro-foundations of discrimination by probing the cognitive processes, social norms, and psychological mechanisms that sustain it. Randomized vignette experiments can manipulate multiple identity dimensions simultaneously, enabling the study of intersectional discrimination and counterfactuals. Behavioral type classification can distinguish between preference-based and belief-based drivers of differential treatment and boost the effectiveness of a policy intervention. Double randomization designs can rigorously estimate the causal effects of interventions on discrimination, mitigating social desirability bias in sensitive domains. And the integration of survey experiments with correspondence studies can link revealed discrimination in the field to individual-level psychological and behavioral measures.

The five field studies presented here illustrate these methods in action, applied to discrimination against migrants in public services (Ecuador) and financial inclusion (Peru), inaccurate beliefs about social policies recipients and inaccurate beliefs about discrimination (Spain), gender-based exclusion from digital finance (Colombia), and workplace harassment (Pakistan). In each case, the behavioral approach reveals aspects of discrimination that traditional methods would miss: the motivational structure underlying differential treatment, the role of social norms in shaping exclusion, the gap between perceived and actual discrimination, and the way gender moderates attributions of responsibility.

Several challenges remain. First, external validity: vignette experiments, however carefully designed, present hypothetical scenarios, and the relationship between vignette responses and real-world behavior requires ongoing validation. Second, measurement of adaptive preferences remains difficult, precisely because those whose preferences have adapted may not report dissatisfaction. Third, the integration of behavioral and structural approaches to discrimination is still in its early stages. Future work should explore how

behavioral micro-foundations interact with institutional and market-level forces to produce systemic discrimination (Bohren et al., 2022).

Finally, the study of intersectional discrimination, where the behavioral approach has perhaps the most to contribute, is still developing the methodological infrastructure it needs. The IDDI represents one step in this direction, but much work remains to be done to develop scalable measures that capture the full complexity of overlapping identities and their interactions with institutional contexts. As Crenshaw (1989) observed more than three decades ago, the people at the intersection are often the least visible. Making them visible to measurement is both a methodological challenge and a moral imperative.

# Bibliographie

Aguilar-Salamanca, C., Fatas, E., & Zuluaga, B. (2026). *Breaking Discriminatory Norms, Building Trust: A Report on Gender Gaps on Financial Inclusion* [BEI Policy Report 2026/PR02,]. Behavioral Economics Institute.

Arrow, K. J. (1973). The theory of discrimination. In O. Ashenfelter & A. Rees (Eds.), *Discrimination in labor markets* (pp. 3–33). Princeton University Press.

Artavia-Mora, L., Fatas, E., Levano, L., & Yamin, P. (2026). *Widespread, but Invisible: A Report on the Financial Exclusion of Migrants in Peru* [BEI Policy Report 2026/PR01,]. Behavioral Economics Institute.

Bardhi, A., Guo, Y., & Strulovici, B. (2024). *Early-career discrimination: Spiraling or self-correcting?* Duke University.

Bartos, V. (2025). Breaking bias: Pathways to reducing discrimination. *CESifo Working Paper*, 11558.

Becker, G. S. (1957). *The economics of discrimination*. University of Chicago Press.

Bertrand, M., & 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.

Bertrand, M., & Duflo, E. (2017). Field experiments on discrimination. In A. V. Banerjee & E. Duflo (Eds.), *Handbook of economic field experiments* (Vol. 1, pp. 309–393). North-Holland.

Bicchieri, C. (2006). *The grammar of society: The nature and dynamics of social norms*. Cambridge University Press.

Bicchieri, C. (2017). *Norms in the wild: How to diagnose, measure, and change social norms*. Oxford University Press.

Bicchieri, C., Fatas, E., Aldama, A., Casas, A., Deshpande, I., Lauro, M., Parilli, C., Spohn, M., Pereira, P., & Wen, R. (2021). In science we (should) trust: Expectations and compliance across nine countries during the COVID-19 pandemic. *PLOS ONE*, 16(6), 0252892.

Blinder, A. S. (1973). Wage discrimination: Reduced form and structural estimates. *Journal of Human Resources*, 8(4), 436–455.

- Block, R., Golder, M., & Golder, S. (2023). Evaluating claims of intersectionality. *The Journal of Politics*, 85(3), 795–808.
- Bohren, J. A., Hull, P., & Imas, A. (2022). Systemic discrimination: Theory and measurement. *National Bureau of Economic Research Working Paper*, 29820.
- Bohren, J. A., Haggag, K., Imas, A., & Pope, D. G. (2025). Inaccurate statistical discrimination: An identification problem. *Review of Economics and Statistics*, 107(3), 605–620.
- Borzino, N., Fatas, E., & Peterle, E. (2023). In transparency we trust an experimental study of reputation, transparency, and signaling. *Journal of Behavioral and Experimental Economics*, 106, 102061.
- Campos-Mercade, P., & Mengel, F. (2024). Non-Bayesian statistical discrimination. *Management Science*, 70(4), 2549–2566.
- Charles, K. K., & Guryan, J. (2008). Prejudice and wages: An empirical assessment of Becker's The Economics of Discrimination. *Journal of Political Economy*, 116(5), 773–809.
- Crenshaw, K. (1989). Demarginalizing the intersection of race and sex: A Black feminist critique of antidiscrimination doctrine, feminist theory and antiracist politics. *University of Chicago Legal Forum*, 1989(1), 139–167.
- Cuesta, J., González-Espinosa, A. C., Hoyos, L. F., Mayorga, J., & Pizano, L. (2024). Social exclusion concepts, measurement, and a global estimate. *PLOS ONE*, 19(2), 0298085.
- Eckel, C. C., Fatas, E., & Kass, M. (2022). Sacrifice: An experiment on the political economy of extreme intergroup punishment. *Journal of Economic Psychology*, 90, 102486.
- Else-Quest, N. M., & Hyde, J. S. (2016). Intersectionality in quantitative psychological research: II. *Methods and Techniques. Psychology of Women Quarterly*, 40(3), 319–336.
- Fatas, E., Nosenzo, D., Sefton, M., & Zizzo, D. J. (2021). A self-funding reward mechanism for tax compliance. *Journal of Economic Psychology*, 86, 102421.
- Fatas, E., Restrepo-Plaza, L. M., Jiménez, N., & Rincón, G. (2021). The behavioral consequences of conflict exposure on risk preferences. In *Oxford Research Encyclopedia of Economics and Finance*. Oxford University Press.
- Fatas, E., & Restrepo-Plaza, L. (2022). When losses can be a gain: A large lab-in-the-field experiment on reference dependent forgiveness in Colombia. *Journal of Economic Psychology*, 88, 102463.
- Fatas, E., & Artavia-Mora, L. (2023). The integration of migrants in modern societies. In L. Artavia-Mora & Z. Khan (Eds.), *Building Behavioral Science in International Development*. Bescy Publisher.

Fatas, E., & Restrepo-Plaza, L. M. (2024). Inequality as a behavioral driver: An inspiring contribution to behavioral political economy by Shaun Hargreaves-Heap. *Review of Behavioral Economics*, 11(2), 235–254.

Fatas, E., Restrepo-Plaza, L., & Banuri, S. (2024). A simple twist of fate: An experiment on election uncertainty and democratic institutions. *Journal of Economic Behavior & Organization*, 228, 106752.

Fatas, E., & Restrepo-Plaza, L. M. (2025). Behaviour and violent conflict exposure. In S.-H. Chuah, R. Hoffmann, & A. Neelim (Eds.), *Elgar Encyclopedia of Behavioural and Experimental Economics*. Edward Elgar Publishing.

Fatas, E., Yamin, P., & Artavia-Mora, L. (2025). Intersectional discrimination. In E. Hoon, C. Swee, R. H., & A. N (Eds.), *Elgar Encyclopedia of Behavioural and Experimental Economics*. Edward Elgar Publishing.

Fatas, E., Restrepo-Plaza, L., & Yamin, P. (2026). *Workplace Violence and Harassment in Pakistan's Healthcare Sector: A Diagnostic Tool Report* [BEI Policy Report 2026/PR04,]. Behavioral Economics Institute.

Fatas, E., Restrepo-Plaza, L., Rojas, A., & Yamin, P. (2026). *Mitigating Discrimination Towards Migrants in Public Services: A Report on a Behavioral Intervention with Civil Servants in Ecuador* [BEI Policy Report 2026/PR05,]. Behavioral Economics Institute.

Fatas, E., Restrepo-Plaza, L., Solaz, H., & Pinazo, J. (2026). *The Behavioral Limits of Inclusion Policies: A Report on a Mixed Methods Study A* [BEI Policy Report 2026/PR03,]. Behavioral Economics Institute.

Fischbacher, U., Gächter, S., & Fehr, E. (2001). Are people conditionally cooperative? Evidence from a public goods experiment. *Economics Letters*, 71(3), 397–404.

FitzGerald, C., & Hurst, S. (2017). Implicit bias in healthcare professionals: A systematic review. *BMC Medical Ethics*, 18(1), 19.

Greenwald, A. G., McGhee, D. E., & Schwartz, J. L. K. (1998). Measuring individual differences in implicit cognition: The Implicit Association Test. *Journal of Personality and Social Psychology*, 74(6), 1464–1480.

Guryan, J., & Charles, K. K. (2013). Taste-based or statistical discrimination: The economics of discrimination returns to its roots. *Economic Journal*, 123(572), 417–432.

Hainmueller, J., Hopkins, D. J., & Yamamoto, T. (2013). Causal inference in conjoint analysis: Understanding multidimensional choices via stated preference experiments. *Political Analysis*, 22(1), 1–30.

Harnois, C. E., Bastos, J. L., & Shariff-Marco, S. (2020). Intersectionality, contextual specificity, and everyday discrimination: Assessing the difficulty associated with identifying a main reason for discrimination among racial/ethnic minority respondents. *Sociological Methods & Research*, 51(3), 983–1015.

- Hoff, K., & Walsh, J. S. (2018). The whys of social exclusion: Insights from behavioral economics. *World Bank Research Observer*, 33(1), 1–33.
- Martin, D., & Marx, P. (2022). A robust test of prejudice for discrimination experiments. *Management Science*, 68(6), 4527–4536.
- Oaxaca, R. (1973). Male–female wage differentials in urban labor markets. *International Economic Review*, 14(3), 693–709.
- Onuchic, P. (2025). *Recent contributions to theories of discrimination*.
- Oswald, F. L., Mitchell, G., Blanton, H., Jaccard, J., & Tetlock, P. E. (2013). Predicting ethnic and racial discrimination: A meta-analysis of IAT criterion studies. *Journal of Personality and Social Psychology*, 105(2), 171–192.
- Phelps, E. S. (1972). The statistical theory of racism and sexism. *American Economic Review*, 62(4), 659–661.
- Pütten, A. M., & Sach, A. (2024). Michael is better than Mehmet: exploring the perils of algorithmic biases and selective adherence to advice from automated decision support systems in hiring. *Frontiers in Psychology*, 15, 1416504.
- Restrepo-Plaza, L., & Fatas, E. (2022). When ingroup favoritism is not the social norm a lab-in-the-field experiment with victims and non-victims of conflict in Colombia. *Journal of Economic Behavior and Organization*, 194, 363–383.
- Restrepo-Plaza, L., & Fatas, E. (2023). Building inclusive institutions in polarized scenarios: Experimental evidence for the peace process and the same sex marriage debate in Colombia. *Constitutional Political Economy*, 34(1), 88–110.
- Rich, J. (2014). What do field experiments of discrimination in markets tell us? A meta-analysis of studies conducted since 2000. *IZA Discussion Paper*, 8584.
- Schaerer, M., Plessis, C., Nguyen, M. H. B., Aert, R. C. M., Tiokhin, L., Lakens, D., Clemente, E. G., Pfeiffer, T., Dreber, A., Johannesson, M., Clark, C. J., & Uhlmann, E. L. (2023). On the trajectory of discrimination: A meta-analysis and forecasting survey capturing 44 years of field experiments on gender and hiring decisions. *Organizational Behavior and Human Decision Processes*, 179, 104280.
- Sen, A. (1999). *Development as freedom*. Alfred A. Knopf.
- Williams, D. R., Yu, Y., Jackson, J. S., & Anderson, N. B. (1997). Racial differences in physical and mental health: Socio-economic status, stress and discrimination. *Journal of Health Psychology*, 2(3), 335–351.

# Notes de bas de page

**1** : The Implicit Association Test (IAT), developed by Greenwald, McGhee, and Schwartz (1998), measures differential association between concepts and attributes. Its predictive validity for discriminatory behavior remains debated; see Oswald et al. (2013) and Greenwald et al. (2009).[↵](#)

**2** : See Borzino et al (2023) for network identities, Eckel et al (2022) for intergroup conflict induced by artificially created identities, Fatas et al (2021) for laboratory participants acting as taxpayers, or Fatas et al (2024) as an example of laboratory methods used to create different political actors identities.[↵](#)

**3** : For a detailed treatment of the IDDI methodology, see Fatas, Yamin, and Artavia-Mora (2025).[↵](#)

**4** : Martin and Marx (2022) develop a test that identifies prejudice regardless of private learning, solving the identification problem noted by Bohren, Haggag, Imas, and Pope (2025).[↵](#)

**5** : Bardhi, Guo, and Strulovici (2024) show how rational agents can fall into learning traps where discrimination persists because decision-makers never accumulate enough evidence to update their beliefs.[↵](#)

**6** : As our recent work in Latin America and Asia, using mini-podcasts, radio shows, and social media, shows. Unfortunately, there is not enough space to present all these interventions in this short paper.[↵](#)