

# Measuring and Mitigating Racial Disparities in Tax Audits \*

Hadi Elzayn      Evelyn Smith      Thomas Hertz      Cameron Guage  
Arun Ramesh      Robin Fisher      Daniel E. Ho      Jacob Goldin

July 19, 2024

## Abstract

Tax authorities around the world rely on audits to detect underreported tax liabilities and to verify that taxpayers qualify for the benefits they claim. We study differences in Internal Revenue Service audit rates between Black and non-Black taxpayers. Because neither we nor the IRS observe taxpayer race, we propose and employ a novel partial identification strategy to estimate these differences. Despite race-blind audit selection, we find that Black taxpayers are audited at 2.9 to 4.7 times the rate of non-Black taxpayers. An important driver of the disparity is differing audit rates by race among taxpayers claiming the Earned Income Tax Credit (EITC). Using counterfactual audit selection models to explore why the disparity arises, we find that maximizing the detection of underreported taxes would not lead to Black EITC claimants being audited at higher rates. Rather, the audit disparity among EITC claimants stems in large part from a policy decision to prioritize detecting overclaims of refundable credits over other forms of noncompliance. Modifying the audit selection algorithm to target total underreported taxes while holding fixed the number of audited EITC claimants would reduce the share of audited taxpayers who are Black, and would lead to more audits focused on accurate reporting of business income and deductions; fewer audits focused on the eligibility of claimed dependents; higher per-audit costs; and more detected noncompliance.

---

\*The views presented here are those of the authors and do not necessarily represent the position of the Treasury Department. The Treasury Department reviewed these results to ensure compliance with statutory prohibitions on the disclosure of taxpayer information and with policy prohibitions on the disclosure of confidential information related to IRS audit processes. For helpful comments, we thank Emily Black, Edith Brashares, Dorothy Brown, Jonathan Choi, Geoffrey Gee, Robert Gillette, Alissa Graff, John Guyton, Anne Herlache, Jim Hines, Janet Holtzblatt, Hilary Hoynes, Tatiana Homonoff, Barry Johnson, Larry Katz, Elaine Maag, Michael Morse, Julian Nyarko, Nina Olson, Derek Ouyang, Claire Lazar Reich, Kit Rodolfa, Dan Rosenbaum, Joel Slemrod, David Splinter, Megan Stevenson, Alex Turk, Caroline Weber, the staff members of the Joint Committee on Taxation who contributed to a consolidated set of comments, and seminar participants. We are grateful for financial support from the Hoffman-Yee Research Grant for Stanford's Institute for Human-Centered Artificial Intelligence and from Arnold Ventures.

Elzayn: First author; Stanford University. Smith: University of Michigan. Hertz: Internal Revenue Service. Guage: Stanford. Ramesh: UCLA. Fisher: U.S. Treasury Department. Ho: Equal co-supervising author; Stanford University. Goldin: Equal co-supervising author; University of Chicago, American Bar Foundation, and NBER; Corresponding author: jsgoldin@uchicago.edu.

# 1 Introduction

In recent decades, U.S. policymakers have increasingly relied on the income tax system to implement a host of social programs; for example, the Earned Income Tax Credit (EITC) has replaced welfare as the largest cash-based safety net program in the United States. The Internal Revenue Service (IRS) administers these programs and is also charged with ensuring that individuals meet their taxpaying obligations. Like tax authorities around the world, it relies on audits to detect underreporting of tax liabilities and to verify that taxpayers qualify for the benefits they claim.

The task of selecting which taxpayers to audit is partly a prediction exercise—which taxpayers have underreported tax obligations that an audit would uncover?—and partly a policy determination—what type of underreporting should be predicted and pursued? Evidence from many domains in which data-driven algorithms are used to allocate enforcement resources suggests that both aspects of this process may inadvertently reinforce disadvantages against historically marginalized groups (Angwin et al., 2016; Buolamwini and Gebru, 2018; Obermeyer et al., 2019). Such concerns are particularly acute for tax audits, which can exacerbate financial strain for the lowest income taxpayers – whose tax refunds are typically frozen while an audit is in place – and can dissuade individuals from participating in safety net programs for which they qualify (Guyton et al., 2018; National Taxpayer Advocate, 2019a).

In this paper, we investigate racial disparities in the selection of taxpayers for audit, focusing on differences in the selection of Black and non-Black taxpayers. Because the IRS does not collect information about taxpayers’ race, identifying differences in audit rates by race is itself a significant challenge. Researchers have developed a range of tools for imputing race from other observed characteristics, but such approaches can lead to biased estimates for the parameters of interest unless restrictive assumptions are satisfied (Chen et al., 2019; Knox et al., 2022). A second methodological challenge is that even if we could observe taxpayer race, the fact that underreporting is observed only for those returns that were

selected for audit (the so-called selective labels problem) (Kleinberg et al., 2018a) makes it difficult to understand why disparities emerge or which policies would lessen them.

We investigate this topic using comprehensive administrative microdata on approximately 148 million tax returns and 780,000 audits.<sup>1</sup> To circumvent the selective labels problem, we also leverage nearly 72,000 audits of randomly selected taxpayers to investigate the effects of counterfactual audit selection policies. To address the problem of missing race, we use Bayesian Improved First Name and Surname Geocoding (BIFSG) to impute race based on taxpayers’ names and census block groups (Imai and Khanna, 2016; Voicu, 2018). We then propose and implement a novel approach for bounding the true audit disparity by race from the (imperfectly measured) BIFSG proxy. By individually matching a subset of the tax data to self-identified race data from other administrative sources, we provide evidence that the assumptions underlying our bounding approach are satisfied in practice.

We report a number of new results about racial disparities in tax audits. First, we estimate that the audit rate for returns filed by Black taxpayers is between 0.81 and 1.34 percentage points higher than the audit rate for non-Black taxpayers. This (unconditional) disparity is substantial when compared to the base audit rate of 0.54% for the overall U.S. population. In relative terms, our estimates imply that Black taxpayers are audited at between 2.9 to 4.7 times the rate of non-Black taxpayers.

Second, we find that an important contributor to the racial audit disparity is a difference in audit rates among Black and non-Black taxpayers who claim the EITC. Others have speculated that Black taxpayers may be audited at higher rates because they are more likely than non-Black taxpayers to claim the EITC, and EITC claimants (of any race) are audited at higher rates than most other taxpayers (Bloomquist, 2019; Kiel and Fresques, 2019). However, we find that this channel does not fully explain the disparity we observe. Rather, we estimate a substantial disparity in the selection of taxpayers for audit *within* the population of EITC claimants: Black taxpayers claiming the EITC are between 2.9 and

---

<sup>1</sup>We primarily focus on tax year 2014 – the most recent year for which audit outcome data was complete at the time of our analysis. We find similar results for other years spanning the 2010–2018 time period.

4.4 times as likely to be audited as non-Black EITC claimants. In contrast, we observe a much smaller, though still statistically significant, difference in audit rates between Black and non-Black taxpayers who do not claim the EITC.

Third, we explore the factors contributing to the observed audit disparity among taxpayers claiming the EITC. Because EITC audit selection is largely automated, and because the IRS does not observe race, the disparity is unlikely to be driven by disparate treatment of Black and non-Black taxpayers.<sup>2</sup> At the same time, we find that the disparity cannot be explained by group-level differences in total dollars of tax underreporting: we estimate that Black EITC claimants are audited at higher rates than non-Black EITC claimants within each decile of under-reported taxes.

To better understand the source of the audit disparity in light of this finding, we simulate counterfactual audit selection algorithms for EITC claimants, using data from audits of randomly selected tax returns. We explore potential explanations related to both aspects of the audit selection process described above: predicting which taxpayers have underreported taxes that an audit would uncover and the policy decision concerning which types of underreporting to predict and pursue (the algorithmic objective).

Beginning with the algorithmic objective, recall that taxpayers may underreport taxes by misreporting their taxable income and/or by overstating their eligibility for tax-administered benefits, often in the form of refundable credits. Although details of the IRS audit selection algorithm are confidential, the agency publicly reports that the model underlying its primary program for auditing EITC claimants is designed to detect noncompliance from this latter source of underreporting — i.e., overclaimed refundable credits. Using our simulated audit algorithms, we find that when audits of EITC claimants are allocated based on this objective, Black taxpayers are audited at higher rates than non-Black taxpayers. In contrast, if audits were allocated based on the objective of maximizing total dollars of detected underreporting — from any source of noncompliance — we find that Black EITC claimants would be audited

---

<sup>2</sup>By disparate treatment, we mean different audit rates for taxpayers who file identical tax returns and differ only with respect to their race.



at *lower* rates than non-Black EITC claimants.

Why does the algorithmic objective shape the distribution of audits by race? Among taxpayers claiming the EITC, the highest underreporting is associated with taxpayers who have substantial independent contractor or other business income, but who underreport this income to such a degree that they appear to qualify for the EITC. Among this group, we find that those with the most underreporting are disproportionately non-Black. Despite their high levels of noncompliance, audit algorithms focused exclusively on refundable credits underprioritize these taxpayers because only a portion of their underreporting takes the form of overclaimed credits – the rest stems from underreported taxable income. At the same time, we find that EITC returns filed by Black taxpayers appear more likely to violate the legal requirements for claiming children, which translates into larger average dollar adjustments to the refundable credits claimed on their returns. As a result, audit algorithms that are exclusively focused on this form of noncompliance select Black taxpayers at higher rates.

With any potential change to the audit selection algorithm, it is important to consider the consequences for the operational process through which audits are conducted. Along these lines, we show that the objective of the audit selection algorithm shapes the composition of audited issues, with substantial downstream implications for audit costs and detected underreporting. In this way, our results highlight an important connection between racial audit disparities and ongoing policy debates about the proper level of IRS funding (e.g., Boning et al., 2023).

Distinct from the objective of the audit selection algorithm, we also explore whether the architecture of the prediction model employed as part of the audit selection process contributes to the observed disparity. We find evidence that it does: although Black EITC claimants tend to overclaim refundable credits at higher rates, the observed audit disparity is larger than what these differences in overclaiming would imply. The confidential nature of the IRS audit selection process limits our ability in this paper to explore the details of the actual predictive model used by IRS to allocate EITC audits; however, we provide

suggestive evidence that differences by race in the predictive power of the available features—such as proxies for child custody or parental status—may contribute to higher audit rates for Black taxpayers. Subject to the same caveat, we provide suggestive evidence that it may be possible to modify the predictive model to lessen the racial audit disparity without substantially reducing the amount of overclaimed refundable credits detected on audit.

Our results contribute to an empirical literature that studies the distributional effects of tax policy by race, with seminal contributions by Moran and Whitford (1996) and Brown (2021), among others.<sup>3</sup> However, the unavailability of administrative microdata with information about race has limited the questions that prior studies could address and the alternative policies they could evaluate (Bearer-Friend, 2019; Brown, 2021; Dean, 2021).<sup>4</sup> We build on this literature by linking imputed race estimates with administrative data on tax returns and audits. Doing so allows us to produce the most direct evidence to date on longstanding questions about racial disparities in the administration of the U.S. income tax.

A second contribution of our paper is to introduce a novel partial identification approach for conducting algorithmic disparity assessments with respect to a protected class, when that class is unobserved to the researcher. This challenge arises in a wide range of settings, including voting rights (Imai and Khanna, 2016), regulatory policy (CFPB, 2014; Anson-Dwamena et al., 2021; Haas et al., 2019), and industry (Alao et al., 2021; Andrus et al., 2021). For example, although many U.S. agencies are required by federal law to conduct disparity assessments of the algorithmic decision tools they employ, protected characteristics like race are often missing from administrative records (G.A.O., 2020; Exec. Order. 13985, 2021).

---

<sup>3</sup>Much of this literature focuses on racial differences in the benefits from substantive tax provisions like the mortgage interest deduction (Brown, 2009, 2018), the EITC (Brown, 2005; Hardy et al., 2021), and the Child Tax Credit (Collyer et al., 2019; Goldin and Micheltore, 2022).

<sup>4</sup>Prior analyses have yielded some suggestive evidence, however. For example, Bloomquist (2019) studies regional bias in IRS audits using estimated county-level data, and notes that the ten most heavily audited counties were predominantly comprised of Black taxpayers, whereas the ten least heavily audited counties were predominantly non-Black. In addition, prior research has linked tax data with Census records on self-reported race to study primarily non-tax outcomes (e.g., Chetty et al., 2020); however, various legal and institutional constraints currently limit our ability to apply this approach to our setting.

Our approach relies on weaker conditions than those required to point-identify differences in outcomes by unobserved protected class (e.g., Chen et al., 2019; Fong and Tyler, 2021), while still yielding bounds that may be informative for policy.<sup>5</sup>

Finally, our results relate to a growing literature studying how the choice of outcome to be predicted by an algorithm shapes the distributional properties of procedures based upon that algorithm (Barocas and Selbst, 2016; Kleinberg et al., 2018b; Passi and Barocas, 2019). For example, Obermeyer et al. (2019) link racial disparities in a health care setting to an algorithm that is trained to predict health care expenditures rather than direct health outcomes. With respect to tax audits, Black et al. (2022) study how the choice between regression and classification prediction tasks shapes the distribution of audits by taxpayer income. Related to this literature, our results highlight how policy decisions to prioritize certain forms of legal noncompliance over others can shift the distribution of enforcement burdens.

The paper proceeds as follows. Section 2 provides background on the U.S. tax system and taxpayer audits. Section 3 describes our empirical strategy. Section 4 describes our data. Section 5 provides results relating to estimated race probabilities and statistical bias of our proposed estimators. Section 6 estimates differences in audit rates between Black and non-Black taxpayers. Section 7 investigates the source of the observed audit disparity. An Online Appendix contains proofs and additional results.

## 2 Institutional Background

This section provides background regarding the U.S. individual income tax, taxpayer audits, and the EITC.

---

<sup>5</sup>Kallus et al. (2021) also consider a partial identification approach to estimating disparity when the protected characteristic must be imputed. Unlike our approach, their bounds cover all joint distributions consistent with the observed marginals. In our setting, the Kallus et al. bounds are largely uninformative as to the magnitude or even direction of the audit rate disparity; they cannot rule out Black taxpayers facing either higher or lower audit rates than non-Black taxpayers. Our approach requires additional structure, but the payoff to that structure is a significantly more informative estimate when our assumptions hold.

## **2.1 The U.S. Income Tax**

Most U.S. citizens, as well as some non-citizens, are required to file an income tax return each year, on which they calculate and report their tax liability based on their income as well as any deductions or credits for which they qualify. Unmarried taxpayers file individual returns, whereas most taxpayers who are married file a joint return with their spouse. Taxpayers with children or other dependents may claim them on their own return to qualify for various tax benefits. The vast majority (over 95%) of taxpayers prepare and file their returns with the help of a professional tax preparer or using guided tax preparation software.

## **2.2 IRS Audits**

The IRS is the federal agency responsible for promoting and enforcing compliance with the tax law. One channel through which it does so is by employing taxpayer audits. Taxpayers selected for audit are required to provide additional information to the IRS or otherwise verify the accuracy of the tax liability or refund reported on their tax return. Audits may occur by mail (“correspondence examinations”) or through in-person (or virtual) meetings with IRS employees (“field” or “office” examinations). In recent years, approximately 70% of audits have been conducted through correspondence. Audits of this form tend to focus on a small number of issues and require a response, with substantiation. If the IRS does not receive a response by the due date, it will generally disallow the claimed item and issue the taxpayer with a notice of deficiency. Correspondence audits are substantially cheaper than other forms of audits for the IRS to conduct. At the same time, correspondence audits can be particularly burdensome for lower-income households, who may face additional barriers to understanding the audit notice, acquiring the required documents, or obtaining expert assistance (G.A.O., 2016; National Taxpayer Advocate, 2021).

If an audit results in an adjustment to the (net) tax reported on a taxpayer’s return, the taxpayer is responsible for remitting the difference to the IRS and may face additional penalties, as well as, in rare cases, criminal sanctions. Audits may occur pre- or post-

refund; in the former case, refunds are not issued until after the audit is resolved. Hence, taxpayers who fail to respond to a pre-refund audit typically forego the tax benefits they claimed on their return. Taxpayers who disagree with the results of an audit may appeal the determination with the IRS office of appeals and/or in federal court.

At a high level, audits can be categorized into two groups: research and operational. Research audits are conducted through the National Research Program (NRP), which consists of a stratified random sample of the tax filing population. NRP audits seek to estimate the correctness of the whole return via a close to line-by-line examination. In part because they are so intensive, research audits constitute a small minority of the audits that the IRS performs each year. For example, about 2% of audited returns for tax year 2014 were selected through NRP.

We refer to audits that are not research audits as “operational audits.” Operational audits constitute the vast majority of audits that the IRS performs. Tax returns are selected for an operational audit through a wide variety of processes, the details of which are kept confidential. These processes can range from simple decision rules to manual examination to prioritization based on model-estimated risk scores.<sup>6</sup> Different audit programs use different processes to select which returns to audit. For some programs, tax returns are ranked in a manner that focuses on total noncompliance, aggregated across issues on the return; selected returns are then “classified” by IRS examiners to identify the most promising potential noncompliance issues on which to focus. Other audit programs are focused more narrowly on a particular issues or set of issues; returns audited through such programs are selected based on criteria that relate to the specific noncompliance issues of focus.<sup>7</sup>

To facilitate our research, the IRS shared information on operational audit selection processes with members of our research team; however, IRS policy limits our ability to

---

<sup>6</sup>We observe all training data and the full set of taxpayer features for EITC audits – our focus below – with the exception of audit referrals from whistleblowers or law enforcement (which are present in a very small share of EITC audits).

<sup>7</sup>The set of issues upon which an audit focuses may constrain the process through which the audit is conducted; for example, audits of claimed business deductions may be more amenable to correspondence audit than audits of unreported business income.

disclose information about these processes that would allow taxpayers to manipulate their risk of selection, such as tax return characteristics that may drive audit selection. The federal government itself has publicly disclosed some inputs into audit selection, such as the use of child custody and child birth records to assess whether a taxpayer is eligible to claim a child for a particular credit (G.A.O., 2015), and has also made clear that certain taxpayer characteristics are not taken into account, such as race or where the taxpayer lives (Kiel and Fresques, 2019).

Distinct from formal audits, the IRS operates several programs through which it screens submitted returns for potential identity theft or makes adjustments to taxpayers' submitted returns based on information reported to it by third parties, the detection of math errors (including returns claiming benefits for which the taxpayer does not qualify based on observable characteristics such as the taxpayer's age or reported income), or other factors. In other cases, the IRS will flag a potential compliance issue on a submitted return through a "soft notice" or other process for the taxpayer to resolve the issue without undergoing a formal audit, such as certain situations in which multiple taxpayers claim the same child for the same year. Thus, many of the returns with "smoking gun" evidence of non-compliance are addressed outside of the formal audit process.

## **2.3 The Earned Income Tax Credit**

The Earned Income Tax Credit (EITC) is a tax credit designed to support low- and middle-income taxpayers with earnings from work. The credit is refundable, meaning that taxpayers receive a payment or "refund" from the government if it brings their tax liability for the year below zero. Today, the EITC constitutes the largest cash-based safety net program in the United States, benefiting approximately 31 million households at an annual budgetary cost of \$64 billion (IRS, 2023a).

The amount of EITC for which a taxpayer qualifies is based on the taxpayer's income and family size. The credit amount initially increases with income for lower-income taxpayers,

plateaus, and then decreases with income once a taxpayer’s income surpasses a specified threshold. The maximum EITC amount varies based on the number of children a taxpayer claims; in 2014 (the year that most of our analysis focuses upon), the maximum EITC ranged from \$496 for taxpayers without children to \$6,143 for taxpayers with three or more children, and was available to taxpayers with combined incomes of up to \$52,247 if married, or \$46,997 if single. Taxpayers without income from work do not qualify for any EITC amount. Approximately 19% of taxpayers claimed the EITC in 2014, with an average credit of \$2,122 among claimants (IRS, 2016).

To claim a child for the EITC, the child must satisfy several eligibility tests with respect to the taxpayer. In particular, the taxpayer must be related to the child through one of a specified set of relationships (e.g., the child’s parent, stepparent, grandparent, aunt, uncle, or sibling) and must reside with the child for over half of the tax year. In addition, the child must generally be below the age of 19, or below the age of 24 if the child is a full-time student. In cases in which two or more taxpayers qualify to claim a single child, a series of “tie-breaker” rules govern which claim takes priority.

Non-compliance with the EITC rules has been a persistent subject of policy concern. Estimates of the EITC improper payment rate hover around 25% (U.S. Treasury Department, 2022), which many observers attribute to the complicated rules governing eligibility for the credit (Holtzblatt and McCubbin, 2004; National Taxpayer Advocate, 2019b). Federal law requiring agencies to measure improper payments effectively imposes a minimal number of research audits of EITC returns that IRS must conduct, but such rules do not directly constrain IRS’s ability to adjust the number of EITC returns selected for operational audit (OMB, 2021). In addition to the EITC, the federal government designates two other refundable credits — the American Opportunity Tax Credit and Additional Child Tax Credit — as high-priority programs that are susceptible to significant improper payments.

In recent years, approximately 40% or more of individual taxpayer audits have been

focused on returns claiming the EITC (Congressional Research Service, 2022).<sup>8</sup> EITC returns are selected for audit through a variety of enforcement programs; as such, EITC claimants may be audited in a manner that is narrowly focused on their eligibility for the EITC (e.g., their eligibility to claim a particular dependent) or in a manner that investigates other potential forms of noncompliance (e.g., the accuracy of their reported income or their eligibility to claim a business deduction). The vast majority (94% in 2014) of audits of EITC claimants are correspondence examinations and approximately two-thirds occur pre-refund (see Appendix Table A.1).

Approximately three-quarters of audited EITC returns in 2014 were selected through the Dependent Database (DDb) program, which flags returns based on a set of rules and heuristics as well as various proprietary risk scores. The features that are used to calculate these proprietary risk scores are drawn from taxpayers' filed returns as well as other administrative data about taxpayers or their children available to IRS.<sup>9</sup> For DDb audits, the selection step of the process is automated without manual review by human examiners (IRS, 2011).

### 3 Empirical Framework

In this section, we provide results relating to the identification of disparities in an outcome (like audits) with respect to a characteristic (like race) that cannot be directly observed by the researcher.

---

<sup>8</sup>The high EITC audit rate originates from a 1997 deal between House Republicans and the Clinton administration to preserve EITC funding in exchange for heightened enforcement (Johnston, 2000).

<sup>9</sup>These administrative data sources include state-provided data on child custody determinations that is compiled by the Department of Health and Human Services, child birth records from the Social Security Administration, and prisoner data from the Bureau of Justice Statistics.



### 3.1 Basic Notation

Tax returns are indexed by  $i$ , and have observable characteristics  $X_i$ . We use  $B_i \in \{0, 1\}$  to indicate whether the primary filer on tax return  $i$  is Black, and  $Y_i \in \{0, 1\}$  to indicate whether the return is audited. The audit rate for Black taxpayers is  $Y^B = \mathbb{E}[Y|B = 1]$ , and the audit rate for non-Black taxpayers is  $Y^{NB} = \mathbb{E}[Y|B = 0]$ .

Our goal is to estimate the audit disparity with respect to Black taxpayers,  $D$ , which we define as the difference in audit rates between Black and non-Black taxpayers:

$$D = Y^B - Y^{NB} = \mathbb{E}[Y|B = 1] - \mathbb{E}[Y|B = 0] \quad (\text{Audit Disparity})$$

An important barrier to studying differences in audit rates by race is that neither we nor the IRS observe taxpayer race. To overcome this challenge, we first estimate the probability that a taxpayer is Black using a subset of characteristics we do observe, and second, use the resulting race probabilities, along with administrative data on audits, to estimate differences in audit rates by race. That is, our approach is to:

1. Estimate  $b_i = \Pr[B_i = 1|Z_i]$ , where  $Z_i \subseteq X_i$  is a subset of  $i$ 's observable characteristics.
2. Use estimated  $b_i$  and observed  $Y_i$  to estimate  $D$ .

In the remainder of this section, we describe these two steps in additional detail.

### 3.2 Imputing Race

To impute race, we apply Bayesian Improved Surname Geocoding, which uses name and geolocation to probabilistically infer race (Imai and Khanna, 2016). This method has been widely applied in academic studies and is recommended when race is missing by the National Academy of Medicine (Nerenz et al., 2009). Recent work has shown that first names are more informative than surnames for identifying Black individuals (Voicu, 2018), so we incorporate first name information as well, applying Bayesian Improved First Name Surname Geocoding

(BIFSG). The method is “naive” in the sense that it assumes that first name, surname, and geography are independent after conditioning on race:

$$\Pr[F, S, G|B] = \Pr[F|B] \Pr[S|B] \Pr[G|B]$$

where  $F$  is first name,  $S$  is surname, and  $G$  is geography. Using Bayes’ rule, this assumption implies

$$\Pr[B = 1|F, S, G] = \frac{\Pr[F|B = 1] \Pr[S|B = 1] \Pr[G|B = 1] \Pr[B = 1]}{\sum_{j=0}^1 \Pr[F|B = j] \Pr[S|B = j] \Pr[G|B = j] \Pr[B = j]}, \quad (1)$$

and similarly for  $\Pr[B = 0|F, S, G]$ . See Appendix B.1 for a formal derivation. Estimating these terms by name and geography yields individual-level race probabilities. Because audits occur at the level of the tax return, we estimate a single race probability per return, focusing on the primary filer in cases of joint returns by married spouses.

The independence assumption underlying BIFSG is strong, and is likely violated in practice (e.g., Greengard and Gelman, 2023). Still, prior research has found that the method performs well across a range of settings. Below, we validate the performance of BIFSG race probabilities as an input to our disparity estimators using a subset of tax records matched to non-IRS administrative data containing taxpayer race. We also verify the robustness of our results to alternative imputation methods, including an approach that sidesteps the independence assumption by predicting race solely on the basis of geographic information.

### 3.3 Estimating Disparity using Imputed Race

After estimating the probability that each taxpayer is Black, we next consider how to use those estimated probabilities to identify the difference in audit rates by race. We consider two estimators: the *probabilistic disparity estimator* and the *linear disparity estimator*. We

characterize the bias of each and provide conditions under which the two estimators bound the true audit disparity.

The probabilistic estimator calculates average audit rates by race by weighting each taxpayer's contribution to the average audit rate by the probability that the taxpayer is or is not Black. Formally, given estimated race probability  $b_i$  and audit status  $Y_i$ , we define the probabilistic audit rate estimators as

$$\hat{Y}_p^B = \frac{\sum_i b_i Y_i}{\sum_i b_i} \quad \hat{Y}_p^{NB} = \frac{\sum_i (1 - b_i) Y_i}{\sum_i (1 - b_i)}$$

where B and NB refer to the estimated audit rates among Black and non-Black taxpayers, respectively. The probabilistic disparity estimator,  $\hat{D}_p$ , is the difference in the probabilistic audit rate estimates for these groups:

$$\hat{D}_p = \hat{Y}_p^B - \hat{Y}_p^{NB} = \frac{\sum_i b_i Y_i}{\sum_i b_i} - \frac{\sum_i (1 - b_i) Y_i}{\sum_i (1 - b_i)}$$

The second estimator we consider is the linear disparity estimator,  $\hat{D}_l$ . Consider the regression of  $Y$  on  $b$ :

$$Y = \alpha + \beta b + \eta.$$

The linear disparity estimator corresponds to the estimated coefficient on  $b$  in this regression:

$$\hat{D}_l = \hat{\beta} = \frac{\sum_i (Y_i - \bar{Y})(b_i - \bar{b})}{\sum_i (b_i - \bar{b})^2}$$

where  $\bar{Y}$  and  $\bar{b}$  denote the sample averages of  $Y$  and  $b$ . The corresponding linear estimators of  $Y^B$  and  $Y^{NB}$  are given by  $\hat{Y}_l^B = \hat{\alpha} + \hat{\beta}$  and  $\hat{Y}_l^{NB} = \hat{\alpha}$ .

The following proposition characterizes the asymptotic bias of the disparity estimators.

**Proposition 1.** Suppose that  $b$  is a taxpayer's probability of being Black given some

observable characteristics  $Z$ , so that  $b = \Pr[B = 1|Z]$ . Define  $D_p$  as the asymptotic limit of the probabilistic disparity estimator,  $\widehat{D}_p$ , and  $D_l$  as the asymptotic limit of the linear disparity estimator,  $\widehat{D}_l$ . Then:

1.

$$D_p = D - \frac{\mathbb{E}[\text{Cov}(Y, B|b)]}{\text{Var}(B)} \quad (1.1)$$

2.

$$D_l = D + \frac{\mathbb{E}[\text{Cov}(Y, b|B)]}{\text{Var}(b)} \quad (1.2)$$

3. Suppose  $\mathbb{E}[\text{Cov}(Y, B|b)] \geq 0$  and  $\mathbb{E}[\text{Cov}(Y, b|B)] \geq 0$ . Then

$$D_p \leq D \leq D_l \quad (1.3)$$

4. Suppose  $\mathbb{E}[\text{Cov}(Y, B|b)] \leq 0$  and  $\mathbb{E}[\text{Cov}(Y, b|B)] \leq 0$ . Then

$$D_l \leq D \leq D_p \quad (1.4)$$

Appendix B provides a proof of this proposition, derives similar results for the audit rate (level) estimators, and characterizes the estimators' asymptotic distribution.<sup>10</sup>

For the probabilistic disparity estimator to be consistent (Proposition 1.1), it must be that any association between race and audits is mediated through predicted race. In our setting, this condition would be violated if Black taxpayers were selected for audit at different

---

<sup>10</sup>In some applications, we consider estimates of the *conditional* audit disparity with respect to a subset of taxpayer characteristics  $x \in X$ :  $E[Y|B = 1, x \in X] - E[Y|B = 0, x \in X]$ . Proposition 1 extends naturally to such applications, with the key covariance terms replaced by  $\mathbb{E}[\text{Cov}(Y, b|B, x \in X) | x \in X]$  and  $\mathbb{E}[\text{Cov}(Y, B|b, x \in X) | x \in X]$ , respectively. In other applications, we consider weighted versions of these estimators. We discuss both of these cases in Appendix B.

rates than non-Black taxpayers with identical names and living in identical neighborhoods.<sup>11</sup>

The linear disparity estimator requires a different assumption for consistency, namely that there be no residual association between predicted race and audits after conditioning on the taxpayer’s actual race (Proposition 1.2). This exclusion restriction would be violated if some of the information used to predict race—e.g., name—is associated with audits through channels other than race, such as socioeconomic status.

Propositions 1.3 and 1.4 highlight a useful implication of Proposition 1.1 and 1.2: the linear and probabilistic disparity estimators asymptotically bound the true disparity when  $E[\text{Cov}(Y, b|B)]$  and  $E[\text{Cov}(Y, B|b)]$  share the same sign.<sup>12</sup> In our setting, using geography and name to predict race, both of these covariance terms are likely to be positive. For example, there are well-documented differences in marital patterns by race, with rates of marriage among Black households below those of other groups (Aughinbaugh et al., 2013). In addition, unmarried taxpayers may be associated with higher rates of EITC audit risk due to the residency and relationship tests that govern the EITC qualifying child rules (Leibel et al., 2020). Hence, to the extent that race predictions derived from name and geography do not fully capture racial differences in household structure, it is likely that some residual correlation between audits and race would remain, so that  $E[\text{Cov}(Y, B|b)] > 0$ . At the same time, some research shows that Black Americans with more distinctively Black names have lower socioeconomic outcomes (Fryer and Levitt, 2004; Cook et al., 2016). To the extent that audit rates are declining in income for some portions of the income distribution (Black et al., 2022), this suggests that, even after conditioning on race, taxpayers with higher predicted probabilities of being Black may be audited at higher rates than taxpayers with lower predicted probabilities of being Black, suggesting  $E[\text{Cov}(Y, b|B)] > 0$ . Below, we

---

<sup>11</sup>Proposition 1.1 is related to a result in Chen et al. (2019) and Kallus et al. (2021), which substitutes  $\mathbb{E}[\text{Cov}(Y, B|Z)]$  for  $\mathbb{E}[\text{Cov}(Y, B|b)]$  in our expression, similar to conditioning on combinations of covariates in lieu of the propensity score. The difference is significant in practice, since it may not be practical to calculate  $\text{Cov}(Y, B|Z)$  when  $Z$  takes on many values, such as in the common circumstance in which race is imputed from name and geography.

<sup>12</sup>Proposition 3 in Appendix B establishes that the Black and non-Black audit rate levels are similarly bounded by the linear and probabilistic estimators under the same conditions.

provide empirical support for these conditions using an auxiliary data set in which race is observed for a subset of taxpayers.

Finally, Proposition 1 requires that the individual race probability estimates are perfectly calibrated,  $b = Pr[B = 1|Z]$  for each  $Z$ ; in Appendix B, we derive the bias of the disparity estimators when this assumption fails (see Proposition 2). One method for researchers to assess the calibration of the individual race probabilities is by calculating those probabilities for a subset of the data for which individual-level race information is available. To the extent the race probabilities are found to be miscalibrated, Proposition 2 shows that a simple linear correction using the available race labels can remove much of the resulting bias (see Appendix B.5 for details). Below, we implement this approach to assess and re-calibrate our BIFSG estimates.

## 4 Data

This section describes the IRS data relating to audits and other tax variables as well as the data we use to impute taxpayer race.

### 4.1 Tax Data

We begin with comprehensive, administrative, and anonymized IRS data from approximately 148 million individual income tax returns with valid social security numbers for 2014. We primarily focus on tax year 2014 because it is the most recent year for which the vast majority of audits were complete and available to us at the time of analysis. For each return, we observe the amount and sources of reported income, deductions, and credits claimed. We also observe information returns for each taxpayer, such as employer-reported wages on Form W-2, and other administrative records, such as Social Security Administration data on gender and year of birth.

Among the returns in our data, there were 780,627 operational audits, constituting 0.53%

of returns filed for the year. We also use data on the research audits conducted under the NRP, which, as described in Section 2, are selected using a stratified random sample of taxpayers. Between 2010 and 2014, there were between approximately 13,500 and 15,500 NRP audits per year. To increase the precision of our analyses that use NRP data, we pool the 71,878 returns selected for NRP audit between 2010 and 2014.

For each filed return, we observe whether the return was selected for audit.<sup>13</sup> In addition, among audited returns, we observe the amount, if any, of the IRS-imposed adjustment to the originally filed return. Throughout, we report quantities in 2014 dollars, inflation-adjusting the NRP returns from prior years.

## 4.2 Race Data

As described above, the IRS does not collect data on taxpayer race, either directly via tax returns or indirectly via merging tax data with administrative data on race from other agencies. Therefore, we rely on a BIFSG approach to estimate the probability that a taxpayer self-reports as Black based on the first name, last name, and location of residence reported on the taxpayer’s return. The taxpayer’s location was measured at the level of the Census Block Group, the smallest geographic unit with racial composition reported by the Census, which typically contain 600-3,000 individuals. Data on the joint distribution of first names and race were obtained from Loan Application Registers under the Home Mortgage Disclosure Act from 2007-2010, following Tzioumis (2018); data on the joint distribution of last names and race were obtained from the 2010 Decennial Census Surname File (U.S. Census Bureau, 2021); and data on the racial make-up of Census block groups from the American Community Survey 5-year estimates (2010-2014).<sup>14</sup> Information

---

<sup>13</sup>More precisely, we observe whether an audit for the return was completed prior to 2023. Nearly all audits are completed within six years of a return being filed.

<sup>14</sup>For purposes of training and validating BIFSG, we follow the racial and ethnic classifications provided in the first name and surname files; this entails treating individuals who report their ethnicity to be Hispanic as non-Black, and treating individuals who report themselves to be multi-racial as non-Black. Note that because of how the Census-derived inputs to BIFSG are measured, the method is technically trained to predict the householder’s report of the taxpayer’s race rather than the taxpayer’s own self-report. The steps

regarding the availability of these characteristics in our sample is described in Appendix Table A.2. We are able to estimate race probabilities based on all three attributes (first name, last name, and geolocation) for 73% of tax year 2014 returns. For the remaining returns, we use the available subset of these attributes to impute race.

To assess the validity of our estimated race probabilities and the statistical bias of our audit disparity estimators, we obtained data on self-reported race for a subset of our sample by matching taxpayers to publicly available voter registration records from North Carolina. North Carolina required all registered voters to report race until 1993, after which reporting became optional. Taxpayer and voter records were matched using name and address, which resulted in a 47% unique match rate and  $\sim 2.5$ M matched records.<sup>15</sup> In some of the calibration exercises using this data, we re-weight North Carolina taxpayers to match the overall U.S. population on observable demographic characteristics. Appendix C provides additional details regarding the data match and the construction of these weights.

## 5 Race Estimate Calibration and Statistical Bias Assessment

This section presents results relating to the calibration of our taxpayer race estimates and statistical bias of the audit disparity estimators on which we rely.

### 5.1 Taxpayer Race Estimates

As described above, we use BIFSG to estimate the probability that a taxpayer is Black based on the taxpayer’s first name, last name, and geography. Figure 1 summarizes the results of this exercise. The left panel of the figure presents the distribution of estimated race

---

of this analysis requiring non-anonymized taxpayer information were conducted by Treasury economists, with the (anonymized) results provided to the other members of our research team.

<sup>15</sup>Matching took place in a fire-walled environment, separately from the main analysis, under the direction of the Treasury Office of Tax Analysis. Information relating to voting and political party were excluded from all analyses.



probabilities. The distribution is bi-modal, with 4.4% of taxpayers having 90% or higher predicted probability of being Black and 77.0% of taxpayers having 10% or lower predicted probability of being Black. The mean prediction is 12.2%, which corresponds closely to the 12.2% of the overall U.S. population that was estimated to be Black by the U.S. Census in 2014 (ACS, 2014).<sup>16</sup>

The right panel of Figure 1 assesses the calibration of the estimated race probabilities. It uses the matched North Carolina data to compare the true probability that a taxpayer identifies as Black with the BIFSG-predicted probability that the taxpayer does so. The figure shows that the predicted race probabilities are generally monotonic in true self-reported race and generally track the 45-degree line. We observe a similar pattern for the sub-population of taxpayers claiming the EITC, although here we observe some evidence that BIFSG may under-estimate the probability a taxpayer is Black; we explore several re-calibration methods below to address this issue. Overall, to the extent our matched North Carolina sample is representative of the population, this analysis suggests that the BIFSG-derived race estimates constitute a reasonably accurate approach for estimating taxpayer race in this setting (see Appendix Table A.3 and Appendix Figure A.1 for additional detail). Below, we consider several robustness checks that employ alternative methods for calculating race probability estimates and obtain qualitatively similar results.

## 5.2 Assessing Bias of the Audit Disparity Estimators

In this section we use the North Carolina data to shed light on the statistical bias of the audit disparity estimators described above.

To visualize the parameters shaping the bias of the disparity estimates from Proposition 1, Figure 2 bins the taxpayers from the North Carolina data set based on their estimated probability of being Black, and plots the fraction of Black and non-Black taxpayers audited

---

<sup>16</sup>The relatively small share of individuals estimated to be Black with very high probability may reflect a limitation of the BIFSG methodology, which outputs a probability near one only for a taxpayer living in a very high Black-share neighborhood *and* with very distinctively Black first and last names. We investigate the robustness of our results to alternative race imputation methods below.

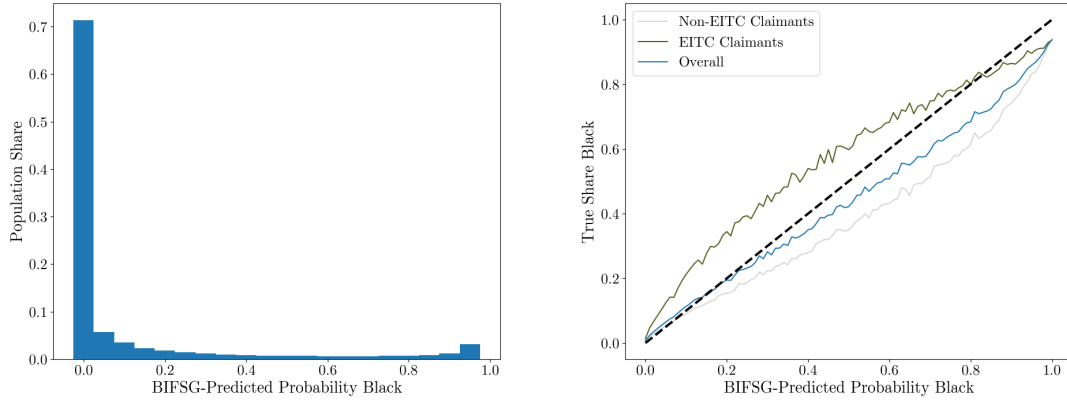
within each bin. The figure shows a positive residual correlation between audit probability and being Black after conditioning on the estimated race probability. From Proposition 1.1, this implies that the probabilistic disparity estimate is downward-biased, i.e.,  $E[\text{Cov}(Y, B|b)] > 0$ . At the same time, the figure suggests an upward-sloping audit rate in predicted race, for both Black and non-Black taxpayers, after conditioning on self-reported race, or  $E[\text{Cov}(Y, b|B)] > 0$ . From Proposition 1.2, this implies that the linear disparity estimator is upward-biased.

Appendix Table A.4 reports a more formal test for the sign of the key covariance terms. To estimate  $E[\text{Cov}(Y, b|B)]$ , we use the North Carolina data to directly calculate the covariance between audits and predicted race probabilities separately for Black and non-Black taxpayers. We aggregate these estimated covariances into an estimate of  $E[\text{Cov}(Y, b|B)]$  by weighting each race-specific covariance by the estimated proportion of all taxpayers that are Black or non-Black, respectively. In similar fashion, we estimate  $E[\text{Cov}(Y, B|b)]$  by calculating the sample covariance between audits and self-reported race separately for each estimated race probability percentile, and then aggregate based on the share of taxpayers in each race probability percentile. For both  $E[\text{Cov}(Y, b|B)]$  and  $E[\text{Cov}(Y, B|b)]$ , we reject the null hypothesis that the parameter is less than or equal to 0 with  $p < 0.01$ .

We obtain similar results when we re-weight the North Carolina data to match the U.S. population on a range of observable characteristics (Column 2 of Appendix Table A.4) and when we restrict the analysis to EITC claimants (Columns 3 and 4). In contrast, for the non-EITC population (Columns 5 and 6), we are unable to sign the second of these covariance terms with statistical significance.

We interpret the results in this subsection to support the hypothesis that  $E[\text{Cov}(Y, B|b)] > 0$  and  $E[\text{Cov}(Y, b|B)] > 0$  for EITC taxpayers and the overall population, and therefore, that the probabilistic and linear disparity estimators bound the true audit rate disparity for these populations.

Figure 1: Distribution and Calibration of Estimated Race Probabilities



*Notes:* Left: Nationwide histogram of BIFSG-predicted probability that a taxpayer is Black (non-Hispanic). The mean prediction is 12.2%. Right: The figure shows the calibration of the BIFSG imputations for the taxpayers in the matched North Carolina data set. Taxpayers are split into groups based on their predicted probability of being Black (discretized into 100 bins 1 percentage point wide). The predicted probability of being Black is on the  $x$ -axis; the  $y$ -axis represents the true proportion of each group that is Black according to self-reported race observed in the North Carolina matched sample, re-weighted to be representative of the overall United States (see Appendix C for details). A perfectly calibrated predictor would fall exactly on the 45-degree line, shown as the black dotted line. The figure shows overall calibration in blue as well as calibration among EITC claimants (dark green) and non-EITC claimants (light green).

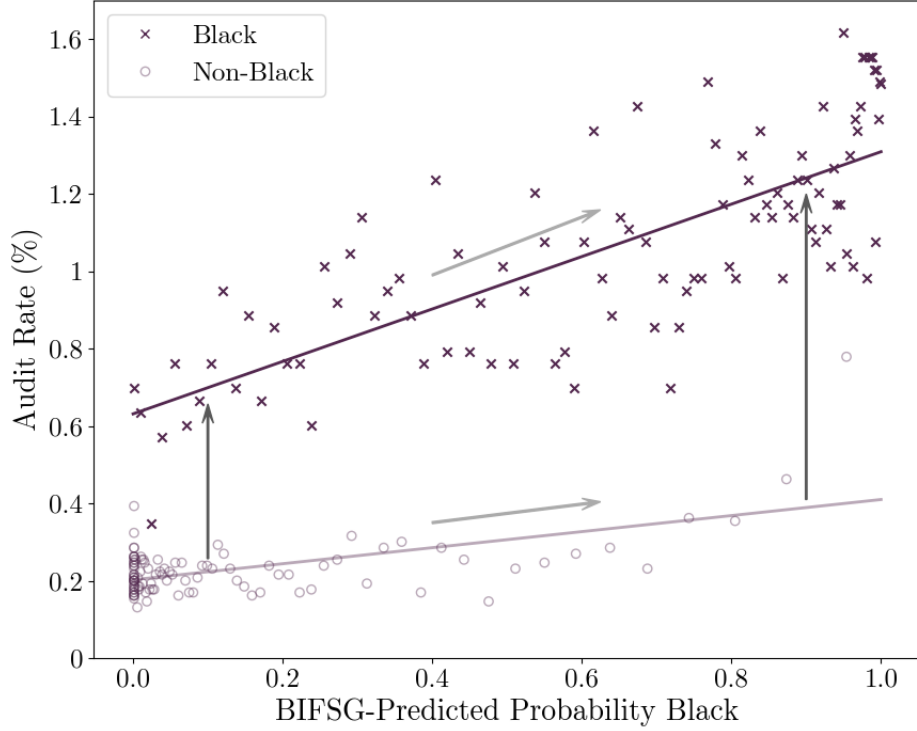
## 6 Audit Disparity Results

In this section, we report our estimates of the difference in audit rates between Black and non-Black taxpayers. We begin with the overall population of U.S. taxpayers before turning to EITC claimants.

Figure 3 presents our main findings concerning racial audit disparities for the population of U.S. taxpayers. The left panel plots the mean audit rate against the binned estimated probability that a taxpayer is Black. The upward-sloping relationship shown in the figure suggests that Black taxpayers are audited at a higher rate than non-Black taxpayers.

The right panel of Figure 3 depicts the estimated audit rates among Black and non-Black taxpayers, respectively, obtained from the probabilistic and linear estimators. Both estimators imply that Black taxpayers were audited at a higher rate than non-Black taxpayers. In particular, the probabilistic estimator implies a racial audit disparity of 0.81 percentage points: 1.24% of Black taxpayers were audited, compared to 0.43% of non-Black taxpayers. These audit rates are precisely estimated: the 95% confidence interval on the

Figure 2: Audit Rate by Predicted and Self-Reported Race



*Notes:* The figure shows the relationship between audit incidence (y axis) and BIFSG-predicted probability that a taxpayer is Black (x axis) for taxpayers filing returns for tax year 2014. Audit rates are plotted separately for Black and non-Black taxpayers in the North Carolina matched sample. Black and non-Black taxpayers are each grouped into 100 equal-sized bins, with Black taxpayers indicated by dark purple x's and non-Black taxpayers indicated by light purple circles. The average vertical distance between the x's and circles (illustrated by the dark gray arrows) provides an estimate of the sign of  $E[\text{Cov}(Y, B|b)]$ . The average slope of the group-specific best-fit lines (illustrated by the light gray arrows) provides an estimate of the sign of  $E[\text{Cov}(Y, b|B)]$ .

probabilistic disparity estimate ranges from 0.81 to 0.82 percentage points.<sup>17</sup> As expected, the linear estimator implies an even larger racial audit disparity, of 1.34 percentage points, and is also precisely estimated. Because the conditions for Proposition 1.3 appear satisfied in our setting, we interpret the probabilistic and linear disparity estimates as bounds on the true racial audit disparity. Thus, our results suggest that Black taxpayers were audited

<sup>17</sup>This confidence interval reflects sampling uncertainty in the outcome model, in the sense that even the universe of 2014 taxpayers may not perfectly reflect the underlying data generating process, but abstracts from uncertainty in the construction of the BIFSG-based probability estimates. We obtain slightly less precise disparity estimates after accounting for this source of uncertainty, following the dual-bootstrap method proposed by Lu et al. (2024) (Appendix Table A.10 and Appendix Figure A.12).

at between 2.9 and 4.7 times the rate of non-Black taxpayers.<sup>18</sup>

Figure 4 plots estimated audit rates by income and race. Black taxpayers appear more heavily audited throughout the income distribution.<sup>19</sup> Notably, the difference in audit rates appears largest for taxpayers with incomes that potentially qualify for the EITC. To more directly explore the role of the EITC in the observed racial audit disparity, we investigate differences in audit rates by EITC claim status as well as differences in EITC claiming by race. The right panel of Appendix Figure A.5 shows that the audit rate among EITC claimants (of any race) is more than 4 times higher than among non-EITC claimants (1.45 vs. 0.31 percent). In addition, the left panel of the figure shows that the EITC claim rate is increasing in the probability that a taxpayer is Black. Hence, one possibility is that the observed difference in audit rates could be due to EITC claimants being audited at higher rates and Black taxpayers being over-represented among that group.

To assess this hypothesis, we estimate audit disparities by race separately for EITC claimants and non-claimants. If differences in EITC claiming rates by race account for the racial disparity, we would expect the difference in audit rates by race to be relatively small *within* the population of EITC claimants. However, Figure 5 shows this is not the case. Instead, the estimated disparity in audit rates between Black and non-Black EITC claimants is substantially larger in percentage point terms (between 1.96 and 2.90 p.p.) than the estimated disparity for the full population (between 0.81 pp and 1.34 p.p.). In contrast, we estimate significantly smaller racial audit disparities among taxpayers not claiming the EITC (between 0.10 and 0.18 p.p.).<sup>20</sup>

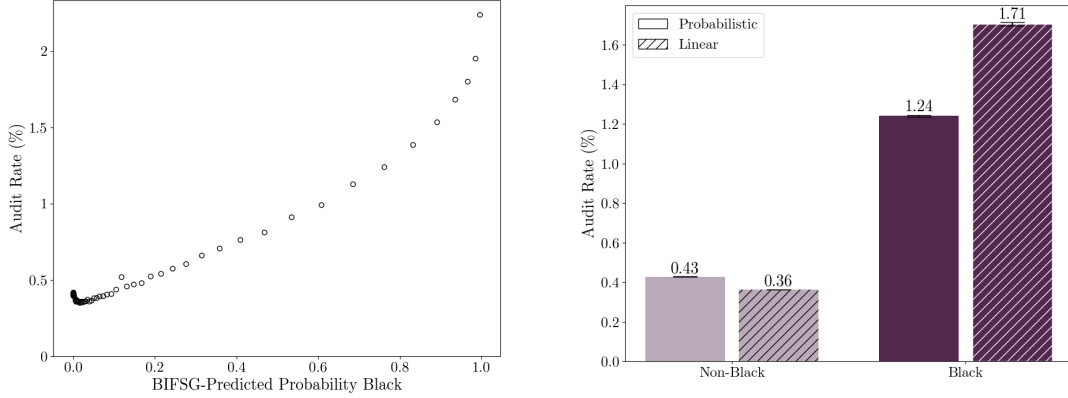
---

<sup>18</sup>We obtain comparable results from using a linear program to calculate maximum and minimum values of disparity that are consistent with the Proposition 1.3 assumptions and the observed joint distribution of  $b$  and  $Y$ ; see Appendix B.9.

<sup>19</sup>For ease of exposition, we focus on the probabilistic estimate of the audit rate; Appendix Figure A.2 shows a similar pattern using the linear estimator. As discussed in Appendix B.8, interpreting the estimators as bounds on a conditional audit disparity requires that the Proposition 1 conditions hold at the subgroup-level. Appendix Figure A.4 reports estimates of the relevant covariance terms by income bin from the NC data. For most bins, the Proposition 1.3 conditions appear satisfied, although the pattern is more consistent for the probabilistic estimator. For higher income taxpayers, the linear disparity estimator may not constitute an upper bound.

<sup>20</sup>These disparities are precisely estimated; refer to Table 1 for standard errors. As discussed in Section 5, we lack empirical evidence that the linear disparity estimator yields an upper bound on the racial audit

Figure 3: Audit Rates by Race



*Notes:* The figure shows the relationship between audits and race among taxpayers filing returns for tax year 2014. Left: Binned scatterplot of audit rate by BIFSG-predicted probability that a taxpayer is Black. Taxpayers have been grouped into 100 equal-sized bins. Right: Estimated audit rates among Black and non-Black taxpayers, calculated using the probabilistic audit rate estimator and the linear audit rate estimator with BIFSG-predicted probabilities. Error bars show the 95% confidence interval based on the distribution of estimates from 100 bootstrapped samples.

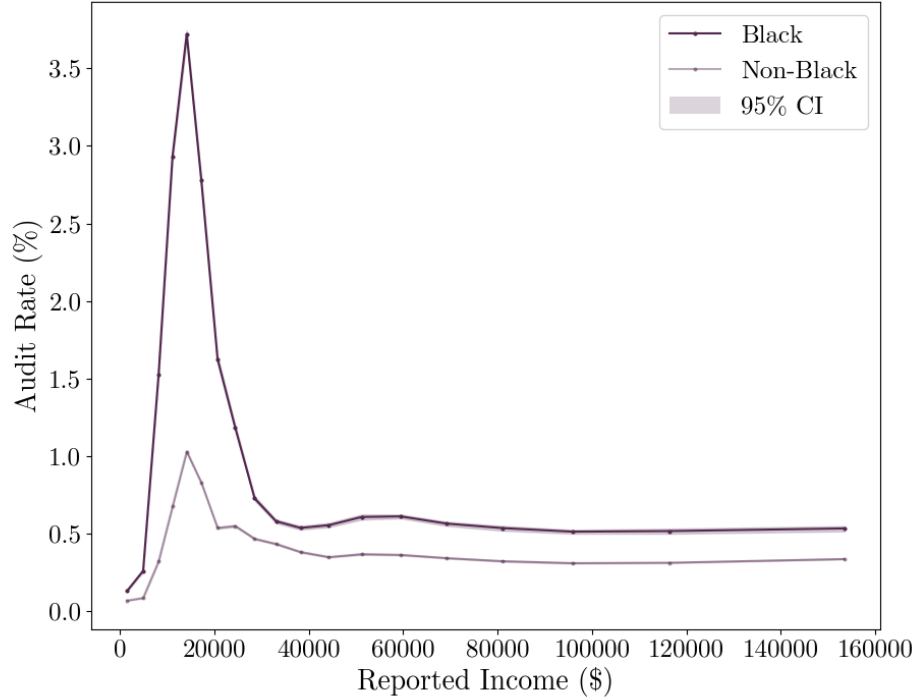
We can formally decompose the overall audit disparity into three components: (1) racial differences in the audit rate among EITC claimants; (2) racial differences in the audit rate among EITC non-claimants; and (3) racial differences in the rate at which taxpayers claim the EITC, scaled by differences in the audit rate for EITC versus non-EITC returns:

$$\begin{aligned}
 Y^B - Y^{NB} = & \underbrace{(Y_C^B - Y_C) C_B - (Y_C^{NB} - Y_C) C_{NB}}_{(1)} \\
 & + \underbrace{(Y_{NC}^B - Y_{NC}) (1 - C_B) - (Y_{NC}^{NB} - Y_{NC}) (1 - C_{NB})}_{(2)} + \underbrace{(C_B - C_{NB}) (Y_C - Y_{NC})}_{(3)}
 \end{aligned}$$

where  $C_B$  and  $C_{NB}$  denote the respective probabilities that Black and non-Black taxpayers claim the EITC;  $Y_C^B$ ,  $Y_C^{NB}$ , and  $Y_C$  denote the respective audit rates for Black, non-Black, and all EITC claimants; and similarly for  $Y_{NC}^B$ ,  $Y_{NC}^{NB}$ , and  $Y_{NC}$  with respect to EITC non-claimants. We estimate that racial differences in the audit rate within EITC claimants

disparity for EITC non-claimants. In our linked North Carolina data set containing self-reported race, the true racial audit disparity for this group is slightly larger than, but similar in magnitude to, the disparity obtained from the linear disparity estimator (see Appendix Table A.7).

Figure 4: Audit Rates by Income and Race



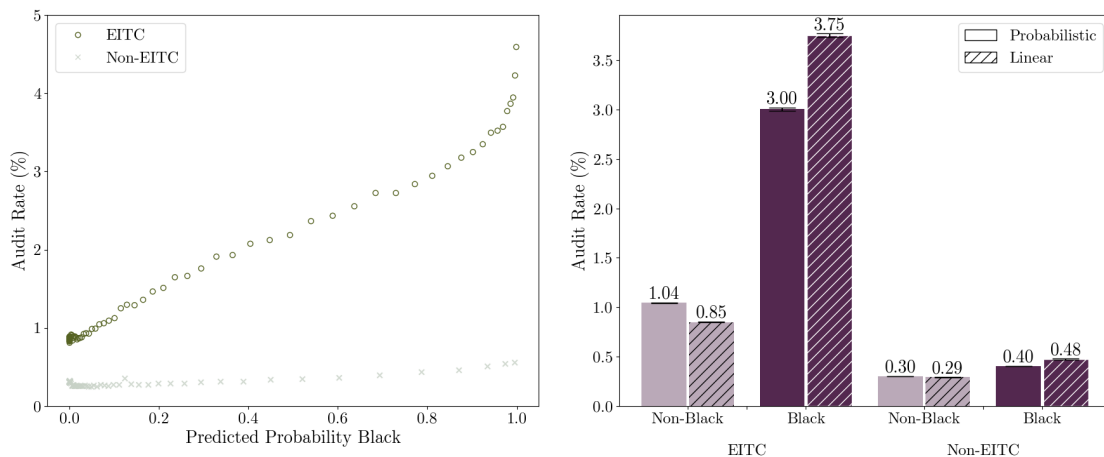
*Notes:* The figure shows the estimated audit rate by income among Black and non-Black taxpayers filing returns for tax year 2014. Income is measured according to the adjusted gross income (AGI) reported on the taxpayer’s return (i.e., prior to audit adjustments). Binned audit rates by race are determined using the probabilistic estimator; Appendix Figure A.2 reports the corresponding analysis with the linear estimator. The sample is limited to taxpayers reporting non-negative AGI. Taxpayers have been grouped into 20 equal-sized bins, based on their AGI. The shaded area around each line shows the 95% confidence interval based on the distribution of estimates from 100 bootstrapped samples. To facilitate presentation, the x-axis is limited to bins with mean reported AGI under \$200,000; a version of the figure with all income percentiles is presented in Appendix Figure A.3.

contribute between 70% and 73% of the total disparity, with the remainder primarily due to racial differences in the rate that taxpayers claim the EITC (20% to 21%) and a smaller portion due to racial differences in the audit rate among taxpayers not claiming the EITC (7% to 8%).<sup>21</sup>

We next explore the type of audits from which the disparity arises. Table 1 shows that audit disparities appear to be largely driven by differences in the selection of correspondence

<sup>21</sup>Appendix D provides additional detail, as well as other potential decompositions. The relative importance of these three components is sensitive to the decomposition considered and the choice of reference group, but a consistent finding is that differences in audit rates among EITC claimants are an important contributor to overall disparity—generating at least 32% and up to 83% of the total difference in audit rates.

Figure 5: Audit Rates by Race and EITC Claiming



*Notes:* The figure shows the relationship between audits and race among taxpayers filing returns for tax year 2014, broken out by whether a taxpayer claims the EITC in that year. Left: Binned scatterplot of audit rate by BIFSG-predicted probability Black by EITC claim status, with EITC claimants and non-claimants each grouped into 100 equal-sized bins based on their estimated probability of being Black. EITC claimants are represented by dark green dots and non-claimants by light gray x's. Right: Estimated audit rate by race and EITC claim status, calculated using the probabilistic audit rate estimator and the linear audit rate estimator with BIFSG-predicted probabilities. Error bars show the 95% confidence interval based on the distribution of estimates from 100 bootstrapped samples.

audits, whereas Black and non-Black taxpayers appear to be selected for field and office audits at roughly similar rates.<sup>22</sup> We observe disparities in both pre-refund and post-refund audits, although the magnitude is larger in the former category than in the latter even once we limit consideration to correspondence audits. Although audit disparities are largely concentrated among EITC claimants, correspondence audits appear to drive the smaller disparity we observe among EITC non-claimants as well. Among EITC claimants, 78.5% of the observed audit disparity is attributable to audits conducted through the DDb program (Appendix Table A.11).

We next explore heterogeneity in audit disparities within distinct groups of EITC claimants. Figure 6 reveals significant absolute and relative disparities by race among unmarried EITC claimants, particularly unmarried men. Strikingly, among unmarried EITC claimants with dependents, the audit rate for Black men is over 4 percentage points

<sup>22</sup>The financial and time costs of being audited differ across these audit categories, as do the average amounts of back taxes, penalties, and interest that are imposed. Appendix Tables A.8 and A.9 provide back-of-the-envelope estimates for how audit burdens vary by race, accounting for these factors.



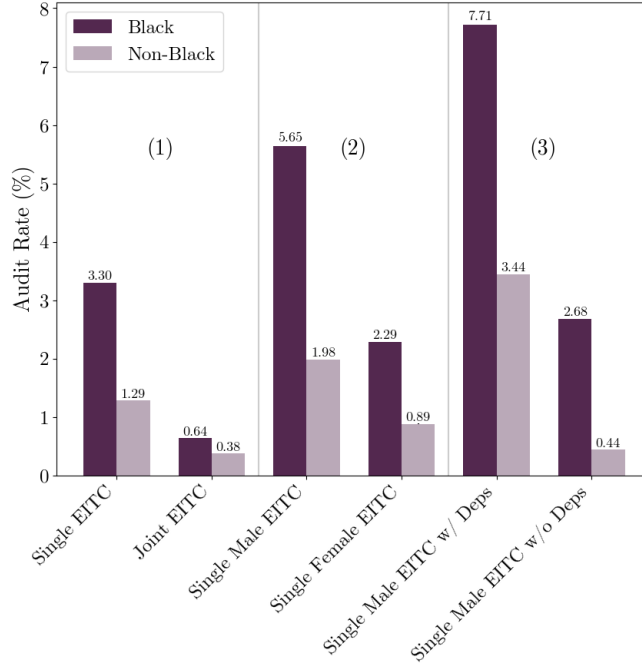
Table 1: Updated Estimated Audit Rate Disparity

	Any Audit (1)	Audit Timing		Audit Type	
		Pre Refund (2)	Post Refund (3)	Correspondence (4)	Field/ Office (5)
Panel A: Full Population					
Probabilistic	0.813 (0.003)	0.569 (0.002)	0.244 (0.002)	0.804 (0.003)	0.010 (0.001)
Linear	1.345 (0.004)	0.941 (0.003)	0.403 (0.002)	1.328 (0.004)	0.016 (0.001)
Mean Audit Rate	0.526	0.243	0.283	0.433	0.094
N (Millions)	148.3	148.3	148.3	148.3	148.3
Panel B: EITC Population					
Probabilistic	1.950 (0.008)	1.475 (0.007)	0.475 (0.004)	1.943 (0.008)	0.007 (0.001)
Linear	2.885 (0.009)	2.182 (0.008)	0.703 (0.005)	2.875 (0.009)	0.010 (0.002)
Mean Audit Rate	1.443	0.956	0.487	1.356	0.087
N (Millions)	28.3	28.3	28.3	28.3	28.3
Panel C: Non-EITC Population					
Probabilistic	0.102 (0.002)	0.005 (0.001)	0.097 (0.002)	0.088 (0.002)	0.013 (0.001)
Linear	0.180 (0.003)	0.009 (0.001)	0.171 (0.002)	0.156 (0.002)	0.024 (0.001)
Mean Audit Rate	0.310	0.075	0.235	0.215	0.096
N (Millions)	120.0	120.0	120.0	120.0	120.0

*Notes:* The table reports probabilistic and linear estimates of the difference in audit rates between Black and non-Black taxpayers filing income tax returns for tax year 2014. Units are percentage points (0-100). The category of audit considered varies across columns; for example, the results in column (4) show the estimated difference in the rate that Black versus non-Black taxpayers are selected for a correspondence audit. Panel A includes all taxpayers, whereas Panels B and C restrict the analysis to EITC claimants and non-claimants, respectively. Standard errors, reported in parentheses, are calculated from the asymptotic distributions described in Appendix B.3. See Appendix Tables A.5 and A.6 for estimates of audit rate levels by race.

larger than the audit rate for non-Black men, and both are an order of magnitude larger than the audit rate for the overall U.S. population. In contrast, we observe smaller racial audit disparities among joint filers, unmarried women, and unmarried men who do not claim dependents, although the ratio of audit rates among Black to non-Black taxpayers

Figure 6: Audit Rate Disparities by EITC Subgroup



*Notes:* The figure shows the estimated audit rate among the specified subgroups of Black and non-Black EITC claimants. Conditional audit rates by race are calculated using the probabilistic audit rate estimator applied to BIFSG-predicted probabilities that a taxpayer is Black. Panel (1) splits EITC claimants by single vs joint filers; (2) splits single EITC claimants by taxpayer gender; and (3) splits single men claiming the EITC by whether they claim dependents. A similar analysis, corresponding to the linear estimator, is presented in Appendix Figure A.6

remains substantial among these groups.<sup>23</sup>

We conduct a number of analyses to assess the robustness of our results to alternative approaches for estimating taxpayer race, reported in Appendix Table A.12. First, recall that 27% of the taxpayers in our sample are missing one or more of the variables used to impute race; Column 1 re-estimates the racial audit disparity after excluding this group. Second, as an alternative to BIFSG, in Column 2 we re-estimate predicted race using additional individual characteristics and relaxing the naive Bayes independence assumption with Gibbs sampling (see Appendix B.6 for details). Third, our BIFSG race probability estimates are derived from samples that are designed to be representative of the overall U.S. population,

<sup>23</sup>Appendix Figure A.7 reports estimates of the conditional covariance terms from Proposition 1.3 for each subgroup in Figure 6. In each case, the terms are estimated to have the expected sign, with the exception of the term associated with joint filers claiming the EITC. Hence, the linear estimator may not identify an upper bound on disparity for that group.

not the subset of the population that files taxes or claims the EITC. As discussed in Section 3, such mis-calibration can shape the bias of our disparity estimators. In Appendix B, we formally derive this bias and use the result to re-calibrate our race probability estimates, treating the North Carolina data as ground truth (Columns 3 and 4).<sup>24</sup> Fourth, we can avoid imposing BIFSG’s conditional independence assumption if race is imputed based on geography alone rather than on geography and name. Column 5 replicates our main results using this simpler proxy.<sup>25</sup> Across columns, the results are largely unchanged from our baseline approach.

A different potential concern with our analysis is that the matched North Carolina data that we use to validate our identifying assumptions may not be representative of the national population, even after re-weighting. To assess this possibility, we matched our sample to voter registration data in six other states in which race is recorded and estimated the key covariance terms from Proposition 1. Across states, we obtain substantially similar results (Appendix Figure A.13).<sup>26</sup>

Finally, our results thus far have focused on returns for tax year 2014. To confirm the patterns we observe are not limited to that year, we estimate disparity among all taxpayers (Appendix Figure A.14) and among EITC claimants (Appendix Figure A.15) for tax years 2010, 2012, 2016, and 2018. In each case, we obtain comparable results to those from 2014.

---

<sup>24</sup>Appendix Figures A.8 and A.9 contain results from the analogous exercise applied to the income bins and demographic groups reported in Figures 4 and 6.

<sup>25</sup>For additional detail on the geography-based estimates, see Appendix Figures A.10 and A.11 and Appendix Table A.13

<sup>26</sup>A limitation of the matched data from states other than North Carolina, and the reason we do not rely on them for our main analysis, is that they are drawn from voting records for 2023; as expected, this leads to a much lower match rate. Separately, a recent working paper independently confirms the existence of the audit rate disparity in a matched national sample of residents in certain tax-subsidized apartments (Derby et al., 2024).

## 7 What Causes the Racial Audit Disparity Among EITC Claimants?

Because the racial audit disparity appears concentrated among EITC claimants, our remaining analyses focus on this population of taxpayers. As discussed above, we are confident that the disparity for this group is not due to disparate treatment in audit selection because the vast majority of EITC returns are selected for audit based on automated processes, and these processes do not include race as an input. In this section, we explore how group-level differences in taxpayer characteristics, as well as choices related to the design of the audit selection algorithm, might contribute to the observed disparity. To simplify exposition, in this section we rely primarily on the probabilistic estimator; we obtain qualitatively similar results using the linear disparity estimator (reported in Appendix A).

### 7.1 Differences in Underreporting

We first investigate whether the observed audit disparity can be explained by differences in the distribution of underreporting between Black and non-Black EITC claimants. By underreporting, we mean the difference between a taxpayer's correct income tax obligations for a tax year (which may be negative in the case of a taxpayer qualifying for refundable credits) and the tax obligations reported on the taxpayer's return. For example, underreporting may arise from reporting too little income, too many deductions, or from claiming a credit for which the taxpayer does not qualify. Underreporting may be intentional or inadvertent, and may be due to decisions by either the taxpayer or a tax preparer.

### 7.1.1 Differences in Actual Underreporting

A challenge in studying whether the observed disparity is due to racial differences in the distribution of underreporting is that we observe underreporting only among those taxpayers who were selected for audit. We can circumvent this obstacle by combining non-random operational audits with data from randomly selected NRP audits. That is, using Bayes rule, we can express the audit rate for taxpayers of race  $j$  and underreporting amount  $k$  as:

$$\Pr[Y = 1|B = j, K = k] = \Pr(K = k|Y = 1, B = j) \frac{\Pr(Y = 1|B = j)}{\Pr(K = k|B = j)}.$$

From this equation, we can estimate the racial audit disparity at a given level of underreporting using: the estimated audit rates by race for EITC claimants from Figure 5 for  $\Pr(Y = 1|B = j)$ ; the NRP data to estimate  $\Pr(K = k|B = j)$ ; and the detected underreporting from the operational audit results to estimate  $\Pr(K = k|Y = 1, B = j)$ .<sup>27</sup>

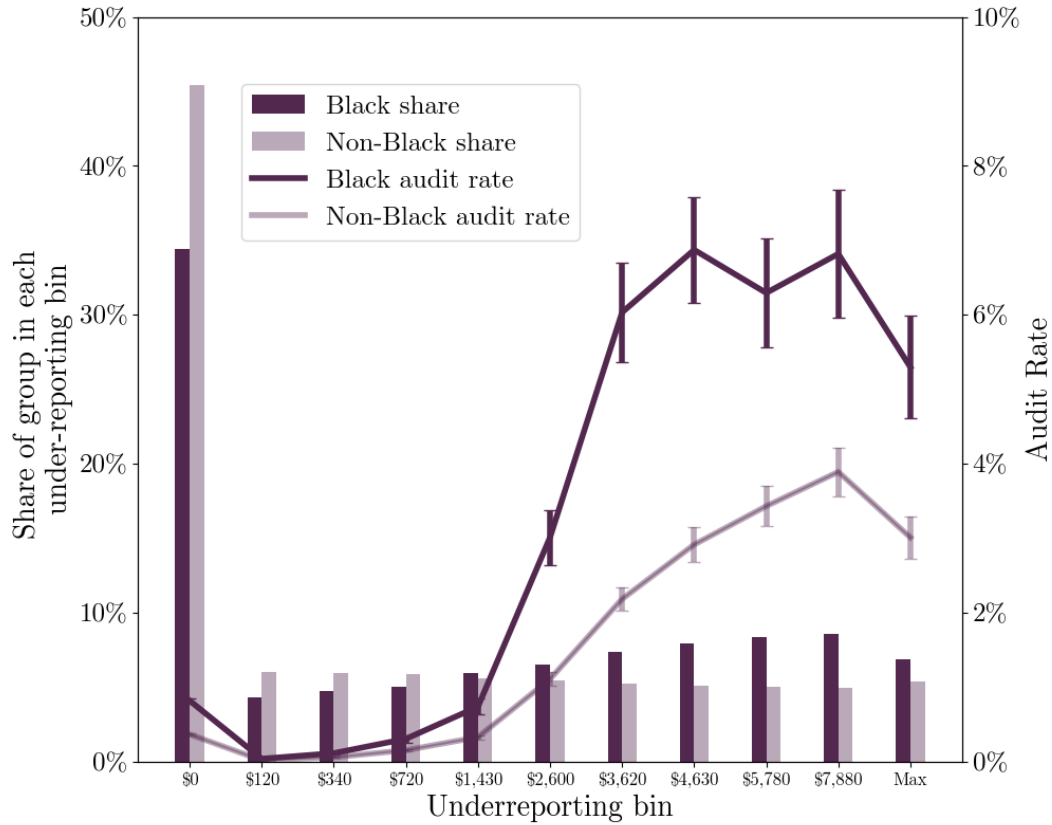
The results of this analysis are presented in Figure 7, with EITC claimants binned according to their underreported taxes. The figure shows that the distribution of underreporting for Black EITC claimants tends to be concentrated at higher values than for non-Black EITC claimants. However, within each underreporting bin, the estimated audit rate for Black taxpayers exceeds that of non-Black taxpayers, and for many bins, the difference is substantial. Figure 7 therefore provides evidence that differences in audit rates remain when comparing Black and non-Black taxpayers with similar levels of underreporting.<sup>28</sup>

---

<sup>27</sup>With respect to the last term, operational audits may not detect all of a taxpayer’s underreporting because, unlike the NRP, operational audits do not assess each issue on the return. In practice, this concern is lessened by the fact that the issues worked by operational audits tend to be those where underreporting is suspected. Conversely, although they are intensive, even NRP audits may miss some underreporting, and differential mis-measurement of underreporting by race could distort our results. For differences between NRP and operational audit scope to affect this analysis, those differences would have to vary by race. In addition, the factors that have been found to induce inaccuracies in NRP audit results are concentrated at the top of the income distribution (Guyton et al., 2021); hence, we do not expect this concern to be important for the EITC population – our focus here. See Appendix B.7 for further details.

<sup>28</sup>A limitation of this analysis is that the relatively small degree of overlap between the EITC NRP sample and the matched NC data set limits our ability to re-calibrate the disparity estimates by underreporting bin, as we did with the subgroup disparity estimates presented in Appendix Figures A.8 and A.9.

Figure 7: Racial Audit Disparity Among EITC Claimants by Underreported Taxes



*Notes:* The figure shows the estimated audit rates for Black and non-Black EITC claimants, respectively, by under-reported taxes. Taxpayers are binned into 11 categories: those with less than \$1 of underreporting, and 10 equal deciles of taxpayers with positive underreporting. Underreporting deciles are defined based on the distribution of underreporting among EITC claimants, as measured by NRP audits. Bin labels on the x-axis reflect the upper dollar limit of each underreporting bin (rounded for confidentiality). Estimated audit rates by race are calculated using the probabilistic estimator. All analyses account for NRP sampling weights. Brackets reflect the estimated 95% confidence interval, derived from bootstrapped standard errors (N=100). The bars show the estimated share of Black and non-Black taxpayers, respectively, that fall into each underreporting bin. A similar analysis, corresponding to the linear estimator, is presented in Appendix Figure A.16.

Figure 7 highlights the relationship between underreporting and audits throughout the underreporting distribution; in contrast, our next analysis zooms in on the extreme right tail of the underreporting distribution. The motivation for doing so is that the IRS selects less than 1.5% of EITC returns for audit; it could be that Black taxpayers are over-represented among EITC claimants with the highest underreporting. If so, selecting from this population could generate the higher observed audit rate for that group.

To explore this possibility, we simulate an “oracle” selection algorithm, which prioritizes returns for audit according to the actual dollar amount of underreporting that would be detected if the return were audited.<sup>29</sup> We treat the set of NRP-audited returns—for which we know true underreporting—to be the population, and pretend that we must select some subset of returns from this population to audit. To implement the oracle, we rank each return in this population based on its under-reported taxes. Then, in descending order of this ranking, we select returns for audit until some pre-specified audit rate has been reached. For each audit rate that we consider, we calculate (1) the sum of detected underreporting over the audited taxpayers, and (2) the racial audit disparity that would be induced by this selection process.<sup>30</sup>

The results of this analysis are shown by the line labeled “Total Underreporting Oracle” in Figure 8. As a benchmark, the Figure also reports total detected underreporting and disparity from audits of tax year 2014 returns claiming the EITC, indicated by the dashed red lines and labeled “Status Quo”. At each audit rate considered, the underreporting oracle selects Black taxpayers at a lower rate than the status quo, and, more strikingly, at a lower rate than non-Black taxpayers. Thus, although Figure 7 shows that the distribution of underreporting for Black EITC claimants tends to be concentrated at higher values than

---

<sup>29</sup>Because the actual amount of underreporting on a return cannot be known before an audit is conducted, such an algorithm is not feasible for the IRS to implement; we consider it as a benchmark before turning to predicted underreporting below.

<sup>30</sup>Throughout, we use “detected underreporting” as shorthand for total recommended adjustments from audit, not accounting for appeals or uncollected tax debts. We report annualized detected underreporting, accounting for NRP sample weights, and adjusting for the fact that our NRP data pools multiple tax years; see Appendix E for details.

for non-Black EITC claimants, these results suggest that the EITC claimants with the very largest underreporting are disproportionately non-Black.

Figure 9 illustrates this contrast more directly. In Section 6, we estimated that the status quo audit rate for non-Black EITC claimants was 1.04% compared to 3.00% for Black EITC claimants. In contrast, if EITC claimants were selected for audit according to their true underreporting, the audit rate for non-Black EITC claimants would be 1.63% compared to 0.74% for Black EITC claimants. We interpret this result, in conjunction with the different audit rates by race within underreporting bins in Figure 7, as evidence that the observed audit disparity cannot be entirely explained by group-level differences in underreporting by race.

### 7.1.2 Differences in Predicted Underreporting

In practice, the actual amount of a taxpayer’s underreporting is unknown at the time that auditing decisions are made. We now ask whether the observed audit disparity is due to the IRS relying on *predicted*, rather than actual, underreporting in selecting which returns to audit.

To study this question, we train a random forest model to predict taxpayer underreporting (in dollars) based on features that the tax authority can observe at the time of the audit decision. The model is trained on NRP returns claiming the EITC, and largely incorporates the same features available to the DIF and DDb programs to select audits of EITC claimants — information reported on tax returns supplemented with additional administrative data available to the IRS.<sup>31</sup> We then simulate selecting taxpayers in descending order of the model’s predictions, until some specified audit rate has been reached. In other words, the approach is the same as the underreporting oracle, except that returns are selected for audit on the basis of predicted, rather than actual, underreporting.

The dark purple line in Figure 8 shows detected underreporting and disparity induced by

---

<sup>31</sup>Additional detail on the predictive model is provided in Appendix E.



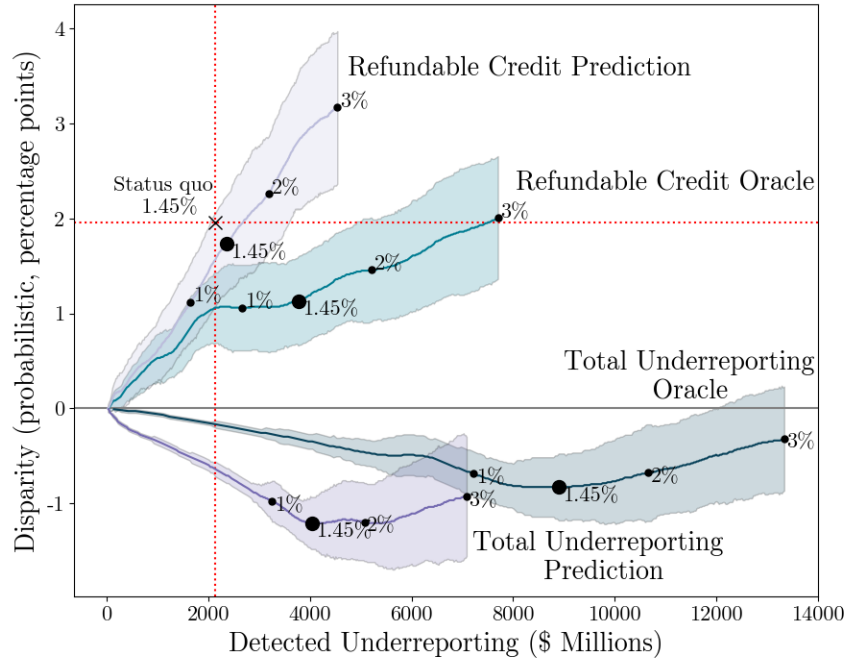
the predicted underreporting algorithm. Unsurprisingly, selecting audits based on predicted underreporting yields significantly less detected underreporting than the oracle at any given audit rate. However, like the underreporting oracle, the underreporting prediction algorithm selects Black taxpayers for audit at lower rates than non-Black taxpayers at each audit rate we consider. At the status quo audit rate, selecting EITC claimants for audit based on predicted underreporting would yield an audit rate of 0.45% for Black taxpayers and 1.71% for non-Black taxpayers. Hence, the fact that the IRS must select audits based on predicted rather than actual noncompliance does not, in itself, appear to explain the observed audit disparity.

## 7.2 Algorithmic Objective: Refundable Credit Overclaims

In this subsection, we explore the possibility that the observed audit disparity arises because the IRS selects among EITC returns for audit based on some objective other than maximizing the detection of underreported taxes. In particular, we consider the effect on disparity of allocating audits based on (1) the amount of underreporting attributable to overclaimed refundable credits, rather than (2) the total amount of underreported taxes (from whatever source).

To do so, we compare the audit selection algorithms described in the prior subsection that focus on total underreporting with algorithms that focus exclusively on refundable credit overclaims. Specifically, we consider two new algorithms (both of which are depicted in Figure 8). The first is a refundable credit oracle, which ranks returns by the sum of actual underreporting (in dollars) attributable to the three refundable credits designated by OMB as high-priority programs susceptible to significant improper payments: the EITC, the Additional Child Tax Credit, and the American Opportunity Tax Credit. As with the total underreporting oracle, the refundable credit oracle is not feasible to implement in practice because the actual amount of overclaimed refundable credits cannot be known with certainty before the audit occurs, but serves as a useful benchmark. The

Figure 8: Detected Underreporting and Disparity by Algorithm



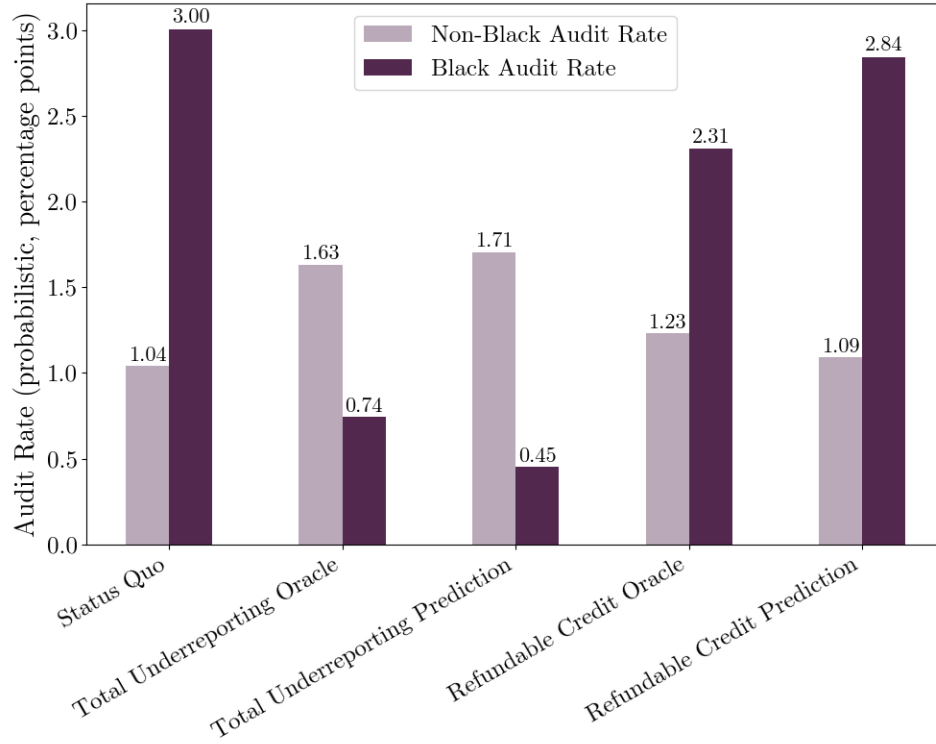
*Notes:* The figure shows estimated disparity ( $y$ -axis) and annualized detected underreporting ( $x$ -axis) under alternative algorithms for selecting audits of EITC claimants and under alternative audit rates (0.1% to 3%). For each algorithm, returns claiming EITC in the NRP sample are first ranked and then selected for audit in descending order of the ranking until the specified audit rate is reached. The ranking varies by algorithm and is based on: total underreporting (dark blue line), predicted total underreporting (dark purple line), underreporting due to overclaimed refundable credits (light blue line), and predicted underreporting due to overclaimed refundable credits (light purple line). The point labeled “Status quo” shows estimated disparity and total underreporting from the 1.45% of EITC returns selected for audit from the population of tax year 2014 returns. The labeled points along each line correspond to the audit rate specified in the label; the audit rate corresponding to the status quo EITC audit rate (1.45%) is denoted by a larger dot. Disparity is calculated using the probabilistic disparity estimator; see Appendix Figure A.17 for results using the linear disparity estimator. Shaded regions around each line correspond to 95% confidence intervals for disparity calculated based on the distribution of estimates from 100 bootstrapped samples. Annualized detected underreporting is the sum of detected underreporting (positive or negative) among returns selected for audit under the specified algorithm, scaled to reflect our use of five years of NRP data. All analyses incorporate NRP sampling weights. The two prediction algorithms are based on random forest regression models. See Appendix E for additional detail.

second new algorithm is based on predicted refundable credit overclaims; it ranks returns by the output of a random forest regressor trained to predict the sum of underreporting attributable to the same three refundable credits. After ranking, both algorithms select returns for audit in descending rank order until a specified audit rate is reached. In other words, the two refundable credit algorithms mirror the two total underreporting algorithms described above; the only difference is that the ranking used to select audits is based on overclaimed refundable credits—actual overclaims in the case of the oracle and predicted overclaims in the case of the prediction algorithm—rather than total underreporting.

As shown in Figure 8, the refundable credit algorithms detect substantially less underreporting than the algorithms focused directly on that objective, indicating that refundable credit overclaims are not the only important source of underreporting among EITC claimants. With respect to disparity, both refundable credit algorithms yield opposite-signed results compared to the algorithms focused on total underreporting: they select Black taxpayers at higher rates than non-Black taxpayers. As highlighted in Figure 9, selecting EITC claimants for audit based on actual refundable credit overclaims while holding the overall audit rate of EITC claimants fixed at the status quo level would lead to an audit rate of 2.31% for Black EITC claimants, compared to 1.23% for non-Black EITC claimants. Similarly, selecting audits according to predicted refundable credit overclaims would lead to an audit rate of 2.84% for Black EITC claimants, compared to 1.09% for non-Black EITC claimants. We interpret these results as evidence that targeting the detection of overclaimed refundable credits leads to a larger share of Black taxpayers being selected for audit compared to targeting the detection of total underreporting.

Does the objective of the audit selection algorithm actually contribute to the audit disparity we document in Section 6? Several pieces of evidence suggest that it does. First, public governmental documents describe the goal of the DDb audit program — the IRS’s primary EITC audit selection tool during our sample period — to be the identification of

Figure 9: Group-Specific Audit Rates by Algorithm



*Notes:* The figure reports estimated audit rates for Black and non-Black EITC claimants that would be induced by the algorithms considered in Figure 8 under an assumed audit rate of 1.45%. The population of tax returns available for audit is based on the NRP sample of taxpayers claiming the EITC; see notes to Figure 8 for additional detail. Status quo refers to the estimated audit rates by race for tax year 2014 returns claiming the EITC, reported in Figure 5. Audit rates by race are estimated using the probabilistic estimator; see Appendix Figure A.18 for results using the linear estimator.

taxpayers who do not meet refundable credit eligibility requirements (G.A.O., 2015).<sup>32</sup> Second, the results in Figure 8 suggest that the predicted refundable credit algorithm serves as a good proxy for the amalgamation of (confidential) IRS programs and algorithms through which EITC audit selection occurs. In particular, the Figure shows that when we select audits based on our predictions of refundable credit overclaiming, where our predictions are based on largely the same features and training data that are available to IRS, we obtain a similar disparity as we observe among the returns actually selected by IRS for audit. Finally, as shown in Appendix Figure A.19, operational audits of EITC returns are strongly associated with predicted refundable credit overclaims; in contrast we observe a much lower association between operational audits and predicted total underreporting. For these reasons, we interpret our results to support the conclusion that the IRS’s choice of algorithmic objective is an important contributor to the observed disparity in EITC audit rates.

The next subsection unpacks this conclusion, exploring why the audit selection objective shapes the distribution of audits by race.

### 7.2.1 Why Algorithmic Objective Shapes Audit Disparity

To understand why the objective of the audit selection algorithm shapes the distribution of audits of EITC claimants by race, we explore racial differences in the distribution of various forms of tax non-compliance. To do so, we draw on the richness of the NRP data to identify specific types of errors present on EITC claimants’ returns, and consider the prevalence of these errors among EITC claimants that would be selected under each of the audit selection algorithms considered in Figure 8. The results are presented in Table 2. As expected, returns selected by the refundable credit oracle tend to have higher refundable credit overclaims than those selected by the total underreporting oracle (\$8,191 vs \$3,909 on average), but smaller total adjustments to the taxpayer’s claimed refund or balance due (\$9,595 vs \$22,578 on

---

<sup>32</sup>Our main results are not limited to DDb-selected audits but Appendix Table A.11 confirms the presence of a disparity for this subset of audited returns.

average).

Table 2: Audit-Selected Tax Returns by Algorithm

	Total Underreporting Oracle	Refundable Credit Oracle	Total Underreporting Prediction	Refundable Credit Prediction
Any Underreporting (%)	100.0	100.0	90.2	89.1
Mean Underreporting (\$)	22,578	9,595	10,164	5,952
Any Refundable Credit Overclaiming (%)	91.4	100.0	81.8	85.5
Mean Refundable Credit Overclaiming (\$)	3,909	8,191	2,174	5,144
Dependent Error Rate (%)	27.1	79.8	3.6	33.1
Head of Household Error Rate (%)	17.7	71.3	4.3	61.8
Any Business Income Underreporting (%)	85.7	27.0	83.9	31.4
Probabilistic Disparity (p.p.)	-0.9	1.1	-1.3	1.8
Linear Disparity (p.p.)	-1.3	1.7	-1.9	2.6

*Notes:* The table reports characteristics of the tax returns that would be selected for audit by each of the algorithms considered in Figure 8 under an assumed audit rate of 1.45%. The population of tax returns available for audit is based on the NRP sample of taxpayers claiming the EITC; see notes to Figure 8 for additional detail. Noncompliance is measured based on NRP audit adjustments. “Dependent Error Rate” refers to the share of tax returns that had one or more dependents living with the taxpayer reduced upon audit. “Head of Household Error Rate” refers to the share of tax returns filing as head of household but found ineligible to do so upon audit. “Any Business Income Underreporting” refers to the share of tax returns with positive underreporting of Schedule C Net Income. The final two rows report estimated disparity; units are percentage points (0-100).

More striking is the different profile of taxpayer errors pursued by the different audit selection algorithms. The vast majority (80%) of returns selected by the refundable credit oracle contained a dependent error – that is, the audit determined that at least one of the dependents claimed on the return was not eligible to be claimed by the taxpayer – compared to only 27% of returns selected by the underreporting oracle. In contrast, 86% of returns selected by the underreporting oracle underreported income from a taxpayer’s

business, compared to only 27% of returns selected by the refundable credit oracle.<sup>33</sup> This pattern arises because eligibility for many refundable credits is linked to children; hence, the detection of erroneously claimed dependents tends to be associated with large reductions in refundable credits. At the same time, the contribution of refundable credit overclaims to total underreporting is necessarily bounded by the maximum credit amount, whereas detecting underreported income can lead to arbitrarily large increases in tax liability. Hence, the largest adjustments to total tax – i.e., those prioritized by the underreporting oracle – primarily stem from returns with large amounts of underreported income.

We next explore whether the type of error on which an algorithm focuses shapes the racial distribution of taxpayers selected for audit. We find evidence that it does. In particular, we estimate that erroneously claimed dependents are more common among Black taxpayers than non-Black taxpayers (Table 3). As a result, audit selection processes that are (implicitly) trained to detect dependent mistakes tend to select Black taxpayers at higher rates.<sup>34</sup>

In contrast, the bottom rows of Table 3 suggest that Black taxpayers are less likely to be selected by audit algorithms focused on total underreporting because they are disproportionately under-represented among those taxpayers with the largest amounts of underreported business income. That is, among EITC claimants are a group of taxpayers who have high (actual) incomes but who underreport their incomes to such a degree that they appear eligible for the EITC. Auditing the business income and deductions claimed by taxpayers in this group yields large upward adjustments to tax liability, which is why they are prioritized by algorithms that focus on total underreporting. Because taxpayers in this group are disproportionately non-Black, audit algorithms focused on total underreporting tend to select Black taxpayers at lower rates.

---

<sup>33</sup>The two prediction algorithms differ from one another in a similar manner.

<sup>34</sup>It is beyond the scope of this paper to explore why child-claiming errors vary by race, but one factor that likely contributes is the lower marriage rate among Black Americans in conjunction with credit eligibility rules that prevent cohabitating individuals from claiming the child of an unmarried partner, even if they contribute to the child’s support (Maag et al., 2016; Goldin and Jurow Kleiman, 2021). Along these lines, Micheltore and Pilkauskas (2022) find that Black children are more likely to reside in families with complex or ambiguous tax filing situations.

Table 3: Types of Tax Noncompliance by Race

	Probabilistic Estimate		Linear Estimate	
	Black Taxpayers	Non-Black Taxpayers	Black Taxpayers	Non-Black Taxpayers
Claims Dependent	73.0	69.4	74.6	69.3
Dependent Error Rate	26.6	16.3	30.8	15.4
Head of Household Error Rate	33.5	19.3	39.1	17.8
Any Business Income	19.0	21.7	17.9	22.0
Any Business Income Underreporting	15.9	18.1	15.0	18.4
Underreported Business Income is Among Top...				
10%	1.03	1.97	0.66	2.06
5%	0.40	1.01	0.17	1.07
1%	0.05	0.21	-0.01	0.23

*Notes:* The table reports estimates of the characteristics of tax returns filed by Black and non-Black taxpayers. Analyses are based on the NRP sample of taxpayers claiming the EITC. All analyses account for NRP sampling weights. “Dependent Error Rate” refers to the share of tax returns that had one or more dependents living with the taxpayer reduced upon audit. “Head of Household Error Rate” refers to the share of tax returns filing as head of household but found ineligible to do so upon audit. Business Income refers to net income reported on Schedule C of a tax return. The final three rows report the shares of taxpayers with business income underreporting within the top 90th, 95th, and 99th percentiles of the distribution of positive business income underreporting.

### 7.2.2 Algorithmic Objective and IRS Operational Considerations

Modifying the objective of the IRS audit selection algorithm could have a range of downstream implications for the agency’s enforcement operations beyond the change in the identity of the audited taxpayers. One channel through which this could occur is a change in the composition of the issues upon which audits focus. As discussed in the prior subsection, switching from an algorithm focused on overclaimed refundable credits to one



focused on total underreporting is likely to shift the focus of examinations away from dependent eligibility issues and toward issues related to the proper reporting of business income and deductions.

What would this change in focus mean from an operational perspective? A likely consequence would be an increase in per-return auditing costs.<sup>35</sup> EITC returns differ from one another in terms of the number of hours required to conduct an audit as well as the training and experience required of the auditor. An important driver of this variation is the presence of business income: on average, EITC returns with substantial business income cost the agency \$369.70 per return to audit, compared with \$23.09 for other EITC returns.<sup>36</sup> By increasing the share of audits of EITC claimants that are focused on issues relating to business income, the change in algorithmic objective would increase average auditing costs by increasing the examiner resources required to conduct each exam.

To quantify the magnitude of this effect, Appendix Figure A.20 plots audit costs by algorithm based on the share of returns with substantial business income that each algorithm selects. Holding the current EITC audit rate fixed, switching from an algorithm focused on refundable credit overclaims to one focused on total underreporting would increase the share of audited returns with substantial business income from 3% to 93%, and would raise EITC examination costs by nearly an order of magnitude (Appendix Table A.15). This result suggests that switching algorithmic objectives may, in practice, require that the IRS reduce the total number of EITC returns it audits, at least in the short-term where the amount of examiner resources available for EITC audits may be relatively fixed.<sup>37</sup>

---

<sup>35</sup>The merits of accounting for auditing costs in selection raises difficult normative considerations that are outside the scope of this project; for example, it may be unfair (or generate perverse incentives) for the IRS to avoid auditing those taxpayers most likely to vigorously contest their assessment. Similarly, the cheapest correspondence audits for the IRS to conduct are those for which a refundable credit is disallowed because the taxpayer does not respond. However, non-response to an audit does not necessarily signal ineligibility for a credit, and is more common among Black taxpayers.

<sup>36</sup>These cost estimates are calculated from operational audits based on the average time logged by IRS employees dealing directly with the case multiplied by the applicable General Schedule payscale given the employee level. The classification of returns into those with and without substantial business income follows the IRS's internal classification of EITC returns into *activity codes* 270 and 271 based on whether the return reports gross business receipts in excess of \$25,000 (see Appendix Table A.14).

<sup>37</sup>To illustrate the potential effects of this type of constraint, Appendix Figure A.21 replicates Figure 8,

At the same time, Appendix Figure A.20 also shows that the increase in detected underreporting from the change in algorithmic objective would substantially exceed the increase in audit costs. Intuitively, EITC returns with business income are more difficult to audit, but yield much higher adjustments on average when an audit is undertaken. As such, the qualitative pattern in Figure 8 is largely unchanged when audits are allocated in a manner that accounts for the higher cost of auditing EITC returns claiming business income (Appendix Figure A.22). We therefore conclude that although prioritizing total underreporting would increase per-return audit costs, reforms to reduce disparity need not conflict with the policy goal of deploying audits to efficiently detect tax noncompliance.<sup>38</sup>

### 7.3 Algorithmic Bias in Refundable Credit Prediction

Thus far we have focused on the difference between the disparity induced by the total underreporting algorithms compared to the disparity induced by the overclaimed refundable credit algorithms. However, Figure 8 also highlights that the observed status quo disparity (1.96 p.p.) is substantially larger than the disparity induced by the refundable credit oracle (1.08 p.p.). This suggests that some of the observed disparity is attributable to factors beyond the choice of algorithmic objective. In other words, the status quo disparity is larger than what we would expect to observe even conditioning on the audit selection objective being the detection of refundable credit overclaims. Additional evidence for this hypothesis is provided in Appendix Figure A.25, which shows that racial disparities in audit rates remain when comparing EITC claimants with similar levels of refundable credit overclaims.

Beyond algorithmic objective, one factor that could contribute to the observed audit holding fixed the share of audits allocated to returns with substantial business income. Doing so raises the share of Black taxpayers selected by the total underreporting algorithms, but not to the level of the refundable credit algorithms or to the status quo disparity.

<sup>38</sup>The analyses in this subsection assume the operational audits from which the cost estimates were derived were focused on the same issues as would be selected under the various algorithms considered. We obtain qualitatively similar results when relying on a more conservative (potentially inflated) measure of auditing costs, derived from the hours that auditors report spending on NRP audits (which cover nearly all potential issues on a return); see Appendix Figures A.23 and A.24.

disparity is if the prediction errors underlying the status quo audit selection algorithm were unevenly distributed by race. A large literature in algorithmic fairness documents how data-driven predictive models may exacerbate ground-truth group-level differences in the outcome being predicted (e.g., Leino et al., 2018; Reich, 2021). Figure 8 provides some initial evidence that this may be occurring in our setting: the disparity induced by the refundable credit prediction model (1.75 p.p.) is similar in magnitude to the status quo disparity and approximately 60% larger than the disparity induced by the refundable credit oracle. Since the refundable credit model shares the same objective and data as the IRS, this raises the possibility that the predictive models used by IRS for pursuing refundable credit overclaims may be amplifying the audit disparity that would emerge due to actual differences by race in the distribution of refundable credit overclaiming and selection along those lines.

In Online Appendix F, we investigate the role of differential prediction errors by race in generating the observed disparity. Beginning with our simulated algorithms, we find that errors in predicting refundable credit overclaims are unevenly distributed by race in a manner that may contribute to higher audit rates for Black taxpayers. Turning to IRS operational models, we document uneven errors in the risk measure used by the IRS’s DDb model to predict the likelihood that a child has been incorrectly claimed. We also investigate missingness of parental information on the birth certificate data that IRS uses to help determine whether the children claimed on tax returns meet the required eligibility criteria. We find that whereas birth certificates are missing maternal information at roughly equal rates by race, the birth certificates of children claimed on the returns of Black taxpayers are substantially more likely to be missing information about the identity of the father. Finally, we provide suggestive evidence that there may be opportunities to reduce audit disparities without substantially degrading accuracy by adjusting the features used to form predictions about overclaimed refundable credits.<sup>39</sup>

---

<sup>39</sup>The Online Appendix explores a number of additional factors apart from prediction errors that could potentially also contribute to the observed audit disparity. As detailed in that appendix, we find no evidence

## 8 Conclusion

In this paper, we found that Black taxpayers are audited at a higher rate than non-Black taxpayers, and that this is primarily due to differences in the audit rate between Black and non-Black EITC claimants. In addition, we found that the audit selection objective for EITC returns contributes to this observed disparity, in conjunction with differences by race in the types of errors that EITC claimants tend to make. In particular, we found that designing audit selection algorithms to maximize the detection of overclaimed refundable credits leads to Black EITC claimants being selected at higher rates, whereas designing such algorithms to maximize the detection of total underreporting (from any source) yields the opposite pattern. We interpret our results to suggest that policymakers seeking to reduce the observed audit disparity should consider reorienting audit selection (at least in part) around the goal of maximizing the detection of total underreporting rather than prioritizing noncompliance from refundable credit overclaims. More generally, while the proper focus of tax enforcement has been debated for decades, our findings shed new light on how competing enforcement priorities shape the distribution of audits by race.

We emphasize that implementing a change in audit selection objectives is not as simple as swapping one predictive model for another. In particular, we found that modifying the algorithmic objective to focus on total underreporting would have important downstream effects on the composition of audited EITC returns – shifting audits away from issues of dependent eligibility and toward issues relating to the accuracy of taxable business income. Audits in the latter category require more resources — both auditor time and expertise — compared to most audits of EITC returns today. Although we found that the increase in detected underreporting would exceed the increase in examiner costs, implementing this

---

that the observed disparity is substantially driven by racial differences in the distribution of EITC claimants’ reported income, household composition, or use of a tax preparer. We also provide evidence that the observed disparity is not driven by racial differences in the distribution of high-information (“smoking gun”) signals of non-compliance, or by the use of “regression” models trained to maximize dollars of detected refundable credit overclaims as opposed to “classification” models trained to maximize the probability that audited returns overclaim a refundable credit by at least some (specified) amount.

type of change in the short-term would likely require that the IRS conduct fewer audits of EITC returns while devoting more resources per audit. Longer term, the IRS may be able to reduce the cost of auditing the issues identified by algorithms focused on total underreporting by conducting a greater share of such audits by correspondence (e.g., mailing requests for documentation of claimed business deductions). The IRS recently announced that it plans to study this possibility (IRS, 2023b).

Although the IRS is responsible for audit selection, some of the factors we have identified as contributing to the observed audit disparity are shaped by forces outside the IRS’s control. For example, Congress sets the rules governing credit eligibility — which may contribute to more mistakes among Black taxpayers due to racial differences in family structure — and IRS funding — which shapes the ability of the agency to allocate resources to complex cases. In contrast, the IRS is not compelled by Congress to prioritize refundable credit overclaims over other forms of tax non-compliance.

We note several limitations to our work. First, we have focused our investigation on audit disparities for Black taxpayers, which has been the subject of significant scholarly and policy interest. Investigating audit disparities for other racial and ethnic groups is an important avenue for future work, with potentially differing causes and opportunities for mitigation.<sup>40</sup> Similarly, although we have focused on audit disparities among EITC claimants, our analysis also suggests (smaller) disparities among higher income taxpayers, for whom we cannot rule out the possibility of disparate treatment (since audit selection of higher income taxpayers may not be fully automated).

Second, we have focused on audit selection algorithms designed to detect underreporting, but policymakers may prioritize other objectives as well, such as deterring non-compliant behavior, avoiding audits of compliant taxpayers, transparency, and promoting a fair distribution of audited returns by income. Relatedly, our analysis of counterfactual audit algorithms does not account for the full set of constraints facing tax

---

<sup>40</sup>See Derby et al. (2024) for a discussion of these issues and some evidence concerning audit disparities with respect to other racial and ethnic groups.

authorities like the IRS, such as the types of compliance issues that can be explored through correspondence audit, or differences in audit response rates or dollars collected depending on whether the audit is pre- versus post-refund. A more complete optimal policy analysis would require accounting for these additional considerations.

Finally, audit selection constitutes only one dimension in which tax administration may differently affect taxpayers by race. Disparities may also exist in how the audit is conducted, and with respect to such processes as collections, appeals, settlements, and guidance (Bearer-Friend, 2021; Book, 2021). The approach described in this paper can serve as a foundation to explore disparities in these areas as well.

## References

- Alao, R., Bogen, M., Miao, J., Mironov, I., and Tannen, J. (2021). How meta is working to assess fairness in relation to race in the u.s. across its products and systems. (Cited on 6)
- Andrus, M., Spitzer, E., Brown, J., and Xiang, A. (2021). What we can’t measure, we can’t understand: Challenges to demographic data procurement in the pursuit of fairness. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pages 249–260. (Cited on 6)
- Angwin, J., Larson, J., Mattu, S., and Kirchner, L. (2016). Machine bias. *ProPublica*, May, 23(2016):139–159. (Cited on 2)
- Anson-Dwamena, R., Pattath, P., and Crow, J. (2021). Imputing missing race and ethnicity data in covid-19 cases. (Cited on 6)
- Aughinbaugh, A., Robles, O., and Sun, H. (2013). Marriage and divorce: Patterns by gender, race, and educational attainment. *Monthly Lab. Rev.*, 136:1. (Cited on 17)
- Barocas, S. and Selbst, A. D. (2016). Big data’s disparate impact. *California law review*, pages 671–732. (Cited on 7)
- Bearer-Friend, J. (2019). Should the IRS Know Your Race? The Challenge of Colorblind Tax Data. *Tax L. Rev.*, 73:1. (Cited on 6)
- Bearer-Friend, J. (2021). Colorblind tax enforcement. *NYU Law Review*. (Cited on 50)
- Black, E., Elzayn, H., Chouldechova, A., Goldin, J., and Ho, D. (2022). Algorithmic fairness and vertical equity: Income fairness with irs tax audit models. In *2022 ACM Conference on Fairness, Accountability, and Transparency*, pages 1479–1503. (Cited on 7, Appendix-78)
- Bloomquist, K. M. (2019). Regional bias in irs audit selection. *Tax Notes*. (Cited on 3, 6)
- Boning, W. C., Hendren, N., Sprung-Keyser, B., and Stuart, E. (2023). A welfare analysis of tax audits across the income distribution. Technical report, National Bureau of Economic Research. (Cited on 5)
- Book, L. (2021). Tax administration and racial justice: The illegal denial of tax based pandemic relief to the nation’s incarcerated. *South Carolina Law Review*, 72. (Cited on 50)
- Brown, D. A. (2005). The tax treatment of children: Separate but unequal. *Emory LJ*, 54:755. (Cited on 6)
- Brown, D. A. (2009). Shades of the american dream. *Wash. UL Rev.*, 87:329. (Cited on 6)
- Brown, D. A. (2018). Homeownership in black and white: The role of tax policy in increasing housing inequity. *U. Mem. L. Rev.*, 49:205. (Cited on 6)

- Brown, D. A. (2021). *The Whiteness of Wealth: How the Tax System Impoverishes Black Americans—And How We Can Fix It*. Crown Publishing Group (NY). (Cited on 6)
- Buolamwini, J. and Gebru, T. (2018). Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on fairness, accountability and transparency*, pages 77–91. PMLR. (Cited on 2)
- CFPB (2014). *Using Publicly Available Information to Proxy for Unidentified Race and Ethnicity: A methodology and assessment*. Consumer Financial Protection Bureau. (Cited on 6)
- Chen, J., Kallus, N., Mao, X., Svacha, G., and Udell, M. (2019). Fairness under unawareness: Assessing disparity when protected class is unobserved. In *Proceedings of the conference on fairness, accountability, and transparency*, pages 339–348. (Cited on 2, 7, 17, Appendix-47)
- Chetty, R., Hendren, N., Jones, M. R., and Porter, S. R. (2020). Race and economic opportunity in the united states: An intergenerational perspective. *The Quarterly Journal of Economics*, 135(2):711–783. (Cited on 6)
- Collyer, S., Harris, D., and Wimer, C. (2019). Left behind: The one-third of children in families who earn too little to get the full child tax credit. (Cited on 6)
- Congressional Research Service (2022). Audits of eite returns: By the numbers. (Cited on 12)
- Cook, L. D., Logan, T. D., and Parman, J. M. (2016). The mortality consequences of distinctively black names. *Explorations in Economic History*, 59:114–125. (Cited on 17)
- Dean, S. A. (2021). *Testimony to House Committee on Ways and Means*. (Cited on 6)
- Derby, E., Dowd, C., and Mortenson, J. (2024). Assessing statistical bias in racial and ethnic disparity estimates using bifsg. *Working Paper*. (Cited on 31, 49)
- Fong, C. and Tyler, M. (2021). Machine learning predictions as regression covariates. *Political Analysis*, 29(4):467–484. (Cited on 7)
- Fryer, R. G. and Levitt, S. D. (2004). The causes and consequences of distinctively black names. *The Quarterly Journal of Economics*, 119(3):767–805. (Cited on 17)
- Goldin, J. and Jurow Kleiman, A. (2021). Whose child is this? improving child-claiming rules in safety-net programs. *Yale Law Journal*, 131:1719. (Cited on 43)
- Goldin, J. and Micheltore, K. (2022). Who benefits from the child tax credit? *National Tax Journal*. (Cited on 6)
- Government Accountability Office (2015). Irs return selection: Wage and investment division should define audit objectives and refine other internal controls. (Cited on 10, 41, Appendix-77)



- Government Accountability Office (2021). Artificial intelligence: An accountability framework for federal agencies and other entities. (Cited on 6)
- Government Accountability Office (GAO) (2017). Refundable Tax Credits: Comprehensive Compliance Strategy and Expanded Use of Data Could Strengthen IRS’s Efforts to Address Noncompliance. (Cited on 8)
- Greengard, P. and Gelman, A. (2023). Bisg: When inferring race or ethnicity, does it matter that people often live near their relatives? *arXiv preprint arXiv:2304.09126*. (Cited on 14)
- Guyton, J., Langetieg, P., Reck, D., Risch, M., and Zucman, G. (2021). Tax evasion at the top of the income distribution: theory and evidence. (Cited on 33)
- Guyton, J., Leibel, K., Manoli, D. S., Patel, A., Payne, M., and Schafer, B. (2018). The effects of eitc correspondence audits on low-income earners. (Cited on 2)
- Haas, A., Elliott, M. N., Dembosky, J. W., Adams, J. L., Wilson-Frederick, S. M., Mallett, J. S., Gaillot, S., Haffer, S. C., and Haviland, A. M. (2019). Imputation of race/ethnicity to enable measurement of hedis performance by race/ethnicity. *Health Services Research*, 54(1):13–23. (Cited on 6)
- Hardy, B., Hokayem, C., and Ziliak, J. P. (2021). Income inequality, race, and the eitc. *Working paper*. (Cited on 6)
- Holtzblatt, J. and McCubbin, J. (2004). Issues affecting low-income filers. *The crisis in tax administration*, 148:148–49. (Cited on 11)
- Imai, K. and Khanna, K. (2016). Improving ecological inference by predicting individual ethnicity from voter registration records. *Political Analysis*, pages 263–272. (Cited on 3, 6, 13)
- Internal Revenue Service (IRS) (2011). Internal revenue manual 4.19.20. (Cited on 12)
- Internal Revenue Service (IRS) (2016). Internal revenue service data book, 2015. (Cited on 11)
- Internal Revenue Service (IRS) (2023a). Commissioner Letter to Chairman Wyden. <https://home.treasury.gov/system/files/136/091823-Wyden-Letter-from-IRS-Commissioner-on-Audit-Disparities.pdf> (Accessed: Feb 2024). (Cited on 49)
- Internal Revenue Service (IRS) (2023b). EITC Fast Facts. <https://www.eitc.irs.gov/partner-toolkit/basic-marketing-communication-materials/eitc-fast-facts/eitc-fast-facts> (Accessed: June 2023). (Cited on 10)
- Johnston, D. C. (April 16, 2000). I.r.s. more likely to audit the poor and not the rich. *New York Times*. (Cited on 12)

- Kallus, N., Mao, X., and Zhou, A. (2021). Assessing algorithmic fairness with unobserved protected class using data combination. *Management Science*. (Cited on 7, 17)
- Kiel, P. and Fresques, H. (2019). Where in The U.S. Are You Most Likely to Be Audited by the IRS? (Cited on 3, 10)
- Kleinberg, J., Lakkaraju, H., Leskovec, J., Ludwig, J., and Mullainathan, S. (2018a). Human decisions and machine predictions. *The quarterly journal of economics*, 133(1):237–293. (Cited on 3)
- Kleinberg, J., Ludwig, J., Mullainathan, S., and Sunstein, C. R. (2018b). Discrimination in the age of algorithms. *Journal of Legal Analysis*, 10:113–174. (Cited on 7)
- Knox, D., Lucas, C., and Cho, W. K. T. (2022). Testing causal theories with learned proxies. *Annual Review of Political Science*, 25(1):null. (Cited on 2)
- Leibel, K., Lin, E., and McCubbin, J. (2020). Social welfare considerations of eitc qualifying child noncompliance. *Treasury Office of Tax Analysis Working Paper*. (Cited on 17)
- Leino, K., Black, E., Fredrikson, M., Sen, S., and Datta, A. (2018). Feature-wise bias amplification. *arXiv preprint arXiv:1812.08999*. (Cited on 47)
- Lu, B., Wan, J., Ouyang, D., Goldin, J., and Ho, D. E. (2024). Quantifying the uncertainty of imputed demographic disparity estimates: The dual-bootstrap. (Cited on 24, Appendix-13, Appendix-35)
- Maag, E., Peters, H. E., and Edelstein, S. (2016). Increasing family complexity and volatility: The difficulty in determining child tax benefits. *Tax Policy Center*. (Cited on 43)
- Micheltore, K. M. and Pilkauskas, N. V. (2022). The earned income tax credit, family complexity, and children’s living arrangements. *RSF: The Russell Sage Foundation Journal of the Social Sciences*, 8(5):143–165. (Cited on 43)
- Moran, B. I. and Whitford, W. (1996). A Black Critique of the Internal Revenue Code. *Wis. L. REv.*, page 751. (Cited on 6)
- National Taxpayer Advocate (2019a). Annual report to congress 2019. (Cited on 2)
- National Taxpayer Advocate (2019b). Report: Making the eitc work for taxpayers and the government. (Cited on 11)
- National Taxpayer Advocate (2021). Annual report to congress 2021. (Cited on 8)
- Nerenz, D. R., McFadden, B., Ulmer, C., et al. (2009). Race, ethnicity, and language data: standardization for health care quality improvement. (Cited on 13)
- Obermeyer, Z., Powers, B., Vogeli, C., and Mullainathan, S. (2019). Dissecting racial bias in an algorithm used to manage the health of populations. *Science*, 366(6464):447–453. (Cited on 2, 7)

- Office of Management and Budget (OMB) (2021). Circular a-123, appendix c: Requirements for payment integrity improvement. (Cited on 11)
- Passi, S. and Barocas, S. (2019). Problem formulation and fairness. In *Proceedings of the conference on fairness, accountability, and transparency*, pages 39–48. (Cited on 7)
- Reich, C. L. (2021). Resolving the disparate impact of uncertainty: Affirmative action vs. affirmative information. *arXiv preprint arXiv:2102.10019*. (Cited on 47)
- Tzioumis, K. (2018). Demographic aspects of first names. *Scientific data*, 5(1):1–9. (Cited on 19, Appendix-29)
- U.S. Census Bureau (2021). Decennial census surname files (2010, 2000).”. <https://www.census.gov/data/developers/data-sets/surnames.html>, Last accessed on 2023-01-12. (Cited on 19)
- U.S. Executive Order 13985 (2021). Exec. order no. 13985 86 fed. reg. 7009, advancing racial equity and support for underserved communities through the federal government. (Cited on 6)
- U.S. Treasury Department (2022). Agency financial report: Fiscal year 2021. (Cited on 11)
- Voicu, I. (2018). Using first name information to improve race and ethnicity classification. *Statistics and Public Policy*, 5(1):1–13. (Cited on 3, 13)

# Online Appendix to Measuring and Assessing Differences in Audit Rates of Black and Non-Black Americans

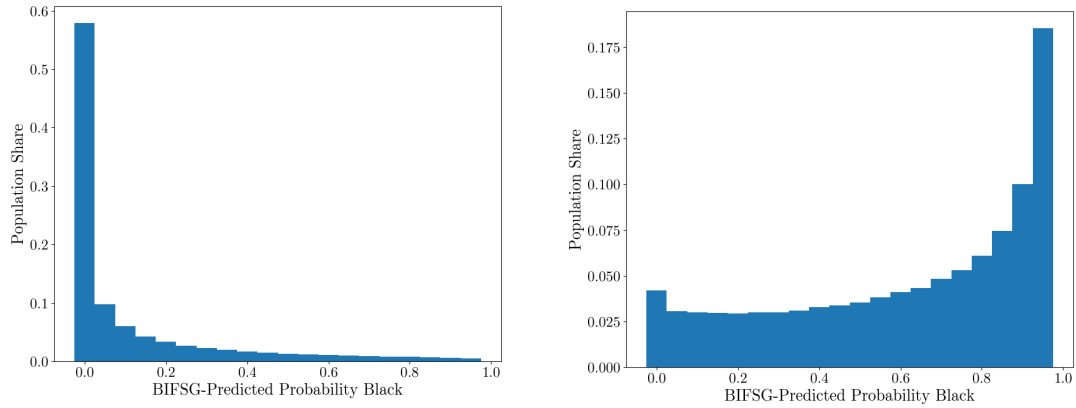
Hadi Elzayn, Evelyn Smith, Thomas Hertz, Cameron Guage, Arun Ramesh, Robin Fisher, Daniel E. Ho, and Jacob Goldin

# List of Appendices

A Additional Tables and Figures	Appendix-2
B Additional Results Relating to Disparity Estimation	Appendix-40
C North Carolina Match and Bias Correction	Appendix-69
D EITC Disparity Decomposition	Appendix-70
E Taxpayer Noncompliance Prediction Model	Appendix-74
F Additional Factors Potentially Contributing to Audit Disparity	Appendix-76
G Appendix References	Appendix-91

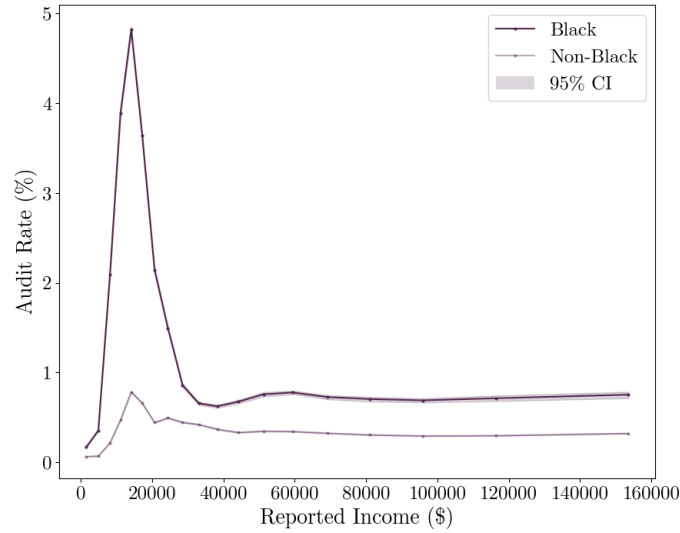
## A Additional Tables and Figures

Figure A.1: Distribution of Race Imputations for Known Black and Non-Black Taxpayers



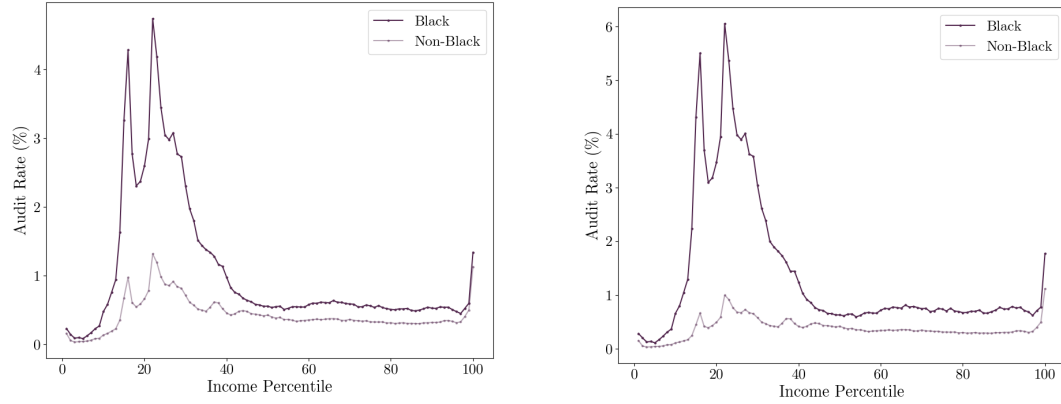
*Notes:* Left: The figure shows the distribution of BIFSG-predicted probabilities that a taxpayer is Black (non-Hispanic) for non-Black taxpayers in our matched North Carolina sample. Right: The figure shows the distribution of BIFSG-predicted probabilities that a taxpayer is Black (non-Hispanic) for Black (non-Hispanic) taxpayers in our matched North Carolina sample.

Figure A.2: Audit Rate Disparity by Income (Linear Estimator)



*Notes:* The figure shows the estimated audit rate by income among Black and non-Black taxpayers filing for tax year 2014 using the linear estimator. Income is measured according to the Adjusted Gross Income (AGI) reported on the taxpayer's return (i.e., prior to audit adjustments). The sample is limited to taxpayers reporting non-negative AGI. Taxpayers are grouped into 20 equal-sized bins. The shaded area around each line shows the 95% confidence interval based on the distribution of estimates from 100 bootstrapped samples. . To facilitate presentation, the x-axis is limited to bins with mean reported AGI under \$200,000; a version of the figure with all income percentiles is presented in Appendix Figure A.3.

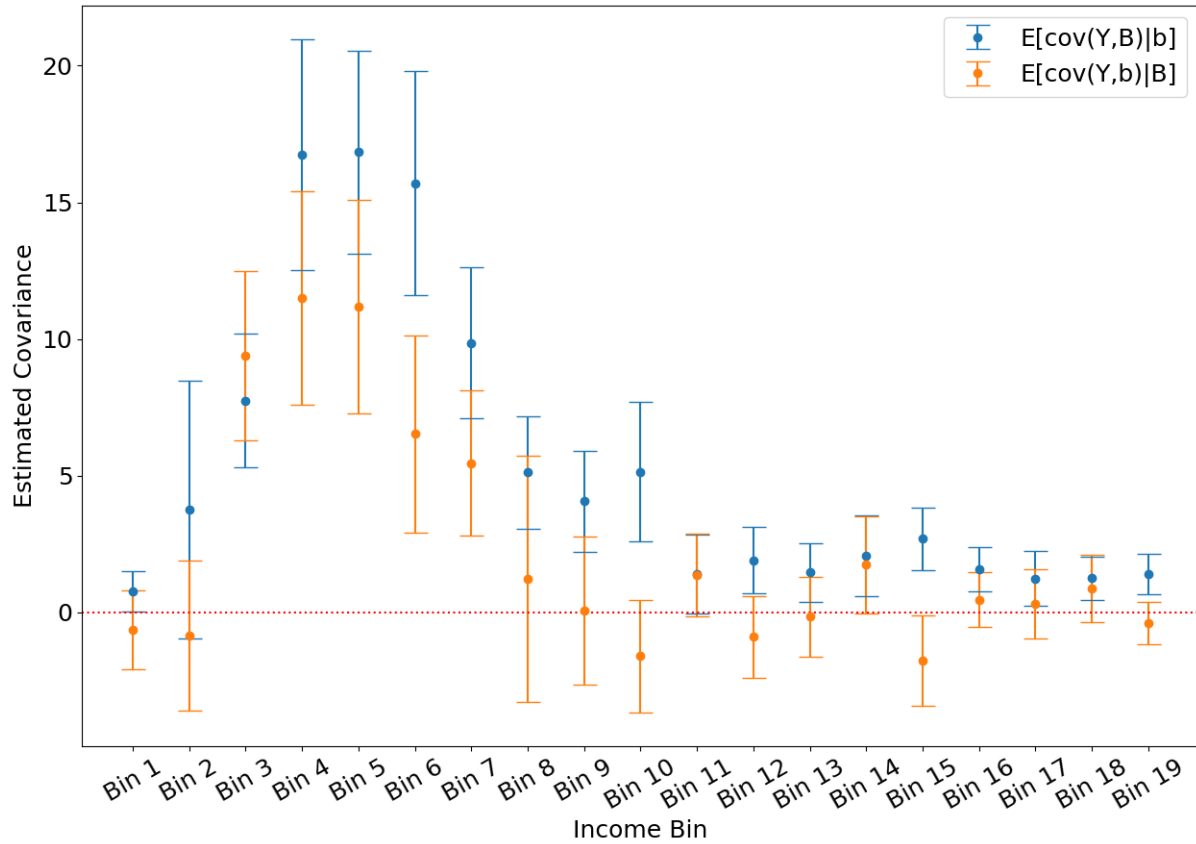
Figure A.3: Audit Rate Disparity by Income (All Income Percentiles)



*Notes:* Left: The figure shows the probabilistic estimate of audit rate by income percentile among Black and non-Black taxpayers filing returns for tax year 2014. Income is measured according to the adjusted gross income (AGI) reported on the taxpayer's return (i.e., prior to audit adjustments). The sample is limited to taxpayers reporting non-negative AGI. Taxpayers have been grouped into 100 equal-sized bins, based on their AGI. Right: The figure shows the corresponding analysis to the left panel, but estimating binned audit rates using the linear estimator instead of the probabilistic estimator.

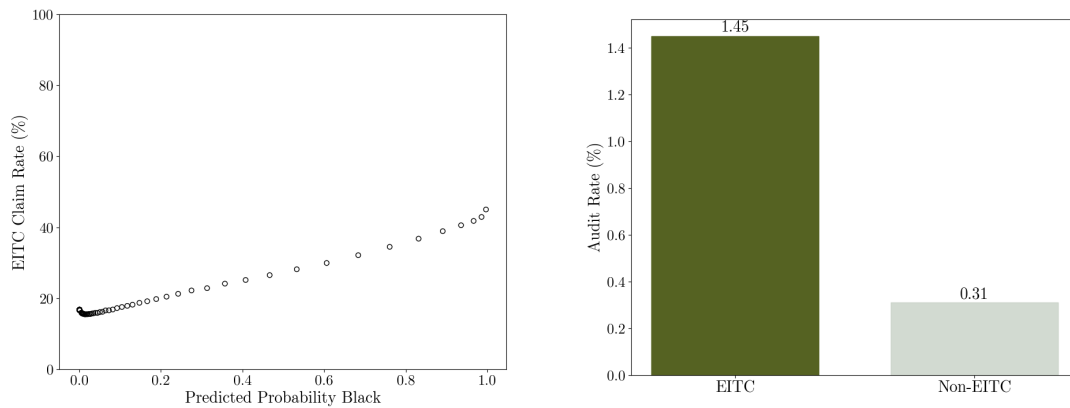


Figure A.4: Covariance Condition Estimates for Income Bins



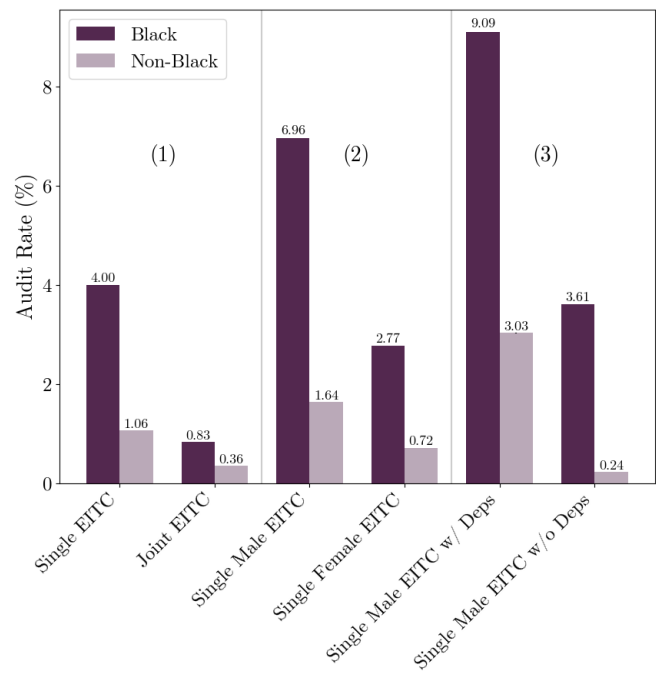
*Notes:* The figure displays the estimated covariance between audits and self-reported race, conditional on estimated race, as well as the estimated covariance between audits and estimated race, conditional on self-reported race, for each income bin reported in Figure 4. The estimates are calculated from the matched sample of North Carolina taxpayers, using the weights described in Appendix C. Brackets denote 95% confidence intervals.

Figure A.5: Audit Rates by EITC Claim Status



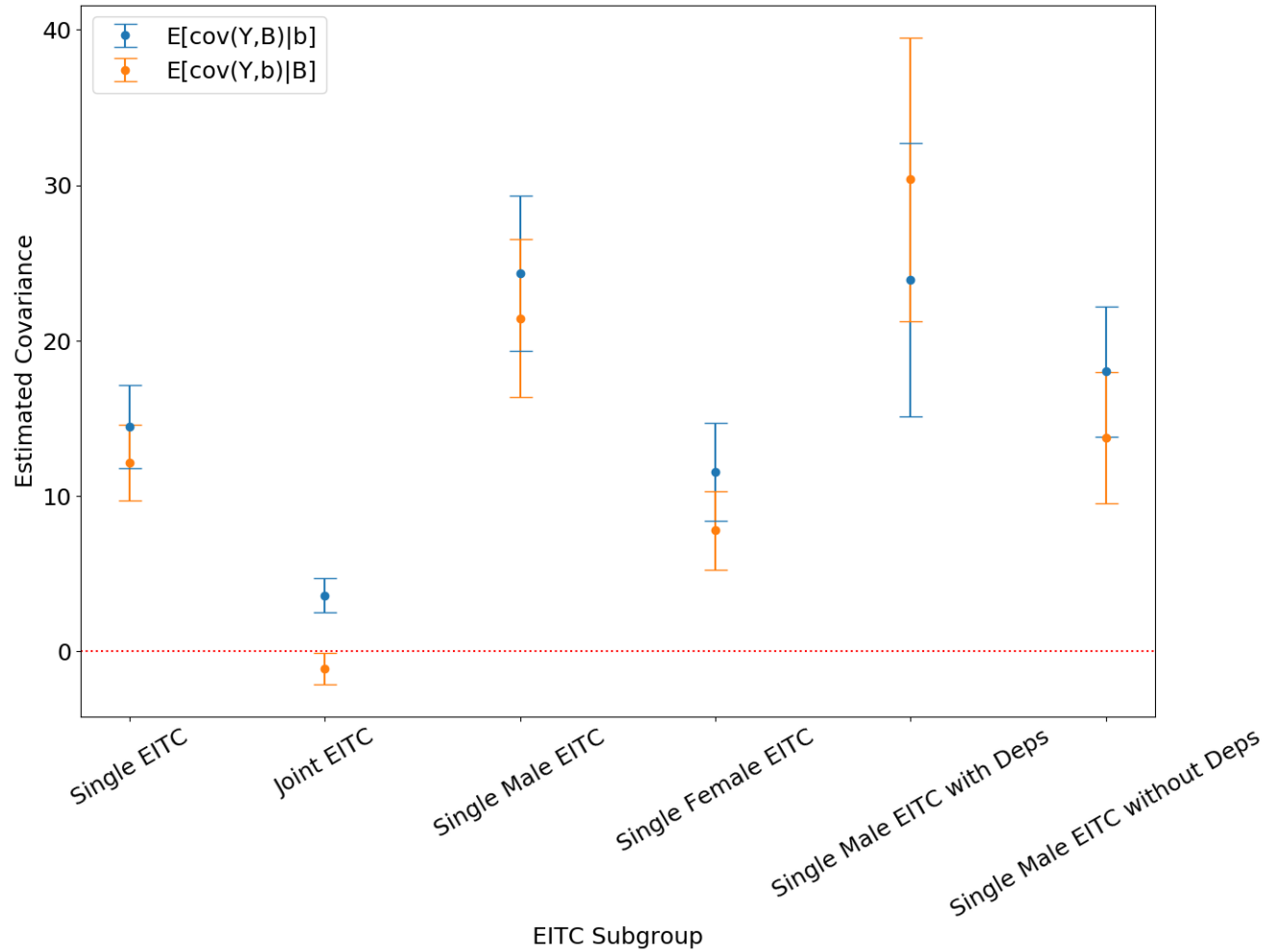
*Notes:* The figure shows the relationship between audits and EITC claim status among taxpayers (of any race) filing returns for tax year 2014. Left: Binned scatterplot of EITC claim rate by BIFSG-predicted probability that a taxpayer is Black. Taxpayers have been grouped into 100 equal-sized bins. Right: Audit rates among EITC claimants and non-EITC claimants.

Figure A.6: Audit Rate Disparities by EITC Subgroup (Linear Estimator)



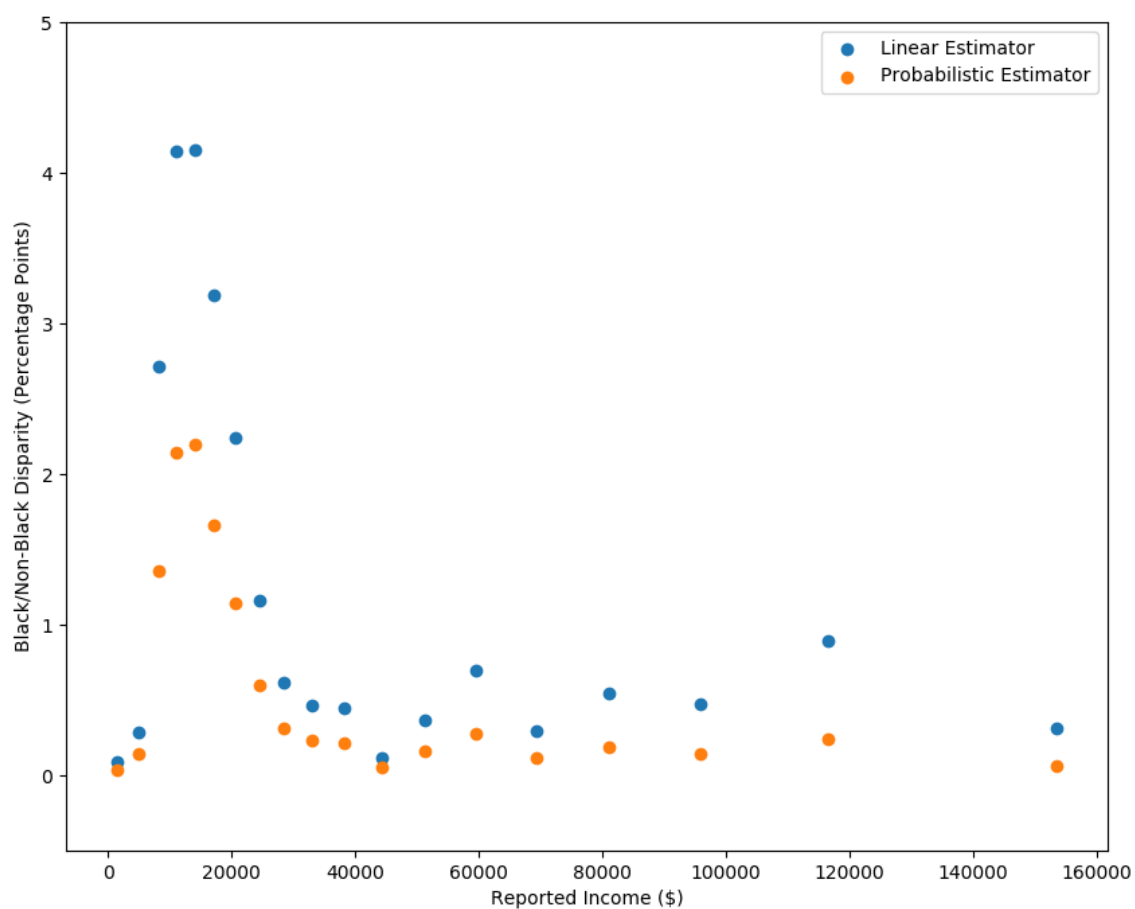
*Notes:* The figure shows the estimated audit rate among the specified subgroups of Black and non-Black taxpayers. Conditional audit rates by race are calculated using the linear audit rate estimator applied to BIFSG-predicted probabilities that a taxpayer is Black. Panel (1) splits EITC claimants by single vs joint filers; (2) splits single EITC claimants by taxpayer gender; and (3) splits single men claiming the EITC by whether they claim dependents.

Figure A.7: Covariance Condition Estimates for EITC Subgroups



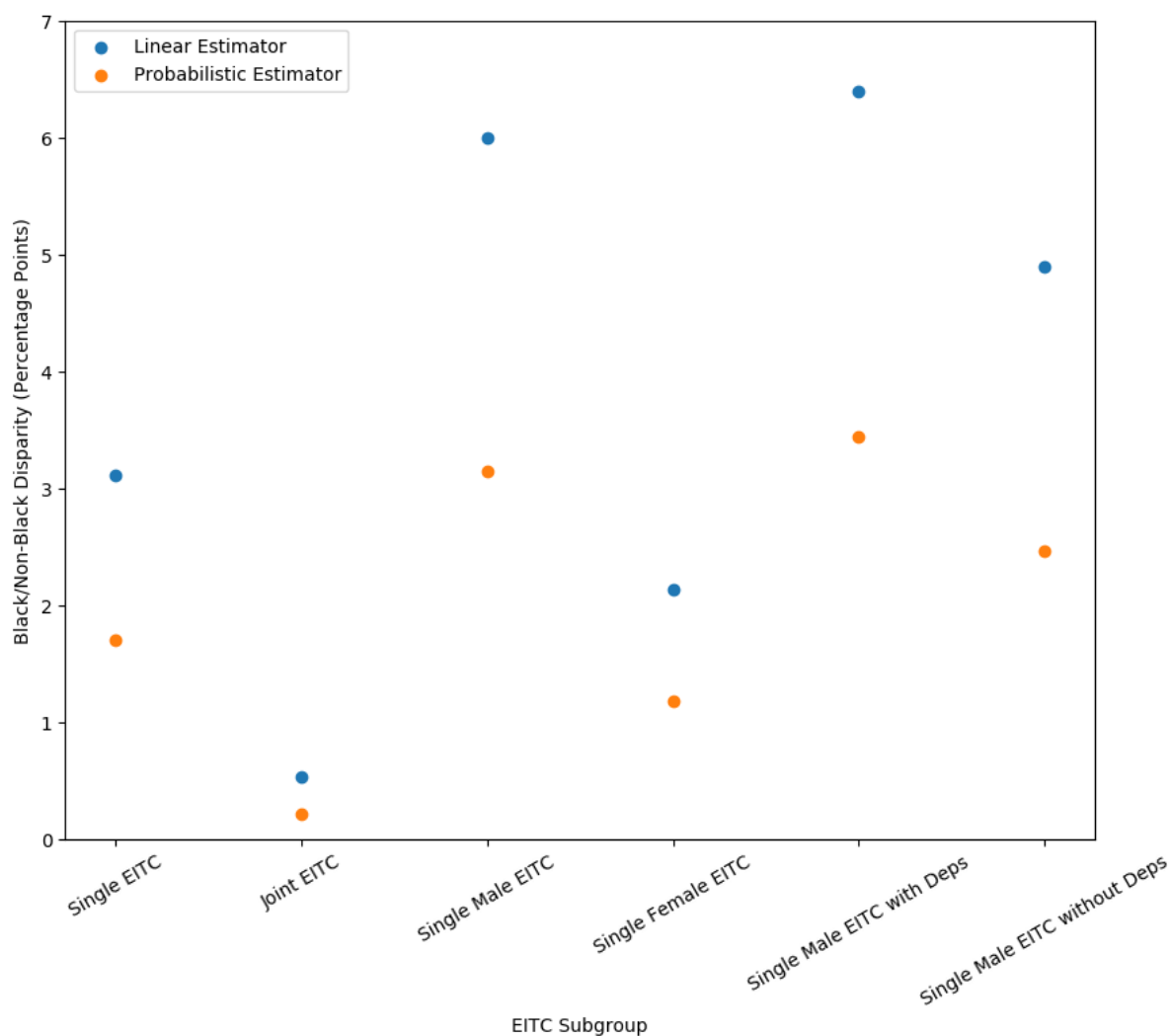
*Notes:* The figure displays the estimated covariance between audits and self-reported race, conditional on estimated race, as well as the estimated covariance between audits and estimated race, conditional on self-reported race, for each income bin reported in Figure 6. The estimates are calculated from the matched sample of North Carolina taxpayers, using the weights described in Appendix C. Brackets denote 95% confidence intervals.

Figure A.8: Disparity Estimates Within Income Bins Using Re-Calibrated BIFSG Probabilities



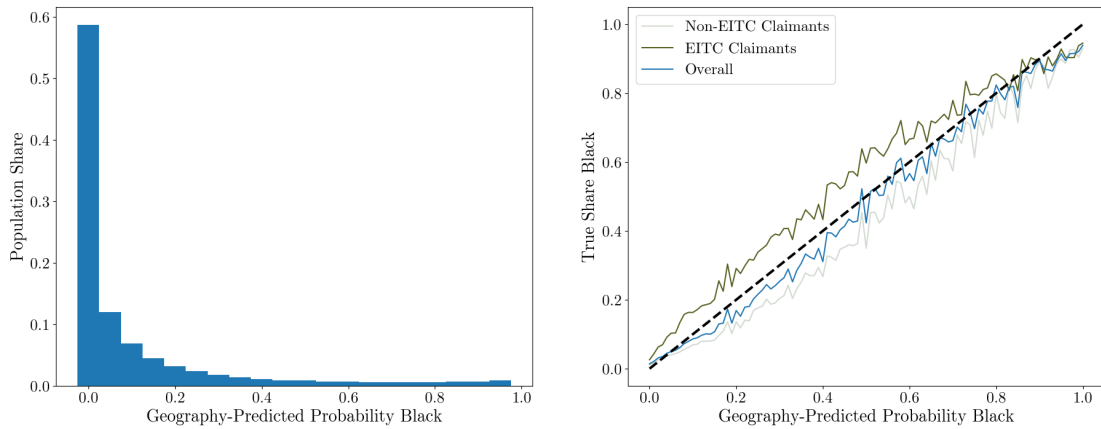
*Notes:* The figure reports linear and probabilistic audit rate disparity estimates in each of the income bins considered in Figure 4 using BIFSG scores recalibrated by the ground truth sample of North Carolina taxpayers in each bin. See Appendix Section B.5 for details on obtaining the recalibrated proxy. The recalibration exercise incorporates the North Carolina weights described in Appendix C.

Figure A.9: Disparity Estimates Within EITC Subgroups Using Re-Calibrated BIFSG Probabilities



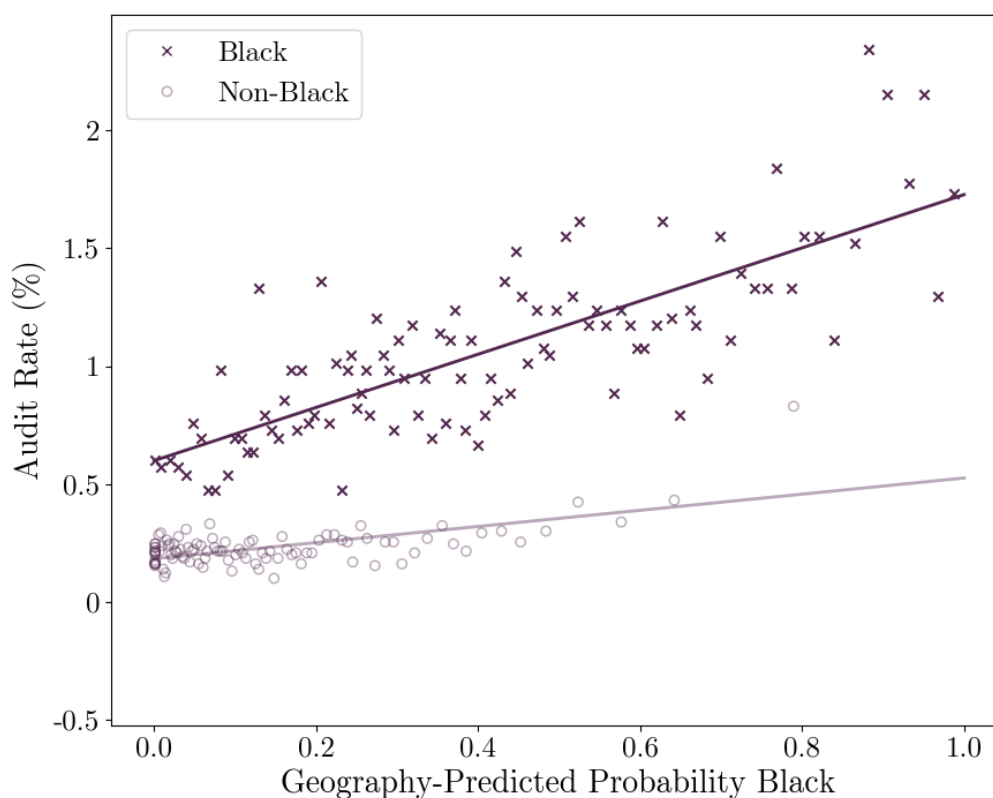
*Notes:* The figure shows probabilistic and linear audit rate disparity estimates corresponding to each of the subgroups considered in Figure 6 using BIFSG scores recalibrated by the ground truth sample of North Carolina taxpayers in each bin. See Appendix Section B.5 for details on obtaining the recalibrated proxy. The recalibration exercise incorporates the North Carolina weights described in Appendix C.

Figure A.10: Distribution and Calibration of Geography-Based Race Imputations



*Notes:* The figure replicates Figure 1 using estimated race probabilities based only on the location of a taxpayer's residence; the probability that a taxpayer is Black is set equal to the fraction of their CBG that is composed of Black residents. Left: Histogram of geography-predicted probabilities that a taxpayer is Black. The mean prediction is 11.8%. Right: The figure shows the calibration of the geographic predictions in the matched North Carolina dataset. Taxpayers are split into groups based on their predicted probability of being Black (discretized into 100 bins 1 percentage point wide). The predicted probability of being Black is on the  $x$ -axis; the  $y$ -axis represents the true proportion of each group that is Black according to the ground-truth race observed in the North Carolina matched sample, re-weighted to be representative of the overall United States (see Appendix C for details). A perfectly calibrated predictor would fall on the 45-degree line, shown as the black dotted line. The figure shows overall calibration in blue as well as calibration among EITC claimants (dark green) and non-EITC claimants (light green).

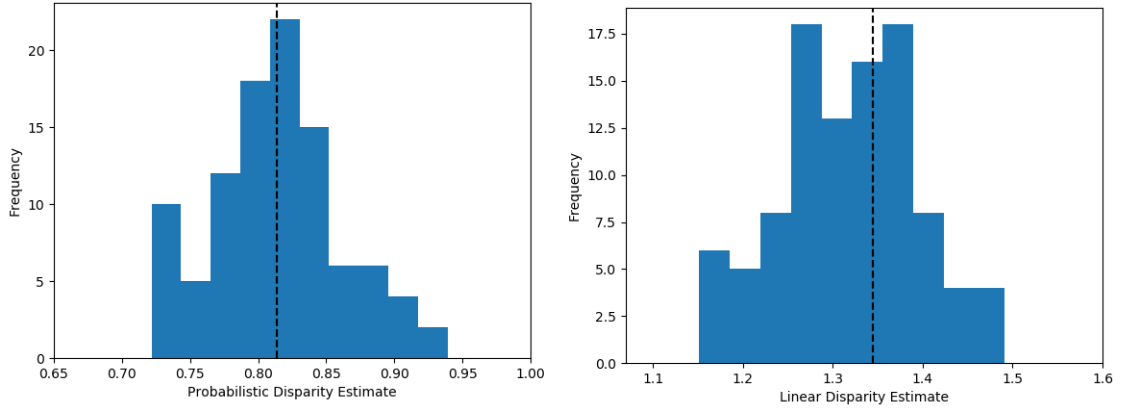
Figure A.11: Audit Rate by Geography-Based Race Imputation and Self-Reported Race



*Notes:* The figure replicates Figure 2 using estimated race probabilities based only on the location of a taxpayer's residence; the probability that a taxpayer is Black is set equal to the fraction of their CBG that is composed of Black residents. The figure shows the relationship between audit incidence and BIG-predicted probability that a taxpayer is Black for taxpayers filing returns for tax year 2014. Audit rates are plotted separately for Black and non-Black taxpayers in the North Carolina matched sample. Black and non-Black taxpayers are each grouped into 100 equal-sized bins, with Black taxpayers indicated by dark purple x's and non-Black taxpayers indicated by light purple circles.

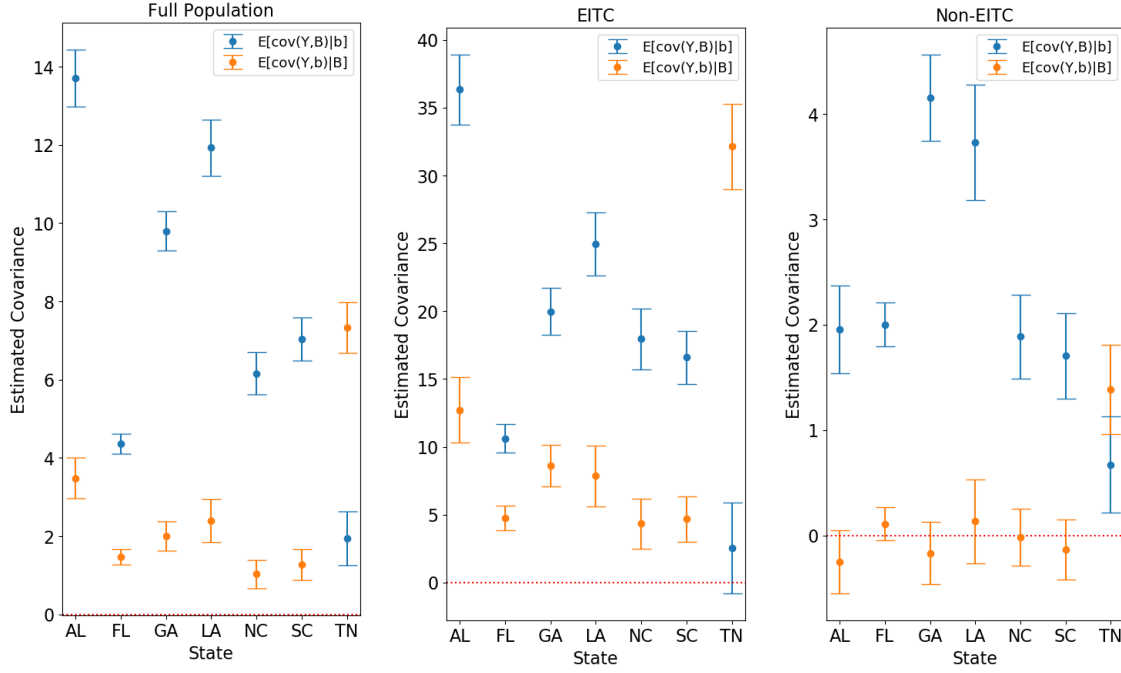


Figure A.12: Distribution of Disparity Estimates from Dual-Bootstrap



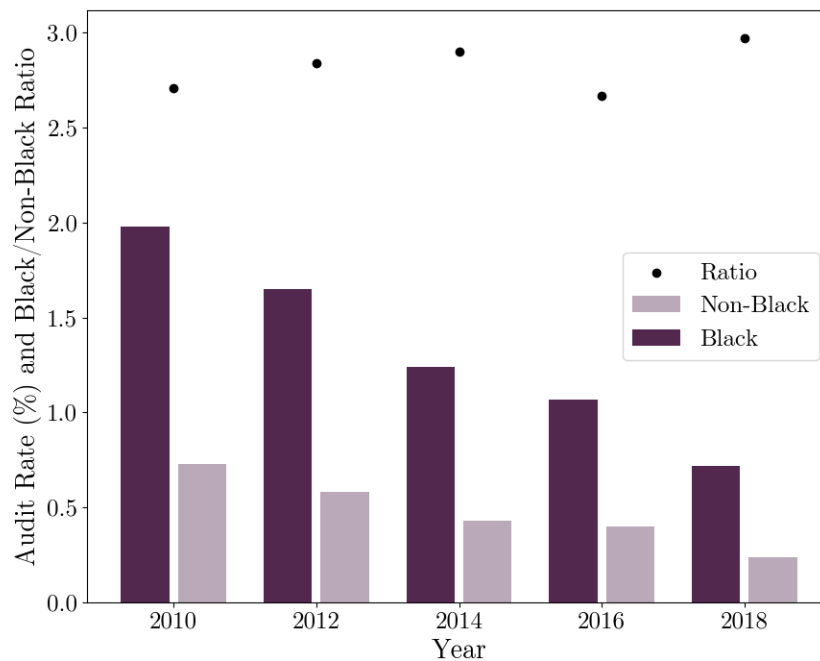
*Notes:* The figure shows the distribution of probabilistic (left panel) and linear (right panel) disparity estimates produced by our implementation of the Lu et al. (2024) dual-bootstrap procedure. To implement the procedure, we resampled first names and geographic information, but did not resample surnames (since the Census surname data we use corresponds to the full population rather than a sample) to generate 100 sets of BIFSG posteriors. We then estimate disparity with each set of BIFSG posteriors, resampling our taxpayer population on each iteration. Units are percentage points (0-100). The dashed line represents the disparity point estimate presented in Table 1.

Figure A.13: Residual Covariance Estimates by State



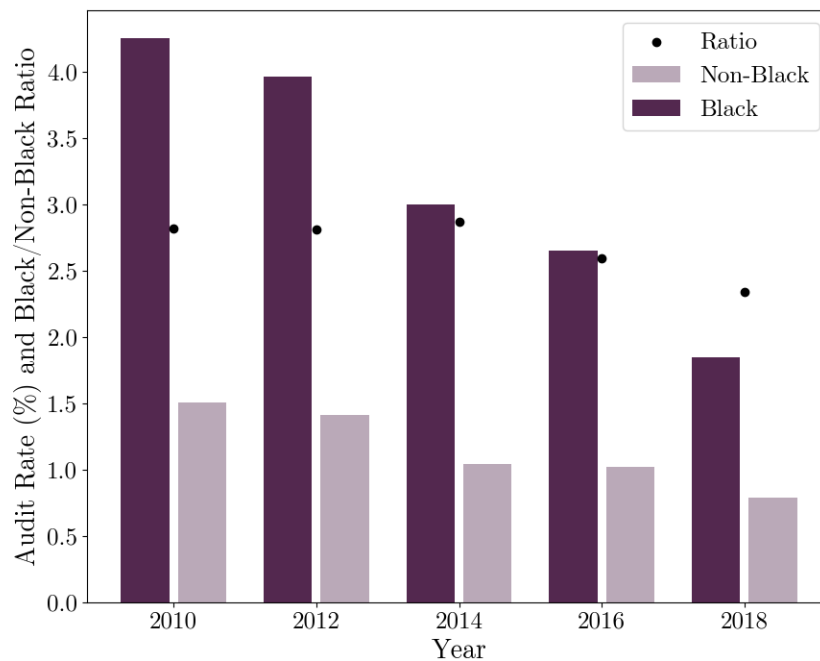
Notes: Left: The figure shows estimates of the covariance between audits ( $Y$ ) and self-reported race ( $B$ ), conditional on estimated race ( $b$ ), as well as the estimated covariance between audits and estimated race, conditional on self-reported race (corresponding to the terms  $\mathbb{E}[\text{Cov}(Y, B|b)]$  and  $\mathbb{E}[\text{Cov}(Y, b|B)]$  respectively) for taxpayers that match to 2023 L2-collected voter data in seven states (match rates are displayed in parentheses): Alabama (37.18%), Florida (30.01%), Georgia (30.94%), Louisiana (38.60%), North Carolina (16.50%), South Carolina (36.21%), and Tennessee (22.63%). The match procedure we use for these states is the same as the procedure we use to match our main North Carolina data set, described in Appendix C. Error bars represent 95% confidence intervals. The displayed estimates and underlying standard errors are multiplied by  $10^4$ . The panels show results for the full population, EITC claimants, and EITC non-claimants, respectively.

Figure A.14: Estimated Audit Rate Disparity by Year



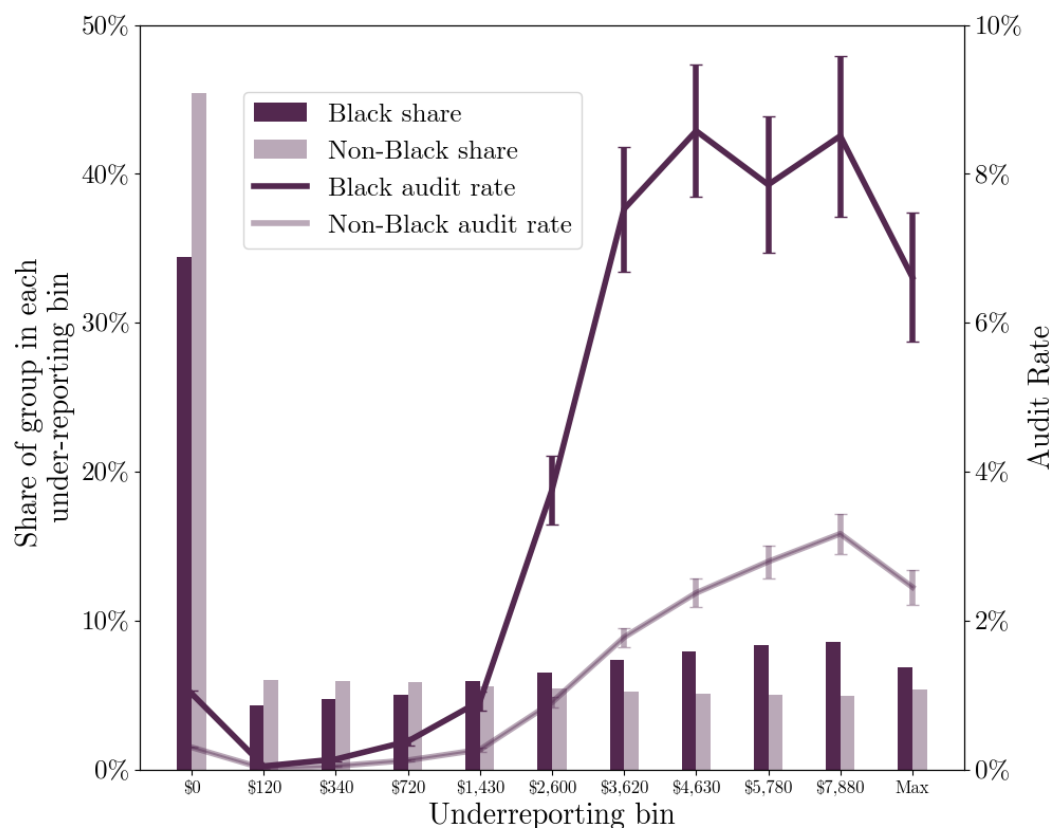
*Notes:* The figure reports the estimated audit rates among Black and non-Black taxpayers for tax years 2010, 2012, 2014, 2016, and 2018, calculated using the probabilistic audit rate estimator applied to BIFSG-predicted probabilities (calculated using the data sources described in Section 4.2). “Ratio” refers to the ratio of the estimated Black audit rate to the estimated non-Black audit rate.

Figure A.15: Estimated Audit Rate Disparity Among EITC Claimants by Year



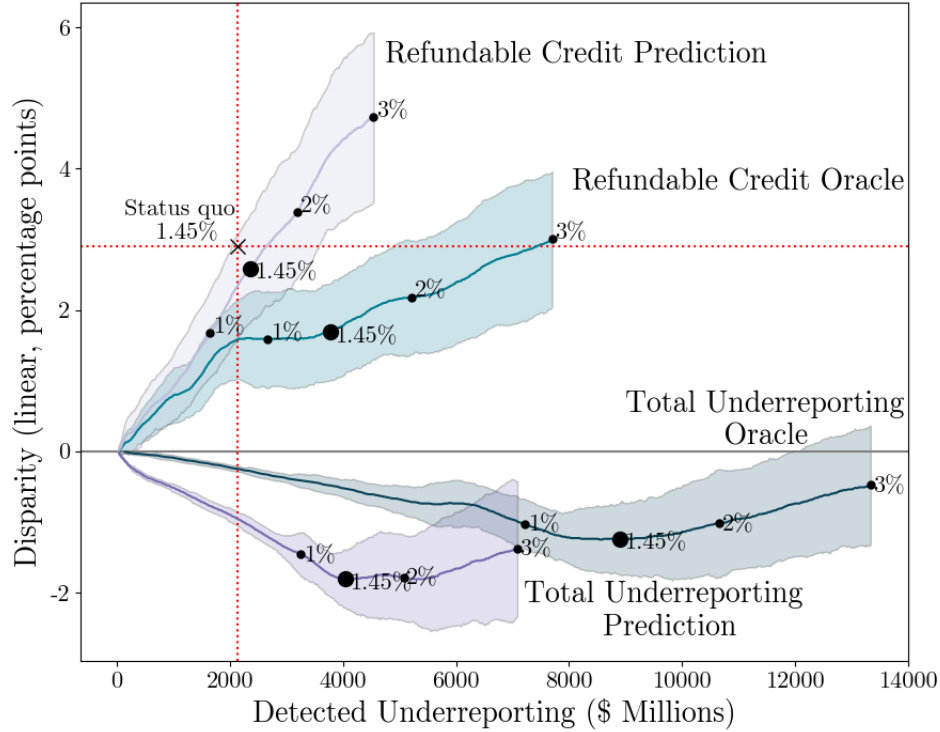
*Notes:* The figure reports the estimated audit rates among Black and non-Black EITC claimants for tax years 2010, 2012, 2014, 2016, and 2018, calculated using the probabilistic audit rate estimator applied to BIFSG-predicted probabilities (calculated using the data sources described in Section 4.2). “Ratio” refers to the ratio of the estimated Black audit rate to the estimated non-Black audit rate.

Figure A.16: Racial Audit Disparity Among EITC Claimants by Underreported Taxes (Linear Estimator)



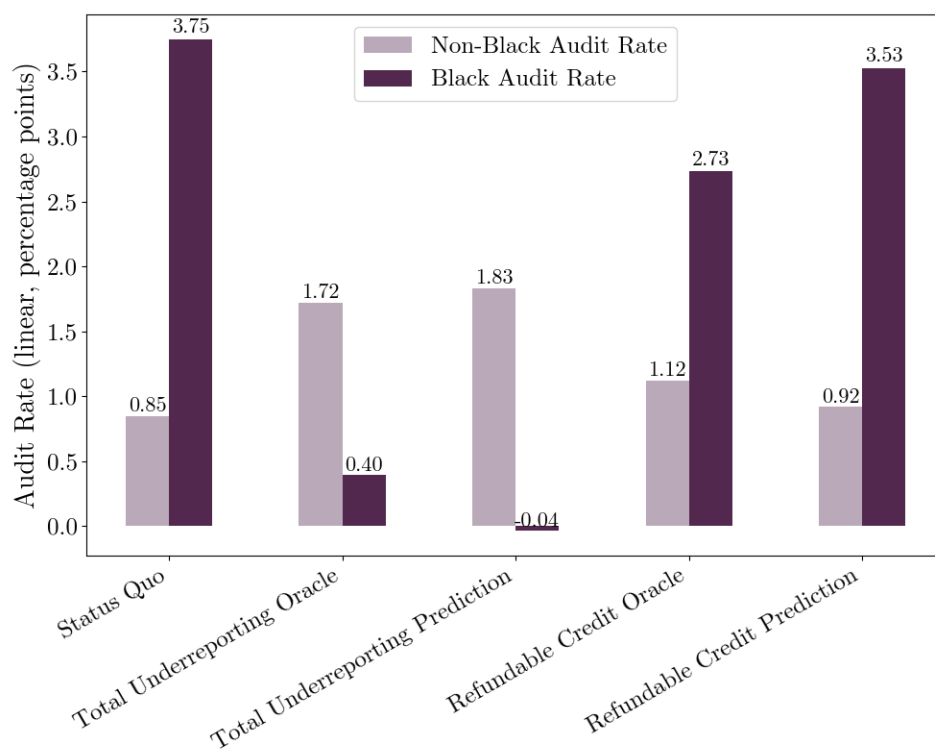
*Notes:* The figure shows the estimated audit rates for Black and non-Black EITC claimants, respectively, by under-reported taxes. Taxpayers are binned into 11 categories: those with less than \$1 of under-reporting, and 10 equal deciles of taxpayers with positive under-reporting. Underreporting deciles are defined based on the distribution of underreporting among EITC claimants, as measured by NRP audits. Estimated audit rates by race are calculated using the linear estimator and the method described in Section 6 of the main text. All analyses account for NRP sampling weights. Brackets reflect the estimated 95% confidence interval, derived from bootstrapped standard errors (N=100). The bars show the estimated share of Black and non-Black taxpayers, respectively, that fall into each under-reporting bin.

Figure A.17: Detected Underreporting and Disparity by Algorithm (Linear Estimator)



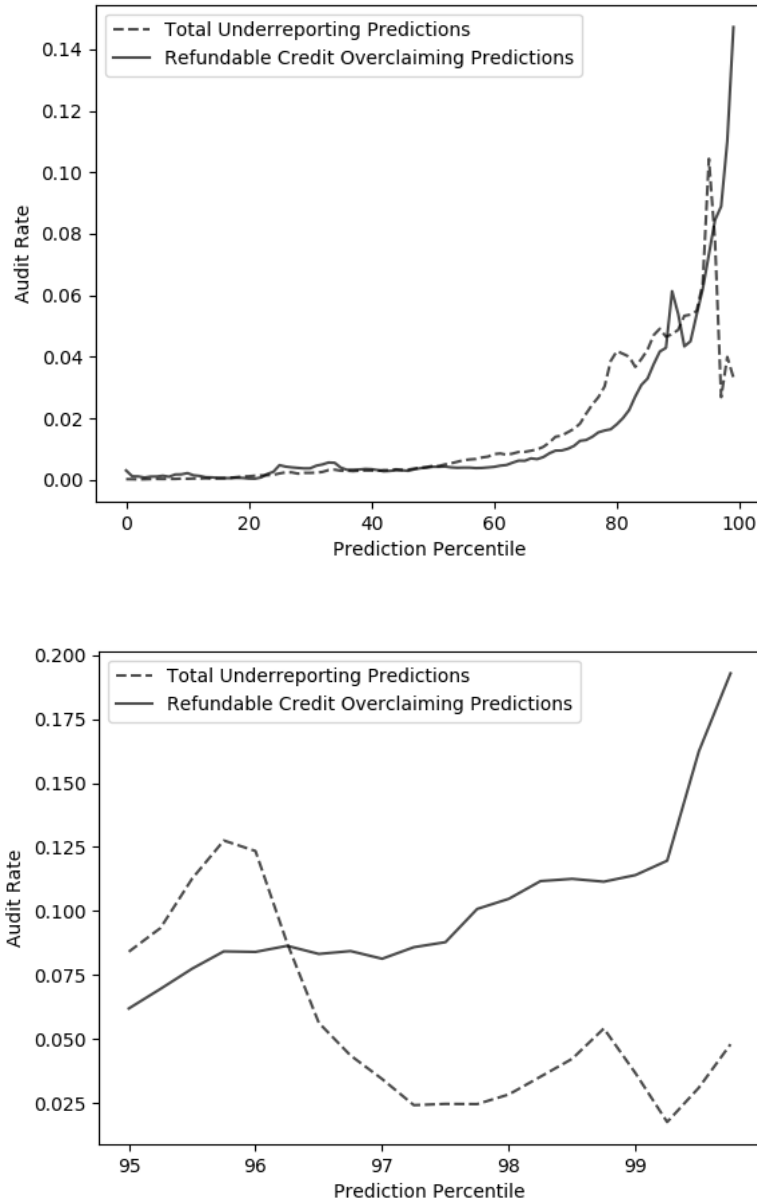
*Notes:* The figure replicates Figure 8, but with disparity calculated based on the linear disparity estimator. The figure shows the implied difference in audit rates between Black and non-Black taxpayers ( $y$ -axis) and annualized detected underreporting ( $x$ -axis) under alternative algorithms for selecting audits of EITC claimants and under alternative audit rates. The trajectories correspond to the total underreporting oracle (dark blue), total underreporting prediction (dark purple), refundable credit oracle (light blue), and refundable credit prediction (light purple) algorithms. The labeled points along each trajectory represent estimated detected underreporting and disparity for the specified algorithm at the audit rate specified in the label. For each algorithm, the audit rates considered range from 0.1% to 3%; the audit rate corresponding to the status quo (1.45%) is denoted by a larger dot. The total underreporting prediction algorithm is based on a random forest regressor trained to predict total underreporting. The refundable credit prediction algorithm is based on a random forest regressor trained to predict total adjustments to EITC, CTC, and AOTC amounts. The total underreporting oracle selects returns in descending order of true underreporting. The refundable credit oracle selects returns in descending order of true EITC, CTC, and AOTC overclaiming. Annualized detected underreporting is calculated as the total amount of adjustments (positive or negative) imposed on returns selected for audit under the specified audit selection algorithm. Detected underreporting and disparity estimates are constructed using the full set NRP EITC claimants from 2010-14. Details relating to the predictive models and associated estimates are contained in Appendix E. The point labeled “Status quo” shows the estimated disparity in audit rates between Black and non-Black taxpayers and detected underreporting from the 1.45% of EITC returns selected for audit from the population of tax year 2014 returns. All analyses incorporate NRP sampling weights. Bars around each trajectory represent 95% confidence intervals around disparity estimates; they are calculated based on the distribution of estimates from 100 bootstrapped samples from the full set of NRP EITC claimants; see Appendix E for details.

Figure A.18: Group-Specific Audit Rates by Algorithm (Linear)



*Notes:* The figure replicates Figure 9, but with audit rates calculated using the linear estimator. It reports estimated audit rates for Black and non-Black EITC claimants that would be induced by the algorithms considered in Figure A.17 under an assumed audit rate of 1.45%. The population of tax returns available for audit is based on the NRP sample of taxpayers claiming the EITC; see notes to Figure A.17 for additional detail. Status quo refers to the estimated audit rates by race for tax year 2014 returns claiming the EITC, reported in Figure 5.

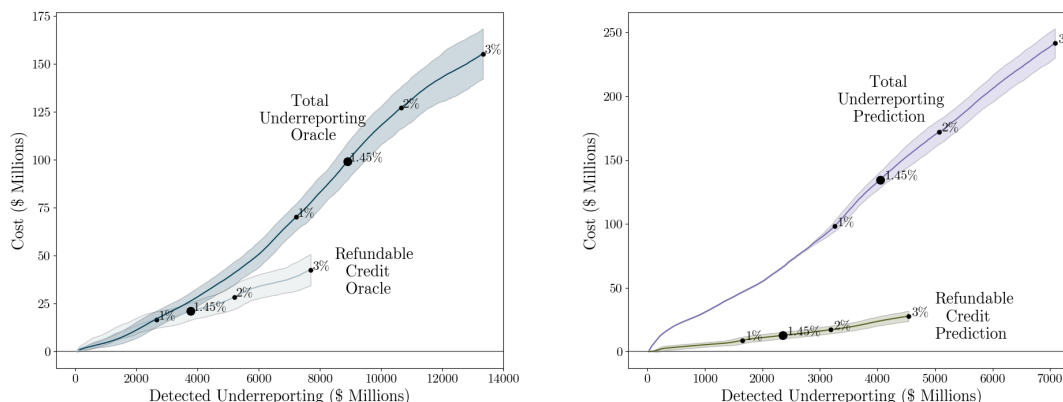
Figure A.19: Audit Rates by Prediction Percentiles



*Notes:* This figure displays the relationship between operational audit rates and predictions from the Total Underreporting and Refundable Credit Overclaiming Prediction algorithms, applied to the population of EITC claimants in 2014. For tractability, the versions of the Total Underreporting and Refundable Credit Overclaiming algorithms applied for the analysis in this figure are simplified in the following sense: We apply our NRP-based models to non-NRP taxpayers by training models with the forty most important features from both the total underreporting prediction and refundable credit overclaiming prediction models respectively and using them to obtain predictions. We validate that these simplified models yield similar results to the full-featured models in the NRP data set. In the top panel, we divide taxpayers into 100 equally-sized bins separately for each model and compute the audit rate separately for each bin. In the bottom panel, we divide taxpayers into 400 equally-sized bins separately for each model and filter to the top 20, computing audit rates separately for each bin. In this panel we divide the bin numbers on the x-axis by 4 to report results for each quarter of a percentile from 95 through 100.

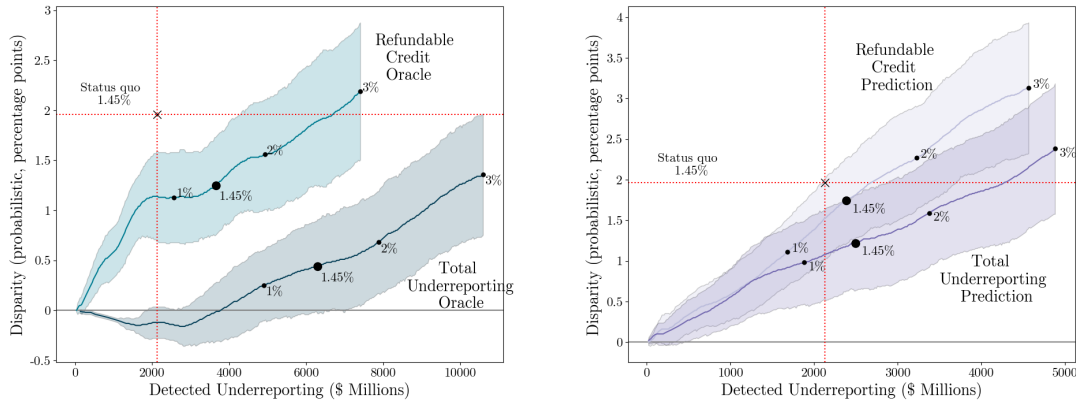


Figure A.20: Detected Underreporting and Cost by Algorithm



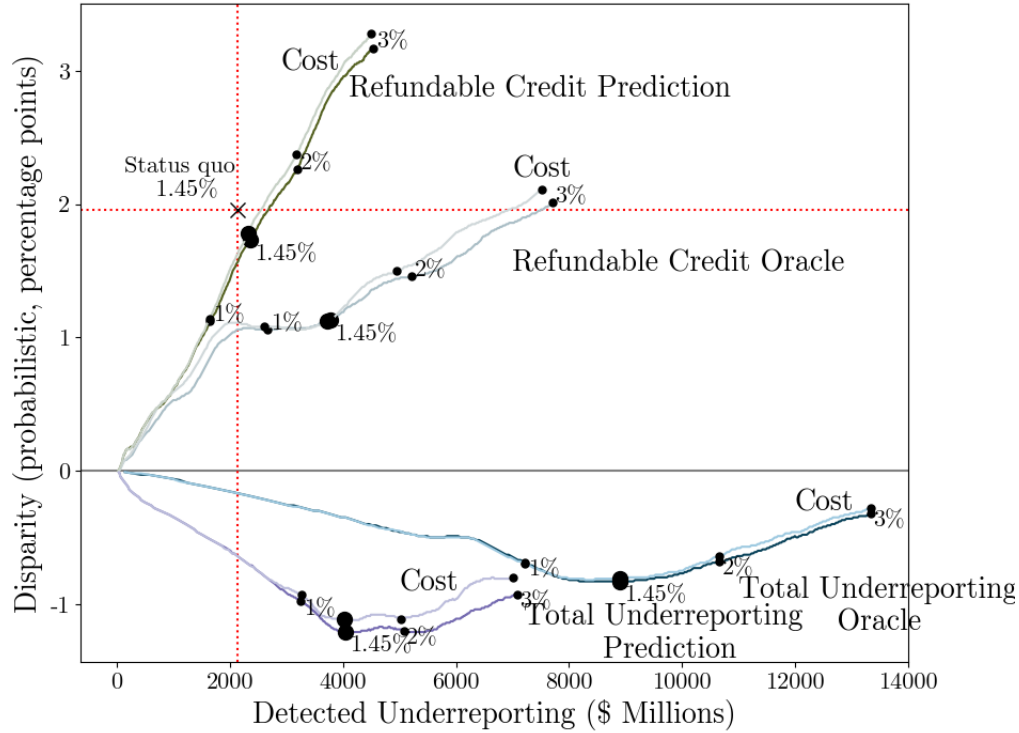
*Notes:* The figure shows the annualized cost ( $y$ -axis) and annualized detected underreporting ( $x$ -axis) under alternative algorithms for selecting audits of EITC claimants and under alternative audit rates. Predictive models are trained and evaluated on the set of NRP EITC claimants from 2010-14; see Appendix E for details. The left panel displays trajectories for the total underreporting oracle (dark blue) and the refundable credit oracle (light blue) algorithms. The right panel displays trajectories for the total underreporting prediction (dark purple) and the refundable credit prediction (light purple) algorithms. The labeled points along each trajectory represent estimated detected underreporting and cost for the specified algorithm at the audit rate specified in the label. The audit rates considered range from 0.1% to 3%. The audit rate corresponding to the status quo (1.45%) is denoted by a larger dot. The total underreporting prediction algorithm is based on a random forest regressor trained to predict total underreporting. The refundable credit prediction algorithm is based on a random forest regressor trained to predict total adjustments to EITC, CTC, and AOTC amounts. The total underreporting oracle selects returns in descending order of true underreporting. The refundable credit oracle selects returns in descending order of true EITC, CTC, and AOTC overclaiming. Annualized cost is calculated as the total cost incurred to audit the returns selected under the specified audit selection algorithm, scaled to reflect our use of five years of NRP data. The cost estimates are calculated from operational audits based on the average time logged by IRS employees dealing directly with the case, multiplied by the applicable General Schedule payscale based on the employee level. Returns without substantial business income (not reporting gross business receipts in excess of \$25,000) are classified into activity code 270 and are assigned a cost of \$23.09, the winsorized average cost of tax returns in activity code 271 in our sample. Returns with substantial business income (reporting gross business receipts in excess of \$25,000) are classified into activity code 271 and are assigned a cost of \$369.70, the winsorized average cost of tax returns in activity code 271 in our sample. Annualized detected underreporting is calculated as the total detected underreporting (positive or negative) imposed on returns selected for audit under the specified audit selection algorithm, scaled to reflect our use of five years of NRP data. All analyses incorporate NRP sampling weights. Bars around each trajectory represent 95% confidence intervals around cost estimates; they are calculated based on the distribution of estimates from 100 bootstrapped samples from the full set of NRP EITC claimants; see Appendix E for details.

Figure A.21: Detected Underreporting and Disparity by Algorithm (Constrained Models)



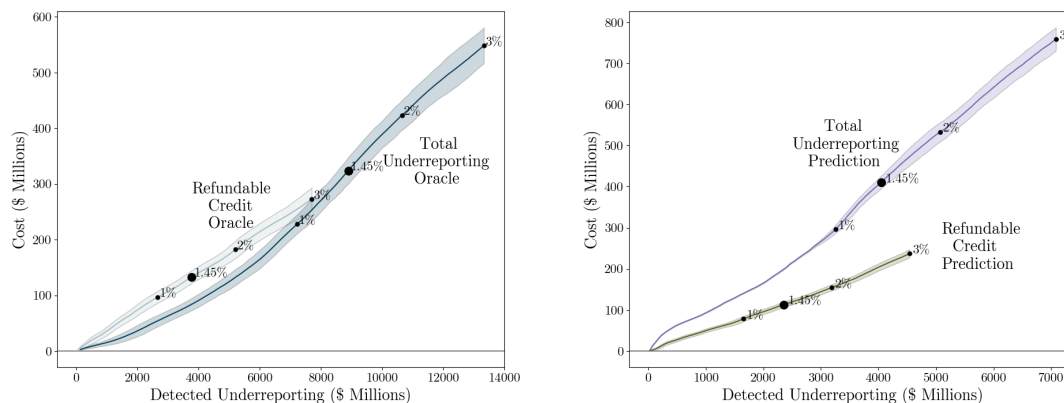
*Notes:* The figure shows the implied difference in audit rates between Black and non-Black taxpayers ( $y$ -axis) and annualized detected underreporting ( $x$ -axis) under alternative algorithms for selecting audits of EITC claimants and under alternative audit rates, where each algorithm's allocation of audits between EITC business and non-business activity codes is constrained to match the status quo allocation. The left panel shows trajectories corresponding to the total underreporting oracle (dark blue) and the refundable credit oracle (light blue) algorithms. The right panel shows trajectories corresponding to the total underreporting prediction (dark purple) and refundable credit prediction (light purple) algorithms. The labeled points along each trajectory represent estimated detected underreporting and disparity for the specified algorithm at the audit rate specified in the label. For each algorithm, the audit rates considered range from 0.1% to 3%; the audit rate corresponding to the status quo (1.45%) is denoted by a larger dot. The total underreporting prediction algorithm is based on a random forest regressor trained to predict total underreporting. The refundable credit prediction algorithm is based on a random forest regressor trained to predict total adjustments to EITC, CTC, and AOTC amounts. The total underreporting oracle selects returns in descending order of true underreporting. The refundable credit oracle selects returns in descending order of true EITC, CTC, and AOTC overclaiming. Disparity is calculated from the probabilistic disparity estimator applied to BIFSG-derived probabilities that a taxpayer is Black. Annualized detected underreporting is calculated as the total amount of adjustments (positive or negative) imposed on returns selected for audit under the specified audit selection algorithm. Detected underreporting and disparity estimates are constructed using the full set NRP EITC claimants from 2010-14. Details relating to the predictive models and associated estimates are contained in Appendix E. The point labeled "Status quo" shows the estimated disparity in audit rates between Black and non-Black taxpayers and detected underreporting from the 1.45% of EITC returns selected for audit from the population of tax year 2014 returns. All analyses incorporate NRP sampling weights. Bars around each trajectory represent 95% confidence intervals around disparity estimates; they are calculated based on the distribution of estimates from 100 bootstrapped samples from the full set of NRP EITC claimants; see Appendix E for details.

Figure A.22: Allocating Audits Based on Reward Net of Audit Costs



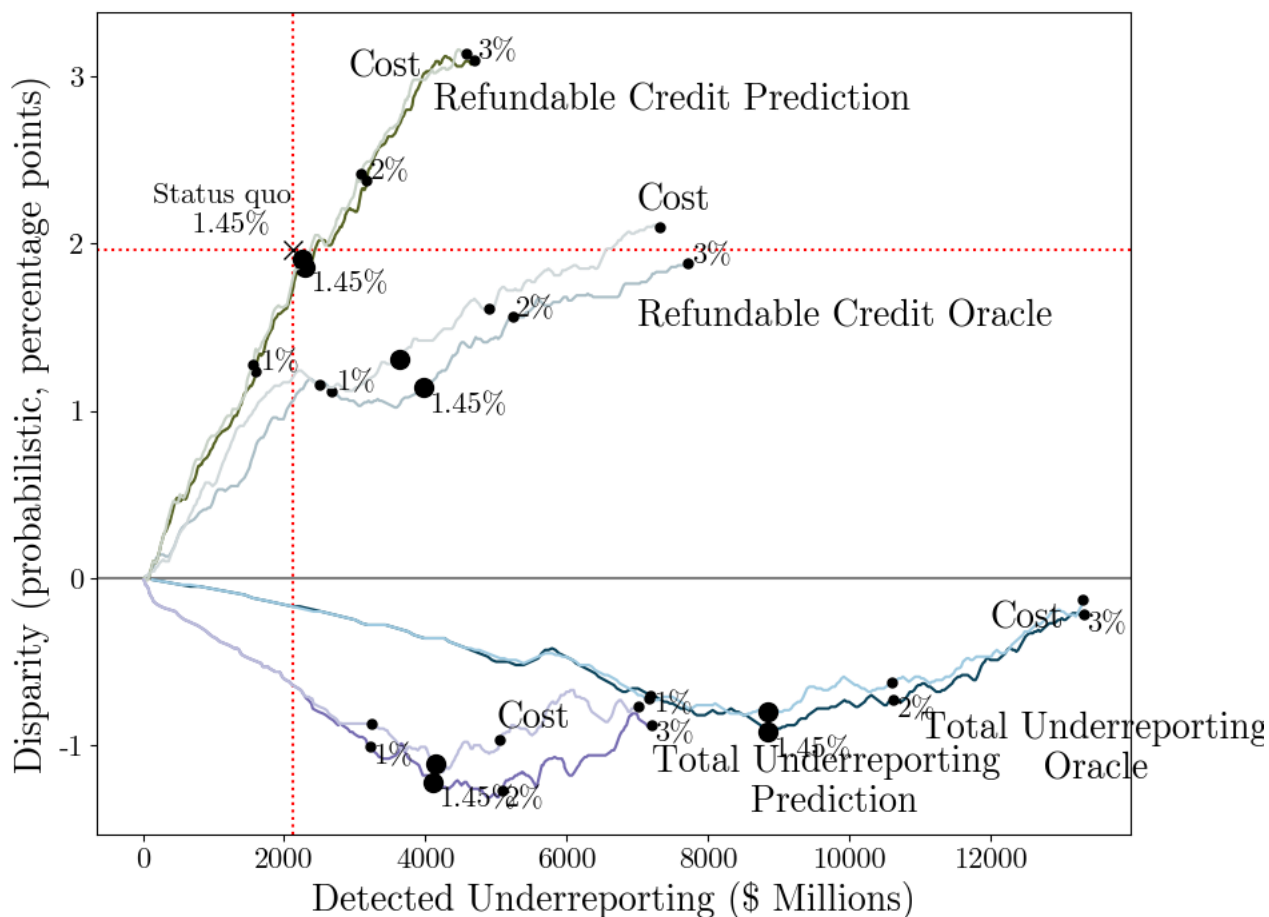
*Notes:* The figure shows the implied difference in audit rates between Black and non-Black taxpayers ( $y$ -axis) and annualized detected underreporting ( $x$ -axis) under alternative assumptions about whether returns are selected for audit based on detected reward, whether total underreporting or refundable credit overclaiming (gross of audit costs, as in our other analyses), or based on detected reward minus expected audit costs. Underreporting/refundable credit overclaiming is based on either the oracle or the random forest regressor prediction algorithm, as specified. Audit costs are measured at the activity code level, using data on the time spent on audit examination and the salary grade of the examiner, and abstracting from non-salary costs associated with the enforcement process, such as appeals, litigation, and collections, or fixed costs, such as overhead. Using this approach, the average cost of auditing an EITC non-business return (activity code 270) is \$23.09, whereas the average cost of auditing an EITC business return (activity code 271) is \$369.70. The labeled points along each trajectory represent estimated detected underreporting and disparity for the specified algorithm at the audit rate specified in the label. Disparity is calculated from the probabilistic disparity estimator applied to BIFSG-derived probabilities that a taxpayer is Black. For each algorithm, the audit rates considered range from 0.1% to 3%; the audit rate corresponding to the status quo (1.45%) is denoted by a larger dot. Annualized detected underreporting is calculated as the total amount of adjustments (positive or negative) imposed on returns selected for audit under the specified audit selection algorithm. Detected underreporting and disparity estimates are constructed using the full set NRP EITC claimants from 2010-14. Details relating to the predictive models and associated estimates are contained in Appendix E. The point labeled “Status quo” shows the estimated disparity in audit rates between Black and non-Black taxpayers and detected underreporting from the 1.45% of EITC returns selected for audit from the population of tax year 2014 returns. All analyses incorporate NRP sampling weights.

Figure A.23: Detected Underreporting and Cost by Algorithm (NRP-Derived Audit Costs)



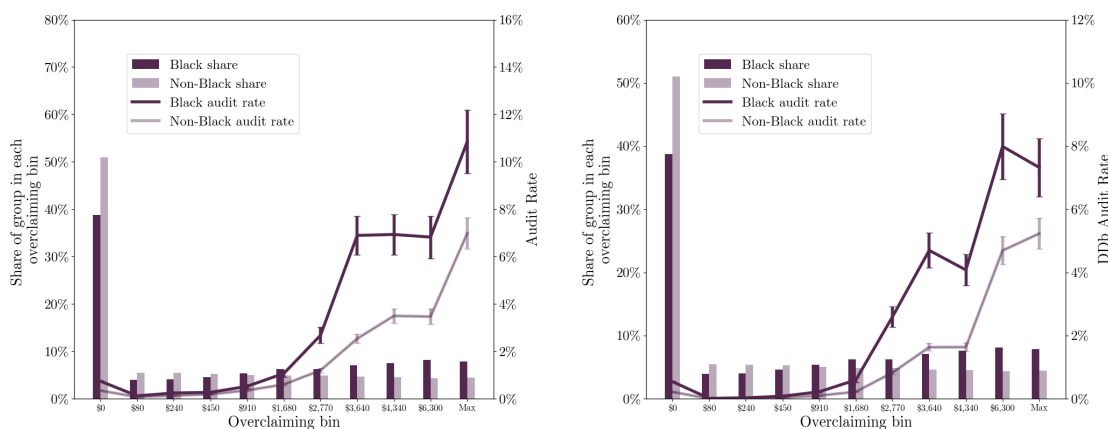
*Notes:* The figure replicates Appendix Figure A.20 using an alternative measure of costs derived from the number of hours reported by examiners conducting NRP examinations. Audit costs are measured at the activity code level, with the average cost of auditing an EITC non-business return (activity code 270) being \$263.51 (compared to \$23.09 under our original cost measure), and the average cost of auditing an EITC business return (activity code 271) being \$1,110.90 (compared to \$369.70 under our original cost measure). The figure shows the annualized cost ( $y$ -axis) and annualized detected underreporting ( $x$ -axis) under alternative algorithms for selecting audits of EITC claimants and under alternative audit rates. Predictive models are trained and evaluated on the set of NRP EITC claimants from 2010-14; see Appendix E for details. The left panel displays trajectories for the total underreporting oracle (dark blue) and the refundable credit oracle (light blue). The right panel displays trajectories for the total underreporting prediction model (dark purple), and the refundable credit prediction model (light purple). The labeled points along each trajectory represent estimated detected underreporting and cost for the specified algorithm at the audit rate specified in the label. The audit rates considered range from 0.1% to 3%. The audit rate corresponding to the status quo (1.45%) is denoted by a larger dot. The total underreporting prediction algorithm is based on a random forest regressor trained to predict total underreporting. The refundable credit prediction algorithm is based on a random forest regressor trained to predict total adjustments to EITC, CTC, and AOTC amounts. The total underreporting oracle selects returns in descending order of true underreporting. The refundable credit oracle selects returns in descending order of true EITC, CTC, and AOTC overclaiming. Annualized cost is calculated as the total cost incurred to audit the returns selected under the specified audit selection algorithm, scaled to reflect our use of five years of NRP data. Annualized detected underreporting is calculated as the total detected underreporting (positive or negative) imposed on returns selected for audit under the specified audit selection algorithm, scaled to reflect our use of five years of NRP data. All analyses incorporate NRP sampling weights. Bars around each trajectory represent 95% confidence intervals around cost estimates; they are calculated based on the distribution of estimates from 100 bootstrapped samples from the full set of NRP EITC claimants; see Appendix E for details.

Figure A.24: Allocating Audits Based on Reward Net of Audit Costs (NRP-Derived Audit Costs)



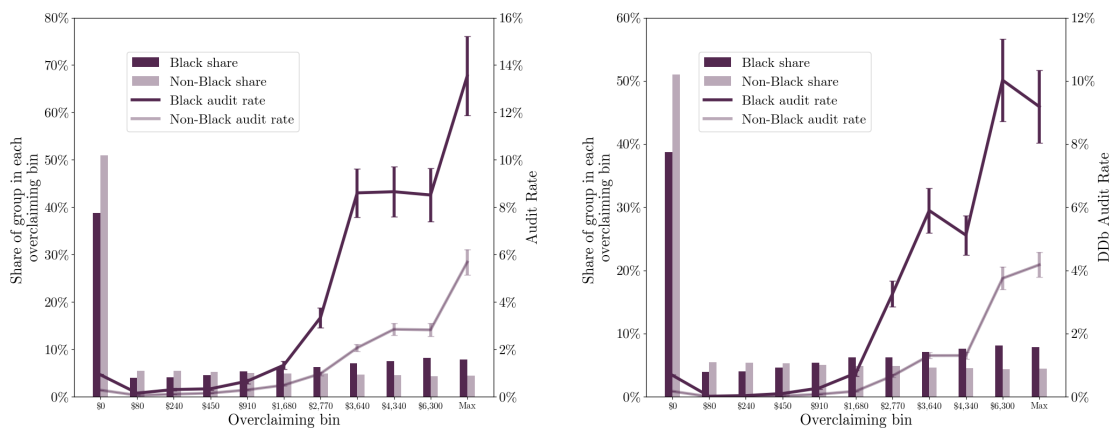
*Notes:* The figure replicates Appendix Figure A.22 using an alternative measure of costs derived from the number of hours reported by examiners conducting NRP examinations. Audit costs are measured at the activity code level, with the average cost of auditing an EITC non-business return (activity code 270) being \$263.51 (compared to \$23.09 under our original cost measure), and the average cost of auditing an EITC business return (activity code 271) being \$1,110.90 (compared to \$369.70 under our original cost measure). The figure shows the implied difference in audit rates between Black and non-Black taxpayers (*y*-axis) and annualized detected underreporting (*x*-axis) under alternative assumptions about whether returns are selected for audit based on detected reward, whether total underreporting or refundable credit overclaiming (gross of audit costs, as in our other analyses), or based on detected reward minus expected audit costs. Underreporting/refundable credit overclaiming is based on either the oracle or the random forest regressor prediction algorithm, as specified. The labeled points along each trajectory represent estimated detected underreporting and disparity for the specified algorithm at the audit rate specified in the label. Disparity is calculated from the probabilistic disparity estimator applied to BIFSG-derived probabilities that a taxpayer is Black. For each algorithm, the audit rates considered range from 0.1% to 3%; the audit rate corresponding to the status quo (1.45%) is denoted by a larger dot. Annualized detected underreporting is calculated as the total amount of adjustments (positive or negative) imposed on returns selected for audit under the specified audit selection algorithm. Detected underreporting and disparity estimates are constructed using the full set NRP EITC claimants from 2010-14. Details relating to the predictive models and associated estimates are contained in Appendix E. The point labeled “Status quo” shows the estimated disparity in audit rates between Black and non-Black taxpayers and detected underreporting from the 1.45% of EITC returns selected for audit from the population of tax year 2014 returns. All analyses incorporate NRP sampling weights.

Figure A.25: Racial Audit Disparity Among EITC Claimants by Overclaimed Refundable Credits



*Notes:* The figure shows two analogs of Figure 7. Each panel groups taxpayers into bins along the x-axis based on refundable credit overclaiming. For the left panel, the y-axis reflects any type of audit; for the right panel, the y-axis reflects DDb audits only. For the left panel, audit-imposed adjustments to refundable credits are used to measure a taxpayer's overclaimed refundable credits, conditional on being selected for operational audit, by race. For the right panel, total audit adjustments from DDb audits are used as a proxy for a taxpayer's overclaimed refundable credits, conditional on being selected for operational audit, by race. In both panels, taxpayers are binned into 11 categories: those with less than \$1 of overclaimed refundable credits, and 10 equal deciles of taxpayers with positive overclaimed refundable credits. Overclaiming deciles are defined based on the distribution of overclaiming among EITC claimants, as measured by NRP audits. Bin labels on the x-axis reflect the upper dollar limit of each overclaiming bin (rounded for confidentiality). Estimated audit rates by race are calculated using the probabilistic estimator. All analyses account for NRP sampling weights. Brackets reflect the estimated 95% confidence interval, derived from bootstrapped standard errors (N=100). The bars show the estimated share of Black and non-Black taxpayers, respectively, that fall into each overclaiming bin. A similar analysis, corresponding to the linear estimator, is presented in Appendix Figure A.26.

Figure A.26: Racial Audit Disparity Among EITC Claimants by Overclaimed Refundable Credits (Linear Estimator)



Notes: The figure replicates the analysis in Figure A.25, using the linear estimator to calculate audit rates by race.

Table A.1: EITC Audit Frequency by Timing and Type

	Correspondence	Field/Office	All Audit Types
Pre-Refund	270,940 (66.4%)	0 (0%)	270,940 (66.4%)
Post-Refund	112,689 (27.6%)	24,361 (6.0%)	137,050 (33.6%)
All Audit Times	383,629 (94.0%)	24,361 (6.0%)	407,990 (100%)

*Notes:* The table reports the frequency of audits of 2014 tax returns claiming the EITC by audit timing (whether the audit occurred pre- or post-refund) and by audit type (whether the audit was conducted by correspondence or as a field or office examination). Percentages (reported in parentheses) reflect the share of all audits of the specified taxpayer population that fall into the specified audit category.



Table A.2: Coverage of BIFSG Features Among 2014 Taxpayers

Case	Count	First Name	Last Name	CBG	Share of Total
1	107,624,714	X	X	X	72.6%
2	10,087,515	X	X		6.8%
3	10,455,708	X		X	7.1%
4	14,981,324		X	X	10.1%
5	2,572,849			X	1.7%
6	1,431,541		X		1.0%
7	903,311	X			0.6%
8	248,356				0.2%
Total	148,305,318				100%

*Notes:* The table shows the availability of the data used to calculate race probabilities for primary filers on tax year 2014 returns. The distribution of race by first name is tabulated from mortgage applications, following Tzioumis (2018); it is missing for names not among the 4,250 most common names in that data. The distribution of race by last names is tabulated from 2010 Census data and includes the 162,253 most common surnames. The distribution of race by Census Block Group (CBG) is tabulated from the Census 2014 5-Year American Community Survey and covers all CBGs. CBG data is missing for taxpayers who cannot be reliably geo-coded to a specific CBG. In our main analysis, taxpayers in row 8 are assigned a BIFSG-predicted probability Black based on the national average share of the population that is Black.

Table A.3: Calibration Metrics for BIFSG Predictions.

Metric	Full Population		EITC Population	
	Imputed (1)	Recalibrated (2)	Imputed (3)	Recalibrated (4)
Area Under ROC Curve	0.9048	0.9048	0.9038	0.9038
Panel A: 50% Threshold				
False Positive	0.0880	0.0632	0.1119	0.1181
True Positive	0.6804	0.6066	0.7237	0.7377
False Negative	0.3196	0.3934	0.2763	0.2623
True Negative	0.9120	0.9368	0.8881	0.8819
Precision	0.6529	0.7002	0.8002	0.7946
Recall	0.6804	0.6066	0.7237	0.7377
Accuracy	0.8667	0.8722	0.8252	0.8268
Panel B: 75% Threshold				
False Positive	0.0338	0.0112	0.0504	0.0504
True Positive	0.4740	0.2822	0.5169	0.5170
False Negative	0.5260	0.7178	0.4831	0.4830
True Negative	0.9662	0.9888	0.9496	0.9496
Precision	0.7733	0.8598	0.8641	0.8641
Recall	0.4740	0.2822	0.5169	0.5170
Accuracy	0.8699	0.8506	0.7842	0.7842
Panel C: 90% Threshold				
False Positive	0.0114	-	0.0216	0.0191
True Positive	0.2855	-	0.3180	0.2946
False Negative	0.7145	-	0.6820	0.7054
True Negative	0.9886	-	0.9784	0.9809
Precision	0.8588	-	0.9013	0.9053
Recall	0.2855	-	0.3180	0.2946
Accuracy	0.8511	-	0.7259	0.7184

*Notes:* The table characterizes the predictive power of BIFSG in the task of predicting whether a taxpayer self-reports as Black as measured in the North Carolina data. A positive label refers to the individual self-reporting as Black and a negative label refers to the individual self-reporting as non-Black. We can then characterize performance using standard metrics from the machine learning literature. We use columns to demarcate different versions of the final predictions and comparing against different populations: Columns 1 and 3 correspond to the standard BIFSG score, while Columns 2 and 4 correspond to the re-calibrated BIFSG score (described in Section B.5); and Columns 1 and 2 are evaluations against the full population, while in Columns 3 and 4 are evaluations against only the EITC claimant population. We use rows to demarcate each error metric. The first error metric we consider, the Area Under the Receiver Operator Characteristic (ROC) curve, requires as input only the probabilistic predictions and labels; the other metrics, which are classification-based, require discrete label choices. We thus convert BIFSG scores into predicted labels of Black/non-Black via thresholding, i.e. labeling all observations with predicted probability Black of  $t$  or greater as Black and all others as non-Black. We consider thresholds at 50%, 75%, and 90%, demarcated by Panels A-C.

Table A.4: Residual Covariance Estimates

	Full Population		EITC		Non-EITC	
E[cov(Y,B) b]	5.76*** (0.20)	5.05*** (0.23)	14.68*** (0.81)	13.39*** (1.03)	1.65*** (0.14)	1.20*** (0.14)
E[cov(Y,b) B]	2.03*** (0.14)	2.28*** (0.24)	8.76*** (0.65)	9.62*** (0.97)	0.00 (0.10)	-0.28 (0.18)
Weighted	No	Yes	No	Yes	No	Yes
N	1,613,130		277,064		1,336,062	

*Notes:* The table displays the estimated covariance between audits and self-reported race, conditional on estimated race, as well as the estimated covariance between audits and estimated race, conditional on self-reported race. The estimates are for the matched sample of North Carolina taxpayers for the full population (columns 1 and 2), EITC claimants (columns 3 and 4), and non-EITC claimants (columns 5 and 6). Columns 2, 4, and 6 are re-weighted to be representative of the U.S. population, using the weights described in Appendix C. Standard errors are displayed in parentheses. The displayed estimates and standard errors are multiplied by  $10^4$ . Stars correspond to p-values derived from two-sided hypothesis tests. \* :  $P < .10$ ; \*\* :  $P < .05$ ; \*\*\* :  $P < .01$ .

Table A.5: Estimated Audit Rate by Race (Probabilistic)

	Any Audit (1)	Audit Timing		Audit Type	
		Pre Refund (2)	Post Refund (3)	Correspondence (4)	Field/ Office (5)
Panel A: Full Population					
Black	1.241 (0.003)	0.743 (0.002)	0.497 (0.002)	1.139 (0.003)	0.103 (0.001)
Non-Black	0.427 (0.001)	0.174 ( $< 0.001$ )	0.253 ( $< 0.001$ )	0.335 ( $< 0.001$ )	0.093 ( $< 0.001$ )
N (Millions)	148.3	148.3	148.3	148.3	148.3
Panel B: EITC Population					
Black	2.989 (0.006)	2.126 (0.005)	0.863 (0.003)	2.897 (0.006)	0.092 (0.001)
Non-Black	1.039 (0.002)	0.651 (0.001)	0.389 (0.001)	0.954 (0.002)	0.085 (0.001)
N (Millions)	28.3	28.3	28.3	28.3	28.3
Panel C: Non-EITC Population					
Black	0.401 (0.002)	0.080 (0.001)	0.322 (0.001)	0.294 (0.001)	0.108 (0.001)
Non-Black	0.300 ( $< 0.001$ )	0.074 ( $< 0.001$ )	0.225 ( $< 0.001$ )	0.206 ( $< 0.001$ )	0.094 ( $< 0.001$ )
N (Millions)	120.0	120.0	120.0	120.0	120.0

*Notes:* The table reports estimates of the audit rate (in group-specific levels, not differences between groups) for Black and non-Black taxpayers filing income tax returns for tax year 2014. Units are percentage points (0-100). All estimates are based on the probabilistic audit rate estimator. The category of audit considered varies across columns; for example, the results in column (4) show the estimated rate at which Black and non-Black taxpayers are selected for correspondence audit. Panel A includes all taxpayers, whereas Panels B and C restrict the analysis to EITC claimants and non-claimants, respectively. Standard errors, reported in parentheses, correspond to the standard deviation of the distribution of estimates from 100 bootstrapped samples.

Table A.6: Estimated Audit Rate by Race (Linear)

	Any Audit (1)	Audit Timing		Audit Type	
		Pre Refund (2)	Post Refund (3)	Correspondence (4)	Field/ Office (5)
Panel A: Full Population					
Black	1.707 (0.004)	1.070 (0.003)	0.637 (0.002)	1.599 (0.004)	0.108 (0.001)
Non-Black	0.363 (0.001)	0.129 ( $< 0.001$ )	0.234 ( $< 0.001$ )	0.271 (0.001)	0.092 ( $< 0.001$ )
N (Millions)	148.3	148.3	148.3	148.3	148.3
Panel B: EITC Population					
Black	3.732 (0.009)	2.688 (0.007)	1.044 (0.005)	3.637 (0.009)	0.095 (0.001)
Non-Black	0.846 (0.002)	0.505 (0.002)	0.342 (0.001)	0.762 (0.002)	0.085 (0.001)
N (Millions)	28.3	28.3	28.3	28.3	28.3
Panel C: Non-EITC Population					
Black	0.471 (0.003)	0.083 (0.001)	0.388 (0.002)	0.355 (0.002)	0.117 (0.001)
Non-Black	0.292 ( $< 0.001$ )	0.074 ( $< 0.001$ )	0.218 ( $< 0.001$ )	0.199 ( $< 0.001$ )	0.093 ( $< 0.001$ )
N (Millions)	120.0	120.0	120.0	120.0	120.0

*Notes:* The table replicates Appendix Table A.5, but using the linear audit rate estimator instead of the probabilistic audit rate estimator.

Table A.7: Ground-Truth Disparities in Matched North Carolina Sample

Estimator	Full Population (1)	EITC (2)	Non-EITC (3)
Unweighted	0.835	1.849	0.204
Re-weighted	1.284	2.393	0.265
N	1,613,124	277,064	1,336,060

*Notes:* The table shows audit rate disparities within the matched North Carolina sample, where we observe self-reported race. Units are percentage points (0-100). The Black/non-Black audit disparity is shown for the full population (column 1), EITC claimants (column 2), and non-EITC claimants (column 3). The first row computes audit rate disparities directly using the data, while the second row re-weights the data to be representative of the full population using the North Carolina weights described in Appendix C.

Table A.8: Audit Burden Disparity Estimates (Probabilistic)

	Correspondence	Office	Field	Overall
Audit Rate (pp)	0.43	0.06	0.03	0.53
Share of Total Audits	0.82	0.12	0.06	1.00
Taxpayer Hours	30	38	34	31.18
Taxpayer Compliance Cost	643	1,717	4,431	1,002.98
Penalties and Interest	320.71	1,580.77	6,434.52	846.77
Assessed Taxes	5,252.56	7,130.46	24,960.56	6,694.85
Disparity (Audit Rate, pp)	0.80	0.02	-0.01	0.66
Disparity (Hours)	0.241	0.008	-0.004	0.199
Disparity (Compliance Cost)	5.17	0.36	-0.50	4.26
Disparity (Penalties and Interest)	2.58	0.34	-0.73	2.11
Disparity (Assessed Taxes)	42.20	1.51	-2.83	34.68

*Notes:* The table reports quantities for estimating differences by race in the consequences to taxpayers of being selected for an audit. All results are presented separately by type of audit (Columns 1-3), as well as aggregated based on the share of 2014 audits in each category (Column 4). The first six rows report summary audit statistics. The audit rate (row 1) is the percent of taxpayers that receive the type of audit in question. Taxpayer hours (row 3) and compliance costs (row 4) are taken from estimates reported in Guyton & Hodge (2014). Hours refers to the total hours the taxpayer spends responding to the audit. Compliance cost refers to the sum of out-of-pocket compliance costs and monetized hours spent responding to the audit (accounting for income differences in the cost of taxpayer time), inflation-adjusted from 2009 dollars to 2014 dollars. Penalties and interest (row 5) are net of IRS abatements. The bottom five rows report the estimated disparity in each outcome in rows two through six. The audit rate disparity is our main outcome of interest in the rest of the paper, here calculated using the probabilistic disparity estimator. For columns 1-3, the estimated disparity for other variables is calculated by multiplying the audit rate disparity by the average value of that outcome for the corresponding type of audit. For example, the estimated disparity in hours from a correspondence audit is the product of the correspondence audit disparity (0.80) and average taxpayer hours per correspondence audit (30). The disparity results in column 4 are obtained by taking the weighted average of the disparity results in columns 1-3, using the weights in row 2. These estimates assume that within-audit category, the components of audit burden do not vary by race. In addition, these estimates do not incorporate factors known to be important to taxpayers such as the non-financial stress of experiencing an audit, delays in receiving anticipated refunds, or reductions in refundable credits claims in future years. For context, taxpayers (whether audited or not) spend an average of 0.16 hours complying with audits

Table A.9: Audit Burden Disparity Estimates (Linear)

	Correspondence	Office	Field	Overall
Audit Rate (pp)	0.43	0.06	0.03	0.53
Share of Total Audits	0.82	0.12	0.06	1.00
Taxpayer Hours	30	38	34	31.18
Taxpayer Compliance Cost	643	1,717	4,431	1,002.98
Penalties and Interest	320.71	1,580.77	6,434.52	846.77
Assessed Taxes	5,252.56	7,130.46	24,960.56	6,694.85
Disparity (Audit Rate, pp)	1.33	0.04	-0.02	1.09
Disparity (Hours)	0.398	0.013	-0.006	0.329
Disparity (Compliance Cost)	8.54	0.60	-0.83	7.04
Disparity (Penalties and Interest)	4.26	0.55	-1.20	3.49
Disparity (Assessed Taxes)	69.77	2.50	-4.67	57.34

*Notes:* The table replicates Table A.8, but using the linear disparity estimator instead of the probabilistic disparity estimator to calculate the audit rate disparity.

Table A.10: Dual-Bootstrap Confidence Intervals

Estimator	Full Population (1)	EITC (2)	Non-EITC (3)
Linear	1.345	2.885	0.180
95% CI	(1.169, 1.461)	(2.633, 2.991)	(0.138, 0.230)
Probabilistic	0.813	1.950	0.102
95% CI	(0.729, 0.909)	(1.846, 2.031)	(0.080, 0.135)
N	148,305,318	28,338,472	119,966,846

*Notes:* The table reports the linear and probabilistic disparity estimates along with 95% confidence intervals obtained via our implementation of the Lu et al. (2024) dual-bootstrap procedure. To implement the procedure, we resampled first names and geographic information, but did not resample surnames (since the Census surname data we use corresponds to the full population rather than a sample) to generate 100 sets of BIFSG posteriors. We then compute our outcomes of interest with each set of BIFSG posteriors, resampling our taxpayer population on each iteration. Units are percentage points (0-100). Confidence intervals were obtained based on the distribution of estimates obtained by this procedure (see Appendix Figure A.12).

Table A.11: DDb vs. Non-DDb Audit Disparity Estimates Among EITC Claimants

Estimator	DDb Audit Disparity (1)	Non-DDb Audit Disparity (2)	Share of Disparity Attributable to DDb (3)
Linear	2.265 (0.008)	0.620 (0.005)	78.5
Probabilistic	1.531 (0.007)	0.419 (0.004)	78.5
Mean	1.071	0.372	
N	28,338,472	28,338,472	

*Notes:* Columns 1 and 2 report the estimated audit rate disparities for DDb and non-DDb audits among EITC claimants. Column 3 is calculated as the ratio of the DDb audit disparity in Column 1 to the total audit disparity (the sum of Columns 1 and 2), multiplied by 100. Units are percentage points (0-100). The “Mean” row is the share of the EITC claimant population that is selected for the specified category of audit. In our data, 74.31% of the audits of EITC claimants are selected through the DDb program.



Table A.12: Audit Disparity Robustness Checks

	BIFSG	Gibbs	Re-calibrated Unweighted	Re-calibrated Weighted	Geography- Only
	(1)	(2)	(3)	(4)	(5)
Panel A: Full Population					
Linear	1.27 (0.005)	1.55 (0.04)	1.623 (0.005)	1.543 (0.004)	1.738 (0.005)
Probabilistic	0.811 (0.004)	1.02 (0.03)	0.774 (0.003)	0.735 (0.003)	0.696 (0.003)
N	107,624,714	1,390,219	148,305,318	148,305,318	134,671,876
Panel B: EITC Population					
Linear	3.00 (0.01)	2.82 (0.08)	3.48 (0.01)	3.31 (0.01)	3.500 (0.01)
Probabilistic	2.14 (0.01)	2.03 (0.07)	1.842 (0.008)	1.816 (0.008)	1.623 (0.008)
N	19,357,514	283,055	28,338,472	28,338,472	25,768,606
Panel C: Non-EITC Population					
Linear	0.172 (0.003)	0.165 (0.03)	0.217 (0.003)	0.206 (0.003)	0.352 (0.004)
Probabilistic	0.104 (0.002)	0.100 (0.02)	0.097 (0.002)	0.091 (0.002)	0.128 (0.002)
N	88,267,200	1,107,164	119,966,846	119,966,846	108,903,270

*Notes:* The table shows the estimated audit rate disparity from the linear and probabilistic disparity estimators, under various modifications to our baseline approach. Units are percentage points (0-100). standard errors, reported in parentheses, are calculated from the asymptotic distributions described in Appendix B.3. Column 1 restricts the analysis to the subset of taxpayers for which each of first name, last name, and census block group are available. Column 2 predicts taxpayer race using the Gibbs sampling approach described in Appendix B.6. Column 3 predicts taxpayer race after re-calibrating the race probability estimates using the North Carolina data, as described in Appendix C. Column 4 replicates Column 3, but re-weights the data to be representative of the full population using the North Carolina weights described in Appendix C. Column 5 predicts taxpayer race using only taxpayers' geographic location. Panel A shows results for the full population; Panel B for the EITC population; and Panel C for the non-EITC population. Each displayed disparity estimate is in terms of percentage points and is statistically different from zero ( $p < .01$ ).

Table A.13: Residual Covariance Estimates (Geography-Only)

	Full Population		EITC		Non-EITC	
E[cov(Y,B) b]	7.52*** (0.23)	8.51*** (0.38)	22.58*** (0.94)	24.87*** (1.41)	1.98*** (0.16)	1.64*** (0.28)
E[cov(Y,b) B]	2.16*** (0.12)	2.73*** (0.24)	8.49*** (0.57)	7.79*** (0.79)	0.08 (0.08)	0.53** (0.24)
Weighted	No	Yes	No	Yes	No	Yes
N	1,612,713		277,005		1,335,708	

*Notes:* The table displays the estimated covariance between audits and self-reported race, conditional on geography-only estimated race, as well as the estimated covariance between audits and estimated race, conditional on self-reported race. The estimates are for the matched sample of North Carolina taxpayers for the full population (columns 1 and 2), EITC claimants (columns 3 and 4), and non-EITC claimants (columns 5 and 6). Columns 2, 4, and 6 are re-weighted to be representative of the U.S. population, using the weights described in Appendix C. Standard errors are displayed in parentheses. The displayed estimates and standard errors are multiplied by  $10^4$ . Stars correspond to p-values derived from two-sided hypothesis tests. \* :  $P < .10$ ; \*\* :  $P < .05$ ; \*\*\* :  $P < .01$ .

Table A.14: Audits of EITC Claimants by Presence of Business Income

	No Substantial Business Income	Substantial Business Income
Share of EITC Claimants	0.94	0.07
Share of Audits of EITC Claimants	0.96	0.04
Audit Rate (pp)	1.46	0.96
Disparity (Probabilistic, pp)	1.98	0.70
Disparity (Linear, pp)	2.93	1.15
Black Audit Rate (Probabilistic, pp)	3.02	1.58
Black Audit Rate (Linear, pp)	3.76	1.99
Non-Black Audit Rate (Probabilistic, pp)	1.04	0.89
Non-Black Audit Rate (Linear, pp)	0.84	0.84
N	26,330,090	1,853,037

*Notes:* The table reports descriptive statistics and estimated audit disparity for EITC claimants, based on whether they report substantial business income on their return. Substantial business income is defined based on the IRS's definitions for activity codes 270 and 271, which partition the set of EITC claimants with total positive income below \$200,000. Activity code 271 includes those taxpayers who report substantial business income (i.e., schedule C or schedule F gross receipts in excess of \$25,000); activity code 270 includes the remainder. A very small share of EITC claimants are classified into activity codes other than 270 or 271 because they report total positive income above \$200,000; we exclude them for purposes of this analysis. Disparity estimates are presented for both the probabilistic disparity estimator and linear disparity estimator applied to BIFSG-derived probabilities that a taxpayer is Black. Audit rates by race present group-specific levels (rather than differences between groups) for Black and non-Black taxpayers filing income tax returns for tax year 2014.

Table A.15: Audit Cost and Revenue by Model

	Prediction Models		Oracles	
	Refundable Credit Overclaiming	Total Underreporting	Refundable Credit Overclaiming	Total Underreporting
Share with Substantial Business Income	0.028	0.927	0.091	0.655
Cost (Operational, \$ Millions)	13	134	23	99
Cost (NRP, \$ Millions)	113	408	138	324
Detected Underreporting (\$ Millions)	2,299	4,109	3,976	8,855

*Notes:* The table reports information about the performance of the alternative audit selection algorithms considered in Figure 8 at the status quo (1.45%) audit rate for EITC returns. The first row refers to the share of selected returns that fall into activity code 271 (gross business receipts above \$25,000). Cost refers to the estimated cost to the IRS of conducting the selected audits; it is calculated based on the share of selected returns in activity codes 270 and 271, and the average audit cost per return in those respective categories. The second row reports costs calculated from the average hours reported by IRS examiners dealing directly with the case multiplied by the applicable General Schedule payscale given the employee level. The third row instead calculates costs similarly but based on the number of hours reported by examiners in the much more intensive and comprehensive NRP exams. Costs are annualized to reflect our use of five years of NRP data. Detected underreporting refers to the total amount of underreporting (positive or negative) discovered on returns selected for audit under the specified audit selection algorithm, and is annualized to reflect our use of five years of NRP data.

## B Additional Results Relating to Disparity Estimation

In this section, we provide additional theoretical results relating to our estimation of disparity from BIFSG-derived race probability estimates.

### Contents

B.1	Results Relating to BIFSG Estimator . . . . .	Appendix-40
B.2	Proof of Proposition 1 and Related Results . . . . .	Appendix-41
B.3	Inference in Finite Samples . . . . .	Appendix-48
B.4	Incorporating Sampling Weights into Disparity Estimation .	Appendix-49
B.5	Estimating Disparity from a Recalibrated Race Probability Estimate . . . . .	Appendix-51
B.6	Gibbs Sampling . . . . .	Appendix-54
B.7	Estimating Audit Disparity Conditional on True Underreporting	Appendix-56
B.8	Conditional Disparity Estimators and Decomposition . . . .	Appendix-58
B.9	Tightening Bounds with a Linear Program . . . . .	Appendix-64

### B.1 Results Relating to BIFSG Estimator

To derive Equation (1), use Bayes rule to write:

$$\begin{aligned}\Pr[B|F, S, G] &= \frac{\Pr[F, S, G|B] \Pr[B]}{\Pr[F, S, G]} \\ &= \frac{\Pr[F|B] \Pr[S|B] \Pr[G|B] \Pr[B]}{\Pr[F, S, G]}\end{aligned}$$

where the second equation follows from the “naive” conditional independence assumption underlying the approach. Equation (1) then follows by dividing  $\Pr[B = 1|F, S, G]$  by  $\Pr[B = 0|F, S, G]$ , and using the fact that  $\Pr[B = 1|F, S, G] + \Pr[B = 0|F, S, G] = 1$ .

In the Census data we use to estimate BIFSG scores, we observe  $\Pr[B|S]$  rather than  $\Pr[S|B]$ , and we cannot back out  $\Pr[B|S]$  due to censoring of uncommon surnames. Hence, the actual BIFSG scores we estimate are derived from

$$\Pr[B|F, S, G] = \frac{\Pr[F|B] \Pr[B|S] \Pr[G|B] \Pr[S]}{\Pr[F, S, G]}$$

Dividing  $\Pr[B = 1|F, S, G]$  by  $\Pr[B = 0|F, S, G]$  leads the (unobserved)  $\Pr[S]$  terms to cancel, and following the same procedure as above we obtain:

$$\Pr[B = 1|F, S, G] = \frac{\Pr[F|B = 1] \Pr[B = 1|S] \Pr[G|B = 1]}{\sum_{j=0}^1 \Pr[F|B = j] \Pr[B = j|S] \Pr[G|B = j]}$$

which we use to estimate taxpayer-level race probabilities.

## B.2 Proof of Proposition 1 and Related Results

Recall Proposition 1:

**Proposition 1.** Suppose that  $b$  is a taxpayer's probability of being Black given some observable characteristics  $Z$ , so that  $b = \Pr[B = 1|Z]$ . Define  $D_p$  as the asymptotic limit of the probabilistic disparity estimator,  $\hat{D}_p$ , and  $D_l$  as the asymptotic limit of the linear disparity estimator,  $\hat{D}_l$ . Then:

1.

$$D_p = D - \frac{\mathbb{E}[\text{Cov}(Y, B|b)]}{\text{Var}(B)} \quad (1.1)$$

2.

$$D_l = D + \frac{\mathbb{E}[\text{Cov}(Y, b|B)]}{\text{Var}(b)} \quad (1.2)$$

3. Suppose  $\mathbb{E}[\text{Cov}(Y, B|b)] \geq 0$  and  $\mathbb{E}[\text{Cov}(Y, b|B)] \geq 0$ . Then

$$D_p \leq D \leq D_l \quad (1.3)$$

4. Suppose  $\mathbb{E}[\text{Cov}(Y, B|b)] \leq 0$  and  $\mathbb{E}[\text{Cov}(Y, b|B)] \leq 0$ . Then

$$D_l \leq D \leq D_p \quad (1.4)$$

Proposition 1 follows from the more general Proposition 2, given below. Before stating and proving it, we state and prove a lemma showing that  $D_p = D_l \cdot \frac{\text{Var}(b)}{\mathbb{E}[b](1-\mathbb{E}[b])}$  (under the mild condition that  $b$  be almost surely nontrivial; in practice, observations for which  $b$  is 0 or 1, i.e. ground truth is available, can be analyzed separately).

**Lemma 1.** Suppose that  $0 < b < 1$  almost surely, and that  $\mathbb{E}|Y|$  is finite. Then as sample size grows, the probabilistic estimator converges almost surely to:

$$D_p = D_l \cdot \frac{\text{Var}(b)}{\mathbb{E}[b](1-\mathbb{E}[b])}$$

*Proof.* We can write  $D_p$  as:

$$D_p = \frac{\sum_i b_i Y_i}{\sum_i b_i} - \frac{\sum_i (1-b_i) Y_i}{\sum_i 1-b_i} = \frac{\frac{1}{n} \sum_i b_i Y_i}{\frac{1}{n} \sum_i b_i} - \frac{\frac{1}{n} \sum_i (1-b_i) Y_i}{\frac{1}{n} \sum_i (1-b_i)}$$

For both the numerator and denominator, the strong law of large numbers holds (since  $\mathbb{E}|Y|$  is finite and, since  $0 < b < 1$ ,  $\mathbb{E}|bY|$  also is also finite), so the numerator and denominator of each of the two terms converge almost surely to their expectations. Since  $0 < b < 1$  almost surely, the continuous mapping theorem gives that the ratio of the terms converges to the

ratio of their limits. That is:

$$\left[ \frac{\frac{1}{n} \sum_i b_i Y_i}{\frac{1}{n} \sum_i b_i} - \frac{\frac{1}{n} \sum_i (1 - b_i) Y_i}{\frac{1}{n} \sum_i (1 - b_i)} \right] \xrightarrow[n \rightarrow \infty]{\text{a.s.}} \left[ \frac{\mathbb{E}[bY]}{\mathbb{E}[b]} - \frac{\mathbb{E}[(1 - b)Y]}{\mathbb{E}[1 - b]} \right]$$

Now, simply combining fractions, we note:

$$\begin{aligned} \frac{\mathbb{E}[bY]}{\mathbb{E}[b]} - \frac{\mathbb{E}[(1 - b)Y]}{\mathbb{E}[1 - b]} &= \frac{\mathbb{E}[bY] - \mathbb{E}[b]\mathbb{E}[bY] - \mathbb{E}[b]\mathbb{E}[Y] + \mathbb{E}[b]\mathbb{E}[bY]}{\mathbb{E}[b](1 - \mathbb{E}[b])} \\ &= \frac{\mathbb{E}[bY] - \mathbb{E}[b]\mathbb{E}[Y]}{\mathbb{E}[b](1 - \mathbb{E}[b])} \\ &= \frac{\text{Cov}(Y, b)}{\mathbb{E}[b](1 - \mathbb{E}[b])} \end{aligned}$$

Finally, we recall that  $D_l = \frac{\text{Cov}(Y, b)}{\text{Var}(b)}$  by construction; substituting in  $\text{Cov}(Y, b) = D_l \text{Var}(b)$  yields the result.  $\square$

**Proposition 2.** Suppose that  $b$  is a (potentially imperfectly calibrated) estimate of the probability that a taxpayer is Black, based on some observable characteristics  $Z$ . Let  $\varepsilon = B - b$  denote the error in a taxpayer's predicted race. Define  $D_p$  as the asymptotic limit of the probabilistic disparity estimator,  $\hat{D}_p$ , and  $D_l$  as the asymptotic limit of the linear disparity estimator,  $\hat{D}_l$ . Define  $\mu = \text{Cov}(\mathbb{E}[\eta|b], \mathbb{E}[\varepsilon|b])$ , where  $\eta$  denotes the residual from the linear projection of  $Y$  on  $b$ .

Then:

1.

$$D_l = D \left( 1 + \frac{\text{Cov}(b, \varepsilon)}{\text{Var}(b)} \right) + \frac{\mathbb{E}[\text{Cov}(Y, b|B)]}{\text{Var}(b)}$$

2.

$$D_p = \frac{D \cdot \text{Var}(B) - D_l \cdot \text{Cov}(b, \varepsilon)}{\mathbb{E}[b](1 - \mathbb{E}[b])} - \frac{\mathbb{E}[\text{Cov}(Y, B|b)] + \mu}{\mathbb{E}[b](1 - \mathbb{E}[b])}$$

*Proof of Proposition 2.* Consider the linear projections of  $Y$  on  $b$  and of  $Y$  on  $B$ :

$$Y = \alpha + \beta b + \eta$$

$$Y = \alpha' + \gamma B + \nu$$

By construction,  $\text{Cov}(b, \eta) = \text{Cov}(B, \nu) = 0$ . In addition,  $E[\nu] = 0$ , so

$$\gamma = E[Y|B = 1] - E[Y|B = 0] = D$$

Also, by construction:

$$\gamma \text{Var}(B) = \text{Cov}(Y, B)$$

and similarly,

$$\beta \text{Var}(b) = \text{Cov}(Y, b)$$

Using the law of total covariance, we can write:

$$\text{Cov}(Y, b) = E[\text{Cov}(Y, b|B)] + \text{Cov}(E[Y|B], E[b|B])$$

The latter term can be expanded as:

$$\begin{aligned} \text{Cov}(E[Y|B], E[b|B]) &= \text{Cov}(E[\alpha' + \gamma B + \nu|B], E[B - \varepsilon|B]) \\ &= \text{Cov}(\alpha' + \gamma B + E[\nu|B], B - E[\varepsilon|B]) \\ &= \gamma \text{Var}(B) - \gamma \text{Cov}(B, E[\varepsilon|B]) + \text{Cov}(E[\nu|B], B) - \text{Cov}(E[\nu|B], E[\varepsilon|B]) \\ &= \gamma \text{Var}(B) - \gamma \text{Cov}(B, E[\varepsilon|B]) \end{aligned}$$

where the last equality follows from the fact that since  $B$  is binary,  $\text{Cov}(B, \nu) = 0 \implies E[\nu|B] = 0$  for all  $B$ .

Next, note that

$$\begin{aligned} \text{Cov}(B, E[\varepsilon|B]) &= E[B E[\varepsilon|B]] - E[B] E[E[\varepsilon|B]] \\ &= E[E[B \varepsilon|B]] - E[B] E[E[\varepsilon|B]] \\ &= E[B \varepsilon] - E[B] E[\varepsilon] \\ &= \text{Cov}(B, \varepsilon) \\ &= \text{Cov}(b + \varepsilon, \varepsilon) \\ &= \text{Cov}(b, \varepsilon) + \text{Var}(\varepsilon) \end{aligned}$$

Combining these results, we have:

$$\begin{aligned} \beta \text{Var}(b) &= \text{Cov}(Y, b) \\ &= E[\text{Cov}(Y, b|B)] + \text{Cov}(E[Y|B], E[b|B]) \\ &= E[\text{Cov}(Y, b|B)] + \gamma \text{Var}(B) - \gamma \text{Var}(\varepsilon) - \gamma \text{Cov}(b, \varepsilon) \end{aligned}$$

From the definition of  $\varepsilon$ , we have:

$$\text{Var}(B) = \text{Var}(b) + \text{Var}(\varepsilon) + 2\text{Cov}(b, \varepsilon) \implies \text{Var}(B) - \text{Cov}(b, \varepsilon) - \text{Var}(\varepsilon) = \text{Var}(b) + \text{Cov}(b, \varepsilon)$$

Thus

$$\beta \text{Var}(b) = \gamma[\text{Var}(b) + \text{Cov}(b, \varepsilon)] + E[\text{Cov}(Y, b|B)]$$

and dividing through by  $\text{Var}(b)$  yields part 1 of the proposition.

To prove part 2 of the proposition, again use the law of total covariance:

$$\text{Cov}(Y, B) = E[\text{Cov}(Y, B|b)] + \text{Cov}(E[Y|b], E[B|b])$$

Expanding the second term of the right-hand side of the equation, we have

$$\begin{aligned} \text{Cov}(E[Y|b], E[B|b]) &= \text{Cov}(E[\alpha + \beta b + \eta|b], E[b + \varepsilon|b]) \\ &= \text{Cov}(\alpha + \beta b + E[\eta|b], b + E[\varepsilon|b]) \\ &= \beta \text{Var}(b) + \beta \text{Cov}(b, E[\varepsilon|b]) + \text{Cov}(E[\eta|b], b) + \text{Cov}(E[\eta|b], E[\varepsilon|b]) \end{aligned}$$

Note that:

$$\begin{aligned} \text{Cov}(b, E[\varepsilon|b]) &= E[b E[\varepsilon|b]] - E[b]E[E[\varepsilon|b]] \\ &= E[E[b \varepsilon|b]] - E[b]E[E[\varepsilon|b]] \\ &= E[b \varepsilon] - E[b]E[\varepsilon] \\ &= \text{Cov}(b, \varepsilon) \end{aligned}$$

By the same logic:

$$\text{Cov}(E[\eta|b], b) = \text{Cov}(\eta, b) = 0$$

Define  $\mu := \text{Cov}(E[\eta|b], E[\varepsilon|b])$ . Then collecting results, we have

$$\begin{aligned} \gamma \text{Var}(B) &= \text{Cov}(Y, B) \\ &= E[\text{Cov}(Y, B|b)] + \text{Cov}(E[Y|b], E[B|b]) \\ &= E[\text{Cov}(Y, B|b)] + \beta \text{Var}(b) + \beta \text{Cov}(b, \varepsilon) + \mu \end{aligned}$$

Rearranging, and recalling that  $D_p = \beta \frac{\text{Var}(b)}{\mathbb{E}[b](1-\mathbb{E}[b])}$  from Lemma 1 yields the result. □

Now, we prove Proposition 1 as a consequence of Proposition 2.

*Proof of Proposition 1.* If  $b = \Pr[B = 1|Z] = \mathbb{E}[B|Z]$ , it follows from the definition of  $\varepsilon$  that

$$\begin{aligned} E[\varepsilon|Z] &= E[B|Z] - E[b|Z] \\ &= E[B|Z] - E[E[B|Z]|Z] \\ &= E[B|Z] - E[B|Z] \\ &= 0 \end{aligned}$$



Hence, we can write

$$\begin{aligned}
\text{Cov}(b, \varepsilon) &= E[b \varepsilon] - E[b] E[\varepsilon] \\
&= E[b \varepsilon] \\
&= E[E[b \varepsilon | Z]] \\
&= E[b E[\varepsilon | Z]] \\
&= E[b \cdot 0] \\
&= 0,
\end{aligned}$$

where the third equality follows from the law of iterated expectations, and the fourth from the fact that  $b$  is a function of  $Z$ .

Substituting the fact that  $\text{Cov}(b, \varepsilon) = 0$  into Proposition 2.1, and noting that since  $\mathbb{E}[b] = \mathbb{E}[\mathbb{E}[B|Z]] = \mathbb{E}[B]$ ,

$$\mathbb{E}[b](1 - \mathbb{E}[b]) = \mathbb{E}[B](1 - \mathbb{E}[B]) = \text{Var}(B).$$

yields Proposition 1.2.

Proposition 1.1 follows by again substituting in  $\text{Cov}(b, \varepsilon) = 0$  and noting that  $\mathbb{E}[\varepsilon|b] = 0$  because:

$$\begin{aligned}
E[\varepsilon|b] &= E[E[\varepsilon|b, Z]|b] \\
&= E[E[\varepsilon|Z]|b] \\
&= E[0|b] \\
&= 0,
\end{aligned}$$

where the second equality follows from the fact that  $b$  is a function of  $Z$ .

Finally, once the forms of  $D_l$  and  $D_p$  are established, Proposition 1.3 and 1.4 follow directly when the respective assumptions on the signs of  $\mathbb{E}[\text{Cov}(Y, b|B)]$  and  $\mathbb{E}[\text{Cov}(Y, B|b)]$  are met.  $\square$

The following Proposition extends Proposition 1 to the case in which the estimand of interest is the level of the outcome by group, rather than the difference in the levels of the outcome across groups. In particular, Proposition 3 characterizes the bias of the linear and probabilistic approaches for estimating the audit rate by race (not the audit rate disparity).

**Proposition 3.** (Statistical Bias of Level Estimators). Suppose  $b = \Pr[B|Z]$ . Consider the following estimators:

$$\begin{aligned}
\hat{Y}_p^B &:= \frac{\sum b_i Y_i}{\sum b_i} & \text{and} & & \hat{Y}_p^{NB} &:= \frac{\sum (1 - b_i) Y_i}{\sum (1 - b_i)} \\
\hat{Y}_l^B &:= \hat{\alpha} + \hat{\beta} & \text{and} & & \hat{Y}_l^{NB} &:= \hat{\alpha},
\end{aligned}$$

where  $\hat{\alpha}$  and  $\hat{\beta}$  are the intercept and slope, respectively, from the regression of  $Y$  on  $b$ . Let  $Y_p^B, Y_p^{NB}, Y_l^B, Y_l^{NB}$  be the respective limits the estimators described above converge to. Then:

1.  $Y_l^B$  and  $Y_l^{NB}$  have the following biases relative to the true audit rates  $Y^B$  and  $Y^{NB}$ :

$$Y_l^B - Y^B = (1 - \mathbb{E}[B]) \frac{\mathbb{E}[\text{Cov}(Y, b|B)]}{\text{Var}(b)} \quad \text{and} \quad Y_l^{NB} - Y^{NB} = -\mathbb{E}[B] \frac{\mathbb{E}[\text{Cov}(Y, b|B)]}{\text{Var}(b)}$$

2.  $Y_p^B$  and  $Y_p^{NB}$  have the following biases relative to the true audit rates  $Y^B$  and  $Y^{NB}$ :

$$Y_p^B - Y^B = -\frac{\mathbb{E}[\text{Cov}(Y, B|b)]}{\mathbb{E}[B]} \quad \text{and} \quad Y_p^{NB} - Y^{NB} = \frac{\mathbb{E}[\text{Cov}(Y, B|b)]}{1 - \mathbb{E}[B]}$$

3. Suppose  $\mathbb{E}[\text{Cov}(Y, b|B)] = 0$ . Then:

$$Y_l^B = Y^B \quad \text{and} \quad Y_l^{NB} = Y^{NB}$$

4. Suppose  $\mathbb{E}[\text{Cov}(Y, B|b)] = 0$ . Then:

$$Y_p^B = Y^B \quad \text{and} \quad Y_p^{NB} = Y^{NB}$$

*Proof.* Notice that 3) and 4) follow directly from 1) and 2). For 1): By construction, we have

$$Y = \alpha + \gamma B + \nu$$

From this, we know  $Y^{NB} = \alpha$  and  $Y^B = \alpha + \gamma$ .

Taking expectations, and rearranging:

$$Y^{NB} = \alpha = E[Y] - \gamma E[B]$$

In contrast, our sample estimate of  $Y^{NB}$  from the linear estimator,  $\hat{Y}_l^{NB}$ , is given by:

$$\hat{Y}_l^{NB} = \hat{\alpha} = \bar{Y} - \hat{\beta} \bar{b}$$

which converges to

$$Y_l^{NB} = E[Y] - D_l E[b] = \mathbb{E}[Y] - D_l \mathbb{E}[B],$$

since  $\mathbb{E}[b] = \mathbb{E}[B]$  (because  $\mathbb{E}[b] = \mathbb{E}_Z[\Pr[B = 1|Z]] = \mathbb{E}[B]$ ).

From Proposition 1, we know  $D_l = \gamma + \frac{\mathbb{E}[\text{Cov}(Y, b|B)]}{\text{Var}(b)}$

Substituting this into the above, we have:

$$\begin{aligned} Y_l^{NB} &= E[Y] - \gamma E[B] - E[B] \frac{\mathbb{E}[\text{Cov}(Y, b|B)]}{\text{Var}(b)} \\ &= Y^{NB} - E[B] \frac{\mathbb{E}[\text{Cov}(Y, b|B)]}{\text{Var}(b)} \end{aligned}$$

Turning to  $Y_l^B$ , we have

$$\widehat{Y}_l^B = \widehat{\alpha} + \widehat{\beta}$$

which converges to

$$\begin{aligned} Y_l^B &= Y_l^{NB} + D_l \\ &= \left( \alpha - E[B] \frac{\mathbb{E}[\text{Cov}(Y, b|B)]}{\text{Var}(b)} \right) + D_l \\ &= \alpha - E[B] \frac{\mathbb{E}[\text{Cov}(Y, b|B)]}{\text{Var}(b)} + \gamma + \frac{\mathbb{E}[\text{Cov}(Y, b|B)]}{\text{Var}(b)} \\ &= \alpha + \gamma + (1 - E[B]) \frac{\mathbb{E}[\text{Cov}(Y, b|B)]}{\text{Var}(b)} \\ &= Y^B + (1 - E[B]) \frac{\mathbb{E}[\text{Cov}(Y, b|B)]}{\text{Var}(b)} \end{aligned}$$

We prove 2) in a very similar manner as the related statement is in Chen et al. (2019):  
Note that:

$$Y^B = \mathbb{E}[Y|B=1] = \frac{\mathbb{E}[YB]}{\mathbb{E}[B]} = \frac{\mathbb{E}[\mathbb{E}[YB|b]]}{\mathbb{E}[B]}$$

On the other hand,

$$\widehat{Y}_p^B = \frac{\frac{1}{n} \sum b_i Y_i}{\frac{1}{n} \sum b_i} \longrightarrow \frac{\mathbb{E}[Yb]}{\mathbb{E}[b]} := Y_p^B$$

since the law of large numbers applies to the numerator and the denominator separately (and the boundedness away from the end of the interval guarantees that the limits of the ratio converges to the ratio of the limits).

But  $\mathbb{E}[b] = \mathbb{E}_Z[\Pr[B=1|Z]] = \mathbb{E}[B]$ , and  $\mathbb{E}[Yb] = \mathbb{E}[\mathbb{E}[Yb|b]] = \mathbb{E}[b\mathbb{E}[Y|b]] = \mathbb{E}[\mathbb{E}[B|b]\mathbb{E}[Y|b]]$ , so:

$$Y_p^B - Y^B = \frac{\mathbb{E}[\mathbb{E}[Y|b]\mathbb{E}[B|b]]}{\mathbb{E}[B]} - \frac{\mathbb{E}[\mathbb{E}[YB|b]]}{\mathbb{E}[B]} = -\frac{\mathbb{E}[\text{Cov}(Y, B|b)]}{\mathbb{E}[B]}$$

where the second equality follows from the definition of conditional covariance. This establishes the result for  $Y_p^B$ . To see the analogous result for  $Y_p^{NB}$ , let  $A = 1 - B$  and  $a = b$ , and observe that  $-\mathbb{E}[\text{Cov}(Y, A|a)] = \mathbb{E}[\text{Cov}(Y, B|b)]$ . The result then follows in the same manner as above.

□

### B.3 Inference in Finite Samples

This section characterizes the asymptotic distributions of the linear and probabilistic estimators.

#### B.3.1 Standard Errors of Disparity Estimators

Call  $\hat{D}_l^n$  and  $\hat{D}_p^n$  the empirically-constructed linear and probabilistic estimators using a sample size of  $n$  observations. ( $D_l$  and  $D_p$  as written above are what  $\hat{D}_l^n$  and  $\hat{D}_p^n$  converge to as  $n \rightarrow \infty$ .)

**Lemma 2.** For any fixed dataset, we relate  $\hat{D}_p^n$  and  $\hat{D}_l^n$  as:

$$\hat{D}_p^n = \hat{D}_l^n \cdot \frac{\frac{1}{n} \sum_i b_i^2 - \bar{b}^2}{\bar{b}(1 - \bar{b})}$$

And asymptotically,

$$\hat{D}_p^n \rightarrow \hat{D}_l^n \cdot \frac{\text{Var}(b)}{\mathbb{E}[b](1 - \mathbb{E}[b])}.$$

*Proof.* Notice that:

$$\begin{aligned} \hat{D}_p^n &= \frac{\sum b_i Y_i}{\sum b_i} - \frac{\sum (1 - b_i) Y_i}{\sum 1 - b_i} = \frac{\frac{1}{n} \sum b_i Y_i}{\frac{1}{n} \sum b_i} - \frac{\frac{1}{n} \sum (1 - b_i) Y_i}{\frac{1}{n} \sum (1 - b_i)} \\ &= \frac{\frac{1}{n} \sum b_i Y_i}{\bar{b}} - \frac{\frac{1}{n} \sum (1 - b_i) Y_i}{1 - \bar{b}} \end{aligned}$$

where we use  $\bar{\cdot}$  to indicate the sample average. We can then write:

$$\begin{aligned} &\frac{\frac{1}{n} \sum b_i Y_i}{\frac{1}{n} \sum b_i} - \frac{\frac{1}{n} \sum (1 - b_i) Y_i}{\frac{1}{n} \sum (1 - b_i)} = \frac{\frac{1}{n} \sum b_i Y_i}{\bar{b}} - \frac{\frac{1}{n} \sum (1 - b_i) Y_i}{1 - \bar{b}} \\ &= \frac{\frac{1}{n} \sum b_i Y_i - \frac{\bar{b}}{n} \sum b_i Y_i - \frac{\bar{b}}{n} \sum Y_i + \frac{\bar{b}}{n} \sum b_i Y_i}{\bar{b}(1 - \bar{b})} \\ &= \frac{\frac{1}{n} \sum b_i Y_i - \bar{b} \bar{y}}{\bar{b}(1 - \bar{b})} \end{aligned}$$

Now consider the regression estimator. By definition:

$$D_l^n = \frac{\sum (b_i - \bar{b})(y_i - \bar{y})}{\sum (b_i - \bar{b})^2} = \frac{\sum b_i y_i - \bar{b} \sum y_i - \bar{y} \sum b_i + n \bar{b} \bar{y}}{\sum (b_i - \bar{b})^2} = \frac{\frac{1}{n} \sum b_i Y_i - \bar{b} \bar{y}}{\frac{1}{n} \sum (b_i - \bar{b})^2}$$

But notice the numerator in both terms are the same. That is:

$$\hat{D}_p^n = \frac{\frac{1}{n} \sum (b_i - \bar{b})^2}{\bar{b}(1 - \bar{b})} \hat{D}_l^n = C_n \hat{D}_l^n$$

where  $C_n = \frac{\frac{1}{n} \sum (b_i - \bar{b})^2}{\bar{b}(1 - \bar{b})}$ .

But now recall Slutsky's theorem, which says that if  $A_n, B_n$  are random variables and  $B_n \rightarrow c$  for some constant  $c$ , then  $A_n B_n \rightarrow A_n c$ . In particular,

$$C_n \rightarrow \frac{\text{Var}(b)}{\mathbb{E}[b](1 - \mathbb{E}[b])}.$$

The second half of the lemma follows.  $\square$

The asymptotic distribution of  $\hat{D}_l^n$  is well understood, as it is the OLS estimator.

Given the relationship between  $\hat{D}_p^n$  and  $\hat{D}_l^n$  shown above, it is mechanically true that  $\hat{D}_p^n$  will, under the same conditions, be distributed normally as well. Formally:

**Proposition 4.** The asymptotic distribution of  $D_p^n$  is given by:

$$\frac{\hat{D}_p^n - D_p}{\sqrt{V_l^n \frac{\text{Var}(b)}{\mathbb{E}[b](1 - \mathbb{E}[b])}}} \rightarrow \mathcal{N}(0, 1)$$

where  $D_p = \frac{\text{Cov}(Y, b)}{\mathbb{E}[b](1 - \mathbb{E}[b])}$  and  $V_l^n$  is the variance of  $\hat{D}_l^n$ .

## B.4 Incorporating Sampling Weights into Disparity Estimation

In some of our analyses, we use data which is re-weighted to be representative of the full population of U.S. taxpayers.  $\hat{D}^l$  can be naturally extended to incorporate sample weights via weighted regression. How to extend the probabilistic estimator, however, may be less obvious. We propose the following as the weighted probabilistic estimator  $\hat{D}_{p,w}$ :

$$\hat{Y}_{p,w}^B = \frac{\sum_i \omega_i b_i Y_i}{\sum_i \omega_i b_i}, \quad \hat{Y}_p^{NB} = \frac{\sum_i \omega_i (1 - b_i) Y_i}{\sum_i \omega_i (1 - b_i)}, \quad \hat{D}_{p,w} := \hat{Y}_{p,w}^B - \hat{Y}_{p,w}^{NB}$$

where  $\omega_i$  is a sample weight for observation  $i$ . (Notice that as with  $D_p$ , replacing  $Y_i$  with any other random variable gives an estimator for disparity in said random variable.) This estimator is closely related to the Horwitz-Thompson family of estimators; see Robinson (1982); Berger (1998); Delevoye & Sävje (2020) for prior results regarding convergence and consistency.

What is the purpose of the weighted estimator? The intention behind it is to use the data we have to estimate what  $D_p$  *would* be given a different dataset or distribution. If  $\hat{D}_{p,w}$  provides this fidelity, then it is the ‘correct’ weighted analogue. We show here that  $\hat{D}_{p,w}$  is the ‘correct’ weighted analogue in a sense we make formal below.

We will distinguish between two cases. In the first, we have access to a subset of a finite population of individuals, and are given *replicate weights*. The replicate weight for observation  $i$  corresponds to the number of individuals in the full population that  $i$  represents. In other words, we have some dataset of observations  $\mathcal{D} := \{X_i\}_{i=1}^n$ , but the full dataset which we do not have access to has observations  $\mathcal{D}' := \bigcup_i \{X_i\}_{j=1}^{\omega_i}$ . The hope is that  $\hat{D}_{p,w}$  estimated on  $\mathcal{D}$  corresponds to  $\hat{D}_p$  estimated on  $\mathcal{D}'$ .

**Proposition 5.** Suppose we are in the case of replicate weights and  $\mathcal{D}$  and  $\mathcal{D}'$  are as above. Let  $\widehat{D}_{p,w}|\mathcal{D}$  be estimated over  $\mathcal{D}$  and  $\widehat{D}_p|\mathcal{D}'$  be what would be estimated over  $\mathcal{D}'$ . Then:

$$\widehat{D}_{p,w}|\mathcal{D} = \widehat{D}_p|\mathcal{D}'$$

*Proof.* This follows simply from the linearity of the numerator and denominator of  $\widehat{Y}_{p,w}^B$  and  $\widehat{Y}_{p,w}^{NB}$ . Take  $\widehat{Y}_{p,w}^B$ :

$$\widehat{Y}_{p,w}^B|\mathcal{D} = \frac{\sum_i w_i b_i Y_i}{\sum_i w_i b_i} = \frac{\sum_i \sum_{j=1}^{w_i} b_i Y_i}{\sum_i \sum_{j=1}^{w_i} b_i} = \widehat{Y}_p^B|\mathcal{D}'.$$

$\widehat{Y}_{p,w}^{NB}$  follows similarly and thus too  $\widehat{D}_{p,w}$ . □

Notice that the case of replicate weights corresponds to our analyses in which we use NRP to estimate quantities over the population.

The second case is more general: weights may be not be integers corresponding the number of people represented in some larger dataset, but rather changes of measure intended to capture some other distribution. (For instance, weighting for non-response attempts to map the data from responders to the overall population.) In this setting, we are agnostic to how the weights are generated; instead, we merely assume that they successfully accomplish re-weighting at the level of the sample mean. We make this precise in the following proposition:

**Proposition 6.** Suppose we have data drawn from a distribution  $\mathcal{D}$ ; this data includes both a quantity of interest,  $Y_i$ , as well as sample weights  $\omega_i$  that map  $\mathcal{D}$  to some other distribution  $\mathcal{D}'$  in the following sense:

$$\frac{1}{n} \sum_{i=1}^n \omega_i Q_i \xrightarrow{n \rightarrow \infty} \mathbb{E}_{\mathcal{D}'}[Q],$$

for any random variable  $Q$ . Then:

$$\widehat{D}_{p,w}^n|\mathcal{D} \xrightarrow{n \rightarrow \infty} D_p|\mathcal{D}'.$$

*Proof.* Consider  $\widehat{Y}_{p,w}^B$ . Let  $Q := bY$ . Then by assumption:

$$\frac{1}{n} \left[ \sum_{i=1}^n \omega_i b_i Y_i \right] \xrightarrow{n \rightarrow \infty} \mathbb{E}_{\mathcal{D}'}[bY]$$

Similarly,

$$\frac{1}{n} \sum_{i=1}^n \omega_i b_i \xrightarrow{n \rightarrow \infty} \mathbb{E}_{\mathcal{D}'}[b]$$

But we have that:

$$\hat{Y}_{p,w}^{B,n} = \frac{\sum_i \omega_i b_i Y_i}{\sum_i \omega_i b_i} = \frac{\frac{1}{n} \sum_i \omega_i b_i Y_i}{\frac{1}{n} \sum_i \omega_i b_i} \xrightarrow{n \rightarrow \infty} \frac{\mathbb{E}_{\mathcal{D}'}[bY]}{\mathbb{E}_{\mathcal{D}'}[b]}$$

Proceeding similarly with  $\hat{Y}_{p,w}^{NB,n}$  and taking the difference, we obtain:

$$\hat{D}_{p,w}^n \xrightarrow{n \rightarrow \infty} \frac{\mathbb{E}_{\mathcal{D}'}[bY]}{\mathbb{E}_{\mathcal{D}'}[b]} - \frac{\mathbb{E}_{\mathcal{D}'}[(1-b)Y]}{\mathbb{E}_{\mathcal{D}'}[1-b]} = D_p | \mathcal{D}'$$

□

Notice that the choice of unit weights, i.e.  $\omega_i = 1$  satisfies the assumption of the theorem and recovers the original convergence results. For another example, suppose we have groups A and B in equal number throughout the population, but in our data we obtain twice as many observations from group B as group A. Then it is easy to verify that the choice of weights  $\omega_i = \begin{cases} 2/3 & i \in A \\ 3/4 & i \in B \end{cases}$  would satisfy the assumptions, and thus this choice of weights would allow us to recover  $D_p$  in the population from our data.

## B.5 Estimating Disparity from a Recalibrated Race Probability Estimate

In general, our estimate of  $\Pr[B_i = 1|Z_i]$  may be drawn from a population that differs from the population of interest; for example, among EITC claimants,  $\Pr[B_i = 1|Z_i]$  may be different than the equivalent quantity for the population as a whole. (This is reflected in the imperfect calibration visible in Figure 2.) Proposition 2 shows that systematic deviation like this can bias our estimates. However, given access to a *recalibrated*  $b^*$ , (i.e. the linear projection of  $B$  on to the space of  $b$  and a constant, based on a subset of data with ground truth race labels), we can use this re-calibrated proxy to obtain similar results as those in Proposition 1 by applying Proposition 2.

To see this, suppose access to such a  $b^*$  and note that by construction,  $\text{Cov}(b^*, \varepsilon^*) = 0$ , so Proposition 2 applies. Moreover,  $\mathbb{E}[b^*] = \mathbb{E}[B]$ , so  $\mathbb{E}[b^*](1 - \mathbb{E}[b^*]) = \mathbb{E}[B](1 - \mathbb{E}[B]) = \text{Var}(B)$ . So designating  $D_l^*$  and  $D_p^*$  as the linear and probabilistic estimators, respectively, applied to  $b^*$ , as well as  $\eta^*$  and  $\varepsilon^*$  for the analogues of  $\eta$  and  $\varepsilon$ , Proposition 2 indicates that:

$$D_l^* = D + \frac{\mathbb{E}[\text{Cov}(Y, b^*|B)]}{\sigma_{b^*}^2} \quad (2)$$

and

$$D_p^* = D - \frac{\mathbb{E}[\text{Cov}(Y, B|b^*)] + \text{Cov}(\mathbb{E}[\eta^*|b^*], \mathbb{E}[\varepsilon^*|b^*])}{\text{Var}(B)}. \quad (3)$$

These equations are similar to those of Proposition 1, but there are two potential challenges to overcome before applying the recalibrated proxies in the same manner as our

initial probability estimates. The first potential challenge is that it might be more difficult to reason about the sign of the covariances between outcome and re-calibrated proxy than the original proxy, since this recalibrated proxy could differ in its conditional relationship to  $Y$  given  $B$  and vice versa.

The following Lemma addresses this issue; it shows that the main covariance terms  $\mathbb{E}[\text{Cov}(Y, B|b^*)]$  and  $\mathbb{E}[\text{Cov}(Y, b^*|B)]$  will have the same signs as their non-recalibrated counterparts, under the minor condition that  $\text{Cov}(B, b) > 0$ .

**Lemma 3** (). Suppose that  $b$  is a (possibly mis-calibrated) estimate of the probability that a taxpayer is Black based on some observable characteristics  $Z$  and  $b^*$  is the re-calibrated proxy which can be written as an orthogonal projection:

$$B = \mu + \rho b + \varphi,$$

i.e.

$$b^*(b) = \mu + \rho b.$$

Suppose further that  $\text{Cov}(B, b) > 0$ . Then

$$\begin{aligned} \text{sign}(\mathbb{E}[\text{Cov}(Y, B|b)]) &= \text{sign}(\mathbb{E}[\text{Cov}(Y, B|b^*)]) \\ \text{sign}(\mathbb{E}[\text{Cov}(Y, b|B)]) &= \text{sign}(\mathbb{E}[\text{Cov}(Y, b^*|B)]) \end{aligned}$$

*Proof.* We note that

$$\text{Cov}(Y, b^*|B) = \text{Cov}(Y, \mu + \rho b|B) = \rho \text{Cov}(Y, b|B)$$

and

$$\text{Cov}(Y, B|b^*) = \text{Cov}(Y, B|b).$$

Then the signs of  $\mathbb{E}[\text{Cov}(Y, B|b)]$  and  $\mathbb{E}[\text{Cov}(Y, B|b^*)]$  are identical, while the signs of  $\mathbb{E}[\text{Cov}(Y, b|B)]$  and  $\mathbb{E}[\text{Cov}(Y, b^*|B)]$  will agree if  $\rho \geq 0$ . Since  $\rho$  is the coefficient on  $b$  in said regression, it is given by  $\text{Cov}(B, b)/\text{Var}(b)$ , which is positive if and only if  $\text{Cov}(B, b) > 0$ .  $\square$

The second potential difficulty in applying the recalibrated proxy is that the additional covariance term that appears in the expression for  $D_p$ 's asymptotic bias,  $\text{Cov}(\mathbb{E}[\varepsilon^*|b^*], \mathbb{E}[\eta^*|b^*])$ , may also be difficult to reason about on the basis of theory. To interpret this term, note that when  $\mathbb{E}[B|b]$  is a linear function of  $b$ , then (asymptotically) recalibrating via linear regression will recover the true conditional expectation function. By construction in that case,  $\mathbb{E}[\varepsilon^*|b^*]$  will be 0, since  $b^*$  is the CEF. Thus, under linearity, the nuisance covariance term would be 0. Of course, we do not expect *exact* linearity to hold, but we will see in the recalibration exercise below that the CEF is close to linear, so that our estimate of  $\text{Cov}(E[\eta^*|b^*], E[\varepsilon^*|b^*])$  is close to 0 and certainly negligible compared to the other terms.



The upshot of the results in this subsection is that *if* we expect the covariance conditions to hold with our original proxy, we can also expect them to hold for our recalibrated proxy. This in turn allows us to treat disparity estimates arrived at using the recalibrated proxies in the same manner as estimates obtained using the original proxies. We now turn to our empirical approach.

**Empirical Approach** We now describe we apply the aforementioned strategy to re-calibrate the BIFSG-predicted probability Black in the North Carolina dataset. We consider North Carolina as a whole as well as EITC and non-EITC specific approaches.

First, we calculate  $\hat{\rho}$  as the coefficient from regressing an indicator for whether a taxpayer self reports as Black on the BIFSG-predicted probability that a taxpayer is Black. That is, we run the regression:

$$B = \alpha_0 + \rho b + \varphi,$$

with  $\hat{\rho}$  estimated once via ordinary least squares and separately via weighted least squares using the North Carolina weights. We also repeat both estimations separately for EITC taxpayers and non-EITC taxpayers. (The additional weighted/non-weighted and EITC/non-EITC calculations will be repeated throughout; where required, weighted estimates will be computed using the estimators described in B.4 above.) These estimates are reported in the first line of Table B.1.

Next, we assign each individual

$$b_i^* := \hat{\alpha}_0 + \hat{\rho} b_i,$$

and

$$\varepsilon_i^* = B_i - b_i^*;$$

we then estimate  $\widehat{\text{Cov}}(b, \varepsilon^*)$  in the straightforward manner of using sample unweighted and weighted averages and product of  $b$  and  $\varepsilon^*$ , again separating out by EITC status. These estimates are reported in the second line of Table B.1.

The next four lines of Table B.1 are computed in a similar manner. That is, we compute the covariance within a given realization of the conditioning variable (e.g. for the set of non-Black taxpayers,  $B = 0$ ) and then weight these estimates by the estimated share of taxpayers they represent. Importantly, we discretize both  $b_i^*$  and  $b_i$  by rounding to the nearest percentage point in order to create realizations to average over; this approach may introduce some arbitrariness to the analysis, but avoids making parametric assumptions.

Next, we run the regression:

$$Y = \alpha^* + \beta^* b^* + \eta^*$$

and interpret the estimated  $\hat{\beta}^*$ ; that is,  $\hat{\beta}^*$  is the linear estimator of disparity as applied to

the re-calibrated  $b^*$ . We then obtain

$$\hat{\eta}_i^* = Y_i - \hat{\alpha}^* - \hat{\beta}^* b_i^*.$$

We use this to compute the next line of Table B.1 in the following manner: first, we estimate  $\mathbb{E}[\eta^*|b^*]$  and  $\mathbb{E}[\varepsilon^*|b^*]$  by computing the sample averages of  $\eta^*$  and  $\varepsilon^*$  within each discretized  $b^*$  category. We then assign each individual their respective sample averages based on their value of  $b_i^*$ , and then compute the overall covariance estimate over the population using these features.

The next three lines of Table B.1 are computed straightforwardly - i.e.  $\hat{D}$  is based on the ground truth, while  $\hat{D}_l^*$  and  $\hat{D}_p^*$  are computed according to the formulas in equations 2 and 3 above and the appropriate values from previously computed rows of Table B.1.

Table B.1: Estimates from the Re-calibration Exercise

Value	Overall		EITC		Non-EITC	
	NC (1)	Rewighted NC (2)	NC (3)	Rewighted NC (4)	NC (5)	Rewighted NC (6)
$\hat{\rho}$	0.828	0.872	0.923	0.964	0.767	0.802
$\widehat{\text{Cov}}(b, \varepsilon)$	-0.000	-0.000	-0.000	-0.000	0.000	-0.000
$\hat{E}[\text{Cov}(Y, b B)]$	0.000	0.000	0.001	0.001	0.000	-0.000
$\hat{E}[\text{Cov}(Y, B b)]$	0.001	0.001	0.001	0.001	0.000	0.000
$\hat{E}[\text{Cov}(Y, b^* B)]$	0.000	0.000	0.001	0.001	0.000	-0.000
$\hat{E}[\text{Cov}(Y, B b^*)]$	0.001	0.001	0.001	0.001	0.000	0.000
$\widehat{\text{Cov}}(\mathbb{E}[\eta^* b^*], \mathbb{E}[\varepsilon^* b^*])$	0.000	0.000	-0.000	-0.000	0.000	0.000
$\hat{D}_l^*$	0.011	0.016	0.026	0.031	0.002	0.002
$\hat{D}$	0.008	0.012	0.018	0.024	0.002	0.002
$\hat{D}_p^*$	0.005	0.007	0.012	0.017	0.001	0.001

*Notes:* The table details the estimates for the disparities and covariance terms obtained from re-calibration, for the North Carolina dataset (both unweighted (odd columns) and re-weighted (even columns)). The estimates are calculated for the overall population (columns 1 and 2), the EITC population (columns 3 and 4), and the non-EITC population (columns 5 and 6).

## B.6 Gibbs Sampling

In addition to taxpayers' first name, surname, and geographic location, the IRS has access to additional information that may correlate to race. In principle, leveraging such additional information could lead to better estimates of race probabilities and thus of disparity. Additionally, it is possible that a finer breakdown of self-identified race and ethnicity could contain additional information that may affect our disparity estimates. Hence, as an additional robustness check, we leverage income (bucketed into 14 categories) and marital status, abbreviated "MARS" (Single, Married Filing Jointly, or Other) to obtain more accurate race/ethnicity estimates (at the more granular level of Hispanic, non-Hispanic White, non-Hispanic Black, and Other).

To our knowledge, there are no readily available marginal distributions of race/Hispanic probabilities conditional on income or marital status (and said distributions may differ

among taxpayers than the general population); hence, we use Gibbs sampling to obtain approximate probabilities from the IRS' data and BIFSG. Gibbs sampling is a Bayesian algorithm that reduces the problem of sampling from complicated joint distributions to sampling from simpler marginal ones; in this section, we describe in detail this procedure and how we apply it to our setting.

As a starting point, we take the conditional distribution of race and Hispanic origin (RH) given first name, surname, and geography (F, S, and G, respectively, and, collectively, FSG), implied by BIFSG to be correct. We model the joint distribution of  $(RH, FSG, X)$ , where  $X$  represents  $(income, MARS)$ , as a decomposable model with generating components  $\{[RH, F][RH, S][RH, G][RH, X]\}$ . (In other words, we make a similar naive Bayes assumption as in BIFSG, but treating  $X$  as a unit and allowing a more general relationship between income and MARS.) Given this model, we can write the conditional distribution of RH given  $(X, FSG)$  as

$$\Pr(RH|X, FSG) = \text{Multi} \left( n, C\boldsymbol{\theta}_{(i,j)} \frac{\Pr(RH|G)}{\Pr(RH)} \frac{\Pr(RH|F)}{\Pr(RH)} \frac{\Pr(RH|S)}{\Pr(RH)} \right)$$

where  $\boldsymbol{\theta}_{i,j}$  is a vector of probabilities for the RH categories, given  $(X_1 = i, X_2 = j)$ ; that is,  $\boldsymbol{\theta}_{i,j} = \Pr(RH|X_1 = i, X_2 = j)$ . Note also that  $\text{Multi}(n, \mathbf{p})$  represents the multinomial distribution with  $n$  draws and class probabilities  $\mathbf{p}$ , and  $C$  is a normalizing constant.

We estimate the parameters in the model with a Bayesian procedure, so we need a prior on the unknown parameter  $\boldsymbol{\theta}_{(i,j)}$ . We set that to the Dirichlet prior with vector parameter  $\boldsymbol{\alpha}_0 = (1, \dots, 1)$ , denoted  $\boldsymbol{\theta}_{(i,j)} \sim \text{Dir}(\boldsymbol{\alpha}_0)$ . This value for  $\boldsymbol{\alpha}_0$  was chosen to contribute a small amount of information to the model while ensuring that the posterior is well-behaved. Denote the unobserved vector of counts in the RH categories as  $\mathbf{n}_{i,j}$  for  $(X_1 = i, X_2 = j)$ . Given the form of the model,  $\mathbf{n}_{i,j}|X \sim \text{Multi}(n, \boldsymbol{\theta}_{i,j})$ , the Dirichlet distribution was chosen because it is the conjugate prior for the multinomial distribution and  $\boldsymbol{\theta}_{(i,j)}|\mathbf{n}_{i,j} \sim \text{Dir}(\boldsymbol{\alpha}_0 + \mathbf{n}_{i,j})$ . Now we have the full conditional distributions for the unobserved variables in the model.

$$\Pr(RH|X, F = f, S = s, G = g) = \text{Multi} \left( n, C\boldsymbol{\theta}_{(i,j)} \frac{\Pr(RH|g)}{\Pr(RH)} \frac{\Pr(RH|f)}{\Pr(RH)} \frac{\Pr(RH|s)}{\Pr(RH)} \right) \quad (4)$$

and

$$\Pr(\boldsymbol{\theta}_{(i,j)}|\mathbf{n}_{i,j}) = \text{Dir}(\boldsymbol{\alpha}_0 + \mathbf{n}_{i,j}),$$

which make a Gibbs sampling algorithm available for estimation.

An outline of the Gibbs Sampling algorithm used here is provided below. Note the superscript  $(b)$  indexes the iteration number; it is not an exponent.

- Initialization

- For each record, indexed by  $m$ , generate  $RH_m^{(0)}$  from

$$\Pr(RH_m|f_m, s_m, g_m) \sim \text{multi} \left( 1, C \frac{\Pr(RH|g_m)}{\Pr(RH)} \frac{\Pr(RH|f_m)}{\Pr(RH)} \frac{\Pr(RH|s_m)}{\Pr(RH)} \right)$$

where again  $\text{Multi}(n, \mathbf{p})$  represents the multinomial distribution with size  $n$  and probability  $\mathbf{p}$  and  $C$  is a normalizing constant.

– Tabulate  $\mathbf{n}_{i,j}^{(0)} = \sum_{X_m=(i,j)} RH_m^{(0)}$

- Main Loop

– for  $b = 1, \dots, B + b_0$ :

\* generate  $\boldsymbol{\theta}_{i,j}^{(b)} \sim \text{Dir}\left(1, \boldsymbol{\alpha}_0 + \mathbf{n}_{i,j}^{(b-1)}\right)$  for each  $i, j$

\* generate  $RH_m^{(b)}$  as

$$RH_m^{(b)} \sim \text{Multi}\left(1, C\boldsymbol{\theta}_{(i,j)_m}^{(b)} \frac{\Pr(RH|g_m)}{\Pr(RH)} \frac{\Pr(RH|f_m)}{\Pr(RH)} \frac{\Pr(RH|s_m)}{\Pr(RH)}\right)$$

\* tabulate  $\mathbf{n}_{i,j}^{(b)} = \sum_{X_m=(i,j)} RH_m^{(b)}$

This generates a sequence of values  $(\boldsymbol{\theta}_{i,j}^{(b)}, b = 1, \dots, B + b_0)$ . Here,  $b_0$  is called the *burn-in time*. If the initial values, where  $b = 0$ , are far from the center of the posterior distribution, it may take several iterations for the sequence to move toward the mode of the posterior. It can be shown that, after a long enough burn-in time  $b_0$ , the set  $\{\boldsymbol{\theta}_{i,j}^{(b)}, b = 1, \dots, B + b_0\}$  will be a sample from the target distribution, that is, the posterior distribution of  $\boldsymbol{\theta}_{i,j}$ , conditioned on the data. (Technical conditions for this are given in e.g. Geman & Geman (1984).) Then if  $B$  is large,

$$E(\boldsymbol{\theta}_{i,j}|\text{Data}) \approx \frac{1}{B} \sum_{b_0}^{B+b_0} \boldsymbol{\theta}_{i,j}^{(b)}$$

The new probabilities for RH are calculated for each record using Equation 4 above. For more details on the Gibbs sampling technique, see e.g. Casella & George (1992).

Using these probabilities, we re-compute the linear and probabilistic estimators in the same manner as described before; the results are given in Column (2) of Table A.12. These estimates are similar to those calculated with vanilla and re-calibrated BIFSG. Note that in practice, the algorithm described can be computationally expensive; we thus perform the entire procedure on a 1% sample of the population.

## B.7 Estimating Audit Disparity Conditional on True Underreporting

In Section 7.1.1 we describe at a high level how we can combine operational audit data, NRP data, and baseline taxpayer data to estimate the audit rate of taxpayers conditional on a given (binned) amount of underreporting. Here, we provide additional detail. Note that the audit rate of taxpayers of group  $g$  and underreporting  $k$  is simply  $\Pr[Y = 1|B = g, K = k]$ .

By Bayes' rule:

$$\Pr[Y = 1|B = g, K = k] = \frac{\Pr[K = k|Y = 1, B = g] \Pr[Y = 1|B = g]}{\Pr[K = k|B = g]}.$$

Each of the quantities on the right-hand side of the equation includes self-reported race as a conditioning variable. Since we do not have access to self-reported race, we must use our predicted race probability, like in our estimates, to measure these quantities.

We consider each quantity in turn. Consider first  $\Pr[K = k|B = 1]$ , i.e. the probability that a taxpayer has non-compliance  $K$  given that they are Black. In practice, we bin taxpayers' underreporting amounts rather than viewing them as exact figures (as exact repeated amounts of underreporting are rare). Viewing  $K = k$  as membership in the set of taxpayers whose underreporting is in bin  $k$ , we can thus apply either the probabilistic or linear estimator to obtain this quantity:

$$\begin{aligned}\widehat{\Pr}^p[K = k|B = 1] &:= \frac{\sum_{i \in \text{NRP}} b_i \cdot \mathbf{1}[K_i = k]}{\sum_{i \in \text{NRP}} b_i} \\ \widehat{\Pr}^p[K = k|B = 0] &:= \frac{\sum_{i \in \text{NRP}} (1 - b_i) \cdot \mathbf{1}[K_i = k]}{\sum_{i \in \text{NRP}} (1 - b_i)},\end{aligned}$$

where we limit our summation to NRP to ensure that our estimates are representative of the overall population, or the linear estimator, by regressing:

$$\mathbf{1}[K_i = k] = \alpha_k + \beta_k \cdot b_i + \xi_i$$

and taking:

$$\begin{aligned}\widehat{\Pr}^\ell[K = k|B = 1] &:= \hat{\alpha}_k + \hat{\beta}_k \\ \widehat{\Pr}^\ell[K = k|B = 0] &:= \hat{\alpha}_k,\end{aligned}$$

where again the regression is run over EITC claimants in NRP. (As before, we can modify our estimators to take into account sample weights accordingly.)

Next, consider  $\Pr[K = k|Y = 1, B = g]$ . This quantity is just as  $\Pr[K = k|B = g]$ , but limited to taxpayers who were audited; thus, we can again apply the probabilistic and linear estimators, but run them over the operational audit data rather than NRP.

Finally,  $\Pr[Y = 1|B = g]$  is simply the overall audit probability conditional on race, which is the main focus of the paper.

Combining the weighted estimators together, we can write:

$$\widehat{\Pr}^p[Y = 1|K = k, B = 1] := \frac{\frac{\sum_{i \in \text{NRP}} b_i \cdot \mathbf{1}[K_i = k]}{\sum_{i \in \text{NRP}} b_i} \cdot \frac{\sum_i b_i Y_i}{\sum_i b_i}}{\frac{\sum_{i \in \text{OP}} b_i \cdot \mathbf{1}[K_i = k]}{\sum_{i \in \text{OP}} b_i}}$$

and similarly for conditioning on  $B = 0$ . We can similarly combine the linear estimates.

Because we have combined the estimators together, Proposition 1 does not directly apply.

Additionally, estimates are not independent across different underreporting amounts. These factors make the behavior of this estimator more difficult to analyze. Thus, in order to obtain confidence intervals we do not attempt to characterize the standard errors analytically, but instead use the bootstrap. That is, we draw 100 re-samplings (of each dataset) without replacement, and re-compute the estimates for each subsample. We then take the mean of these estimates for each bin  $k$  to be our estimated audit rates, and add/subtract 1.96 times standard error to obtain the confidence intervals. We note that, while we do not have a formal statement about the direction of bias these combined estimators may have and whether the ground truth need lie between the combined probabilistic and linear estimators, the probabilistic estimator tends to produce smaller estimates of within-bin audit rate disparities than the linear estimator does, at least for the bulk of the distribution.

## B.8 Conditional Disparity Estimators and Decomposition

We can generalize the disparity estimators straightforwardly to answer the following question: What is the racial disparity in audits *conditional* on some covariates? Formally, let  $X$  represent a (possibly vector-valued) feature, taking on values  $x \in \mathcal{X}$ . We define the *audit disparity at  $x$*  as:

$$D_x := \mathbb{E}[Y|B = 1, X = x] - \mathbb{E}[Y|B = 0, X = x].$$

To obtain estimates for these quantities, we again apply our linear and probabilistic estimators. That is, we define the probabilistic estimator for audit disparity at  $x$  as:

$$\hat{D}_x^P := \frac{\sum_{i: X_i=x} Y_i b_i}{\sum_{i: X_i=x} b_i} - \frac{\sum_{i: X_i=x} Y_i (1 - b_i)}{\sum_{i: X_i=x} (1 - b_i)},$$

and the linear estimator for audit disparity at  $x$  as the estimated coefficient  $\hat{\beta}_x$  in the regression of  $Y$  on  $b$  over the set of observations  $i$  such that  $X_i = x$ . Formally, that is,

$$\hat{D}_x^l = \frac{\sum_{i: X_i=x} (b_i - \bar{b}_x)(Y_i - \bar{Y}_x)}{\sum_{i: X_i=x} (b_i - \bar{b}_x)^2},$$

where  $\bar{\cdot}_x$  indicates the average taken with respect to observations  $i$  having  $X_i = x$ .

It is easy to see that, assuming that  $\Pr[X_i = x] > 0$ , we will have that:

$$\hat{D}_x^P \xrightarrow{n \rightarrow \infty} \frac{\text{Cov}(Y, b|X = x)}{\mathbb{E}[b|X = x](1 - \mathbb{E}[b|X = x])} \equiv D_x^P$$

by applying the Law of Large Numbers to observations with  $X_i = x$ . Similarly, we can observe that

$$\hat{D}_x^l \xrightarrow{n \rightarrow \infty} \frac{\text{Cov}(Y, b|X = x)}{\text{Var}(b|X = x)} \equiv D_x^l,$$

and that the multiplicative relationship between  $D_x^l$  and  $D_x^p$  holds with  $x$ -specific constants:

$$D_x^p = D_x^l \cdot \frac{\text{Var}(b|X = x)}{\mathbb{E}[b|X = x](1 - \mathbb{E}[b|X = x])}.$$

Accordingly, we can obtain  $x$ -specific covariance conditions for these estimators to write the following theorem:

**Proposition 7** (Conditional Disparity Bounds). Suppose that  $b$  is a taxpayer's probability of being Black given some observable characteristics  $Z$ , so that  $b = \Pr[B = 1|Z, X]$ . Define  $D_x^p$  as the asymptotic limit of the probabilistic disparity estimator at  $X = x$ ,  $\hat{D}_x^p$ , and  $D_x^l$  as the asymptotic limit of the linear disparity estimator at  $X = x$ ,  $\hat{D}_x^l$ . Then:

1.

$$D_x^p = D_x - \frac{\mathbb{E}[\text{Cov}(Y, B|b, X)|X = x]}{\text{Var}(b|X = x)} \quad (7.1)$$

2.

$$D_x^l = D_x + \frac{\mathbb{E}[\text{Cov}(Y, b|B, X)|X = x]}{\text{Var}(B|X = x)} \quad (7.2)$$

3. Suppose  $\mathbb{E}[\text{Cov}(Y, B|b, X)|X = x] \geq 0$  and  $\mathbb{E}[\text{Cov}(Y, b|B, X)|X = x] \geq 0$ . Then

$$D_x^p \leq D \leq D_x^l \quad (7.3)$$

4. Suppose  $\mathbb{E}[\text{Cov}(Y, B|b, X)|X = x] \leq 0$  and  $\mathbb{E}[\text{Cov}(Y, b|B, X)|X = x] \leq 0$ . Then

$$D_x^l \leq D \leq D_x^p \quad (7.4)$$

Importantly, Proposition 7 requires that  $b$  be  $\Pr[B = 1|Z, X]$ . In other words, if we wish to apply it, we must not only have an accurate measurement of the probability a taxpayer is Black given their name and geography, but also including additionally the feature of interest. Moreover, there is again the possibility of bias and noise, and thus the “at- $x$ ” analogue of Proposition 2; a more general formulation and associated proof follow similarly. Thus, in applying this result, recalibration within the values taken on by  $X$  may be important.

Now, Proposition 7 provides an analogue of our bounds for each  $x$ . But we might wish to summarize  $D_x$ , which is a function, into a single number, e.g. by estimating  $\mathbb{E}[D_x]$ . Again, it is easy to see that

$$\frac{1}{n} \sum_x n_x \hat{D}_x^p \xrightarrow{n \rightarrow \infty} \mathbb{E}[D_x^p] \quad \frac{1}{n} \sum_x n_x \hat{D}_x^l \xrightarrow{n \rightarrow \infty} \mathbb{E}[D_x^l],$$

but can we similarly say something about  $\mathbb{E}[D_x^l]$  vs  $\mathbb{E}[D_x]$  vs  $\mathbb{E}[D_x^p]$ ? The following theorem provides both necessary and sufficient conditions; it boils down to asking that the covariance conditions hold *on average*.

**Theorem 1.** Let  $D_x^p$ ,  $D_x^l$ ,  $D_x$  be as above. Let  $\mathcal{P}$  be an arbitrary distribution over  $X$ . Then:

$$\mathbb{E}_{\mathcal{P}}[D_x^p] \leq \mathbb{E}_{\mathcal{P}}[D_x] \leq \mathbb{E}_{\mathcal{P}}[D_x^l]$$

whenever:

1. (Sufficient, but not necessary:)

$$\mathbb{E}[\text{Cov}(Y, b|B, X)] \geq 0, \mathbb{E}[\text{Cov}(Y, B|b, X)] \geq 0 \quad \forall x$$

2. (Necessary and sufficient):

$$\mathbb{E}_{\mathcal{P}} \left[ \frac{\mathbb{E}[\text{Cov}(Y, b|B, X)|X = x]}{\text{Var}(B|X = x)} \right] \geq 0, \mathbb{E}_{\mathcal{P}} \left[ \frac{\mathbb{E}[\text{Cov}(Y, B|b, X)|X = x]}{\text{Var}(b|X = x)} \right] \geq 0$$

The sufficient, but not necessary, condition is easy to reason about because it only requires an assumption about the sign of covariance terms, but of course less likely to hold given its strength. By contrast, the necessary condition is more reasonable in that it allows for the covariance conditions to be negative over some values of  $X$  as long as it is positive over enough of the others; however, ascertaining this requires not merely knowledge of signs, but quantitative variation.

*Proof.* As noted, that

$$D_x = D_x^l - \frac{\mathbb{E}[\text{Cov}(Y, b|B, X)|X = x]}{\text{Var}[b|X = x]}$$

and

$$D_x = D_x^p + \frac{\mathbb{E}[\text{Cov}(Y, B|b, X)|X = x]}{\text{Var}[B|X = x]}$$

follows from the same logic, conditioned on the event  $X = x$ , as used in the proof of Proposition 2. Then:

$$\begin{aligned} \mathbb{E}_{\mathcal{P}}[D_x^l] &= \sum_x \mathcal{P}(x) D_x^p = \sum_x \left[ \mathcal{P}(x) \left[ D_x + \frac{\mathbb{E}[\text{Cov}(Y, b|B, X)|X = x]}{\text{Var}[b|X = x]} \right] \right] \\ &= \mathbb{E}_{\mathcal{P}}[D_x] + \sum_x \mathcal{P}(x) \frac{\mathbb{E}[\text{Cov}(Y, b|B, X)|X = x]}{\text{Var}[b|X = x]} \\ &= \mathbb{E}_{\mathcal{P}}[D_x] + \mathbb{E}_{\mathcal{P}} \left[ \frac{\mathbb{E}[\text{Cov}(Y, b|B, X)|X = x]}{\text{Var}[b|X = x]} \right]. \end{aligned}$$

Thus:

$$\mathbb{E}[\text{Cov}(Y, b|B, X)|X = x] \geq 0 \quad \forall x \implies \mathbb{E}_{\mathcal{P}}[D_x^l] \geq \mathbb{E}_{\mathcal{P}}[D_x],$$



But more weakly, it also holds that:

$$\mathbb{E}_{\mathcal{P}} \left[ \frac{\mathbb{E}[\text{Cov}(Y, b|B, X)|X = x]}{\text{Var}[b|X = x]} \right] \geq 0 \implies \mathbb{E}_{\mathcal{P}}[D'_x] \geq \mathbb{E}_{\mathcal{P}}[D_x].$$

Similarly, we can write that:

$$\mathbb{E}_{\mathcal{P}}[D_x^p] = \mathbb{E}_{\mathcal{P}}[D_x] - \mathbb{E}_{\mathcal{P}} \left[ \frac{\mathbb{E}[\text{Cov}(Y, B|b, X)|X = x]}{\text{Var}[B|X = x]} \right]$$

and so again,

$$\mathbb{E}[\text{Cov}(Y, B|b, X)|X = x] \geq 0 \ \forall X \implies \mathbb{E}_{\mathcal{P}}[D_x^p] \leq \mathbb{E}_{\mathcal{P}}[D_x],$$

and the weaker version

$$\mathbb{E}_{\mathcal{P}} \left[ \frac{\mathbb{E}[\text{Cov}(Y, B|b, X)|X = x]}{\text{Var}[B|X = x]} \right] \geq 0 \implies \mathbb{E}_{\mathcal{P}}[D_x^p] \leq \mathbb{E}_{\mathcal{P}}[D_x]$$

also holds. □

### B.8.1 Using $D_x$ to decompose $D$ into constituent parts

If the distribution over  $X$  is the same for both groups, then  $D = \mathbb{E}[D_x]$ , where the expectation is taken over the population distribution of  $X$ ,  $\mathcal{P}_X$ . Otherwise, there is a *compositional effect*; that is, unless audit rates are constant for all values in  $\mathcal{X}$ , groups will have different total audit rates because individuals tend to have values of  $X$  associated with higher audit rates in one group more than another, even if audit rates are the same for both groups conditional on every  $x$ . Arguably, one might take the point of view that the conditional disparity (either averaged over  $\mathcal{P}_X$  or some other reference distribution, e.g.  $\mathcal{P}_{X|B=1}$ ) is of primary importance, as it captures the portion of disparity stemming from differences in the outcome of interest given the same individual features. We do not take the view that the overall disparity is unimportant for many reasons; for example, even compositional effects can be related to systemic inequality and may thus have relevant policy implications. But it still useful to understand how much of disparity is driven by differing audit rates with the same features rather than compositional effects.

In Proposition 8, we thus decompose disparity into a portion coming from  $\mathbb{E}[D_x]$  and a portion coming from the compositional effects. To understand this decomposition, we define  $\mathcal{P}_g(X)$  to be the conditional distribution  $\Pr[X = x|B = g]$  and, as mentioned,  $\mathcal{P}_X$  to be the unconditional distribution  $\Pr[X = x]$ . We will use  $\mathcal{P}$  for an arbitrary distribution over  $X$ . We also define an operator  $\Delta$  representing the difference of two probability distributions weighted by a function of a random variable. In the discrete case, that is:

$$\Delta_{\mathcal{P}, \mathcal{P}'}(f(x)) := \sum_x f(x) [\mathcal{P}(x) - \mathcal{P}'(x)].$$

With this notation, we can state and prove the following proposition capturing the decomposition with respect to an arbitrary “reference” distribution over  $X$ :

**Proposition 8.** Suppose  $\mathcal{P}$  is an arbitrary distribution over  $X$ . Then overall disparity  $\mathcal{D}$  can be decomposed with respect to  $\mathcal{P}$  as:

$$D = \mathbb{E}_{\mathcal{P}}[D_x] + \Delta_{\mathcal{P}_1, \mathcal{P}}(\mathbb{E}[Y|B = 1, X = x]) - \Delta_{\mathcal{P}_0, \mathcal{P}}(\mathbb{E}[Y|B = 0, X = x]).$$

*Proof.* The law of iterated expectations says that:

$$\mathbb{E}[Y|B = g] = \mathbb{E}_{\mathcal{P}_g} \mathbb{E}[Y|B = g, X = x],$$

so

$$\begin{aligned} D &:= \mathbb{E}[Y|B = 1] - \mathbb{E}[Y|B = 0] \\ &= \mathbb{E}_{\mathcal{P}_1} [\mathbb{E}[Y|B = 1, X = x]] - \mathbb{E}_{\mathcal{P}_0} [\mathbb{E}[Y|B = 0, X = x]]. \end{aligned}$$

Now, adding and subtracting both  $\mathbb{E}_{\mathcal{P}} [\mathbb{E}[Y|B = 1, X = x]]$  and  $\mathbb{E}_{\mathcal{P}} (\mathbb{E}[Y|B = 0, X = x])$  we get:

$$\begin{aligned} D &= \mathbb{E}_{\mathcal{P}} [\mathbb{E}[Y|B = 1, X = x]] - \mathbb{E}_{\mathcal{P}} [\mathbb{E}[Y|B = 0, X = x]] \\ &\quad + \mathbb{E}_{\mathcal{P}_1} [\mathbb{E}[Y|B = 1, X = x]] - \mathbb{E}_{\mathcal{P}} [\mathbb{E}[Y|B = 1, X = x]] \\ &\quad + \mathbb{E}_{\mathcal{P}} [\mathbb{E}[Y|B = 0, X = x]] - \mathbb{E}_{\mathcal{P}_0} [\mathbb{E}[Y|B = 0, X = x]] \end{aligned}$$

The terms in the first line is simply  $\mathbb{E}_{\mathcal{P}}[D_x]$ . For the second term, note that

$$\begin{aligned} &\mathbb{E}_{\mathcal{P}_1} [\mathbb{E}[Y|B = 1, X = x]] - \mathbb{E}_{\mathcal{P}} [\mathbb{E}[Y|B = 1, X = x]] \\ &= \sum_x \mathcal{P}_1(x) \mathbb{E}[Y|B = 1, X = x] - \sum_x \mathcal{P}(x) \mathbb{E}[Y|B = 1, X = x] \\ &= \sum_x \mathbb{E}[Y|B = 1, X = x] [\mathcal{P}_1(x) - \mathcal{P}(x)] \\ &= \Delta_{\mathcal{P}_1, \mathcal{P}}(\mathbb{E}[Y|B = 1, X = x]) \end{aligned}$$

and similarly the third simplifies to  $\Delta_{\mathcal{P}, \mathcal{P}_0}(\mathbb{E}[Y|B = 0, X = x])$ .  $\square$

The most natural choices to use as reference distributions in Proposition 8 are either  $\mathcal{P}_X$ ,  $\mathcal{P}_1$ , or  $\mathcal{P}_0$ . For the latter two, the form becomes slightly simpler:

**Corollary 1.** By noting that  $\Delta_{\mathcal{P}, \mathcal{P}} = 0$ , we can also write that:

$$D = \mathbb{E}_{\mathcal{P}_1}[D_x] + \Delta_{\mathcal{P}_1, \mathcal{P}_0}(Y|B = 0, X = x)$$

and

$$D = \mathbb{E}_{\mathcal{P}_0}[D_x] + \Delta_{\mathcal{P}_0, \mathcal{P}}(Y|B = 1, X = x)$$

Regardless of the reference distribution, Proposition 7 and Corollary 1 show that overall disparity can be decomposed into an average conditional disparity, under some distribution,

and compositional effects that capture the difference of one or both groups' distribution over  $X$  with respect to that reference distribution.

To apply this decomposition, we must know or estimate the constituent components. Non-group-specific distributions, e.g.  $\mathcal{P}(X)$  and  $\mathbb{E}[Y|X]$ , can be simply estimated from the data in the usual manner, but we must also estimate group-specific distributions over  $X$  and group-specific audit rates given  $X$ . We turn to estimation concerns in the following section.

### B.8.2 Estimating decomposition-relevant quantities

We focus on the case where  $\mathcal{X}$  is discrete and has a relatively small number of unique values. In this case, we need to estimate  $D_x$ ,  $\mathbb{E}[Y_x|B = 1, X = x]$ ,  $\mathbb{E}[Y_x|B = 0, X = x]$ ,  $\mathcal{P}_1(x)$ ,  $\mathcal{P}_0(x)$ . Once we obtain these estimates, we combine them in a plug-in manner to estimate the expression in Proposition 8:

$$\begin{aligned} \hat{D} = \sum_x \hat{\mathcal{P}}(x) \hat{D}_x + \sum_x \hat{\mathbb{E}}[Y|B = 1, X = x] \cdot [\hat{\mathcal{P}}_1(x) - \hat{\mathcal{P}}(x)] \\ + \sum_x \hat{\mathbb{E}}[Y|B = 0, X = x] \cdot [\hat{\mathcal{P}}(x) - \hat{\mathcal{P}}_0(x)] \end{aligned} \quad (5)$$

As with the overall disparity estimation and related problems, we can obtain these estimates via either the probabilistic disparity estimator or the linear disparity estimator. For the probabilistic, we estimate  $\hat{D}_x^p$  by computing the probabilistic estimator for individuals with  $X_i = x$ ; we estimate the audit rate conditional at each  $x$  using the levels ; and we estimate the distribution using the probabilistic estimator over the full dataset with an indicator for having  $X_i = x$  as the outcome of interest. That is:

$$\begin{aligned} \hat{D}_x^p &:= \frac{\sum_{i: X_i=x} b_i Y_i}{\sum_{i: X_i=x} b_i} - \frac{\sum_{i: X_i=x} b_i Y_i}{\sum_{i: X_i=x} b_i} \\ \hat{E}[Y|B = 1, X = x] &:= \frac{\sum_{i: X_i=x} b_i Y_i}{\sum_{i: X_i=x} b_i} \\ \hat{E}[Y|B = 0, X = x] &:= \frac{\sum_{i: X_i=x} (1 - b_i) Y_i}{\sum_{i: X_i=x} (1 - b_i)} \\ \hat{\mathcal{P}}_1(x) &:= \frac{\sum_i b_i \mathbf{1}[X_i = x]}{\sum_i b_i} \\ \hat{\mathcal{P}}_0(x) &:= \frac{\sum_i (1 - b_i) \mathbf{1}[X_i = x]}{\sum_i (1 - b_i)}. \end{aligned}$$

For the linear, we estimate, two regressions for each value of  $x \in \mathcal{X}$ :

$$Y_i = \alpha_x + \beta_x b_i$$

over all  $i : X_i = x$ , and

$$\mathbf{1}[X_i = x] = \zeta_x + \xi_x b_i$$

over all  $i$ . We then construct have:

$$\begin{aligned}\widehat{D}_x^p &=: \widehat{\beta}_x \\ \widehat{E}[Y|B=1, X=x] &:= \widehat{\alpha}_x + \widehat{\beta}_x \\ \widehat{E}[Y|B=0, X=x] &:= \widehat{\alpha}_x \\ \widehat{\mathcal{P}}_1(x) &:= \widehat{\zeta}_x + \widehat{\xi}_x \\ \widehat{\mathcal{P}}_0(x) &:= \widehat{\zeta}_x\end{aligned}$$

## B.9 Tightening Bounds with a Linear Program

In this subsection we use a linear program to calculate the maximum and minimum values of disparity that are consistent with the Proposition 1.3 assumptions and the observed joint distribution of  $b$  and  $Y$ . To implement this approach, we discretize the  $b$  distribution into  $n$  mutually exclusive and equal width bins, labeling them with the sample average  $\bar{b}_i$  of taxpayers in the bin. The share of taxpayers in each bin is given by  $p_1, \dots, p_n$ . As in Proposition 1, we assume race probabilities for each bin are perfectly calibrated, so that  $P(B=1|b=b_i) = b_i$  for all  $i$  from 1 to  $n$ . Along with the marginal distribution of  $b$ , we observe the distribution of audits conditional on  $b$ . Denote the audit rate in each bin as  $Y_1, \dots, Y_n$ , and the overall audit rate as  $Y$ . The audit rate among Black taxpayers in bin  $i$  is denoted by  $Y_i^B$  and among non-Black taxpayers as  $Y_i^{NB}$ . Finally, as in Proposition 1.3, we assume the sign of the covariance conditions, i.e. that  $\mathbb{E}[\text{Cov}(Y, b|B)] \geq 0$  and  $\mathbb{E}[\text{Cov}(Y, B|b)] \geq 0$ .

We ask: if we can allocate the  $Y_i$  into  $Y_i^B$  and  $Y_i^{NB}$  for each bin  $i$ , what is the maximum and minimum disparity we could obtain consistent with the results? This is a constrained optimization problem. Disparity is given by:

$$Y^B - Y^{NB} = \sum_i \Pr[b_i|B=1] \cdot Y_i^B - \sum_i \Pr[b_i|B=0] \cdot Y_i^{NB}$$

Noting that

$$b_i Y_i^B + (1 - b_i) Y_i^{NB} = Y_i \implies Y_i^{NB} = \frac{Y_i - b_i Y_i^B}{1 - b_i}$$

we can replace the  $Y_i^{NB}$  into the objective to get:

$$\sum_i \left[ \Pr[b_i|B=1] Y_i^B - \Pr[b_i|B=0] \left( \frac{Y_i - Y_i^B b_i}{1 - b_i} \right) \right]$$

and we may as well drop  $Y_i$  from the objective (since it will not change based on our choice variables) to get:

$$f(\vec{Y}^B) = \sum_i \left[ \Pr[b_i|B=1] Y_i^B + \frac{\Pr[b_i|B=0] b_i}{1 - b_i} Y_i^B \right]$$

Now by Bayes' rule, we can write that:

$$\Pr[b_i|B=1] = \frac{\Pr[B=1|b_i]p_i}{\Pr[B=1]} = \frac{b_i p_i}{\Pr[B=1]} \quad \Pr[b_i|B=0] = \frac{\Pr[B=0|b_i]p_i}{\Pr[B=0]} = \frac{(1-b_i)p_i}{\Pr[B=0]}.$$

If we substitute these into the objective function, we have:

$$\sum_i \left[ \frac{b_i p_i}{\Pr[B=1]} Y_i^B + \frac{(1-b_i)p_i}{\Pr[B=0]} \frac{b_i Y_i^B}{1-b_i} \right] = \sum_i Y_i^B b_i p_i \left[ \frac{1}{\Pr[B=1]} + \frac{1}{\Pr[B=0]} \right]$$

The constant multiplier will also not be affected by our decisions, so we can drop it, and our objectives will be:

$$\max_{\langle Y_i^B \rangle} \sum_i p_i Y_i^B b_i \quad \text{and} \quad \min_{\langle Y_i^B \rangle} \sum_i p_i Y_i^B b_i$$

These are intuitive - to maximize (minimize) disparity, we allocate as much of the audit rate as possible to bins of  $b_i$  that are more likely to contain Black (non-Black) taxpayers, either because  $b_i$  or  $\Pr[b_i]$  is higher (lower) than others.

As for constraints: First, we need  $Y_i^B$  to be in the right range.  $[0, 1]$  is a first pass, but note that this does not guarantee consistency with the data. For instance, if  $Y_i^B$  is 0, even  $Y_i^{NB} = 1$  might not be enough to satisfy that the overall audit rate in the bin matches  $Y_i$ . So in particular, the lowest that  $Y_i^B$  could be is  $\max \left\{ 0, \frac{Y_i - (1-b_i)}{b_i} \right\}$ , while the most that it could be is  $\min \left\{ 1, \frac{Y_i}{b_i} \right\}$ . These numbers can be calculated from the data, so we will simply refer to them as upper bound  $u_i$  and lower bound  $l_i$ .

Second, the first covariance constraint,  $\mathbb{E}[\text{Cov}(Y, B|b)] \geq 0$ .

We note that:

$$\text{Cov}(Y, B|b) = \mathbb{E}[YB|b] - \mathbb{E}[Y|b]\mathbb{E}[B|b].$$

Now,  $\mathbb{E}[Y|b = b_i]$  is simply  $Y_i$ , and  $\mathbb{E}[B|b = b_i]$  is  $b_i$  by assumption of perfect calibration.  $\mathbb{E}[YB|b = b_i] = b_i Y_i^B$ , since  $YB = Y_i^B$  whenever  $B = 1$  and 0 otherwise, and again  $\Pr[B|b = b_i] = b_i$  by perfect calibration.

Then we can write that:

$$\text{Cov}(Y, B|b) = b_i Y_i^B - Y_i b_i \implies \mathbb{E}[\text{Cov}(Y, B|b)] = \sum_i p_i b_i (Y_i^B - Y_i)$$

Finally, the third condition: we need that  $\mathbb{E}[\text{Cov}(Y, b|B)] \geq 0$ . Let's consider just  $B = 1$  first.

$$\begin{aligned} \mathbb{E}[Yb|B=1] - \mathbb{E}[Y|B=1]\mathbb{E}[b|B=1] = \\ \sum_i Y_i^B b_i \Pr[b_i|B=1] - \left( \sum_i \Pr[b_i|B=1] Y_i^B \right) \left( \sum_i \Pr[b_i|B=1] b_i \right) \end{aligned}$$

Define  $\bar{b}_B = \sum_i \Pr[b_i|B=1]b_i$ , which has the interpretation of the average probability of being Black given ground truth race being Black.

Then we can write that the above is:

$$\begin{aligned} & \sum_i Y_i^B b_i \Pr[b_i|B=1] - \left( \sum_i Y_i^B \Pr[b_i|B=1] \right) \left( \sum_i b_i \Pr[b_i|B=1] \right) \\ &= \sum_i Y_i^B b_i \Pr[b_i|B=1] - \sum_i Y_i^B \Pr[b_i|B=1] \bar{b}_B \\ &= \sum_i Y_i^B (b_i - \bar{b}_B) \Pr[b_i|B=1] \end{aligned}$$

Using our Bayes' rule calculation as above, we can write:

$$\sum_i Y_i^B (b_i - \bar{b}_B) \Pr[b_i|B=1] = \frac{1}{\Pr[B]} \sum_i Y_i^B (b_i - \bar{b}_B) b_i p_i$$

For given  $B=0$ , we have that:

$$\begin{aligned} & \sum_i Y_i^{NB} b_i \Pr[b_i|B=0] - \left( \sum_i Y_i^{NB} \Pr[b_i|B=0] \right) \left( \sum_i b_i \Pr[b_i|B=0] \right) \\ &= \sum_i Y_i^{NB} b_i \Pr[b_i|B=0] - \sum_i Y_i^{NB} \Pr[b_i|B=0] \bar{b}_{NB} \\ &= \sum_i Y_i^{NB} (b_i - \bar{b}_{NB}) \Pr[b_i|B=0] \end{aligned}$$

and again applying Bayes' rule we can write the above as:

$$\frac{1}{1 - \Pr[B]} \sum_i Y_i^{NB} (b_i - \bar{b}_{NB}) (1 - b_i) p_i.$$

Then the overall constraint is:

$$\begin{aligned} & \Pr[B] \cdot \frac{1}{\Pr[B]} \sum_i Y_i^B \cdot b_i \cdot (b_i - \bar{b}_B) \cdot p_i \\ &+ (1 - \Pr[B]) \frac{1}{1 - \Pr[B]} \sum_i Y_i^{NB} (1 - b_i) (b_i - \bar{b}_{NB}) p_i \\ &= \sum_i p_i (Y_i^B b_i (b_i - \bar{b}_B) + Y_i^{NB} (1 - b_i) (b_i - \bar{b}_{NB})) \geq 0 \end{aligned}$$

Noting that  $Y_i^{NB} = (Y_i - Y_i^B b_i)/(1 - b_i)$  and factoring out, we can rewrite the last inequality as:

$$\sum_i p_i (Y_i^B b_i (\bar{b}_{NB} - \bar{b}_B) + Y_i (b_i - \bar{b}_{NB})) \geq 0$$

Putting these together, our problem is:

**Program 1** (Maximum Consistent Disparity).

$$\begin{aligned}
\max_{\langle Y_i^B \rangle} \quad & \sum_i p_i Y_i^B b_i \text{ s.t. } Y_i^B \leq u_i \\
& l_i \leq Y_i^B \\
& 0 \leq \sum_i p_i b_i (Y_i^B - Y_i) \\
& 0 \leq \sum_i p_i (Y_i^B b_i (\bar{b}_{NB} - \bar{b}_B) + Y_i (b_i - \bar{b}_{NB}))
\end{aligned}$$

to obtain the maximum disparity consistent with the information, and minimizing the same (or maximizing the negative) to get the minimum disparity.

### B.9.1 Linear Program Results

In Table B.2, we report the information that we observe, including the average probability Black within each bin, the audit rate within each bin, and the share of taxpayers that fall into each bin. Using these along with the assumption of calibration, we can compute the conditional probability of a taxpayer falling into each bin given that they are Black and the upper and lower bounds on  $Y_i^B$ . These together constitute the requisite input to Program 1 and its counterpart minimization. To solve these programs, we use the SciPy (Virtanen et al., 2020) library’s linprog function. The results are given in the following table:

The solutions to the linear programs closely match the disparity estimates obtained by the probabilistic and linear estimators. Compare Column 1 of Table B.3 to Column 1 of Panel A of Table A.6 and Column 2 of Table B.3 to Column (1) of Panel A of Table A.5.

Table B.2: Inputs to Linear Program for Bounding Disparity

Bin	Probability Black	Fraction in Bin	Audit Rate
1	0.0057	0.7130	0.0039
2	0.0722	0.0574	0.0039
3	0.1232	0.0353	0.0047
4	0.1737	0.0233	0.0049
5	0.2241	0.0179	0.0055
6	0.2740	0.0143	0.0060
7	0.3245	0.0119	0.0068
8	0.3743	0.0101	0.0072
9	0.4246	0.0087	0.0077
10	0.4746	0.0079	0.0082
11	0.5248	0.0073	0.0090
12	0.5748	0.0068	0.0095
13	0.6249	0.0065	0.0101
14	0.6750	0.0064	0.0111
15	0.7251	0.0065	0.0118
16	0.7753	0.0068	0.0127
17	0.8255	0.0075	0.0136
18	0.8759	0.0089	0.0149
19	0.9268	0.0123	0.0165
20	0.9820	0.0312	0.0199

*Notes:* The Table reports statistics used to bound disparity through the linear program described in this subsection. Taxpayers have been binned into 20 groups based on their estimated probabilities of being Black, with each group corresponding to the values of  $b$  in a 5 percentage point range. Probability Black denotes the average estimated probability of being Black within the bin. Fraction in Bin denotes the share of the overall taxpayer population in the specified bin.

Table B.3: Maximum and Minimum EITC Audit Disparity Estimates

	Maximum Disparity (1)	Minimum Disparity (2)
Black Audit Rate	1.710	1.241
Non-Black Audit Rate	0.362	0.427
Disparity	1.347	0.813
Audit Rate Ratio	4.7	2.9

*Notes:* The table reports the results of the Maximum Disparity and Minimum Disparity constrained optimization problems described in Section B.9. Units are percentage points (0-100).



## C North Carolina Match and Bias Correction

### C.1 North Carolina Match

In our North Carolina voter registration data (reflecting the state’s voter registration file as of the close of 2020), we observe individuals’ first names, last names, zip codes, residential street addresses, and mailing addresses at time of registration or filing. We match these data to IRS data using these common features according to the following procedure:

1. First, look for exact match on zip code, first name, last name, and full text of residential street address. Remove matched records from both datasets and append matched records to output file.
2. Among unmatched records, look for match on zip code, first four characters of first and last name, and full text of residential street address (after minor data cleaning).
3. Among unmatched records, look for match on zip code, first character of first name, first four characters of last name, and full text of residential street address.
4. Among unmatched records, look for match on zip code, first four characters of first and last name, residential street number, and city.
5. Among unmatched records, look for match on zip code, first character of first name, first four characters of last name, and full text of mailing address.

Using this procedure, we are able to match 2.5 million taxpayer and voter records, or approximately 47% of the population of North Carolina taxpayers for tax year 2014. We use the same procedure to match taxpayers to the 2023 voter registration data from other states reported in Appendix Figure A.13.

### C.2 North Carolina Reweighting

When specified, we use inverse-probability weighting to align the composition of the North Carolina matched sample with that of the full population of tax returns for 2014. The weights are generated from a linear probability model whose binary outcome equals one for records appearing in the IRS-matched North Carolina sample, and whose features are chosen to reflect observable taxpayer characteristics that we would like to align with their US means. These are entered as categorical variables and are fully interacted with one another, resulting in a flexible nonparametric model of the conditional probability of appearing in the North Carolina data. Features include quintiles (as calculated on the full population) of the BIFSG-predicted probability that a taxpayer self reports as black; four activity code groupings<sup>41</sup>; gender; the presence of dependents; joint/non-joint filing status; and whether a taxpayer was audited. The weights are then given by the inverse of these conditional probabilities. The weights were successful in aligning the weighted sample proportions along all included dimensions to within 0.02% of their U.S. population means.

---

<sup>41</sup>Activity codes are grouped as: 270-271 (EITC claimants), 272 (1040 filers without additional schedules or very high income), 273-278 (filers with Schedule C etc. but not very high income), and 279-281 (filers with very high (\$1M) or more income or high (\$ > 250K) with additional schedules).

## D EITC Disparity Decomposition

This appendix provides additional detail regarding the decomposition of the total racial audit disparity into components associated with the disparity within and between EITC claimants presented in Section 6.

As in the main text,  $Y_C^B$  and  $Y_C^{NB}$  refer to the average audit rates for Black and non-Black EITC claimants, respectively;  $Y_{NC}^B$  and  $Y_{NC}^{NB}$  refer to the average audit rates for Black and non-Black EITC non-claimants, respectively; and  $C_B$  and  $C_{NB}$  refer to the respective probabilities that Black and non-Black taxpayers claim the EITC.

Our preferred decomposition, presented in the main text, is given by:

$$Y^B - Y^{NB} = \underbrace{(Y_C^B - Y_C) C_B - (Y_C^{NB} - Y_C) C_{NB}}_{(1)} + \underbrace{(Y_{NC}^B - Y_{NC}) (1 - C_B) - (Y_{NC}^{NB} - Y_{NC}) (1 - C_{NB})}_{(2)} + \underbrace{(C_B - C_{NB}) (Y_C - Y_{NC})}_{(3)}$$

The first component reflects racial differences in the audit rate among EITC claimants: if Black and non-Black EITC claimants were selected at the same rate, it would imply  $Y_C^B = Y_C^{NB} = Y_C$ , so that both terms in (1) would be zero. Similarly, the second component reflects racial differences in the audit rate among non-EITC claimants. The third component reflects compositional differences in the rate at which Black and non-Black taxpayers claim the EITC as well as differences in the audit rate of EITC versus non-EITC returns. This component would be zero if Black and non-Black taxpayers claimed the EITC at equal rates, or if EITC and non-EITC claimants (of any race) were audited at the same rate.

The following table presents estimates of the elements of the decomposition using both the linear and probabilistic audit rate estimators. These estimates are reported in Section 6 of the main text.<sup>42</sup>

An appealing feature of the above decomposition is that it weights the contribution of differences in the EITC and non-EITC audit rate for each racial group based on the share of that racial group claiming the EITC. On the other hand, to the extent that Black and non-Black taxpayers claim the EITC at different rates, there is a sense in which the first two components of the decomposition can be interpreted to reflect compositional differences in EITC claiming as well as differences in the audit rate. We consider two alternative decompositions below, which differ from the above in that they weight the within EITC and within non-EITC components of the disparity based on the EITC claim rates of either Black or non-Black taxpayers.

Using Black taxpayers as the reference group, we can decompose the total disparity as:

$$Y^B - Y^{NB} = \underbrace{(Y_C^B - Y_C^{NB}) C_B}_{(1)} + \underbrace{(Y_{NC}^B - Y_{NC}^{NB}) (1 - C_B)}_{(2)} + \underbrace{(C_B - C_{NB}) (Y_C^{NB} - Y_{NC}^{NB})}_{(3)}$$

---

<sup>42</sup>In finite samples, the individual components of the decomposition may not exactly sum to the estimated total disparity. The percentage contributions we report are calculated by dividing each component estimate by the sum of the three component estimates.

Here, the first component represents the contribution of the disparity within EITC claimants, the second represents the contribution of the disparity within EITC non-claimants, and the third component is the same as with the first decomposition presented, but with the difference in EITC versus non-EITC audit rates evaluated for Black taxpayers.

Using this decomposition with the estimates reported in Appendix Table D.1 implies a larger contribution to the overall disparity from the disparity among EITC claimants, between 78% and 83% depending on whether the probabilistic or linear estimator is used (see Appendix Table D.2).

Alternatively, using non-Black taxpayers as the reference group, we can decompose the total disparity as:

$$Y^B - Y^{NB} = \underbrace{(Y_C^B - Y_C^{NB}) C_{NB}}_{(1)} + \underbrace{(Y_{NC}^B - Y_{NC}^{NB}) (1 - C_{NB})}_{(2)} + \underbrace{(C_B - C_{NB}) (Y_C^B - Y_{NC}^B)}_{(3)}$$

With this decomposition, the estimates reported in Appendix Table D.1 imply a smaller contribution to the overall disparity from the disparity within EITC claimants, and a larger contribution from racial differences in EITC claiming (see Appendix Table D.2). The explanation for this difference is that because non-Black taxpayers claim the EITC at lower rates, using that group as the reference leads to attaching less weight to the disparity among EITC claimants.

Table D.1: Decomposition of Disparity with Resepct to EITC Claiming

	Probabilistic Estimate	Linear Estimate
EITC Claim Rate Among...		
Black Taxpayers	32.43	41.14
Non-Black Taxpayers	17.26	16.06
Audit Rate Among...		
EITC Claimants	1.44	1.44
Black EITC Claimants	2.99	3.73
Non-Black EITC Claimants	1.04	0.85
EITC Non-Claimants	0.31	0.31
Black EITC Non-Claimants	0.40	0.47
Non-Black EITC Non-Claimants	0.30	0.29
Disaprity Contribution From...		
Within EITC Claimants	70%	72%
Within EITC Non-Claimants	9%	8%
Racial Differences in EITC Claiming	21%	20%

*Notes:* The table reports linear and probabilistic estimates of the terms used to calculate the decompositions of disparity described in this appendix section. The final three rows of the table report the contribution of terms (1), (2), and (3) of the preferred disparity decomposition described at the beginning of this Appendix section. Units are percentage points (0-100).

Table D.2: Alternative Disparity Decompositions

Reference Group	Estimator	Within-EITC Contribution	Outside-EITC Contribution	Differences in EITC Claiming Contribution
Black	Probabilistic	78%	8%	14%
Black	Linear	83%	7%	10%
Non-Black	Probabilistic	41%	10%	48%
Non-Black	Linear	32%	10%	57%

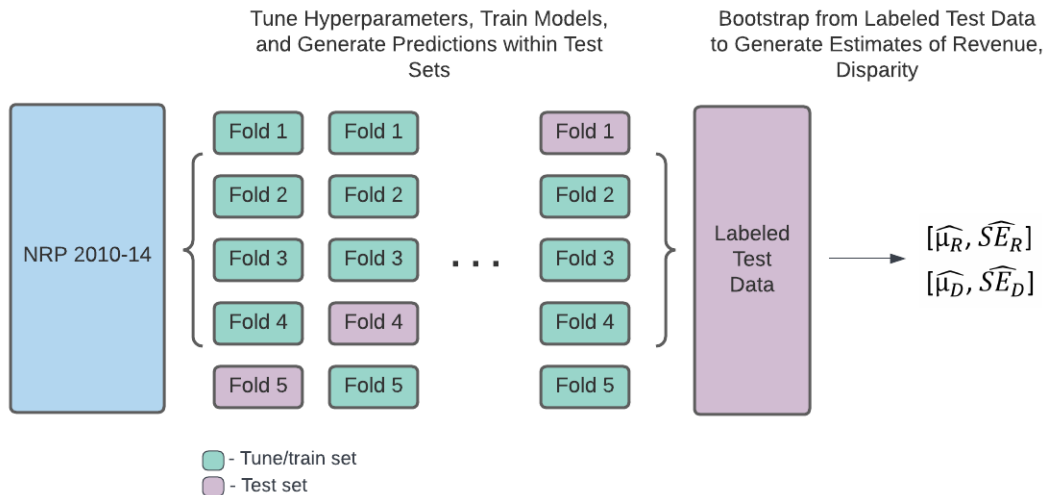
*Notes:* The table reports the decomposition of disparity under the alternative decomposition methods described in this appendix section and under alternative estimation approaches (linear versus probabilistic estimators). Estimates of the terms used to calculate each component of disparity are reported in Appendix Table D.1.

## E Taxpayer Noncompliance Prediction Model

The outcomes of interest  $Y$  of the taxpayer non-compliance random forest models are the dollar amount of adjustment in either total tax liability or refundable credit amount following an audit (for the regression models) and a  $\{0, 1\}$  indicator of whether non-compliance exceeds \$100 (for the classifier models). The inputs to the model are characteristics of the tax return, denoted  $X$ , which include wages and other sources of income, claimed deductions, and flags for whether dependents claimed on the return may violate IRS dependent rules. These features do not include race, gender, age, location, or other demographic variables.

To train and evaluate the models, we first subset the NRP data from tax years 2010-14 to taxpayers claiming the EITC. We then randomly divide the data into 5 folds. We designate 4 of these folds as the training set, and the remaining fold as the test set. We tune the hyperparameters of each model, including the total number of decision trees in each forest, the maximum depth of each decision tree, and the maximum number of features available to each decision tree, using 5-fold cross validation within the training set and random grid search over the space of hyperparameters we consider. We then fit each tuned model on the full set of training data to generate a function  $\hat{Y} = m(X)$  which maps features into predictions, and we apply this model to the test set. We repeat this process 5 times, until each observation in the NRP data has a predicted label. The train-test splits are the same for both the regression and classification models. Figure E.1 provides an exposition of the data flow and model training process.

Figure E.1: Data Flow for Random Forest Models



To obtain estimates of disparity and annualized adjustments at each audit rate, we first bootstrap from the population of labeled test data, and then sort observations within the bootstrapped sample by either the magnitude of their predicted noncompliance (for the regression model) or the predicted likelihood of noncompliance above a \$100 threshold (for

the classification model). Within each sample, annualized adjustments are given by:

$$R_s = \frac{1}{5} \sum_{t=2010}^{2014} \frac{W_t}{\sum_{i=1}^{n_{st}} w_{ist}} \sum_{i=1}^{n_{st}} (a_{ist} w_{ist} r_{ist}) \quad (6)$$

The rightmost sum computes the total weighted audit adjustments across observations in sample  $s$  from tax year  $t$ , where  $a_{ist}$  indicates whether individual  $i$  in sample  $s$  and tax year  $t$  was audited and  $(w_{ist} r_{ist})$  is the weighted adjustment from the audit in 2014 dollars. The term to the left of this sum takes the total sample weights from NRP observations in year  $t$  (denoted  $W_t$ ) over the total sample weight from this year included in sample  $s$ , to account for the fact that each fold only contains a portion of the total population available in each study year. We then sum across each of the 5 study years and divide by 5 to approximate one year of annual adjustments in 2014 dollars. Disparity measures are computed within each sample using both the linear and probabilistic estimators described in Section 3, adjusted to account for NRP sample weights. We take the mean and standard error of these measures across the bootstrapped samples to construct our trajectories and 95% confidence intervals.

Oracle adjustments and disparity calculations are analogous to the random forest calculations, with the exception that the data are sorted using true underlying noncompliance, rather than predicted amount or likelihood of noncompliance.

## F Additional Factors Potentially Contributing to Audit Disparity

Section 7 of the main text provided evidence that the observed audit disparity among EITC claimants stems in part from audit selection algorithms designed to pursue refundable credit overclaims rather than total underreporting, in conjunction with differences by race in the prevalence of different types of noncompliance. This appendix section explores additional factors that may also contribute to the observed racial disparity.

### F.1 Prediction Model Errors

Whereas the refundable credit oracle induces a smaller disparity than what we observe in the operational data, Figure 8 showed that focusing on *predicted* refundable overclaims – where the predictions are obtained from a random forest prediction model trained on the same data that is available to IRS – would lead to a disparity similar to that observed in the operational data. In this subsection, we attempt to better understand why focusing on predicted overclaims can exacerbate the disparity that would be implied by actual differences in refundable credit overclaiming by race.

As a starting point, note that we can express the disparity from a prediction-based algorithm relative to the disparity that would be induced by an oracle in terms of differential false positive and false negative rates by race, where a true positive is defined as a taxpayer’s actual refundable credit overclaiming being in the top share of the refundable credit overclaiming distribution (see Appendix F.5 for details). Appendix Table F.1 reports the frequency of these prediction errors by race for our simulated refundable credit prediction algorithm. Relative to non-Black taxpayers, Black taxpayers experience higher false positive rates and lower false negative rates. Both of these errors push in the direction of higher audit rates for Black relative to non-Black taxpayers.

Of course, the distribution of prediction errors in our simulated algorithm might not translate to the actual audit selection algorithm employed by IRS. Motivated by the important role of child eligibility errors documented in the prior subsection, we next consider errors in the actual measure used by IRS to predict which refundable credit claims are based on ineligible children. This variable is used as part of the DDb audit selection process and is constructed to reflect the likelihood that a child claimed on a return violates the required residency or relationship test with respect to the taxpayer. Appendix Table F.3 links IRS-predicted ineligibility with ground-truth data on child eligibility, as determined through NRP audits. As with the simulated refundable overclaim algorithm explored above, this table also shows that the distribution of both false positive and false negative prediction errors push in the direction of higher audit rates for Black taxpayers.

Which features in the IRS algorithm are responsible for the prediction errors we observe? To maintain the confidentiality of the audit selection process, IRS policy prevents us from publicly disclosing which taxpayer characteristics drive audit selection or how its child eligibility prediction variable is constructed. However, the IRS publicly discloses that it draws on administrative social security records, including birth certificates, to help determine whether a taxpayer claiming a child satisfies the required relationship test



(G.A.O., 2015).<sup>43</sup> Appendix Table F.4 investigates missingness of parental information on the birth certificate data that SSA provides to IRS for EITC claimants; whereas mothers are missing at roughly equal rates, children claimed on the returns of Black taxpayers are substantially more likely to be missing information about the identity of their father on their birth certificate (53.6% vs 36.9%).<sup>44</sup> Because of this missingness, this data source is likely to be less reliable at identifying Black fathers who satisfy the relationship test as compared to non-Black fathers. Consistent with this hypothesis, Figure 6 showed that the racial audit disparity is especially large in percentage point terms among unmarried men claiming children for the EITC.

Finally, it may be that certain model features yield an outsized effect on disparity relative to their importance for accuracy; adding or subtracting features could differentially affect the ability of the model to detect overclaiming for Black and non-Black taxpayers, leading to different shares of each group being selected for audit. In addition, although dropping features generally leads to (weakly) lower model accuracy, some features may be important for minimizing overall MSE but less important for accurately ranking the very top of the overclaiming distribution. To explore this possibility, we identified the top 20 features of the refundable credit overclaiming model in terms of importance and plotted their unconditional correlations with race; six of these features exhibited markedly uneven distributions (Appendix Figure F.1). Re-training the refundable overclaiming prediction model without these six features yielded similar performance at the status quo audit rate but substantially lower disparity compared to the baseline refundable credit prediction algorithm (see Appendix Figure F.2).<sup>45</sup>

The finding that prediction errors are disproportionately concentrated among Black EITC claimants suggests there may be opportunities for policymakers to reduce disparities by modifying the predictive algorithm to improve accuracy, even conditional on the objective of prioritizing the detection of refundable credit overclaims. A full exploration of this possibility would require divulging more details of the EITC audit selection algorithm than is possible under IRS policy, however, so we defer additional exploration of this issue to internal IRS analyses.

## F.2 High-Risk Signals

Section 7 explored whether differences by race in the dollar value of predicted noncompliance contribute to the observed racial audit disparity. A distinct factor that might contribute to the disparity are differences in informational signals that strongly indicate fraud or other errors. With respect to EITC eligibility, for example, this might take the form of multiple taxpayers claiming the same child in the same year (which is not allowed), or a claim by a

---

<sup>43</sup> As discussed above, the IRS also draws on administrative child custody data to predict credit eligibility; prior research has documented substantial inaccuracies in that data, but we lack direct evidence on the distribution of those inaccuracies by race (see National Taxpayer Advocate, 2018).

<sup>44</sup> Fathers are more likely to be listed on a child's birth certificate when the mother is married at the time of birth, and as discussed above, marriage rates tend to be lower among Black parents.

<sup>45</sup> We reiterate that because our refundable credit prediction algorithm represents an approximation to the algorithm underlying actual EITC audit selection, we do not conclude that the IRS could achieve this outcome by excluding the same six features we identified. Rather, this exercise illustrates how the IRS's current approach may not be at the accuracy-disparity frontier.

childless taxpayer whose age falls outside of the statutorily prescribed range for eligibility. In such cases, the IRS might choose to prioritize these likely-noncompliant returns for audit, even above other returns with larger expected (but less certain) adjustments.

For several reasons, it appears unlikely that differences by race in the distribution of these high-risk signals significantly contributes to the observed audit disparity. First, many of the apparent errors that fall into this category (such as claimed children with ages outside of the allowable range) are addressed by the IRS outside of the audit workstream, such as through automatic “math error” adjustments to the taxpayer’s return, or “Automatic Underreporter” corrections that address discrepancies between what the taxpayer reports and information returns that the IRS receives.<sup>46</sup> As such, errors detected through these programs would not be counted as audits in our data. Second, with respect to multiple taxpayers claiming the same child, our disparity estimates are largely unchanged when we exclude from our analysis operational audits that are focused on this issue (Appendix Table F.5).

To further investigate the role of highly predictive (and therefore hard-to-ignore) signals in driving the observed disparity, we consider simple predictive models to identify features, which, on their own, are highly predictive of non-compliance (“smoking guns”). For each of the most important features in the underreporting prediction model, we train a minimalistic model to predict the presence of non-compliance above a \$100 threshold using that feature alone.<sup>47</sup> Appendix Figure F.3 summarizes the performance and disparity of each feature, via these models, on the x-axis and the implied disparity on the y-axis. A “smoking gun” feature that drives disparity would appear as an outlier in terms of both the x and y axes — that is, a model that has much higher performance than others but also much higher disparity. We do not observe a feature that meets this criterion; the apparent outliers in terms of high disparity are near the center of mass of the performance distribution. We interpret this analysis as further evidence against the hypothesis that the disparity is primarily driven by differences in the distribution of high-information signals of noncompliance by race.

### F.3 Regression vs Classification Prediction Task

A related possibility involves the IRS selecting audits to maximize the share of returns with a positive adjustment, rather than total dollars of detected overclaims or underreporting (Black et al., 2022). To investigate how this change to the audit objective shapes disparity, we trained random forest classifier models to predict whether an audited return would yield a positive adjustment of \$100 or higher (see Appendix E for details).<sup>48</sup> Appendix Figure F.5 shows this aspect of the prediction model objective yields an ambiguous effect on disparity, depending on whether the underlying objective is to predict total underreporting

---

<sup>46</sup>Children of any age may be claimed for the EITC if they are “permanently and totally disabled”, but taxpayers claiming these children must indicate on their return that this exception applies to avoid a math error correction.

<sup>47</sup>For continuous features (38 of the 40 we consider) we use a two-layer decision tree, and for discrete features we use a single layer decision tree. A one (two) layer decision tree can segment observations into two (four) categories; these models are thus extremely simple without imposing monotonicity, and are flexible in that they consider all potential splits for each layer.

<sup>48</sup>The classifier threshold of \$100 does not necessarily correspond to the dollar amount used to train any classifier model that the IRS employs. We explore alternative thresholds in Appendix Figure F.4 and find no systematic relationship between disparity and threshold amount.

or refundable credit overclaims. That is, the underreporting classifier yields higher disparity than our baseline underreporting prediction models (which ranks return according to predicted dollars of underreporting), whereas the overclaims classifier yields lower disparity than our baseline overclaims prediction models. Thus, although the question of whether to select audits based on the expected dollars of noncompliance versus the probability of any non-compliance can shape disparity in some settings, it does not appear to explain the audit disparity we observe given IRS’s focus on refundable credit overclaims.

## F.4 Other Group-Level Differences in Taxpayer Characteristics

Without disparate treatment, the observed audit disparity must arise via group-level differences in observable characteristics between Black and non-Black taxpayers. Hence, with sufficient data and model flexibility (i.e., allowing non-linearities and interactions), it must be the case that controlling for all of the inputs available to IRS would eliminate any estimated independent effect of race in a model of audit selection. Rather than focus on which of the observable characteristics driving audit selection induce the disparity, our main approach has been to investigate how the disparity arises from IRS policy goals and differences in the distribution of types of noncompliance that the individual features are predicting.

Still, it may be helpful to understand the role played by several high-salience taxpayer characteristics in driving the disparity, such as the number of dependents a taxpayer claims or whether a taxpayer uses a tax preparer to file her return. Towards that end, in Appendix B.8 we propose a method for estimating the audit disparity after adjusting for certain observable characteristics between groups, and provide conditions under which we can interpret these adjusted disparity estimates as bounds (analogous to Proposition 1). Implementing this approach, we find that the observed disparity remains after adjusting for potential differences by race in the distribution of taxpayer characteristics related to income, household composition (i.e., marital status and number of dependents), and tax preparation method (Appendix Table F.6). These findings appear to undermine some commonly conjectured candidate sources of group-level differences that could drive the observed disparity.

## F.5 Disparity Decomposition with Respect to Underreporting and Prediction Errors

This appendix section derives a simple decomposition (that was referred to in Appendix Section F.1) for an observed disparity in terms of the role of group-level differences in underreporting and errors in which taxpayers are selected. The decomposition also provides insight into the difference in disparity induced by an oracle versus the disparity induced by some other model, such as the refundable credit overclaim prediction model.

Let  $A^O$  indicate whether someone was selected by the oracle. Let  $A^R$  indicate whether someone selected by the other algorithm.

The rate that a taxpayer from group  $j$  is selected by the non-oracle model is given by:

$$\begin{aligned}
p(A^R = 1 | B = j) &= P(A^R = 1 \& A^O = 1 | B = j) + P(A^R = 1 \& A^O = 0 | B = j) \\
&= P(A^O = 1 | B = j) p(A^R = 1 | A^O = 1, B = j) \\
&+ p(A^O = 0 | B = j) p(A^R = 1 | A^O = 0, B = j) \\
&+ p(A^O = 1 | B = j) - p(A^O = 1 | B = j) \\
&= p(A^O = 1 | B = j) - p(A^O = 1 | B = j) (1 - p(A^R = 1 | A^O = 1, B = j)) \\
&+ P(A^O = 0 | B = j) p(A^R = 1 | A^O = 0, B = j) \\
&= p(A^O = 1 | B = j) - p(A^O = 1 | B = j) p(A^R = 0 | A^O = 1, B = j) \\
&+ P(A^O = 0 | B = j) p(A^R = 1 | A^O = 0, B = j) \\
&= c_j - c_j \nu_j + (1 - c_j) \pi_j
\end{aligned}$$

where  $c_j = p(A^O = 1 | B = j)$  is the rate that a taxpayer from group  $j$  would be selected by the oracle;  $\nu_j = p(A^R = 0 | A^O = 1, B = j)$  is the false negative rate (i.e., the probability that a non-compliant taxpayer does not get audited); and  $\pi_j = p(A^R = 1 | A^O = 0, B = j)$  is the false positive rate (i.e., the probability that a compliant taxpayer would be audited, where “compliant” here refers to whether the taxpayer would be selected by the oracle for audit.

Disparity induced by the oracle is:

$$D_O = p(A^O = 1 | B = 1) - p(A^O = 1 | B = 0) = c_B - c_N$$

Disparity induced by the other model is:

$$\begin{aligned}
D_R &= p(A^R = 1 | B = 1) - p(A^R = 1 | B = 0) \\
&= c_B - c_B \nu_B + (1 - c_B) \pi_B - (c_N - c_N \nu_N + (1 - c_N) \pi_N) \\
&= D_O - (c_B \nu_B - c_N \nu_N) + ((1 - c_B) \pi_B - (1 - c_N) \pi_N)
\end{aligned}$$

So the “excess disparity” (relative to oracle) of the other model is:

$$\Delta_R = D_R - D_O = ((1 - c_B) \pi_B - (1 - c_N) \pi_N) - (c_B \nu_B - c_N \nu_N)$$

We can add and subtract terms to obtain:

$$\begin{aligned}
\Delta_R &= c_N \nu_N - c_B \nu_B + (1 - c_B) \pi_B - (1 - c_N) \pi_N \\
&= c_N \nu_N - c_N \nu_B + c_N \nu_B - c_B \nu_B \\
&+ (1 - c_B) \pi_B - (1 - c_N) \pi_B + (1 - c_N) \pi_B - (1 - c_N) \pi_N \\
&= c_N (\nu_N - \nu_B) + \nu_B (c_N - c_B) + \pi_B (c_N - c_B) + (1 - c_N) (\pi_B - \pi_N) \\
&= (1 - c_N) (\pi_B - \pi_N) + c_N (\nu_N - \nu_B) - D_O (\nu_B + \pi_B)
\end{aligned} \tag{7}$$

The three terms correspond to: (1) differences in the false positive rate by group; (2) differences in the false negative rate by group; and (3) to the extent the oracle would induce

a disparity, that disparity will be attenuated by the non-oracle model to the extent that the latter induces accuracy mistakes (of either type).

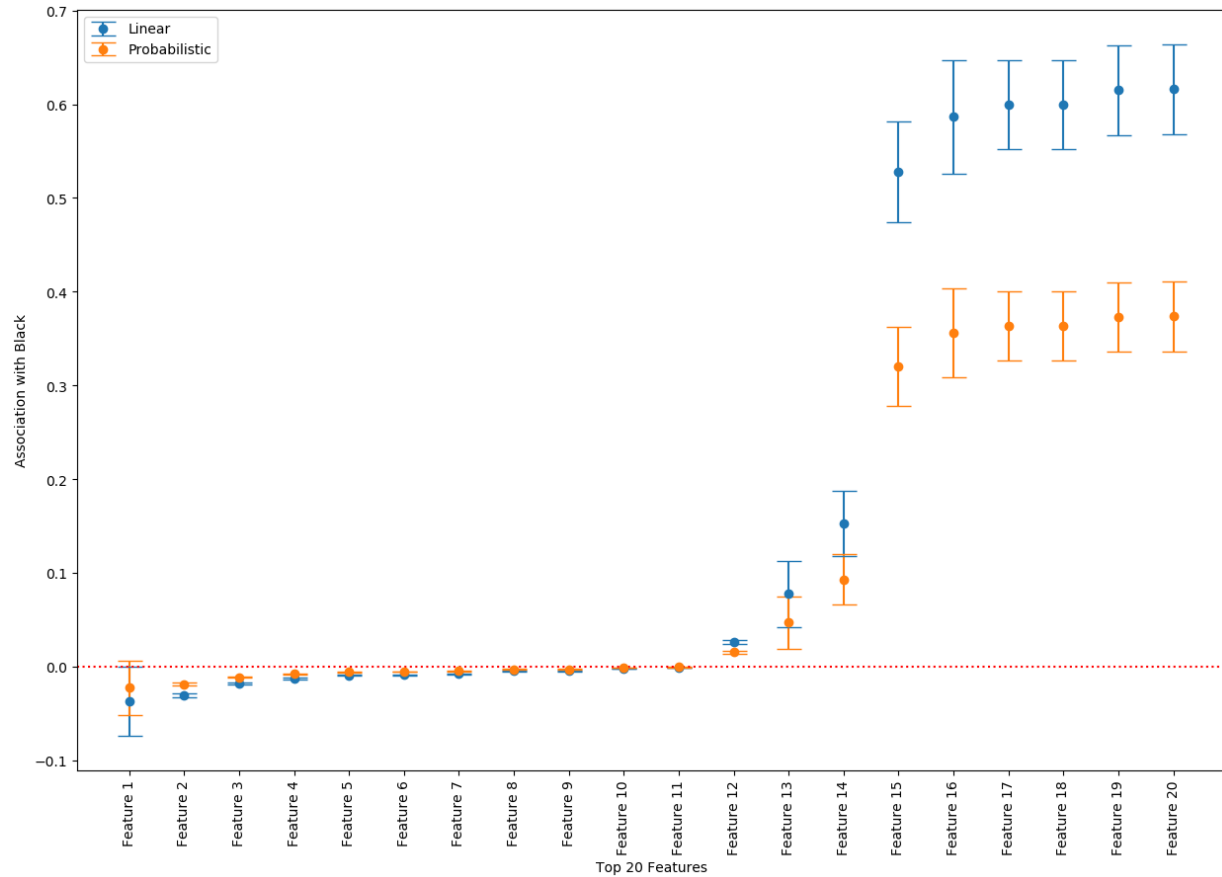
Finally, adding the oracle disparity to both sides of this equation yields an expression for the disparity induced by the non-oracle model:

$$D_R = (1 - c_N) (\pi_B - \pi_N) + c_N (\nu_N - \nu_B) + D_O (1 - \nu_B - \pi_B) \quad (8)$$

Intuitively, the disparity induced by a non-oracle model deviates from the disparity induced by the oracle based on average prediction errors (false positives and false negatives) as well as differences in the prediction errors by race.

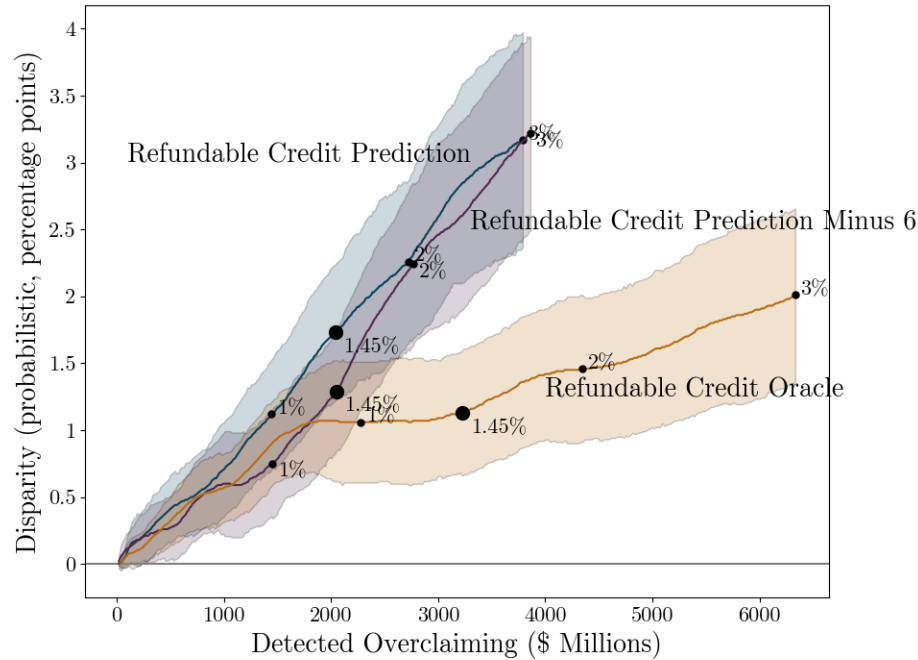
## Figures and Tables for Appendix F

Figure F.1: Distributions by Race for the Important Features of the Refundable Credit Prediction Model



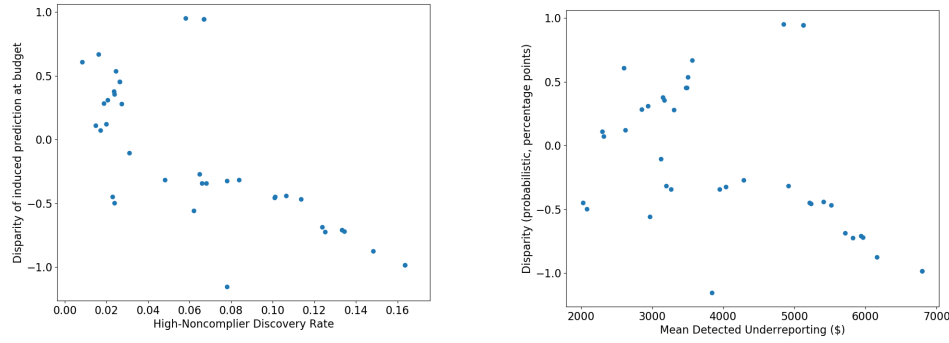
*Notes:* This figure shows the association between the top 20 most important features in the refundable credit model and the estimated probability a taxpayer is Black. Associations are calculated by first applying the linear and probabilistic estimators to each feature and then dividing the output by the standard deviation of the feature. Standard errors are calculated by dividing the standard error of the linear and probabilistic estimators (calculated from the asymptotic distributions described in Appendix B.3) by the standard deviation of the feature. Feature importance scores are computed as the mean decrease in impurity at each node of the decision tree, averaged over trees in the random forest model.

Figure F.2: Detected Underreporting and Disparity by Algorithm (Minus Identified Features)



*Notes:* The figure shows the estimated difference in audit rates between Black and non-Black taxpayers (y-axis) and annualized detected underreporting (x-axis) under alternative algorithms for selecting audits of EITC claimants and under alternative audit rates. Predictive models are trained and evaluated on the set of NRP EITC claimants from 2010-14; see Appendix E for details. The displayed trajectories correspond to the refundable credit prediction regressor (teal), total refundable credit oracle (orange), and refundable credit regressor trained without the 6 features most correlated with race among the twenty most important features of the refundable credit prediction regressor. These correspond to Features 15 through 20 in Appendix Figure F.1. The labeled points along each trajectory represent estimated detected overclaiming and disparity for the specified algorithm at the audit rate specified in the label. The audit rates considered range from 0.1% to 3%. The audit rate corresponding to the status quo (1.45%) is denoted by a larger dot. The refundable credit prediction algorithm is based on a random forest regressor trained to predict overclaiming of adjustments to EITC, CTC, and AOTC amounts. The refundable credit prediction minus 6 is the same, but is trained without the 6 features described above. The refundable credit oracle selects returns in descending order of true EITC, CTC, and AOTC overclaiming. Disparity is calculated using the probabilistic disparity estimator. Annualized detected underreporting is calculated as the total detected underreporting (positive or negative) imposed on returns selected for audit under the specified audit selection algorithm, scaled to reflect our use of five years of NRP data. The point labeled “Status quo” shows estimated disparity and total underreporting from the 1.45% of EITC returns selected for audit from the population of tax year 2014 returns. All analyses incorporate NRP sampling weights. Bars around each trajectory represent 95% confidence intervals around disparity estimates; they are calculated based on the distribution of estimates from 100 bootstrapped samples from the full set of NRP EITC claimants; see Appendix E for details.

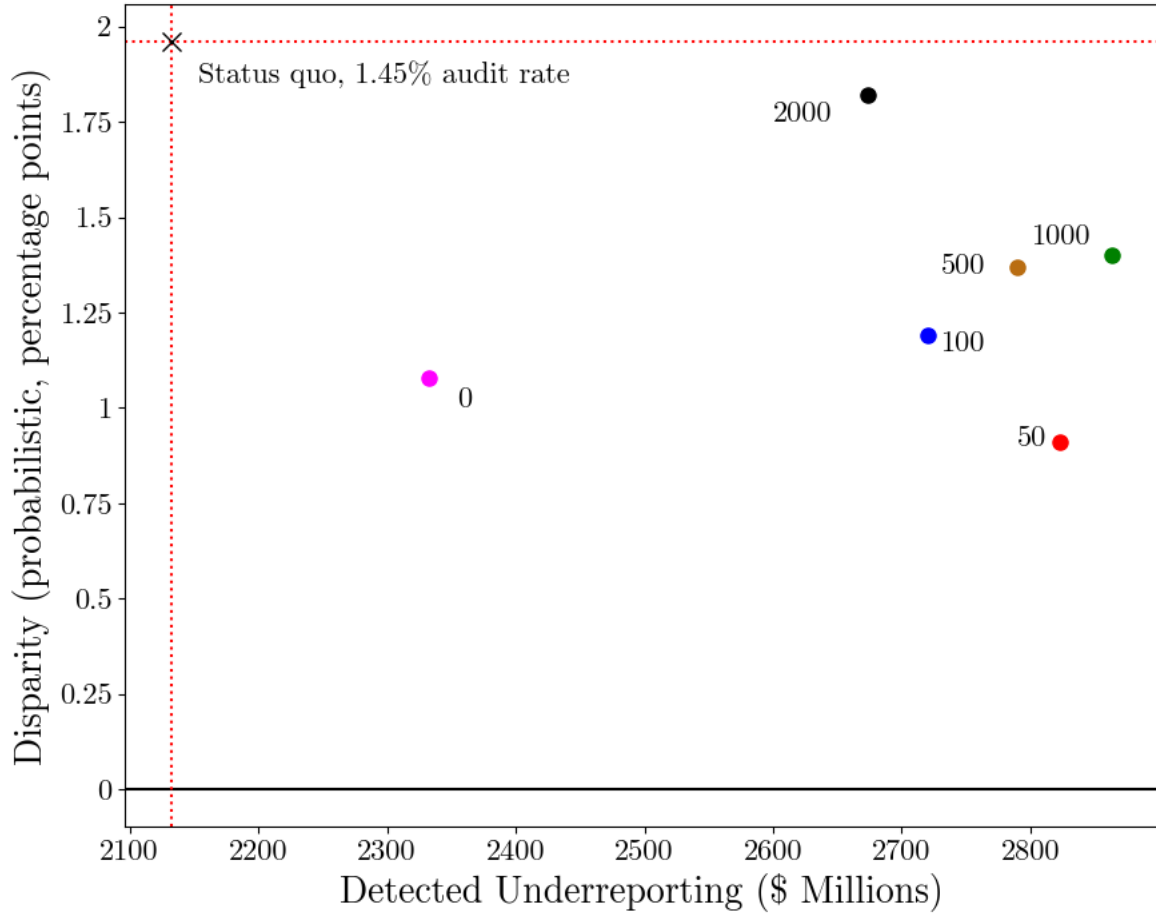
Figure F.3: Performance and Disparity of Single-Feature Models



*Notes:* The figure shows performance and disparity obtained by the single-feature decision trees for predicting total underreporting described in Section 7. Split-points are selected based on the induced impurity of the resulting segmentation. Five instances of the model are trained, each using a different single held-out fold for evaluation and the remaining four folds for training. Model performance is evaluated by selecting the top 1.45% (weighted) EITC NRP returns as ordered by the model predictions. Within each fold, we repeat the selection process 1000 times and permute returns for tie-breaking purposes, and report the average of the specified metric. Each point represents a model trained with a given feature; the x-axis represents performance (right panel: mean non-compliance conditional on selection by the model; left panel: fraction of the top 1.45% of non-compliant taxpayers discovered) and the y-axis represents disparity (as measured by the probabilistic estimator).

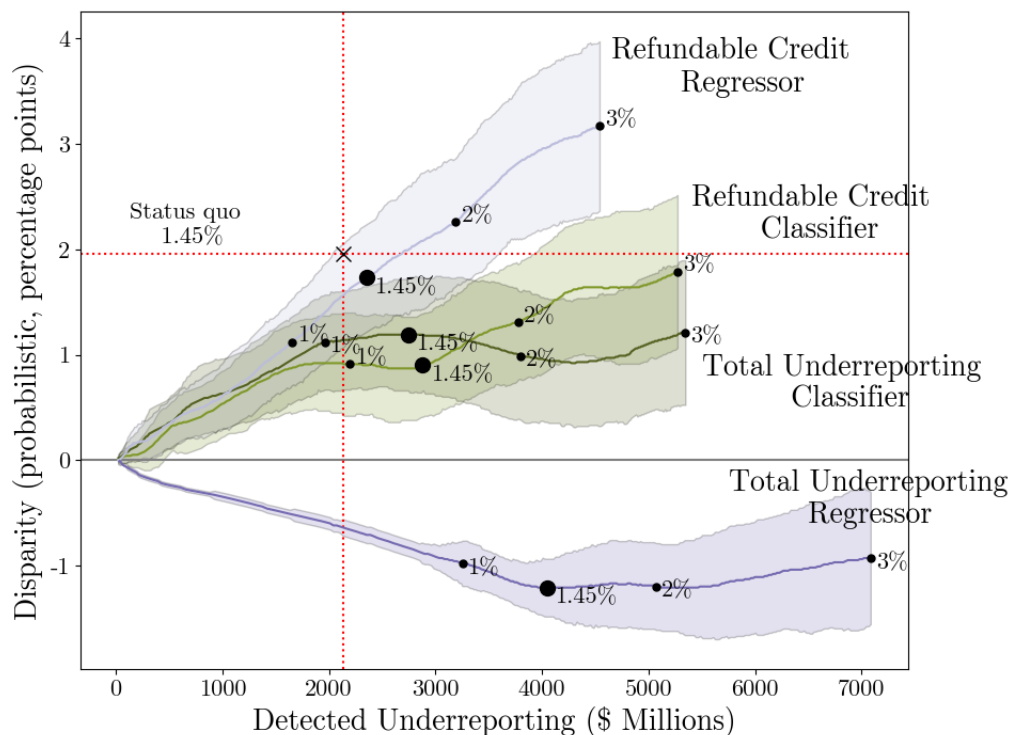


Figure F.4: Detected Underreporting and Disparity by Classifier Threshold



*Notes:* The figure shows the implied difference in audit rates between Black and non-Black taxpayers ( $y$ -axis) and annualized detected underreporting ( $x$ -axis) for random forest classification models trained on alternative dollar thresholds, under the assumption that 1.45% of the EITC population is selected for audit. Each point corresponds to a different classification model, trained to predict whether or not total adjustments exceed the specified dollar threshold. Disparity is calculated from the probabilistic disparity estimator applied to BIFSG-derived probabilities that a taxpayer is Black. Annualized detected underreporting is calculated as the total amount of adjustments (positive or negative) imposed on returns selected for audit under the specified audit selection model. Detected underreporting and disparity estimates for all models are constructed using the full set of NRP EITC observations from 2010-14. Details relating to the predictive models and associated estimates are contained in Appendix E. The point labeled “Status quo” shows the estimated disparity in audit rates between Black and non-Black taxpayers and detected underreporting from the 1.45% of EITC returns selected for audit from the population of tax year 2014 returns. All analyses incorporate NRP sampling weights.

Figure F.5: Detected Underreporting and Disparity by Algorithm (Regression vs. Classification)



*Notes:* The figure shows the estimated difference in audit rates between Black and non-Black taxpayers ( $y$ -axis) and annualized detected underreporting ( $x$ -axis) under alternative algorithms for selecting audits of EITC claimants and under alternative audit rates. Predictive models are trained and evaluated on the set of NRP EITC claimants from 2010-14; see Appendix E for details. The displayed trajectories correspond to the total underreporting classifier (dark green), total underreporting regressor (dark purple), refundable credit classifier (light green), and refundable credit regressor (light purple) algorithms. The labeled points along each trajectory represent estimated detected underreporting and disparity for the specified algorithm at the audit rate specified in the label. The audit rates considered range from 0.1% to 3%. The audit rate corresponding to the status quo (1.45%) is denoted by a larger dot. The total underreporting classifier is based on a random forest model trained to predict whether or not underreporting exceeds \$100. The total underreporting regressor is based on a random forest model trained to predict total underreporting. The refundable credit classifier is based on a random forest model trained to predict whether or not total adjustments to EITC, CTC, and AOTC amounts exceed \$100. The total underreporting regressor is based on a random forest model trained to predict total adjustments to EITC, CTC, and AOTC amounts. Disparity is calculated from the probabilistic disparity estimator applied to BIFSG-derived probabilities that a taxpayer is Black. Annualized detected underreporting is calculated as the total detected underreporting (positive or negative) imposed on returns selected for audit under the specified audit selection algorithm, scaled to reflect our use of five years of NRP data. The point labeled “Status quo” shows estimated disparity and total underreporting from the 1.45% of EITC returns selected for audit from the population of tax year 2014 returns. All analyses incorporate NRP sampling weights. Bars around each trajectory represent 95% confidence intervals around disparity estimates; they are calculated based on the distribution of estimates from 100 bootstrapped samples from the full set of NRP EITC claimants; see Appendix E for details.

Table F.1: Refundable Credit “Excess” Disparity Decomposition (Probabilistic)

	Black	Non-Black
Oracle Selection Rate	2.309	1.230
False Positive Rate	2.236	0.886
False Negative Rate	71.547	82.282
Excess Disparity Contribution from ...		
Scaled Difference in False Positive Rates	1.333	
Scaled Difference in False Negative Rates	0.132	
Attenuation of Oracle Disparity	-0.796	
Total Excess Disparity	0.669	
Oracle Disparity	1.079	
Prediction Model Disparity	1.748	

*Notes:* The table decomposes the excess disparity between the Refundable Credit Prediction Model and the Refundable Credit Oracle, as observed in Figure 8. All analyses uses NRP EITC claimants 2010-14, and incorporate NRP sampling weights. Quantities are percentage points (0-100). All estimates are based on the probabilistic estimator. The first row reports the audit rate by race when the refundable credit oracle selects 1.45% of EITC claimants. The second row reports the fraction of taxpayers who would be selected by the refundable credit prediction algorithm among those who would not be selected by the refundable credit oracle (“false positives”). The third row reports the fraction of taxpayers who would be not be selected by the refundable credit prediction algorithm among those who would be selected by the refundable credit oracle (“false negatives”). The next three rows correspond to the three terms in the excess disparity decomposition (see Equation (7) in Appendix F.5). The fifth row corresponds to the first term in the decomposition,  $(1 - c_N)(\pi_B - \pi_N)$ . The sixth row corresponds to the second term in the decomposition,  $c_N(\nu_N - \nu_B)$ . The seventh row corresponds to the third term in the decomposition,  $D_O(\nu_B + \pi_B)$ . The eighth row corresponds to the “excess” disparity and is the sum of the three previous rows. The ninth row corresponds to the disparity induced by the Refundable Credit Oracle. The final row corresponds to the disparity induced by the Refundable Credit Prediction Model and is the sum of the two previous rows. Table F.2 presents an analog to this analysis using the linear estimator.

Table F.2: Refundable Credit “Excess” Disparity Decomposition (Linear)

	Black	Non-Black
Oracle Selection Rate	2.732	1.121
False Positive Rate	2.772	0.750
False Negative Rate	68.927	83.543
Excess Disparity Contribution from ...		
Scaled Difference in False Positive Rates	1.999	
Scaled Difference in False Negative Rates	0.164	
Attenuation of Oracle Disparity	-1.155	
Total Excess Disparity	1.007	
Oracle Disparity	1.611	
Prediction Model Disparity	2.611	

*Notes:* The table replicates Appendix Table F.1 using the linear estimator. Specifically, the table decomposes the excess disparity between the Refundable Credit Prediction Model and the Refundable Credit Oracle, as observed in Figure 8. All analyses uses NRP EITC claimants 2010-14, and incorporate NRP sampling weights. Quantities are percentage points (0-100). All estimates are based on the linear estimator. The first row reports the audit rate by race when the refundable credit oracle selects 1.45% of EITC claimants. The second row reports the fraction of taxpayers who would be selected by the refundable credit prediction algorithm among those who would not be selected by the refundable credit oracle (“false positives”). The third row reports the fraction of taxpayers who would be not be selected by the refundable credit prediction algorithm among those who would be selected by the refundable credit oracle (“false negatives”). The next three rows correspond to the three terms in the excess disparity decomposition (see Equation (7) in Appendix F.5). The fifth row corresponds to the first term in the decomposition,  $(1 - c_N)(\pi_B - \pi_N)$ . The sixth row corresponds to the second term in the decomposition,  $c_N(\nu_N - \nu_B)$ . The seventh row corresponds to the third term in the decomposition,  $D_O(\nu_B + \pi_B)$ . The eighth row corresponds to the “excess” disparity and is the sum of the three previous rows. The ninth row corresponds to the disparity induced by the Refundable credit Oracle, and the final row corresponds to the disparity induced by the Refundable Credit Prediction Model.

Table F.3: High-Risk Classification in NRP and DDB

Panel A: Black			
	Not High-Risk (NRP)		High-Risk (NRP)
Not High-Risk (DDb)	0.68		0.22
High-Risk (DDb)	0.04		0.07
False Positive Rate			0.05
False Negative Rate			0.76
Panel B: Non-Black			
	Not High-Risk (NRP)		High-Risk (NRP)
Not High-Risk (DDb)	0.79		0.16
High-Risk (DDb)	0.03		0.03
False Positive Rate			0.03
False Negative Rate			0.84

*Notes:* The table displays the estimated distribution of Black (Panel A) and non-Black (Panel B) taxpayers for two tests of high-risk classification, one imputed in the Dependent Database, and one determined by line-by-line audits in NRP. Each cell shows the fraction of the group which was classified as either high risk or not high-risk according to the test result in DDB and high risk or not high-risk according to the test result as verified in the NRP. Estimates are computed using the probabilistic estimator and weighted using NRP weights to be representative of the full taxpayer population of EITC claimants with dependents. We report results for all taxpayers in our sample that have a non-missing high-risk indicator in the NRP, and in cases when such a taxpayer does not have a risk-indicator in DDB, we impute not high-risk. The false positive rate is calculated as the share of Black/non-Black taxpayers that are classified as high-risk by DDB but not high-risk by NRP, divided by the share of Black/non-Black taxpayers that are classified as not high-risk by NRP. The false negative rate is calculated as the share of Black/non-Black taxpayers that are classified as high-risk by NRP but not high-risk by DDB, divided by the share of Black/non-Black taxpayers that are classified as high-risk by NRP.

Table F.4: Missingness of Parental Social Security Numbers for EITC-Claimed Dependents

	Overall	Black	Non-Black
Missing Mother's SSN	15.40	15.02	15.51
Missing Father's SSN	40.56	53.55	36.85

*Notes:* The table reports the rate of missing parental information on birth certificates for EITC-claimed children on returns for tax year 2014. Units are percentage points (0-100). Rates are calculated at the return-level for the overall population (column 1), the population of Black taxpayers (column 2), and the population of non-Black taxpayers (column 3). The last two columns are calculated using BIFSG estimates and the probabilistic estimator. Note that taxpayer race in columns 2 and 3 refers to the estimated race of the taxpayer claiming the child, not the race of the child.

Table F.5: Disparity Estimates (Excluding Duplicate Child Claim Audits)

Estimator	Full Population (1)	EITC (2)	Non-EITC (3)
Linear	1.283 (0.004)	2.792 (0.009)	0.165 (0.003)
Probabilistic	0.776 (0.003)	1.887 (0.008)	0.093 (0.002)
N	148,305,318	28,338,472	119,966,846

*Notes:* The table shows estimated audit rate disparities using both the linear and the probabilistic estimators for audits initiated because the same child was claimed as a dependent on multiple returns. Units are percentage points (0-100). The Black/non-Black audit disparity is shown for the full population (column 1), the EITC population (column 2) as well as the non-EITC population (column 3). Standard errors, reported in parentheses, are calculated from the asymptotic distributions described in Appendix B.3. Each displayed disparity estimate is in terms of percentage points and is statistically different from zero ( $p < .01$ ).

Table F.6: Subgroup Disparity Estimates

Estimator	Overall	Income Category	Family Type	Preparer	Combined
Linear	2.900 (0.009)	2.379 (0.009)	2.413 (0.009)	2.899 (0.009)	2.087 (0.008)
Probabilistic	1.960 (0.008)	1.627 (0.007)	1.636 (0.007)	1.963 (0.008)	1.420 (0.007)

*Notes:* The table reports the conditional disparity as described in Section B.8.2. For each feature considered, we compute the  $D_x^p$  and  $D_x^l$  by applying the probabilistic and linear estimators, respectively, to the set of taxpayers whose value of the feature is  $x$ . Then we estimate bounds on  $\mathbb{E}[D_x]$  by averaging  $D_x^p$  or  $D_x^l$  over  $x$ . The standard error on each bound is the square root of the sum over all subgroups of the standard error of the disparity estimate within the subgroup squared times the share of taxpayers in that subgroup squared. We repeat this for: the overall population (as a reference); Income Category, constructed as a cross-product of total adjusted gross income quintiles and Schedule C Income Status (positive amount, non-positive amount, or none); Family Type, constructed as a cross-product of family status (married, single male, or single female) and number of dependents claimed (0, 1, 2, or 3+); Preparer, i.e. whether the taxpayer prepares their own taxes; and Combined, i.e. the cross-product of all of the above.

## G Appendix References

### References

- Berger, Y. G. (1998). Rate of convergence to normal distribution for the Horvitz-Thompson estimator. *Journal of Statistical Planning and Inference*, 67(2):209–226. (Cited on Appendix-49)
- Black, E., Elzayn, H., Chouldechova, A., Goldin, J., & Ho, D. (2022). Algorithmic fairness and vertical equity: Income fairness with IRS tax audit models. In *2022 ACM Conference on Fairness, Accountability, and Transparency*, pages 1479–1503. (Cited on 7, Appendix-78)
- Casella, G. & George, E. I. (1992). Explaining the Gibbs sampler. *The American Statistician*, 46(3):167–174. (Cited on Appendix-56)
- Chen, J., Kallus, N., Mao, X., Svacha, G., & Udell, M. (2019). Fairness under unawareness: Assessing disparity when protected class is unobserved. In *Proceedings of the Conference on Fairness, Accountability, and Transparency* (pages 339–348). (Cited on 2, 7, 17, Appendix-47)
- Delevoeye, A. & Sävje, F. (2020). Consistency of the Horvitz–Thompson estimator under general sampling and experimental designs. *Journal of Statistical Planning and Inference*, 207:190–197. Elsevier. (Cited on Appendix-49)
- Geman, S. & Geman, D. (1984). Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, PAMI-6(6):721–741. <https://doi.org/10.1109/TPAMI.1984.4767596> (Cited on Appendix-56)
- Government Accountability Office. (2015). IRS return selection: Wage and investment division should define audit objectives and refine other internal controls. (Cited on 10, 41, Appendix-77)
- Guyton, J. & Hodge, R. (2014). The compliance costs of IRS post-filing processes. *IRS Research Bulletin*. (Cited on Appendix-34)
- Lu, B., Wan, J., Ouyang, D., Goldin, J., & Ho, D. E. (2024). Quantifying the uncertainty of imputed demographic disparity estimates: The dual-bootstrap. (Cited on 24, Appendix-13, Appendix-35)
- National Taxpayer Advocate. (2018). Annual report to Congress 2018: Most serious problem 11. (Cited on Appendix-77)
- Robinson, P. M. (1982). On the convergence of the Horvitz-Thompson estimator. *Australian Journal of Statistics*, 24(2):234–238 (Cited on Appendix-49)

- Tzioumis, K. (2018). Demographic aspects of first names. *Scientific Data*, 5(1):1–9 (Cited on 19, Appendix-29)
- Virtanen, P., et al. (2020). SciPy 1.0: Fundamental algorithms for scientific computing in Python. *Nature Methods*, 17:261–272. (Cited on Appendix-67)