# ESTIMATING THE LIKELIHOOD OF ARREST FROM POLICE RECORDS IN PRESENCE OF UNREPORTED CRIMES

BY RICCARDO FOGLIATO[1], ARUN KUMAR KUCHIBHOTLA[2], ZACHARY LIPTON[2], DANIEL NAGIN[2], ALICE XIANG[3], AND ALEXANDRA CHOULDECHOVA[2]

[1]*Amazon Web Services[*]*

[2]*Carnegie Mellon University*

[3]*Sony AI*

Many important policy decisions concerning policing hinge on our understanding of how likely various criminal offenses are to result in arrests. Since many crimes are never reported to law enforcement, estimates based on police records alone must be adjusted to account for the likelihood that each crime would have been reported to the police. In this paper, we present a methodological framework for estimating the likelihood of arrest from police data that incorporates estimates of crime reporting rates computed from a victimization survey. We propose a parametric regression-based two-step estimator that (i) estimates the likelihood of crime reporting using logistic regression with survey weights; and then (ii) applies a second regression step to model the likelihood of arrest. Our empirical analysis focuses on racial disparities in arrests for violent crimes (sex offenses, robbery, aggravated and simple assaults) from 2006–2015 police records from the National Incident Based Reporting System (NIBRS), with estimates of crime reporting obtained using 2003–2020 data from the National Crime Victimization Survey (NCVS). We find that, after adjusting for unreported crimes, the likelihood of arrest computed from police records decreases significantly. We also find that, while incidents with white offenders on average result in arrests more often than those with black offenders, the disparities tend to be small after accounting for crime characteristics and unreported crimes.

**1. Introduction.** Characterizing the likelihood that a criminal offense will result in an arrest is central to multiple lines of criminological research, including crime control, deterrence, and racial disparities (Nagin, 2013; Piquero and Brame, 2008). Analyses of arrests traditionally rely on data collected by law enforcement agencies. The offenses captured by these records, however, represent only a fraction of all crimes that occur. By neglecting the "dark figure of crime" (Skogan, 1974), these analyses inevitably overestimate the underlying arrest rate per crime committed. The overestimation can potentially be severe, as data of criminal victimization reveal that less than half of violent offenses in the US ever become known to law enforcement (Morgan and Truman, 2021).

In order to estimate the likelihood of arrest for all crimes that are committed, police records can be augmented with data on crime reporting from victimization surveys. In the US, the National Crime Victimization Survey (NCVS) collects information on whether respondents experienced a victimization, and whether police were made aware of the offense. The idea of combining victimization data with police records was first proposed by Blumstein and Cohen (1979), who estimated arrest rates for violent offenses in Washington D.C. in the 1970s. However, owing to the limited data available, their approach could not account for variations in crime reporting rates across offense characteristics.
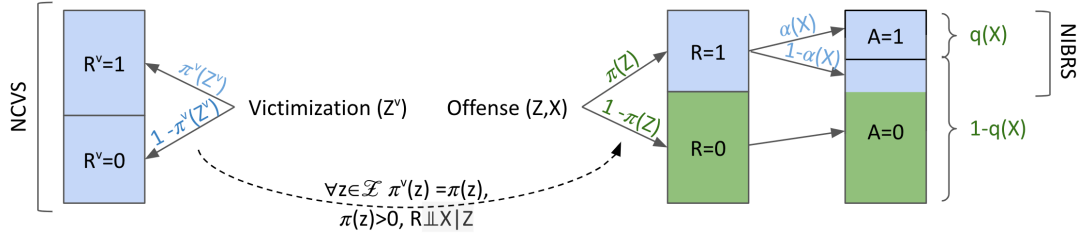
---

FIG 1. *Summary of the proposed methodology. On the left, a victimization with characteristics $Z^v$ is reported to law enforcement ($R^v = 1$) with probability $\pi^v(Z^v)$. NCVS data capture both reported and unreported victimizations. On the right, an offense with characteristics $(Z, X)$ is reported to the police ($R = 1$) with some positive probability $\pi(Z)$. $X$ captures all of the information that is contained in $Z$, i.e., the distribution of $Z|X$ is degenerate. Only reported offenses ($R = 1$) appear in NIBRS data. A reported crime can result in an arrest ($A = 1$) with probability $\alpha(X)$. The target of interest in our work is the conditional probability $q(X)$ that an offense will result in an arrest. Although the estimation of $\pi$ (and consequently of $q$) cannot be pursued solely on police records, assuming that $\pi^v(z) = \pi(z)$ for all $z$'s on the support $\mathcal{Z}$ of $Z$ allows for the estimation of the likelihood of police notification on these data. The rigorous description of this setup is detailed in Section 4.*

The increasing granularity and availability of incident-level crime data released by police agencies through the National Incident Based Reporting System (NIBRS) offers an opportunity to improve this analysis by accounting for crime characteristics. Combined with the information collected through the NCVS, NIBRS records can yield more accurate estimates of arrest rates per crime committed. Unfortunately, it is not possible to link records in NCVS and NIBRS directly. In this paper, we propose statistical methods to estimate the likelihood of arrest on incident-level data from police records while accounting for each crime's likelihood of police notification, computed from victimization survey data. The proposed methodology consists of two simple steps (Figure 1). We first estimate the likelihood of police notification conditional on offense characteristics via logistic regression with survey weights on victimization data. Then, we derive estimators of the total number of offenses, the rate of police notification, and the rate of arrest that leverage the characteristics of crimes in police records and the likelihood that the crime would be reported. The likelihood of arrest conditional on crime characteristics is modeled using logistic regression, and the coefficient estimates are obtained through a two-step estimation approach which accounts for the rates of crime reporting. When fitting logistic regression on crimes with multiple offenders, we handle the data dependence using generalized estimating equations (GEEs) (Liang and Zeger, 1986). We show that the proposed estimators are consistent and asymptotically normal for the target parameters, under a series of assumptions. Although our analytical results rely on the assumption that the models are correctly specified, the model results can be interpreted if this assumption does not hold in practice (Buja et al., 2019a,b; Berk et al., 2019). Although we focus on logistic regression, the proposed framework can be used to show asymptotic normality of any parametric regression model using an analysis similar to the one we conduct.

Our empirical investigation focuses on the assessment of differences in the likelihood of arrest across racial groups on 2006–2015 NIBRS data, with estimates of crime reporting obtained from NCVS data. We focus on violent crimes (sex offenses, robbery, aggravated and simple assaults) because the race of the offender is observed by the victim in the majority of such incidents. By contrast, many property crimes occur without the victim present, and so there is often no opportunity for them to directly observe offender characteristics. Our analysis reveals that on average about one in two violent offenses becomes known to law enforcement and one in five eventually results in arrest. Since the likelihood of crime reporting and of arrest vary with crime characteristics, arrestees do not form a representative sample of all offenders.

In terms of racial disparities, we find that crimes involving black offenders are reported at (marginally) higher rates than those involving white offenders, yet they result in arrests less often. Once crime characteristics are accounted for, the estimated differences in arrests across racial groups tend to be small. We further validate our results through an additional analysis where we employ nonparametric models instead of logistic regression to estimate the likelihood of police notification on survey data. Throughout the discussion one should keep in mind that our empirical findings rely on multiple data-related assumptions, which may not hold true (see limitations in Section 7).

The rest of the paper is organized as follows: Section 2 contains related work. Section 3 describes the data sources and related data processing. Section 4 contains the methodological framework. Section 5 includes the empirical analysis and in Section 6 we present the results of our study. Limitations of our analysis and future work are discussed in Sections 7 and 8 respectively.

**2. Related work.**    The methodology developed in our work is related to the literature on missing data problems (being a case of missingness not at random) (Kang and Schafer, 2007; Little and Rubin, 2019) and on capture-recapture (Petersen, 1896; Lee and Chao, 1994). In both lines of work, the target of key interest is the expected value of a random variable that is only partially observed. Works in this area mainly employ inverse probability weighting methods, e.g., using the Horvitz-Thompson and Hájek estimators (Horvitz and Thompson, 1952; Basu, 2011). In this work, we use the latter estimator to compute the share of unreported crimes and arrest rates from police records. The inclusion probabilities correspond to the estimated likelihood of police notification obtained via logistic regression. This approach can be thought of as a special case of the capture-recapture setting studied in Huggins (1989) when there is only one occasion to recapture. Similar results in the capture-recapture setting are also obtained by Van Der Heijden et al. (2003) and Böhning and Van Der Heijden (2009). However, these works impose distributional assumptions to handle data dependence. Instead, following prior empirical analyses of NIBRS (D'Alessio and Stolzenberg, 2003; Fogliato et al., 2021), we assume independence of the observations. A popular procedure used to fit models in presence of sampling bias is the two-step approach proposed by Heckman (1979). This approach most commonly applies probit regression in the first stage and OLS in the second stage. Our approach instead uses logistic regression in both steps and relies on a different set of assumptions. The design and derivation of our estimation procedure also draws from the literature on survey sampling (Särndal, Swensson and Wretman, 2003) and two-step M-estimation (Newey and McFadden, 1994). The regression analysis of arrests for crimes involving multiple offenders (hence with associated outcomes) via GEEs is inspired by the methods developed in the epidemiological literature (Hubbard et al., 2010). Lastly, our analysis operates under the assumption of covariate shift, i.e., that the distribution of the regressors but not of the outcome may vary between train and test sets (Sugiyama, Krauledat and Müller, 2007). In our setting, regressors and outcomes are represented by victimization and offense characteristics and by whether the crime has been reported respectively, while train and test set correspond to NCVS and NIBRS data respectively. Differently from these works, however, we assume that the posited regression models for the likelihood of crime reporting are well-specified and consequently no adjustments of the loss, such as by reweighting (Byrd and Lipton, 2019), are required.

Our work contributes to the literature on crime control. Estimates of the dark figure of crime and arrest rates for violent offenses have traditionally been obtained either from cross-sectional data and from self-reports of offending behavior, or solely from victimization data. Unlike our study, these analyses generally focus on arrest rates per individual rather than per crime committed. The approach taken by the studies in the first line of work was pioneered

by Blumstein and Cohen (1979), and is in spirit similar to ours. They compute arrest rates as the ratio of the arrest rates measured on police records and of the crime reporting rates on victimization surveys (Blumstein et al., 1986; Blumstein and Cohen, 1987). These studies suffer from one major drawback: As we mentioned in the Introduction, by using aggregate data from police agencies, they cannot account for variations in the likelihood of police notification across crime types, as Blumstein et al. (1986, page 335) also noted. This is relevant to our analysis especially because NIBRS and NCVS may capture populations of offenses with different characteristics, as not all police agencies have adopted NIBRS yet. Our proposed methodology addresses this issue. Within the second line of work, Blumstein et al. (2010) estimate arrest rates for violent crimes on the Rand Second Inmate Survey, a survey of inmates in three US jails conducted in the 1970s. Their estimates of arrest rates are close to ours. In the same study, the authors also assess arrest rates across racial groups and find no evidence of disparities. Two analyses focused on data from the Pathways to Desistance study, a longitudinal investigation of serious juvenile offenders from adolescence to young adulthood (Piquero and Brame, 2008; Brame et al., 2004), similarly do not find racial disparities in arrests. The seeming contrast with our results may be explained by the different nature of the populations of offenders we focus on, or by temporal differences. Within the third and last line of work, Buil-Gil, Medina and Shlomo (2021) estimate how the dark figure of crime varies with crime and neighborhood characteristics on victimization survey data in the UK. They conclude that this figure is associated with the socioeconomic status of the parties involved. Although our analysis does not account for these specific characteristics, our results similarly reveal that the likelihood of crime reporting varies with the demographics of the victim and of the offender. Multiple studies have examined racial disparities in police notification and arrests for violent offenses known to law enforcement. There is evidence that, overall, incidents with black offenders are at least as likely as those with white offenders to be reported to law enforcement (Morgan et al., 2017; Beck and Blumstein, 2018). After accounting for contextual factors, incidents are generally more likely to be reported when one of the parties involved is black (Avakame, Fyfe and McCoy, 1999; Xie and Lauritsen, 2012; Baumer and Lauritsen, 2010; Bachman, 1998; Fisher et al., 2003), although there exist both conflicting and null findings (Baumer, 2002; Dugan, 2003). In our analysis, we find that incidents with black offenders are reported at slightly higher rates than those with white offenders, even conditional on crime characteristics.

There is also mixed evidence concerning the magnitude of differences in arrests across racial groups for crimes known to law enforcement. While some works have concluded that crimes are more likely to result in arrest when the offender is black (Kochel, Wilson and Mastrofski, 2011; Lytle, 2014), multiple analyses focused on violent offenses on NIBRS data have reached a different conclusion (D'Alessio and Stolzenberg, 2003; Pope and Snyder, 2003; Roberts and Lyons, 2009). These works have found that, even after accounting for crime characteristics, black offenders are less likely to be arrested than white offenders for assault and robbery. Differences for rape and homicide were found to be negligible. In our analysis, we find that accounting for unreported crimes reduces the estimated gap in arrest rates for robbery, and the estimated gaps for assaults are close to zero. While most studies have focused on the analysis of incidents with single offenders and victims, Lantz and Wenger (2019) analyze incidents involving violent offenses where white and black individuals offend together. They fit one single regression model and they conclude that white offenders are less likely to be arrested than black offenders. We instead focus on all crimes with multiple offenders, fit separate models for each crime type, and find that the likelihood of arrest is mostly similar across racial groups of offenders. The only exception is robbery, for which arrest appears to be more likely for crimes involving white offenders, regardless of whether unreported crimes are accounted for. Lastly, we note that the results in the literature are likely

susceptible to issues stemming from model misspecification. For instance, Fogliato et al. (2021) showed that model misspecification can impact the magnitude and even the direction of the estimated racial disparities. In this work, we reach analogous conclusions: Our model estimates vary depending on the subsets of crimes considered in the analysis.

**3. Data.** Our empirical analysis leverages data from the National Crime Victimization Survey (NCVS) and the National Incident Based Reporting System (NIBRS). Similarly to past studies on NIBRS (D'Alessio and Stolzenberg, 2003; Fogliato et al., 2021), the main analysis in the paper centers on incidents with one victim and one offender. The data processing to obtain the dataset of crimes involving multiple offenders requires stronger assumptions. This is because NCVS only allows for inference at the level of the incident, while NIBRS contains also offender-level data. We now describe each of the two data sources and related data processing in turn. In this step of the analysis, we wish to identify a set of incidents captured by NCVS and NIBRS that share similar characteristics to ensure that the covariate shift assumption underpinning our analysis plausibly holds.

3.1. *Data on criminal victimization.* The NCVS represents the primary source of information on victimization in the US (Barnett-Ryan, Langton and Planty, 2014). By collecting information on the magnitude and extent of criminal victimization from a nationally representative sample of households, it is designed to complement data from police agencies with an alternative measurement of crime. Survey respondents aged 12 or older are interviewed regarding the criminal victimizations that they experience for nonfatal personal crimes (United States Department of Justice, 2017a). Our analysis focuses on data from interviews conducted between 2003 and 2020, which we obtain from the repository of the Inter-university Consortium for Political and Social Research (ICPSR) (United States Department of Justice, 2021). The data contain information on the stratified, multi-stage cluster sampling design, namely (pseudo-)strata, primary sampling units (PSUs), and observations (incidents) weights for serious crimes.[1] Information about the sampling design and on the construction of the sampling weights can be found in United States Department of Justice (2017a).

In the data, we consider only victimizations that satisfy the following criteria. (i) Incidents need to include an offense of simple assault (excluding verbal threats of assault), aggravated assault, robbery, or rape/sexual assault, which we will refer to as "sex offenses" in the rest of the paper.[2] (ii) We keep only incidents that have occurred within the United States. (iii) Since the NCVS collects information (e.g., demographics) only about the respondent, we drop incidents involving more than one victim. (iv) We consider only incidents with black or white individuals, with the inclusion of Hispanics.[3] In case of incidents involving multiple

---

[1]Similarly to past studies on NCVS (Xie and Lauritsen, 2012; Xie and Baumer, 2019a), our analysis assumes that nonresponse bias is accounted for by the use of survey weights. We acknowledge that in practice this assumption may not hold true.

[2]The Bureau of Justice Statistics (BJS) conflates simple assault with verbal threats of assault in their annual reports. In NIBRS, however, only physical attacks are coded as simple assaults (c.f. page 18 in United States Department of Justice (2019)). In order to align the definitions of simple assaults in NIBRS and NCVS, these crimes are excluded from the analysis. Consequently, statistics based on our proposed taxonomy, which has been chosen to ensure the maximal overlap of offense types between NCVS and NIBRS, will not match those in the reports produced in the BJS reports.

[3]The ethnicity information for victim and offender has been available in NCVS data since 2003 and 2012 respectively. However, the exclusion of Hispanics from NIBRS data is rather challenging because not all agencies report the offender ethnicity information, which was introduced in 2013. One could potentially attempt to identify the law enforcement agencies that generally report such information and consider only data from those agencies, as Roberts and Lyons (2011) have done. That procedure, however, would introduce a geographical bias in our NIBRS sample and thus we do not adopt such an approach.

offenders, we consider only those in which at least one of the offenders belongs to these racial groups. Our final dataset of incidents with single offenders consists of 11145 observations which, when reweighted by the survey weights, correspond to about 40 million crimes. The most frequent types of offense is simple assault (54%, based on survey weights), followed by aggravated assault (24%). Robbery and sex offense are the least frequent types of crime and each of them comprises about 10% of the available observations. The dataset of incidents with one or more offenders comprises 3405 additional observations and in total it corresponds to about 50 million incidents.

The outcome of interest in NCVS is whether the police are aware of the incident, as reported by the NCVS respondent in the survey ($R \in \{0, 1\}$). We consider the likelihood of an incident being reported $\pi^v(Z^v)$ to depend on a set of factors $Z^v$ which include characteristics of the parties involved and contextual factors. In terms of demographics, we account for the age, sex, and race of both victim and offenders. In the analysis of multiple offenders, we consider the sex of the majority of the offenders, and the age of the youngest and of the oldest offenders. We also consider the relationship between victim and offenders (e.g., if they are relatives), whether the victim suffers from a serious or minor injury, and whether the offenders have a firearm or a different weapon. We include two variables corresponding to whether the incident happens during the day and whether it occurs in a public area. To account for geographical variations in the likelihood of police notification, we account for whether the incident took place in a metropolitan statistical area (MSA), the corresponding US Census region in which the incident took place, and the year of the interview. Lastly, we consider whether the offense has been only attempted, the type of crime, and, in case of sex offenses, whether the offense consists of either rape or sexual assault. All variables other than the victim's age and the year are categorical.

3.2. *Crime data from law enforcement agencies.* NIBRS is part of the Federal Bureau of Investigation's Uniform Crime Reporting (UCR) data collection program. Through this program, law enforcement agencies submit detailed data on the characteristics of incidents that are known to them, including information on victims and offenders, and on the nature of the offenses. Our analysis builds on the assumption that when a crime becomes known to law enforcement, it will be recorded in the data released by law enforcement. Our analysis relies on 2006–2015 NIBRS data obtained from the ICPSR repository (United States Department of Justice, 2008a, 2009a, 2010a, 2011a, 2012a, 2013a, 2014a, 2015a, 2016a, 2017b). Note that while we rely on NCVS data from the period 2003-2020, the NIBRS data spans a shorter time period.

For this analysis, we identify incidents with characteristics that are similar to those included in our NCVS dataset. Thus, we apply the following data restrictions. (i) We consider incidents involving crimes of rape and sexual assault (i.e., sex offenses)[4], robbery, aggravated assault, and simple assault. The majority of the incidents (about 99% of cases) contain only one of these offenses. In the analysis of incidents with multiple offenders, we similarly found that in almost all of the incidents the offenders were charged with the same offense. Thus we can reasonably make the simplifying assumption that all offenders involved in the same crime incident commit the same offense. (ii) We keep only data from the 16 states that reported most of their crime data through the NIBRS in this time period. These states are Arkansas, Colorado, Delaware, Idaho, Iowa, Kentucky, Michigan, Montana, New Hampshire, North Dakota, South Carolina, South Dakota, Tennessee, Vermont, Virginia, and West

---

[4]In the category of rape and sexual assault we include crimes of forcible rape, forcible sodomy, sexual assault with an object, and forcible fondling. We do not consider statutory rape and incest because such offenses are unlikely to be reported by NCVS respondents in the interviews. The definition of rape in the UCR was revised in 2013 to also include male victims and female offenders.

Virginia. This exclusion makes our sample representative of a population that is well defined, i.e., the crimes that have become known to police and reported by agencies in the 16 states considered. (iii) We drop incidents that involve more than one victim and, for the analysis of incidents with single offenders, we also drop those that involve more than one offender. (iv) We account only for incidents where the races of victims and offenders are either black or white, including Hispanics. We observe that, based on the data of ethnicity that are available, about 90% of the offenders of Hispanic origin present in our sample are classified as whites. (v) We drop incidents that are cleared by exceptional means due to the death of the offenders or because the offender is in the custody of another jurisdiction.[5] We consider the remaining incidents that are cleared by exceptional means, namely those for which a juvenile offender was not taken into custody, prosecution was declined, or the victim refused to cooperate, as having no arrest. (vi) Lastly, to align our sample with the population of NCVS respondents, we drop incidents with victims aged 11 or younger. Our final samples of offenses involving only individual and one or more offenders consist of approximately 3.3 million and 4.9 million offenses respectively. As in the NCVS data, most of the offenses are simple assault (about 70%) and aggravated assault (about 17%).

The outcome of interest in our analysis is whether the incident results in the arrest of the offender (denoted $A \in \{0, 1\}$). In order to estimate the likelihood that a crime becomes known to the police for each incident in this dataset, $\pi(Z)$, we process the features in the data to obtain a set of crimes characteristics $Z$ that is analogous to those captured in our final NCVS dataset ($Z^v$). NIBRS also contains additional information that can be used in estimating the likelihood of arrest. In our application, we estimate the likelihood of the crime resulting in the offender's arrest, $q(X)$, based on crimes characteristics, $X$. In our analysis of incidents with individual offenders, $X$ includes not only all the variables that are present in $Z$, but also information about the state where the crime occurred, the size of the police force in the agency, and the number of police officers per capita. This additional police agency data is obtained from police employee datasets downloaded from ICPSR (United States Department of Justice, 2008b, 2009b, 2010b, 2011b, 2012b, 2013b, 2014b, 2015b, 2016b, 2017c). In the analysis of incidents with multiple offenders, $Z$ captures aggregate information about the incident, e.g., the age of the youngest offender. $X$ contains variables measured both at the level of the incident and of the individual offender, e.g., the age of the individual offender.

---

**Algorithm 1** Estimation strategy on NCVS and NIBRS

| | |
|---|---|
| $\hat{\gamma} \leftarrow \text{solve} \sum_{i=1}^{N^v} w_i I_i h^v(R_i^v, Z_i^v; \gamma) = 0$ | ▷ Likelihood of police notification $\pi^v(Z^v; \gamma)$ on NCVS data |
| $\hat{N} \leftarrow \sum_{i=1}^{N} R_i / \pi(Z_i; \hat{\gamma})$ | ▷ Total number of offenses $N$ on NIBRS data |
| $\hat{\pi}^* \leftarrow \sum_{i=1}^{N} R_i / \hat{N}, \ \hat{q}^* \leftarrow \sum_{i=1}^{N} A_i / \hat{N}$ | ▷ Rate of police notification $\pi^*$ and arrest $q^*$ on NIBRS data |
| $\hat{\theta} \leftarrow \text{solve} \sum_{i=1}^{N} R_i h(A_i, Z_i, X_i; \theta, \hat{\gamma}) = 0$ | ▷ Likelihood of arrest $q(X; \theta)$ on NIBRS data |

---

**4. Methods.** In this section, we present the statistical methodology behind our empirical analysis. Algorithm 1 describes the proposed estimation approach. We first estimate $\pi^v(Z^v; \gamma)$, the likelihood of police notification conditional on crime characteristics via survey-weighted logistic regression. We then introduce the offense data setup and describe the assumptions underpinning our inference strategy. Next, we review estimators of NIBRS summary statistics, namely the total number of offenses $N$, the rate of police notification $\pi^*$,

---

[5]The excluded incidents represent less than 1% of all offenses, so their inclusion is unlikely to change the conclusions of our analysis. In addition, there are not large differences in clearance by exceptional means across racial groups; see the results in Section A of the Appendix in Fogliato et al. (2021).

and the arrest rate for all crimes committed $q^*$. Lastly, we estimate $q(X; \theta)$, the likelihood of arrest conditional on crime characteristics via logistic regression on NIBRS data. Under regularity conditions, all of the estimators that we present are asymptotically normal. Detailed derivations and proofs of the results are deferred to the Appendix.

4.1. *Crime reporting on NCVS.* Consider the finite population of criminal victimizations in the US, denoted by $V^{N^v} = \{(Z_i^v, R_i^v)\}_{i=1}^{N^v}$. This can be viewed as an i.i.d. sample $(Z^v, R^v) \sim P^v$, where $Z^v = (Z(1), \ldots, Z(d_z)) \in \mathcal{Z}$ indicates the victimization's characteristics and $R^v \in \{0, 1\}$ is the indicator of whether the victimization becomes known to law enforcement. We have access to the NCVS survey sample of size $n^v$, which is drawn from $V^{N^v}$ under some probability sampling design $\psi$. Let the random variable $I_i = 1$ if $i^{th}$ observation is included in this sample, and $I_i = 0$ otherwise. The sample has an associated set of sampling weights $\{w_i : 1 \le i \le N^v, I_i = 1\}$, which are commonly thought as representing the number of units that each sampled observation represents in the larger finite population (Lohr, 2007).

We model the conditional probability of police notification, $\pi^v(Z^v)$, via logistic regression. That is, we take $\pi^v(z; \gamma) := 1/(1 + e^{-\gamma^T z})$ to describe $\mathbb{P}_{P^v}(R^v = 1 | Z = z)$, for $\gamma \in \Gamma$ for some compact set $\Gamma \subset \mathbb{R}^{d_z}$. The superpopulation target parameter $\gamma_0 \in \text{Int}(\Gamma)$ is defined by the moment condition $\mathbb{E}_{P^v}[h^v(R^v, Z^v; \gamma)] = 0$ where $h^v(R^v, Z^v; \gamma) := (R^v - \pi^v(Z^v; \gamma))Z^v$. The design-based estimator $\hat{\gamma}$ of $\gamma_0$ is the solution to the estimating equation $\sum_{i=1}^{N^v} w_i I_i h^v(R_i^v, Z_i^v; \gamma) = 0$ (Lumley and Scott, 2017). Under certain regularity conditions, $(\Sigma^v)^{-1/2} \sqrt{n^v}(\hat{\gamma} - \gamma_0) \xrightarrow{d} \mathcal{N}(0, I_{d_z})$ where $\Sigma^v$ is a positive definite matrix.

4.2. *NIBRS setup.* Let $O^N = \{(X_i, Z_i, R_i, A_i)\}_{i=1}^N$ denote the sample of all offenses committed, which is assumed to be an i.i.d. sample of $(X, Z, R, A) \sim P$. In this part of the analysis, we consider crimes with only one offender, so the independence assumption is likely to hold (but see the longer discussion in Section 7). Let $X = (X(1), \ldots, X(d_x)) \in \mathcal{X}$ and $Z = (Z(1), \ldots, Z(d_z)) \in \mathcal{Z}$ indicate incident characteristics. Let $R \in \{0, 1\}$ and $A \in \{0, 1\}$ indicate whether the offense is known to the police ($R = 1$) and whether it results in an arrest ($A = 1$) respectively. Given that an offense can result in an arrest only if it is known to the police, we assume that $R = 0$ implies $A = 0$. Note that police-recorded data contain only offenses that have been reported, i.e., those for which $R = 1$. We denote with $\mathbb{E}$ the expectation over $P$.

In order to estimate parameters of interest on the entire population using solely the observations for which $R = 1$, we will make use of the following set of assumptions.

**A.1** $\forall z \in \mathcal{Z}$, $\pi^v(z; \gamma_0) = \mathbb{P}_{P^v}(R^v = 1 | Z^v = z)$ for $\gamma_0 \in \Gamma$ where $\Gamma$ is a compact set.
**A.2** $\forall (x, z) \in \mathcal{X} \times \mathcal{Z}$, $\mathbb{P}(R = 1 | X = x, Z = z) = \mathbb{P}(R = 1 | Z = z)$.
**A.3** $\forall z \in \mathcal{Z}$, $\mathbb{P}(R = 1 | Z = z) = \mathbb{P}_{P^v}(R^v = 1 | Z^v = z)$.
**A.4** $\|X\|_\infty < M$ and $\|Z\|_\infty < M$ for some $M > 0$.

A.1 states that the parametric model $\pi^v(z; \gamma_0)$ is correctly specified for $\mathbb{P}_{P^v}(R^v = 1 | Z^v = z)$. We empirically assess this assumption by comparing the logistic regression model with a nonparametric approach, and find small differences in the estimates produced by the two methods for three of the four offense types considered. A.2 states that $R$ is independent of $X$ after conditioning on $Z$. In our empirical analysis, we study the likelihood of arrest per crime committed, $q$, as a function of only $X$ (see 4.4) because $X$ contains at least as much information as $Z$. In other words, $X$ includes more refined details about the incident such as specific geographical information and characteristics about each of the offenders within an incident (hence the distribution of $Z|X$ is degenerate). These characteristics are available in NIBRS but not in NCVS. However, through A.2 we assume that this additional

information is not relevant to the estimation of the distribution of $R|Z$. A.2 may be violated if $Z$ did not capture, for instance, variations in reporting rates across police agencies, but $X$ did. A.3 allows us to compare the probability of police notification in NIBRS and NCVS. This assumption casts our learning problem into the covariate shift setting. Together, A.1 and A.3 imply that $\pi^v(z;\gamma_0) = \mathbb{P}(R = 1|Z = z)$. Thus, in what follows we drop "$v$" from the superscript of $\pi(Z;\gamma_0)$. A.4 ensures that functions of these random variables will have finite moments. This assumption clearly holds true in our application. Lastly, A.1 and A.4 imply that $\pi(z;\gamma) > (1 + e^{\sqrt{d}M \sup_{\gamma \in \Gamma} \|\gamma\|})^{-1} > 0$ for all $z \in \mathcal{Z}$ and $\gamma \in \Gamma$, an assumption that is known as *positivity* in the causal inference literature (Little and Rubin, 2019). In our setting, it rules out the possibility that there exist offenses with certain characteristics that will never be reported to the police. Under conditions A.1–A.4, we can establish the consistency of the estimators we propose in the next section. By imposing additional assumptions on the rate of growth of $n^v$, $N^v$, and $N$, we can derive their asymptotic distributions as well.[6]

4.3. *Crime reporting and arrest rates on NIBRS.* There are three targets of key interest on NIBRS. First, the total number of offenses, which can be estimated using the Horvitz-Thompson estimator $\hat{N} := \sum_{i=1}^N R_i/\pi(Z_i;\hat{\gamma})$. Second, the expected rate of police notification $\pi^* := \mathbb{E}[R]$. Third, the arrest rate $q^* := \mathbb{E}[A]$. These rates are estimated by $\hat{\pi}^* := \sum_{i=1}^N R_i/\hat{N}$ and $\hat{q}^* := \sum_{i=1}^N A_i/\hat{N}$ respectively. Under the assumption that $\hat{\gamma}$ is consistent for $\gamma_0$ and asymptotically normal, as well as some regularity conditions, these estimators are asymptotically normal. The critical step in the derivation of the limiting distributions is to leverage the fact that $\mathbb{E}[R/\pi(Z;\gamma_0)] = 1$ in order to rewrite the unconditional expectations with respect to the event $\{R = 1\}$. This make possible the estimation based only on the sample we have access to.

4.4. *Conditional probability of arrest on NIBRS via logistic regression.* We model the probability of arrest conditional on the covariates $\mathbb{E}[A|X]$ using logistic regression; i.e., we consider $q(x;\theta) := 1/(1 + e^{-\theta^T x})$ where $\theta \in \Theta$ for a compact set $\Theta \subset \mathbb{R}^{d_x}$. The parameter $\theta_0 \in \text{Int}(\Theta)$ is defined by the following moment condition

$$\text{(1)} \qquad \mathbb{E}\left[(A - q(X;\theta))X\right] = 0.$$

Since $A = 0$ whenever $R = 0$, it follows that $\mathbb{E}[AX] = \mathbb{E}[RAX]$. Then, under Assumptions A.1–A.3, the moment condition (1) can be rewritten as

$$G(\theta,\gamma_0) := \mathbb{E}\left[\left(A - \frac{q(X;\theta)}{\pi(Z;\gamma_0)}\right)XR\right] = 0.$$

Thus, in practice, we compute the estimator $\hat{\theta}$ of $\theta_0$ by solving the following estimating equation

$$\text{(2)} \qquad \hat{G}_{\mathcal{N}}(\theta,\hat{\gamma}) := \frac{1}{N}\sum_{i=1}^N R_i h(A_i, Z_i, X_i; \theta, \hat{\gamma}) = 0,$$

where $h(A_i, Z_i, X_i; \theta, \gamma) := [A_i - q(X_i;\theta)/\pi(Z_i;\gamma)]X_i$. The estimate $\hat{\theta}$ of $\theta_0$ can be found using iteratively (re-)weighted least squares. Under Assumptions A.1–A.4, together with the consistency and asymptotic normality of $\hat{\gamma}$ as an estimator of $\gamma_0$, then $\Sigma^{-1/2}\sqrt{n}(\hat{\theta} - \theta_0) \xrightarrow{d} \mathcal{N}(0, I_{d_x})$ as $n \to \infty$ where $\Sigma$ is positive definite.

---

[6]Note that no assumption on $n := \mathbb{E}[R]$ is needed here because, unlike the survey sampling setting, the i.i.d. assumption on the data of offenses guarantees that $n$ and $N$ will grow at the same rate. Since the sample sizes of the NCVS and NIBRS samples are of comparable magnitude, we assume that $\lim_{n,n^v \to \infty} n/n^v = \kappa = O(1)$. The results readily generalize to the cases where $n \ll n^v$ or $n \gg n^v$.

4.5. *Conditional probability of arrest on NIBRS via GEEs.*   In the previous sections we considered only crimes with a single offender. We now consider estimating the conditional probability of arrest from police records on crimes where one *or more* offenders are involved. When multiple offenders act together, the i.i.d. assumption on $O^N$ clearly doesn't hold. We denote the sample of crime incidents, where each incident represents a cluster containing observations corresponding to the individual offenders, occurring in the US with $O^N = \{(K_i, \mathbf{X}_i, Z_i, R_i, \mathbf{A}_i)\}_{i=1}^N$, which is an i.i.d. sample of $(K, \mathbf{X}, Z, R, \mathbf{A}) \sim P$ with $K \in \mathbb{Z}_+$. Here $\mathbf{X}_i$ is a matrix of dimension $K_i \times d_x$ whose $k^{th}$ column corresponds to $\mathbf{X}_{ik}$, the characteristics relative to the offense committed by the $k^{th}$ offender in the $i^{th}$ incident. The vector $\mathbf{A}_i = (\mathbf{A}_{i1}, \dots, \mathbf{A}_{iK_i})^T$ indicates whether each of the offenders in the $i^{th}$ incident are arrested ($\mathbf{A}_{ij} = 1$ for $1 \leq j \leq K_i$) or not. $R_i$ indicates whether the $i^{th}$ incident is known to the police, and $Z_i$ represents characteristics of the same incident. Note that $Z_i$ contains information that is shared across all offenders within the same incident (e.g., location), while $\mathbf{X}_i$ may also include covariates that are specific to the individual offender (e.g., demographics of that offender).

Despite the likely positive correlation across outcomes within the same incident, the logistic regression coefficient estimates discussed in the previous section remain consistent for the target parameters. However, the estimates of their asymptotic variance need to be adjusted (Fitzmaurice, Laird and Rotnitzky, 1993). To account for the correlation in the variance estimation and to increase efficiency, we employ generalized estimating equations (GEEs) (Liang and Zeger, 1986). We assume that

**A.5** $\mathbb{E}[\mathbf{A}_{ij}|\mathbf{X}_i] = q(\mathbf{X}_{ij}; \theta_0)$ where $q(\mathbf{X}_{ij}; \theta_0) := (1 + e^{-\theta_0^T \mathbf{X}_{ij}})^{-1}$ for $1 \leq i \leq N$, $1 \leq j \leq K_i$, $\theta_0 \in \text{Int}(\Theta)$.

According to this assumption, the probability of arrest for an individual does not depend on incident's characteristics related to their co-offenders; see Fitzmaurice et al. (2008, Section 3.2) for a longer discussion of this assumption. To model the covariance across outcomes, we define the matrix $W_i(\theta, \alpha) := W(\mathbf{X}_i, Z_i; \theta, \alpha) = D_i(\theta)^{1/2} C_i(\alpha) D_i(\theta)^{1/2}$ for $1 \leq i \leq N$, where $D_i(\theta)$ is a diagonal matrix of dimension $K_i \times K_i$ whose $k^{th}$ diagonal entry corresponds to $q(\mathbf{X}_{ik}; \theta)(1 - q(\mathbf{X}_{ik}; \theta))$. $C_i(\alpha)$ is the so-called "exchangeable working correlation" matrix, which has dimension $K_i \times K_i$ with 1 on the diagonal and any $\alpha \in [-1, 1]$ elsewhere (Liang and Zeger, 1986).

The estimator $\hat{\theta} \in \text{Int}(\Theta)$ solves the following generalized estimating equation

$$\frac{1}{N} \sum_{i=1}^N R_i h_{GEE}(\mathbf{X}_i, Z_i, \mathbf{A}_i; \theta, \hat{\alpha}) := \frac{1}{N} \sum_{i=1}^N R_i \mathbf{X}_i D_i(\theta) W_i(\theta, \hat{\alpha})^{-1} \left( \mathbf{A}_i - \frac{\mathbf{q}_i(\theta)}{\pi(Z_i; \hat{\gamma})} \right) = 0,$$

where $\mathbf{q}_i(\theta) = (q(\mathbf{X}_{i1}; \theta), \dots, q(\mathbf{X}_{iK_i}; \theta))^T$ and $\hat{\alpha}$ is a consistent estimator of $\alpha_0$, the true correlation parameter, given $\theta$. Then, under certain certain conditions, $\Sigma_{GEE}^{-1/2} \sqrt{n}(\hat{\theta} - \theta_0) \xrightarrow{d} \mathcal{N}(0, I_{d_x})$ as $N \to \infty$ where $\Sigma_{GEE}$ is a positive definite matrix.
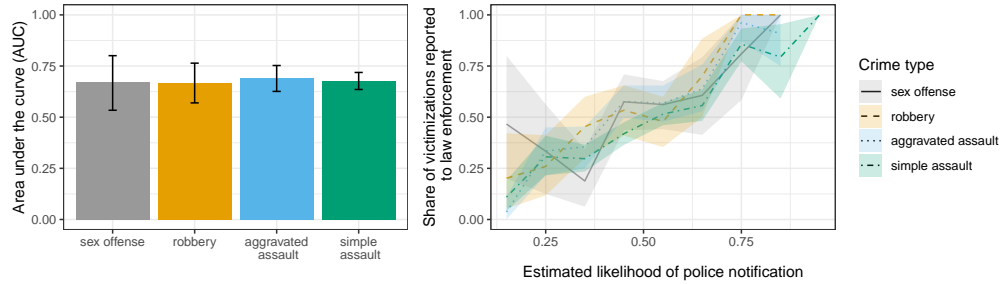
FIG 2. *Area under the curve (AUC) and calibration for the logistic regression model with survey weights that estimates the likelihood of police notification $\pi^v(Z^v)$ on 2003–2020 NCVS data, obtained via cross-validation for the four types of crimes considered. Error bars and shaded regions indicate 95% confidence intervals for the mean.*

**5. Empirical strategy.** Our empirical analysis leverages the analytical framework presented in Section 4. The code used for data processing and the analysis is available at github.com/ricfog/arrests-with-unreported-crimes.

5.1. *Missing data.* Both NCVS and NIBRS data contain a small share of missing values in certain crime characteristics. In NCVS data, a handful of the other variables are missing in less than 5% of the cases. The only exceptions are the MSA information and offender's age in crimes of robbery, which are missing in about one fourth of the observations. In NIBRS data, a few of the variables are missing for less than 5% of the observations. Past work on NCVS and NIBRS has assumed the data to be either missing completely at random (MCAR) or missing at random (MAR) (D'Alessio and Stolzenberg, 2003; Xie and Lynch, 2017; Fogliato et al., 2021). We assume the data to be MAR and impute the missing values via multiple imputation by chained equations (MICE) using predictive mean matching for numerical data, and multinomial and logistic regression for categorical data (Azur et al., 2011). We repeat the same procedure for both crimes involving individual and multiple offenders. It is common practice to create multiple imputed datasets and then pool the estimates computed from each (Graham, Olchowski and Gilreath, 2007). However, our preliminary analysis showed that multiple imputed datasets yielded similar estimates, likely because only a small proportion of observations have missing values. Due to the considerable computational costs of our workflow, we decided to use only one imputed dataset. Consequently, in the downstream inference we will not account for the uncertainty arising from the imputation procedure, which we anticipate to be negligible.

5.2. *Estimation on NCVS.* We now detail the key steps for the analysis of incidents with one offender. An analogous procedure is carried out for incidents with multiple offenders with similar results, so we omit the details for brevity.

We begin by describing the process for using the 2003–2020 NCVS victimization data to estimate the likelihood that an offense becomes known to law enforcement $\pi$ conditionally on its characteristics. We first split the data into two groups stratifying by outcome and year. One subset, which comprises one fifth of the data, is used for model selection. The validation subset, which consists of the rest of the data, is used to estimate the chosen model that will be employed for inference on NIBRS. In terms of model selection, we evaluate several regression models that include interactions between the crime types and the other regressors, which reflect different modeling choices made by past studies (Xie and Lauritsen, 2012; Xie and Baumer, 2019a; Baumer, 2002). This step is warranted by the diversity of feature sets and

model choices adopted in regression analyses by these works. We assess model performance via model calibration and area under the curve (AUC) separately for each type of offense using cross-validation. Both evaluations account for survey weights. The predictions appear to be well calibrated across the various models for all crime types considered other than sex offenses. The AUCs produced by the predictions of the various models are similar, and all are in the range 0.55–0.7 across offense types and models. We proceed with the logistic regression model without interactions between features and crime types, a model that was also considered in past work on NCVS similar to ours (Xie and Lauritsen, 2012).[7]

Next, we fit the selected model on the validation set and compute the variance of the coefficient estimates using the information on pseudo-strata and pseudo-primary sampling unit (PSU) information contained in the data. We report the cross-validated AUC and calibration of this model on this second set of observations in Figure 2. Despite the wide confidence intervals, the model predictions are well calibrated and the AUCs are above 0.6 across all four types of offenses.

Simultaneously, we conduct a sensitivity analysis to compare the estimates produced by the logistic regression with those obtained from a SuperLearner, which represents a more flexible approach (Van der Laan, Polley and Hubbard, 2007; Polley and Van Der Laan, 2010). For this purpose, we train a logistic regression model, a logistic Lasso model (Tibshirani, 1996), a multilayer perceptron with one hidden layer, a Naive Bayes classifier, and a random forest model (Breiman, 2001). We tune the hyperparameters of each of these models separately via cross-validation. To account for the survey weights in training, we resample the observations within each data fold selecting an observation with probability proportional to the corresponding survey weight. We select the set of hyperparameters that yield the highest average AUC across the four crime types. On the second subset of the data, we select by cross-validation the weights that correspond to the convex combination of the predictions produced by these models achieving the best predictive performance. As we describe in Section 6, the estimates of the likelihood of police notification generated by this nonparametric model are close to those obtained through the logistic regression approach.

5.3. *Estimation on NIBRS.* In the next stage, we use the results from the weighted logistic regression analysis on NCVS data. These results help us calculate the likelihood that the police will be notified about each individual incident in the NIBRS dataset. Particularly low values of the assessed probabilities would represent a potential violation of the positivity assumption, which would skew our estimates. Accordingly, we examine the predictions generated by the two models across crime types. The smallest detected values range between 0.05 for sex offenses to above 0.1 for the other offense types. Thus, the results of our analysis on NIBRS will not be overly influenced by a few outliers. We first estimate the total number of crimes $N$, the rate of police notification $\pi^*$, and the arrest rate $q^*$, along with their corresponding variances. Then, we follow the procedure described in Section 4.4 to estimate the likelihood of arrest conditional on covariates, $q(X; \theta_0)$, for incidents with individual offenders via logistic regression.

We additionally perform a number of robustness checks to assess the sensitivity of the downstream estimates to the modeling assumptions. First, we repeat these analyses using estimates of the likelihood of police notification obtained using the SuperLearner in place of the weighted logistic regression. We also investigate whether the logistic regression for $q(X; \theta_0)$ may be misspecified (provided that the model for $\pi(Z; \gamma_0)$ is correct), using "focal slope" model diagnostics proposed by Buja et al. (2019b). Using these graphical tools,

---

[7]We have conducted an additional analysis employing the same model with an interaction between offender's and victim's races. The results from this analysis and those that we report in the paper are similar.

we analyze how the coefficient estimates (specifically, offender's race) change when fitting the regression model on various configurations of the regressor distributions. To implement the reweighting procedure, we proceed as follows. We first construct a grid of five evenly spaced values for the numeric features, and use the grid values of $\{0, 1\}$ for the binary features. For each regressor separately, we split the observations into groups based on the grid's cell center that is closest to each observation's feature value in absolute distance. For each feature-grid cell pair, we then obtain 50 estimates of the logistic regression coefficients by bootstrap resampling 10,000 observations from the given group. We conclude the discussion by presenting our final analysis of incidents with one or more offenders which employs the GEEs framework described in Section 4.5.

**6. Results.** This section is organized as follows. We first empirically demonstrate how our approach, by virtue of accounting for covariate shift, strictly improves upon prior analyses such as Blumstein and Cohen (1979). Then, we assess racial disparities in the rates of police notification, and disparities in arrest rates on all crimes and on only those known to law enforcement. This first part of the analysis focuses on crime incidents with individual offenders. Next, we present the regression results for the estimation of the likelihood of arrest conditional on crime characteristics via GEEs. When describing results as statistically significant, we apply a significance level of 0.01.

6.1. *The necessity of accounting for covariate shift.* In their approach, Blumstein and Cohen (1979) assume that all incidents in NIBRS are equally likely to be reported to law enforcement once we condition on the offender's race and the crime type. Their method naively estimates the number of actual crimes underlying the reported crimes in NIBRS by applying the (fixed) ratio of actual crimes to reported crimes in NCVS. However, because reporting rates may vary according to crime characteristics, and because of potential covariate shift between NCVS and NIBRS in the distribution of crime characteristics, these estimates may be significantly biased.

Both phenomena are observed in the data. As previously discussed, covariate shift arises in part due to the different geographical coverage of NIBRS and NCVS. For example, 35% of the offenses of simple assault known to law enforcement in 2006–2015 NCVS data occur in the southern regions of the US, compared to 60% of those in the NIBRS. As another example, in about 60% of the offenses of aggravated assault recorded in NCVS the offender is known to the victim. By contrast, this occurs in 85% of the cases in NIBRS data. The coefficient estimates produced by the logistic regression fitted on NCVS data reveal that the likelihood of reporting varies across most of the crime characteristics considered by our analysis, often quite substantially (see the Appendix). Unlike Blumstein and Cohen (1979), we account for these variations in our analysis.

6.2. *Racial disparities in crime reporting.* In the available NIBRS data, 59% of all offenders are white. We estimate that the NIBRS data capture 44% (standard error=5%) and 48% (5%) of all violent offenses committed by white and black offenders, respectively. Equivalently, slightly more than half of the crimes that occur in the jurisdictions covered by NIBRS are not reported to law enforcement for both racial groups. The lower reporting rates for white offenders relative to black offenders indicates a (not statistically significant) marginal overrepresentation of black offenders in the data recorded by police agencies compared to their representation in the larger population of offenders. More specifically, we estimate that 61% of all crimes that occur are committed by white offenders.

Table 1 shows the breakdown of the rates of police notification by offense types and offenders' racial groups. Sex offenses are the least likely to be reported to police, with only one

TABLE 1
Summary Statistics (2006–2015 NIBRS Data): Estimating Unreported Crimes Using the Likelihood
of Police Notification Computed from NCVS Data.

| Variable | sex offense | robbery | aggravated assault | simple assault |
|---|---|---|---|---|
| % police notification | | | | |
| • black offenders | 21% (19%) | 55% (6%) | 63% (4%) | 47% (5%) |
| • white offenders | 19% (17%) | 51% (7%) | 60% (4%) | 45% (4%) |
| % arrests (reported crimes) | | | | |
| • black offenders | 22% (<1%) | 17% (<1%) | 40% (<1%) | 38% (<1%) |
| • white offenders | 24% (<1%) | 33% (<1%) | 56% (<1%) | 50% (<1%) |
| % arrests (all crimes) | | | | |
| • black offenders | 5% (20%) | 9% (3%) | 25% (3%) | 18% (4%) |
| • white offenders | 5% (19%) | 17% (5%) | 34% (4%) | 22% (4%) |

Notes: Standard errors are reported within parentheses. The summary statistics are computed on incidents with only one offender. The likelihood of police notification is estimated on NCVS data via logistic regression with survey weights using the methodology described in Section 5.

in four incidents being reported compared to one in two for the other crime types. Since these crimes are unlikely to be reported to law enforcement, the estimates suffer from large variance. We find that crimes with black offenders are associated with higher rates of reporting than those with white offenders across all offense types. However, the regression model fitted on NCVS data reveals that there is only a weak and not statistically significant association between reporting and the offender's racial group once other crime characteristics are taken into account.

It is possible that the logistic regression model fitted on NCVS data is misspecified. Thus, we conduct an analysis of the reporting rates where the likelihood of police notification is estimated via the SuperLearner. By comparing the estimates produced by the two models on NIBRS data, we find that the estimates are close for the crimes of robbery and assaults. The estimates of arrest rates obtained using the two models are virtually identical for robbery and assault, even when we condition on the offender's racial group (the SuperLearner estimates are within a 2% difference from those in Table 1). For sex offenses, on the other hand, the logistic regression tends to underestimate the likelihood of reporting compared to the SuperLearner. This underestimation is quite substantial. The resulting rates of police notification based on the SuperLearner are larger than those produced by the logistic regression.

6.3. *Racial disparities in arrest rates.* We now turn to the estimation of arrest rates. Overall, 49% (standard error<1%) of the offenses known to law enforcement involving white offenders resulted in arrest, compared to 37% (<1%) of those involving black offenders. Table 1 reveals that arrest rates are similar across racial groups for sex offenses, while robbery and assault incidents white offenders result in arrest considerably more often than those with black offenders. Past works on NIBRS have reached qualitatively similar conclusions (D'Alessio and Stolzenberg, 2003; Lantz and Wenger, 2019). Despite the lower crime reporting rates, crimes with white offenders remain more likely to result in arrests than those with black offenders once unreported crimes are accounted for. Overall, arrest rates for crimes are 21% (7%) for white offenders and 17% (5%) for black offenders. Table 1 shows that arrest rates are higher for white offenders in case of assaults and robbery, and are comparable across racial groups in case sex offenses. As in the observed police data, these rates greatly vary across offense types: Arrests occur in about one in twenty sex offenses and one in five simple assaults. The sensitivity analysis via the SuperLearner produces quantitatively similar results except for sex offenses.

TABLE 2
*Regression Results (2006–2015 NIBRS Data): Assessing Racial Differences in Arrest Likelihood for Single-Offender Incidents, Using 2003–2020 NCVS Data for Police Notification Estimates.*

| Variable | Sex offense | Robbery | Aggravated assault | Simple assault |
|---|---|---|---|---|
| Intercept | 0.05 (0.01)*** | 0.48 (0.08)*** | 0.46 (0.06)*** | 0.20 (0.02)*** |
| Age of offender | 1.01 (0.00)*** | 1.01 (0.00)*** | 1.01 (0.00)*** | 1.00 (0.00) |
| Off. is male | 0.99 (0.08) | 0.87 (0.05)* | 0.89 (0.04)* | 0.90 (0.05). |
| Off. is white | 1.00 (0.07) | 1.23 (0.06)*** | 1.03 (0.05) | 1.01 (0.06) |
| Age of victim | 1.00 (0.00) | 1.01 (0.00)*** | 1.01 (0.00)*** | 1.01 (0.00)*** |
| Victim is male | 0.84 (0.06)** | 0.87 (0.04)** | 0.96 (0.04) | 0.93 (0.05) |
| Victim is white | 0.83 (0.07)* | 1.00 (0.06) | 1.05 (0.06) | 1.02 (0.07) |
| Off. is acquaintance | 0.87 (0.06)* | 1.34 (0.07)*** | 0.98 (0.05) | 0.67 (0.04)*** |
| Off. is family member | 1.20 (0.14) | 2.12 (0.19)*** | 1.54 (0.12)*** | 1.24 (0.12)* |
| Off. is intimate partner | 1.33 (0.13)** | 2.30 (0.18)*** | 1.92 (0.13)*** | 1.58 (0.13)*** |
| Minor injury | 1.70 (0.11)*** | 1.24 (0.06)*** | 1.45 (0.07)*** | 1.85 (0.10)*** |
| Serious injury | 2.74 (0.28)*** | 1.99 (0.12)*** | 2.19 (0.14)*** | |
| During day | 0.90 (0.05). | 1.24 (0.05)*** | 0.95 (0.04) | 0.93 (0.04) |
| Private location | 1.33 (0.10)*** | 1.08 (0.05). | 1.40 (0.07)*** | 1.34 (0.08)*** |
| Firearm present | 1.03 (0.16) | 0.97 (0.10) | 0.99 (0.10) | |
| Other weapon present | 0.90 (0.13) | 0.94 (0.10) | 0.90 (0.08) | 0.85 (0.10) |
| Multiple offenses | 1.90 (0.03)*** | 1.58 (0.04)*** | 1.14 (0.01)*** | 1.15 (0.01)*** |
| Offense only attempted | 0.86 (0.12) | 0.98 (0.09) | | |
| MSA, central city | 0.68 (0.05)*** | 0.89 (0.05). | 0.94 (0.05) | 0.92 (0.06) |
| MSA, not central city | 0.84 (0.07)* | 1.02 (0.06) | 1.10 (0.06) | 1.02 (0.07) |
| Nb. of officers per 1000 capita (ORI) | 1.00 (0.00)*** | 0.99 (0.00)*** | 0.99 (0.00)*** | 1.00 (0.00)*** |
| Log population served (ORI) | 0.96 (0.00)*** | 0.83 (0.01)*** | 0.88 (0.00)*** | 0.91 (0.00)*** |

Significance codes: $p < 0.001$ '\*\*\*', $p < 0.01$ '\*\*', $p < 0.05$ '\*', $p < 0.1$ '.'.

Notes: The table shows the odds ratios of the logistic regression coefficients for $q$, the likelihood of arrest that accounts for unreported crimes. The model is fitted on 2006–2015 NIBRS data and uses the estimates of $\pi$ obtained from 2003–2020 NCVS data. Standard errors are reported inside parentheses. Significance codes correspond to the p-values ($p$) of Wald tests to assess the statistical significance of the odds ratios. Year- and state-level fixed effects are included in the regression model but are omitted from the table. "ORI" stands for "originating agency identifier", a regressor whose value is specific to that law enforcement agency.

6.4. *Racial disparities in the likelihood of arrest accounting for crime characteristics.* We estimate the likelihood of arrest conditional on crimes characteristics via the two-step logistic regression detailed in Section 4. The resulting odds ratios of the coefficient estimates are reported in Table 2. In case of robbery offenses, we find that there is a positive and statistically significant association between whether the offender is white and the likelihood that the incident results in arrest. Provided that our model is correctly specified, these results would indicate that white offenders are more likely to be arrested for robbery than black offenders, ceteris paribus. The estimates of this coefficient for the other types of crimes are close to zero and not statistically significant. Thus, the estimated disparities disappear once we account for crime characteristics other than the offender's race.

One outstanding concern is that our logistic regression model estimated on NCVS data may not accurately capture the location-specific patterns in crime reporting existing in the data (e.g., due to omitted variable bias or modeling misspecification). For example, by studying restricted-use NCVS data Baumer (2002) and Xie and Lauritsen (2012) report significant variations in crime reporting rates across neighborhoods. Although the available data do not allow us to analyze reporting rates at the level of the individual law enforcement agencies, we can still assess whether regional patterns are accounted for by employing a flexible modeling approach. Thus, we run the two-step regression analysis using the estimates of the SuperLearner in place of those from the logistic regression on NCVS. The estimates of the offender's race coefficients produced by this approach are close to those presented in Table 2.

We also analyze how the odds ratios of the offender's race coefficient estimates vary under various configurations of the covariates distributions through the focal slope model diagnostics described in Section 5. The results of the diagnostics are reported in the Appendix. We find that both the magnitude and the direction of the coefficients estimates vary with the characteristics of the crimes. The most notable pattern is the change in the association between the likelihood of arrest and the offender's racial group when either only black or white victims are considered. This suggests an interaction between the two covariates. For example, in case of assaults we observe that the estimates association between the offender being white and the likelihood of arrest is close to zero when the victims being considered are white individuals, but it is large and positive in case of black victims. This suggests that, ceteris paribus, white offenders may be more likely to be arrested than black offenders only when they commit interracial crimes. For sex offenses, however, we find that the association is negative in case of crimes with black victims and close to zero otherwise.

The estimates in Table 2 can be compared with those obtained from a logistic regression model fitted directly on NIBRS data without accounting for unreported crimes (see the Appendix). The association between the offender being white and arrests estimated by the models without adjustments for unreported crimes is stronger (and positive) compared to the estimates in Table 2 in case of robberies and assaults. We focus on two other examples of differences between coefficient estimates that stand out. First, the regression model that accounts for unreported crimes estimates a stronger positive association between the victim being injured (vs. no injury) and the occurrence of an arrest. This pattern could be explained by the fact that incidents without injuries are less likely to be reported to law enforcement (see the Appendix). Second, the logistic regression without adjustments estimates a negative and strong association between the presence of a firearm (vs. no weapon) and arrests. This association disappears once unreported crimes are accounted for, again potentially because incidents with firearms tend to be more likely to be reported.

6.5. *Racial disparities in arrests for incidents with multiple offenders.* The estimates of the crime reporting rates for incidents involving more than one offender are similar to those for crimes with individual offenders in Table 1. However, we find that arrest rates computed solely on police-recorded data for these incidents are substantially lower in case of aggravated assaults (by more than 10%), followed by robbery and simple assaults (within a 5% difference). By contrast, arrest rates for sex offenses with multiple offenders are marginally higher than those of crimes with individual offenders. Arrest rates shrink proportionally within racial groups once unreported crimes are taken into account. We continue to observe that white offenders are arrested more often than black offenders across all crime types. The only exception is robbery for which the reduction is limited to white offenders but it is not large enough to reverse the sign of the disparity.

We first fit the two-step regression model using GEEs on only incidents with multiple offenders. The model estimates that, conditionally on other crime characteristics, white offenders face a higher likelihood of arrest than black offenders across all offense types. We next fit the same model specification on incidents with both single and multiple offenders. In doing so, we need to keep in mind that only about one in ten incidents of assault and sex offenses are committed by multiple offenders, and these incidents generally have few offenders. An exception is represented by robbery for which half of the incidents involve multiple offenders. We find that white offenders are associated with a higher likelihood of arrest in case of robbery (estimate is 0.2 with standard error equal to 0.04), while the other estimated associations are virtually zero (full results in Table 5).

**7. Limitations.** Our empirical analysis relies on a series of assumptions about the modeling and data that may not hold in reality. One key limitation of the modeling is that the assumed independence across incidents may be violated. For example, when the same offender is part of multiple separate crime incidents, arrest outcomes become correlated. Given identifying offender-level information, we could, in principle, correct the variance estimates to account for this dependence (Andrews and Monahan, 1992; White, 2014), but such data is not available.

Another limitation of our modeling approach is the potential misspecification of the regression function used in estimating the likelihood of police notification. Through the sensitivity analysis in Section 6, we have shown that employing a more flexible classifier on NCVS data yields results that are similar to those of the logistic regression for most offense types. Leveraging the model diagnostics, we have shown that the logistic regression model fitted on NIBRS was misspecified. Consequently, the coefficient estimates require careful interpretation; see Buja et al. (2019b) and Berk et al. (2019) for detailed treatments of this topic, and Fogliato et al. (2021) for a discussion of the limitations of similar approaches.

Certain variation in reporting rates may also fail to be captured specifically due to omitted variable bias. For example, the NCVS data that are used in our analysis contain little information on the geographical location. Obtaining and incorporating this information may influence the results (Baumer, 2002; Xie and Lauritsen, 2012; Xie and Baumer, 2019b). Although analyses of the restricted-use NCVS data would overcome some of these issues, relevant pieces of information, such as the specific location of the victimization, may simply be missing from the data (Cernat et al., 2021).

Even more importantly, our analysis suffers from limitations related to the nature of the data. These limitations are not unique to our study; they have been discussed in a plethora of criminological works. Firstly, the recorded data may be of poor quality. With respect to survey data, measurement errors arising from sampling design, data collection, victims' recollection of the events and untruthful reporting affect the quality of the data. What victims report in the survey may not always coincide with the same information that is recorded in police data.

Information in NIBRS may not always accurately reflect the characteristics of the crime incident. In this work, for example, we have observed that NCVS respondents were far more likely to report serious injuries in case of sex offenses than what was recorded in NIBRS data. This pattern is unlikely to be explained solely by differences in the underlying populations. Overall, police data can be seen as an artefact of a manipulation process (Richardson, Schultz and Crawford, 2019). It is also possible that instances of wrongful arrests, which we do not consider in the analysis, may be present in the data (Loeffler, Hyatt and Ridgeway, 2019).

In addition to issues of data quality, the data are missing certain information that we hypothesize being relevant to our analysis. For example, we included Hispanics in the analysis because, as ethnicity information is not always recorded (and when recorded it can be imprecise), this population could not be entirely excluded from the sample. However, there is evidence that this ethnic group may be characterized by unique offending and reporting behaviors (Steffensmeier et al., 2011; Roberts and Lyons, 2011; Rennison, 2010).

One further limitation of our analysis concerns the matching of offense categories between the NCVS and the NIBRS, that do not perfectly map. However, even if the definitions were to fully overlap, the type of offense that is reported by the victim may not correspond to the coding of the same offense done by law enforcement. This potential issue may affect mainly simple assaults, which represent the least serious type of crime considered in this analysis. We also do not consider incidents where the victim does not personally see the offender, which represent a minimal share of all incidents reported by NCVS respondents. Thus, together with the fact that not all reported crime may be recorded, this implies that our estimates of the arrest rates represent upper bounds of the true quantities.

**8. Discussion.** In this work, we have proposed estimators of the rates of police notification and of arrest for nonfatal violent crime on NIBRS that leverage data of unreported crimes from NCVS. These estimators are consistent and asymptotically normal under some assumptions. Our empirical investigation of racial disparities revealed that incidents are marginally more likely to be reported to the police when the offender is black. However, in cases of assaults and robbery, crimes with black offenders are generally less likely to result in arrests. These differences are small after accounting for crime characteristics. Additionally, the model diagnostics showed that the direction of these disparities varies with crime characteristics.

We envision three directions in which the proposed methodology can be further developed. First, we could employ nonparametric methods in place of the logistic regression model. In this work, we obtained asymptotic normality for a two-step estimation approach where logistic regression was used in both steps. Nonparametric approaches that yield similar convergence rates could be applied in the first step. For example, the method of kernels introduced in Racine and Li (2004) can handle both categorical and continuous data. Although our empirical analysis did not uncover significant differences in the estimates of the likelihood of police notification produced by parametric and nonparametric models, the latter is more flexible and thus may be more suitable in certain applications.

Secondly, mixed effects logistic regression models could be employed for the estimation of the likelihood of arrest. This represents a modeling approach often used in the social sciences. In this work, we have employed a model that does not account for city- or agency-level effects, which may drive many of the disparities, e.g., see the results of Fogliato et al. (2021). Third, we assume covariate shift between NCVS and NIBRS. Future work could use a reweighted loss to adjust for the shift in the two datasets.

Our study opens multiple avenues of research in the criminology field as well. Despite a longstanding interest in the "dark figure of crime" (Skogan, 1977), how to accurately characterize this figure remains challenging and not well understood. Our methodology represents one way through which it can be described and its magnitude be assessed. It would be interesting to compare results obtained through our methodology with those from the simulation-based approach proposed by Buil-Gil, Moretti and Langton (2021). Future work may also leverage information about the socioeconomic status of the victim, which is available in NCVS data, and of the characteristics of the population in the police agency, which can be obtained from auxiliary data sources and merged with NIBRS data. These aspects were not considered in our work.

Similar to past studies on NCVS and NIBRS, the results described in this paper build on several assumptions. Some of these assumptions may be violated. We hope that, over time, police records will become more accurate and comprehensive, and that detailed information about incidents will be made available through NIBRS, allowing for improved analyses.

APPENDIX A: ADDITIONAL RESULTS

This section contains additional results. Table 3 shows the odds ratios of the coefficients estimates of the logistic regression model run on NCVS data. Table 4 presents the results of the regression analysis targeting the likelihood of arrest for crimes known to law enforcement with individual offenders. Figure 3 shows the predicted likelihood of crime reporting $\pi$ produced by logistic regression and SuperLearner for each observation in NCVS data. Lastly, Figure 4 shows the focal slope model diagnostics.

TABLE 3
*Regression Results (2003–2020 NCVS Data): Estimating Police Notification Likelihood for Single-Offender Incidents Using Logistic Regression with Survey Weights.*

| Variable | Odds ratio estimate |
|---|---|
| Age of off. 12-14 | 1.36 (0.43) |
| Age of off. 15-17 | 2.17 (0.58)** |
| Age of off. 18-20 | 2.00 (0.62)* |
| Age of off. 21-29 | 2.96 (0.84)*** |
| Age of off. 30+ | 2.66 (0.72)*** |
| Off. is male | 0.87 (0.08) |
| Off. is white | 0.95 (0.08) |
| Age of victim | 1.01 (0.00)*** |
| Victim is male | 0.90 (0.07) |
| Victim is white | 0.85 (0.09) |
| Off. is acquaintance | 0.71 (0.06)*** |
| Off. is family member | 0.94 (0.13) |
| Off. is intimate partner | 0.95 (0.12) |
| Minor injury | 1.41 (0.11)*** |
| Serious injury | 2.83 (0.35)*** |
| During day | 0.96 (0.07) |
| Private location | 1.37 (0.12)*** |
| Firearm present | 1.39 (0.27). |
| Other weapon present | 0.81 (0.14) |
| Offense only attempted | 0.82 (0.14) |
| MSA, central city | 0.90 (0.09) |
| MSA, not central city | 1.03 (0.10) |
| Crime is robbery | 0.85 (0.14) |
| Crime is sex offense | 0.23 (0.05)*** |
| Crime is simple assault | 0.57 (0.10)** |

Significance codes: $p < 0.001$ '***', $p < 0.01$ '**', $p < 0.05$ '*', $p < 0.1$ '.'.

TABLE 4

*Regression Results: Assessing Racial Disparities in Arrest Likelihood Based on Known Incidents to Police Agencies, Using Crime Characteristics $\alpha(X)$.*

| Variable | Sex offense | Robbery | Aggravated assault | Simple assault |
|---|---|---|---|---|
| Intercept | 0.42 (0.03)*** | 1.53 (0.16)*** | 1.55 (0.05)*** | 0.93 (0.01)*** |
| Age of offender | 1.01 (0.00)*** | 1.01 (0.00)*** | 1.00 (0.00)*** | 0.99 (0.00)*** |
| Off. is male | 1.15 (0.04)*** | 0.93 (0.03)* | 0.94 (0.01)*** | 0.98 (0.00)*** |
| Off. is white | 1.04 (0.02)* | 1.33 (0.03)*** | 1.09 (0.01)*** | 1.05 (0.00)*** |
| Age of victim | 0.98 (0.00)*** | 1.00 (0.00)* | 1.00 (0.00)*** | 1.01 (0.00)*** |
| Victim is male | 0.90 (0.02)*** | 0.93 (0.02)*** | 1.01 (0.01) | 0.99 (0.00)*** |
| Victim is white | 0.94 (0.02)*** | 1.12 (0.02)*** | 1.21 (0.01)*** | 1.22 (0.01)*** |
| Off. is acquaintance | 1.24 (0.02)*** | 1.80 (0.04)*** | 1.25 (0.01)*** | 0.85 (0.00)*** |
| Off. is family member | 1.41 (0.03)*** | 2.67 (0.15)*** | 1.98 (0.02)*** | 1.58 (0.01)*** |
| Off. is intimate partner | 1.56 (0.03)*** | 2.91 (0.11)*** | 2.49 (0.02)*** | 2.04 (0.01)*** |
| Minor injury | 1.35 (0.02)*** | 1.02 (0.02) | 1.27 (0.01)*** | 1.76 (0.00)*** |
| Serious injury | 1.31 (0.03)*** | 1.26 (0.04)*** | 1.38 (0.01)*** | |
| During day | 0.94 (0.01)*** | 1.35 (0.02)*** | 0.98 (0.01)** | 0.98 (0.00)*** |
| Private location | 0.98 (0.01)* | 0.89 (0.02)*** | 1.23 (0.01)*** | 1.07 (0.00)*** |
| Firearm present | 0.74 (0.04)*** | 0.79 (0.02)*** | 0.82 (0.01)*** | |
| Other weapon present | 1.08 (0.03)*** | 1.07 (0.02)** | 1.04 (0.01)*** | 0.99 (0.01)* |
| Multiple offenses | 2.37 (0.05)*** | 1.82 (0.06)*** | 1.20 (0.01)*** | 1.22 (0.01)*** |
| Offense only attempted | 1.01 (0.03) | 1.12 (0.03)*** | | |
| MSA, central city | 0.68 (0.01)*** | 0.90 (0.03)** | 0.97 (0.01)** | 0.96 (0.00)*** |
| MSA, not central city | 0.77 (0.01)*** | 0.99 (0.03) | 1.11 (0.01)*** | 0.99 (0.00)** |
| Nb. of officers per 1000 capita (ORI) | 1.00 (0.00)*** | 0.99 (0.00)*** | 0.99 (0.00)*** | 1.00 (0.00)*** |
| Log population served (ORI) | 0.94 (0.00)*** | 0.80 (0.01)*** | 0.84 (0.00)*** | 0.86 (0.00)*** |

Significance codes: $p < 0.001$ '***', $p < 0.01$ '**', $p < 0.05$ '*', $p < 0.1$ '.'.

Notes: The table shows the odds ratios of the logistic regression coefficients for $\alpha$, the likelihood of arrest for incidents known to police agencies, fitted on the NIBRS data considered in the analysis. Standard errors are reported inside parentheses. Significance codes correspond to the p-values ($p$) of Wald tests to assess the statistical significance of the odds ratios. Year- and state-level fixed effects are included in the regression model but are omitted from the table.

TABLE 5

*Regression Results (NIBRS Data): Assessing Racial Disparities in Arrest Likelihood for Incidents with One or More Offenders, Using NCVS Data for Police Notification Estimates.*

| Variable | Sex offense | Robbery | Aggravated assault | Simple assault |
|---|---|---|---|---|
| Intercept | 0.07 (0.01)*** | 0.92 (0.12) | 0.40 (0.05)*** | 0.18 (0.02)*** |
| Age of offender | 1.01 (0.00)*** | 1.01 (0.00)*** | 1.01 (0.00)*** | 1.00 (0.00)* |
| Off. is male | 0.88 (0.06). | 0.82 (0.03)*** | 0.91 (0.04)* | 1.00 (0.06) |
| Off. is white | 1.00 (0.06) | 1.18 (0.04)*** | 1.04 (0.04) | 1.00 (0.05) |
| Age of victim | 0.99 (0.00)** | 1.00 (0.00)*** | 1.01 (0.00)*** | 1.01 (0.00)*** |
| Victim is male | 0.83 (0.06)** | 0.77 (0.03)*** | 0.76 (0.03)*** | 0.77 (0.04)*** |
| Victim is white | 0.78 (0.06)** | 0.93 (0.05) | 0.96 (0.05) | 0.95 (0.06) |
| Off. is known | 0.90 (0.06) | 1.01 (0.04) | 1.12 (0.05)** | 0.97 (0.05) |
| Minor injury | 1.64 (0.11)*** | 1.21 (0.05)*** | 1.50 (0.07)*** | 1.89 (0.10)*** |
| Serious injury | 2.83 (0.29)*** | 1.95 (0.10)*** | 2.30 (0.13)*** | |
| During day | 0.95 (0.05) | 1.30 (0.05)*** | 1.00 (0.04) | 0.94 (0.04) |
| Private location | 1.47 (0.10)*** | 1.32 (0.05)*** | 1.70 (0.07)*** | 1.73 (0.09)*** |
| Firearm present | 1.05 (0.14) | 1.08 (0.09) | 0.97 (0.08) | |
| Other weapon present | 1.05 (0.13) | 1.03 (0.09) | 1.00 (0.08) | 0.96 (0.10) |
| Multiple offenses | 1.84 (0.05)*** | 1.60 (0.03)*** | 1.12 (0.01)*** | 1.14 (0.01)*** |
| Offense only attempted | 0.84 (0.11) | 0.91 (0.08) | | |
| MSA, central city | 0.60 (0.05)*** | 0.76 (0.04)*** | 0.85 (0.04)** | 0.83 (0.05)** |
| MSA, not central city | 0.82 (0.07)* | 0.94 (0.05) | 1.08 (0.06) | 1.01 (0.07) |
| Nb. of officers per 1000 capita (ORI) | 1.00 (0.00)*** | 0.99 (0.00)*** | 0.99 (0.00)*** | 1.00 (0.00)*** |
| Log population served (ORI) | 0.95 (0.00)*** | 0.81 (0.00)*** | 0.88 (0.00)*** | 0.91 (0.00)*** |

Significance codes: $p < 0.001$ '***', $p < 0.01$ '**', $p < 0.05$ '*', $p < 0.1$ '.'.

Notes: The table shows the odds ratios of the regression coefficients for $q$, the likelihood of arrest that accounts for unreported crimes, estimated via generalized estimating equations (GEEs). The model is fitted on NIBRS data and uses the estimates of the likelihood of crime reporting $\pi$ obtained from NCVS data. Standard errors are reported inside parentheses. Significance codes correspond to the p-values ($p$) of Wald tests to assess the statistical significance of the odds ratios. Year- and state-level fixed effects are included in the regression model but are omitted from the table.
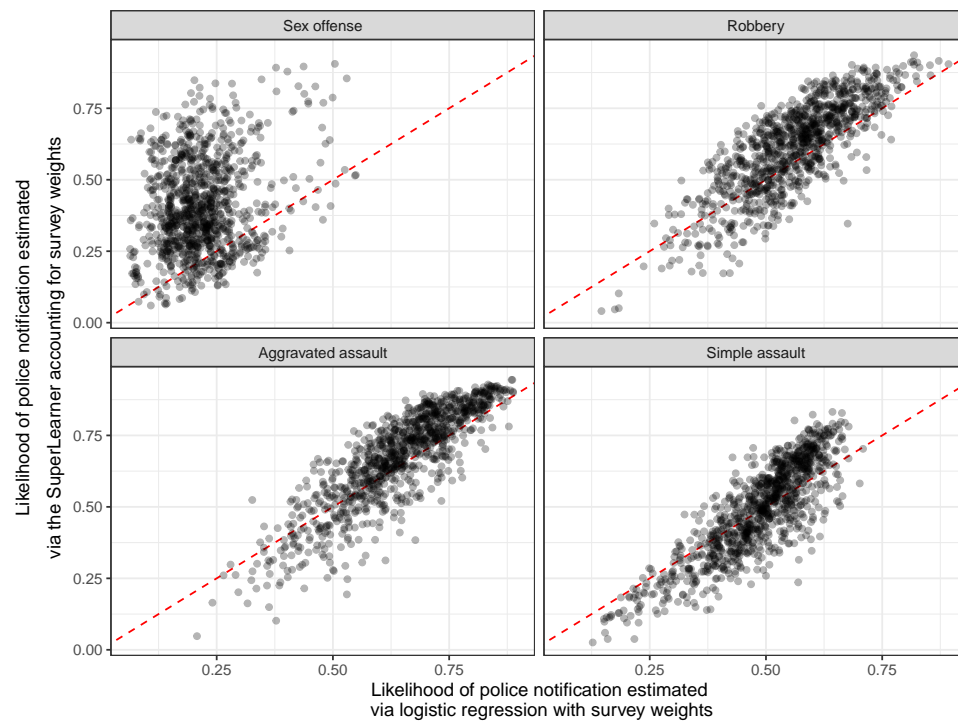
FIG 3. *Estimates of the likelihood of police notification π for observations in 2006–2015 NIBRS data produced by the logistic regression (horizontal axis) and by the SuperLearner (vertical axis) fitted on 2003–2020 NCVS data. For visualization purposes, we show the estimates relative to 1000 randomly sampled observations for each crime type.*
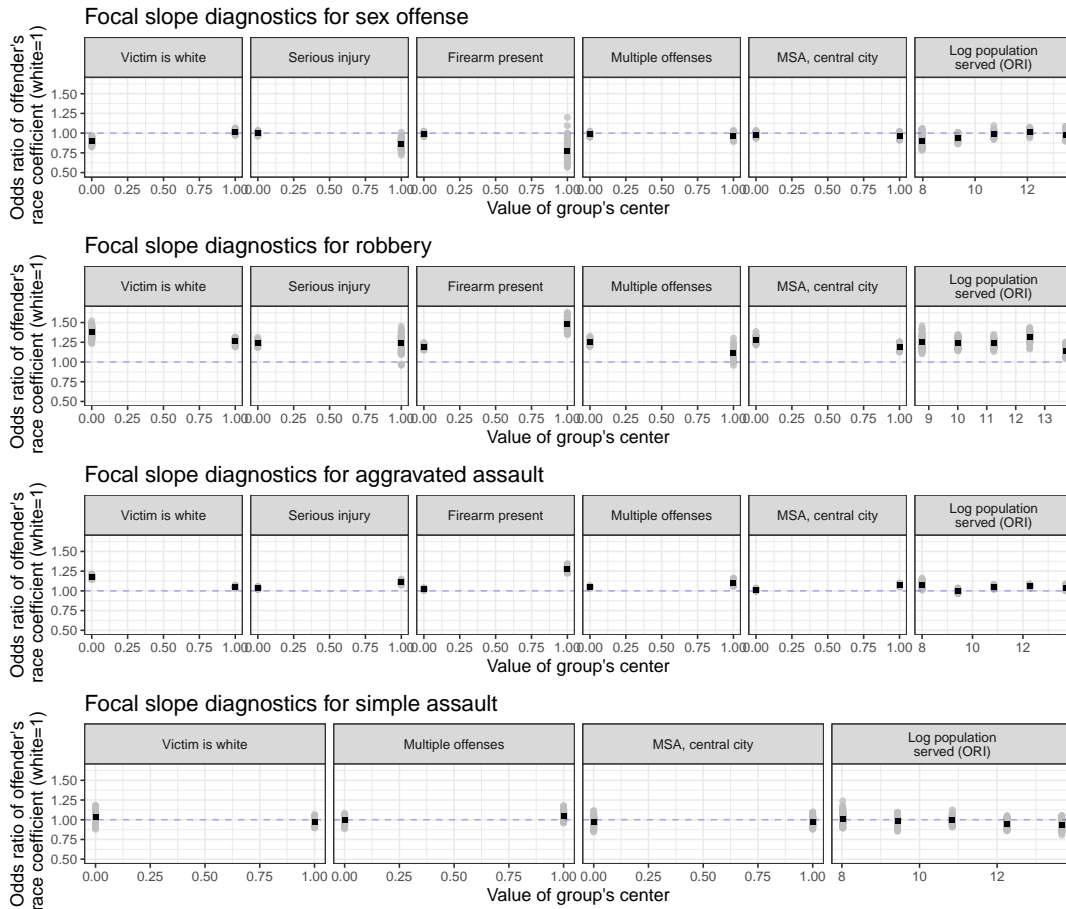
FIG 4. "Focal slope" model diagnostic for the logistic regression model for $q$, the likelihood of arrest that accounts for unreported crimes, on 2006–2015 NIBRS data whose odds ratios of the coefficients estimates are presented in Table 2. Only incidents with one offender are considered. Methodological details are described in Section 5. The grey points correspond to the coefficient estimates relative to the offender's race (white=1) obtained by fitting the logistic regression model on each of 100 bootstrapped datasets, for each variable (panel's title) and variable's grid value (value on the grid, horizontal axis). The black dots correspond to the means of such estimates. We observe that the size and sign of the values of the black dots vary across the range of the regressors. This suggests the presence of interactions between race and the regressors, which in turn indicates that our modeling approach is misspecified.

## APPENDIX B: DETAILS AND PROOFS

This section contains the proofs of the results presented in Section 4. In Section B.1, we present the proof relative to the consistency and asymptotic normality of the coefficient estimates for the logistic regression parameters obtained on survey data (Proposition 1). Then, in Section B.2 we present the asymptotic properties of the estimators of the total number of offenses $N$, expected rate of police notification $\pi^*$, and expected rate of arrest $q^*$ (Lemma 1, Propositions 2, 3, and 4). Lastly, we describe the results for estimation via the two-step logistic regression (Propositions 5 and 6).

**B.1. Estimation on NCVS.** We provide some additional details on the framework presented in Section 4.1 before turning to the proof of Proposition 1. As a reminder, our aim is to make inference on superpopulation parameters. This differs from the finite population framework, for which logistic regression parameter estimation has been studied by Binder (1983). In the following, the subscripts $P^v$ and $\psi$ in the probability $\mathbb{P}$ and expectation $\mathbb{E}$ operators denote superpopulation and sampling design, respectively.

Formally, let the superpopulation target parameter $\gamma_0 \in \mathrm{Int}(\Gamma)$ be defined by the following moment condition

$$\mathbb{E}_{P^v}[h^v(R^v, Z^v; \gamma)] = 0,$$

where $h^v(R^v, Z^v; \gamma) := (R^v - \pi^v(Z^v; \gamma))Z^v$. The parameter $\tilde{\gamma} \in \mathrm{Int}(\Gamma)$ and the design-based estimator $\hat{\gamma} \in \Gamma$ are the solutions to (Lumley and Scott, 2017),

$$(3) \qquad \sum_{i=1}^{N^v} h^v(R_i^v, Z_i^v; \gamma) = 0, \text{ and}$$

$$(4) \qquad \sum_{i=1}^{N^v} w_i I_i h^v(R_i^v, Z_i^v; \gamma) = 0$$

respectively. The estimating equation (3) is unbiased for $\gamma_0$. Conditionally on the finite population $V^{N^v}$, equation (4) is unbiased for $\tilde{\gamma}$ provided that $\mathbb{E}_\psi[I_i w_i] = 1$ for $i = 1, \ldots, N^v$. Since we only have access to the observations for which $I_i = 1$, our estimation will be based on the estimating equation (4). In the presence of endogenous or informative sampling, the estimate $\hat{\gamma}$ obtained by solving (4) may not coincide with the one we would obtain by solving the unweighted estimating equation (Solon, Haider and Wooldridge, 2015).

In order to establish the asymptotic properties of the estimator $\hat{\gamma}$ obtained by solving equation (4) on the sample $V^{N^v}$, we assume that the observations we have access to are sampled from a finite number of strata with known size. Thus, for each stratum, our survey data represent a sample of the finite population belonging to that stratum, which in turn represents an i.i.d. sample of the superpopulation distribution specific to that stratum. The following proposition borrows the setup from Theorem 1.3.9 in Fuller (2011) and leverages the results of Rubin-Bleuer and Kratina (2005).

PROPOSITION 1. Consider an increasing sequence of finite populations where the $N^v$-th population has size $N^v$ and consists of $H \in \mathbb{Z}_+$ strata. The $h$-th stratum is formed by the $N_h^v$ observations $\mathcal{F}_{N^v h} = \{(Z_{N^v hi}^v, R_{N^v hi}^v)\}_{i=1}^{N_h^v}$ which represent an i.i.d. sample of $(Z_h^v, R_h^v) \sim P_h^v$, for $h = 1, \ldots, H$, where $P_h^v$ is the distribution of the superpopulation of the specific stratum. Assume that $\|Z_h^v\|_\infty < M$ for some $M > 0$ and $h = 1, \ldots, H$. For the $h$-th stratum, we have access to a sample of observations that are drawn from $\mathcal{F}_{N^v h}$ according to some sampling design $\psi_{N^v h}$ and let $I_{N^v hi} = 1$ if the $i$-th observation is selected, and $I_{N^v hi} = 0$

otherwise. Let $\{w_{N^v hi}\}_{i=1}^{N_h^v}$ indicate the set of weights associated with the sample in the $h$-th stratum where $w_{N^v hi} := \mathbb{E}_{\psi_{N^v}}[I_{N^v hi}]^{-1}$, and assume that $\max_{h,i} w_{N^v hi} < K$ for some $K > 0$. We denote with $n_{N^v h}^v$ the (expected or fixed) sample size of the $h$-th stratum, with $n_{N^v}^v := \sum_{h=1}^H n_{N^v h}^v$ the size of the entire survey sample, and with $\lambda := \lim_{N^v \to \infty} n_{N^v}^v / N^v$ the limit of the size of the surveyed population compared to the entire finite population. Consider a sequence of stratified samples that is selected such that $N_h^v \to \infty$, $n_{N^v h}^v \to \infty$, and $\lim_{N^v \to \infty} N_h^v / N^v = \lim_{N^v \to \infty} n_{N^v h}^v / n_{N^v}^v = \beta_h \in (0,1]$, for $h = 1, \ldots, H$. The parameters $\gamma_0$ and $\tilde{\gamma}_{N^v}$, and the estimator $\hat{\gamma}_{N^v}$, with $\gamma_0, \tilde{\gamma}_{N^v} \in \text{Int}(\Gamma)$ and $\hat{\gamma}_{N^v} \in \Gamma$, solve respectively

$$\sum_{h=1}^H \beta_h \mathbb{E}_{P_h^v}[\mathbb{E}_{\psi_{N^v h}}[h^v(R_{N^v hi}^v, Z_{N^v hi}^v; \gamma)]] = 0,$$

$$G_{N^v}^v(\gamma) := \frac{1}{N^v} \sum_{h=1}^H \sum_{i=1}^{N_h^v} h^v(R_{N^v hi}^v, Z_{N^v hi}^v; \gamma) = 0,$$

$$\hat{G}_{N^v}^v(\gamma) := \frac{1}{n_{N^v}^v} \sum_{h=1}^H \sum_{i=1}^{N_h^v} w_{N^v hi} I_{N^v hi} h^v(R_{N^v hi}^v, Z_{N^v hi}^v; \gamma) = 0.$$

Assume that, conditionally on the finite population,

$$\sqrt{n_{N^v}^v} \hat{G}_{N^v}^v(\gamma_{N^v}) \xrightarrow{d} \mathcal{N}\left(0, \sum_{h=1}^H \beta_h \Xi_h^f\right)$$

as $n_{N^v}^v \to \infty$ and, in addition, for $\gamma \in \Gamma$,

$$(5) \qquad \lim_{N^v \to \infty} \frac{1}{n_{N^v}^v} \sum_{h=1}^H \sum_{i=1}^{N_h} I_{N^v hi} w_{N^v hi} \nabla_\gamma h^v(Z_{N^v hi}, R_{N^v hi}; \gamma) = J^v(\gamma)$$

where the positive definite covariance matrices $\Xi_h^f$, for $h = 1, \ldots, H$, and $J^v(\gamma)$ are non-stochastic in the population, and $J^v := J^v(\gamma_0) = \lim_{N_h^v \to \infty} \nabla G_{N^v}^v(\gamma_0)$. Then

$$(\Sigma^v)^{-1/2} \sqrt{n_{N^v}^v} (\hat{\gamma}_{N^v} - \gamma_0) \xrightarrow{d} \mathcal{N}(0, I_d)$$

as $n_{N^v}^v \to \infty$ where $\Sigma^v := (J^v)^{-1}[\sum_{h=1}^H \beta_h(\mathbb{E}_{P_h^v} \Xi_h^f + \lambda \Xi_h^s)](J^v)^{-1}$ with the following matrices

$$\Xi_h^s := \text{Var}_{P_h^v}\left(h^v(R_h^v, Z_h^v; \gamma_0)\right),$$

$$J := \sum_{h=1}^H \beta_h \mathbb{E}_{P_h^v}\left[\nabla_\gamma h^v(R_h^v, Z_h^v; \gamma_0)(\nabla_\gamma h^v(R_h^v, Z_h^v; \gamma_0))^T\right].$$

PROOF. To simplify the notation, we will drop "$v$" and "$N^v$" from most of the subscripts and superscripts. Consistency and asymptotic normality of $\hat{\gamma}$ for $\gamma_0$ follow from Theorem 6.1 of Rubin-Bleuer and Kratina (2005), which relies on the following five Assumptions.

**C.1** $G_N(\gamma_0) \xrightarrow{p} 0$ as $N \to \infty$.

**C.2** There is a compact neighborhood $\Gamma$ of $\gamma_0$ on which with probability one all $G_N(\gamma)$ are continuously differentiable and $\nabla_\gamma G_N(\gamma)$ converge uniformly in $\gamma$ to a nonstochastic limit $J^v(\gamma)$ that is nonsingular at $\gamma_0$.

**C.3** $\sqrt{N} G_N(\gamma_0) \xrightarrow{d} \mathcal{N}(0, \sum_{h=1}^H \beta_h \Xi_h^s)$ as $N \to \infty$.

**C.4** Conditionally on the finite population, there is a compact neighborhood $\Gamma$ of $\gamma_0$ on which $\nabla_\gamma \hat{G}_N(\gamma)$ converge uniformly in the design probability to limit that is nonstochastic in the design probability and coincides with $J^v(\gamma)$ at $\gamma_0$ almost surely.

**C.5** Conditionally on the finite population, $\sqrt{n}\hat{G}_N(\gamma_N) \xrightarrow{d} \mathcal{N}(0, \sum_{h=1}^H \beta_h \Xi_h^f)$ as $n \to \infty$ where the covariance matrices $\Xi_h^f$ are nonstochastic in the superpopulation.

Note that C.1 is implied by C.3. To show that C.3 holds, we can prove that the Lindeberg condition is satisfied and then apply the central limit theorem (proposition 2.27 in Van der Vaart (2000)). For any $\epsilon > 0$,

$$
\frac{1}{N}\sum_{h=1}^H\sum_{i=1}^{N_h}\mathbb{E}_{P_h}[\|(R_{hi} - \pi(Z_{hi};\gamma_0))Z_{hi}\|^2\, \mathbb{1}(\|(R_{hi} - \pi(Z_{hi};\gamma_0))Z_{hi}\| > \epsilon\sqrt{N})]
$$

(6)

$$
< \frac{1}{N}\sum_{h=1}^H\sum_{i=1}^{N_h} dM^2 \mathbb{1}(\sqrt{d}M > \epsilon\sqrt{N})]
$$

where we have used the fact that $\|Z_{hi}\| \le \sqrt{d}\|Z_{hi}\|_\infty < \sqrt{d}M$ and $|R_{hi} - \pi(Z_{hi};\gamma_0)|^2 \le 1$. Then $\lim_{N\to\infty}\mathbb{1}(\sqrt{d}M > \epsilon\sqrt{N}) = 0$, and thus the RHS of (6) converges to 0. In addition,

$$
\lim_{N\to\infty}\frac{1}{N}\sum_{i=1}^H\sum_{i=1}^{N_h}\mathrm{Var}_{P_h}((R_{hi} - \pi(Z_{hi};\gamma_0))Z_{hi}) = \lim_{N\to\infty}\sum_{h=1}^H\frac{N_h}{N}\Xi_h^s = \sum_{h=1}^H\beta_h\Xi_h^s
$$

where the equality follows from the fact that observations are identically distributed within strata, and $\Xi_h^s$ represents the covariance matrix for stratum $h$. Condition C.3 follows from an application of the central limit theorem.

In order to show that C.2 holds, it suffices to prove that for any random vector $\gamma_N \in \Gamma$ converging in probability to $\gamma_0$, $\nabla_\gamma G_N(\gamma_N) \xrightarrow{p} J^v(\gamma_0)$ for some nonstochastic limit $J := J(\theta_0)$ (Theorem 1 in Iséki (1957)). We can decompose $\nabla_\gamma G_N(\gamma_N)$ as follows

$$
\frac{1}{N}\sum_{h=1}^H\sum_{i=1}^{N_h}\frac{e^{-\gamma_0 Z_{hi}}}{(1 + e^{-\gamma_0 Z_{hi}})^2}Z_{hi}Z_{hi}^T + \frac{1}{N}\sum_{h=1}^H\sum_{i=1}^{N_h}\left[\frac{e^{-\gamma_N Z_{hi}}}{(1 + e^{-\gamma_N^T Z_{hi}})^2} - \frac{e^{-\gamma_0 Z_{hi}}}{(1 + e^{-\gamma_0^T Z_{hi}})^2}\right]Z_{hi}Z_{hi}^T
$$

where the first term converges in probability to

$$
J^v(\gamma_0) := \sum_{h=1}^H\beta_h\mathbb{E}_{P_h}\left[e^{-\gamma_0^T Z_h}(1 + e^{-\gamma_0^T Z_h})^{-2}Z_h Z_h^T\right].
$$

The second term can be rewritten as

(7)
$$
\frac{1}{N}\sum_{h=1}^H\sum_{i=1}^{N_h}\left[\frac{e^{\gamma_0^T Z_{hi}}(1 - e^{(\gamma_N - \gamma_0)^T Z_{hi}}) + e^{-\gamma_0^T Z_{hi}}(1 - e^{(\gamma_0 - \gamma_N)^T Z_{hi}})}{(1 + e^{-\gamma_N^T Z_{hi}})(1 + e^{\gamma_N^T Z_{hi}})(1 + e^{-\gamma_0^T Z_{hi}})(1 + e^{\gamma_0^T Z_{hi}})}\right]
$$

which can be upper bounded by $e^{\|\gamma_0\|\sqrt{d_z}M}(e^{\sqrt{d_z}M\|\gamma_N - \gamma_0\|} - 1)$. By the continuous mapping theorem this bound is $o_p(1)$. It follows that that C.2 is verified.

Condition C.5 follows from the Assumptions. For a discussions of the specific conditions needed under various sampling designs, see Section 3.5 of Thompson (1997). Similarly, condition C.4 follows from (5).

The result then follows from Theorem 6.1 of Rubin-Bleuer and Kratina (2005).

$\square$

**B.2. Estimation on NIBRS.** Throughout the proofs, we will use the following lemma.

LEMMA 1. Let $f : \mathcal{X} \mapsto \mathbb{R}$. Assume that A.1–A.3 hold. Then

$$\mathbb{E}[f(X)] = \mathbb{E}\left[\frac{f(X)}{\pi(Z;\gamma_0)}\bigg|R=1\right]\pi^*.$$

PROOF. We can show that

$$\mathbb{E}[f(X)] = \mathbb{E}\left[f(X)\frac{R}{\mathbb{P}(R=1|Z,X)}\right] = \mathbb{E}\left[f(X)\frac{R}{\mathbb{P}(R=1|Z)}\right]$$

where the first equality follows from the law of iterated expectations, while the second follows from A.2. Now, thanks to A.3 and A.1 we obtain that

$$\mathbb{E}\left[f(X)\frac{R}{\mathbb{P}(R=1|Z)}\right] = \mathbb{E}\left[f(X)\frac{R}{\pi(Z;\gamma_0)}\right].$$

The result follows. $\qquad\square$

Then we can derive the asymptotic properties of the estimator $\hat{N}$.

PROPOSITION 2. Consider the conditions of Proposition 1 to be satisfied, and Assumptions A.1–A.4 to hold. Then

(8) $$V_N^{-1/2}\sqrt{n}(\hat{N}/N - 1) \xrightarrow{d} \mathcal{N}(0,1)$$

as $N \to \infty$ where

(9) $$V_N := (\pi^*)^2\left[\mathbb{E}\left[\frac{1-\pi(Z;\gamma_0)}{\pi(Z;\gamma_0)^2}\bigg|R=1\right] + \kappa W^T\Sigma^v W\right].$$

with $W := \mathbb{E}[e^{-Z^T\gamma_0}\pi(Z;\gamma_0)^{-1}Z|R=1]$.

PROOF. Consider the following first-order Taylor expansion of $\hat{N}/N - 1$

$$\frac{1}{N}\sum_{i=1}^{N}\left(\frac{R_i}{\pi(Z_i;\gamma_0)}-1\right)$$

$$- (\hat{\gamma}-\gamma_0)^T\frac{1}{N}\sum_{i=1}^{N}R_i e^{-Z_i^T\gamma_0}Z_i + (\hat{\gamma}-\gamma_0)^T\frac{1}{N}\sum_{i=1}^{N}R_i e^{-Z_i^T\tilde{\gamma}}Z_i Z_i^T(\hat{\gamma}-\gamma_0)$$

where $\tilde{\gamma}$ is a vector between $\hat{\gamma}$ and $\gamma_0$. Using the Cauchy-Schwarz inequality together with A.4, we obtain $((\tilde{\gamma}-\gamma_0)^T Z_i)^2 \leq \|Z_i\|^2\|\hat{\gamma}-\gamma_0\|^2 < \sqrt{d_z}M\|\hat{\gamma}-\gamma_0\|^2$. Thus, we can rewrite $\sqrt{n}(\hat{N}/N - 1)$ as

(10) $$\frac{\sqrt{n}}{\sqrt{N}}\frac{1}{\sqrt{N}}\sum_{i=1}^{N}\left(\frac{R_i}{\pi(Z_i;\gamma_0)}-1\right)$$

$$- \sqrt{n^v}(\hat{\gamma}-\gamma_0)^T\sqrt{\kappa}\frac{1}{N}\sum_{i=1}^{N}Z_i R_i e^{-Z_i^T\gamma_0} + \sqrt{\kappa}O_p(\sqrt{n^v}\|\hat{\gamma}-\gamma_0\|^2).$$

The last term in (10) can be rewritten as $O_p(\|\sqrt{n^v}(\hat{\gamma}-\gamma_0)\|^2/\sqrt{n^v}) = O_p(1/\sqrt{n^v})$ thanks to Proposition 1. The first term is a sum of i.i.d. random variables that are bounded by A.4

and thus it is asymptotically normal with mean 0 thanks to the central limit theorem. By A.4, $N^{-1}\sum_{i=1}^{N} R_i Z_i e^{-Z_i^T \gamma_0}$ in the second term is an average of i.i.d. bounded random variables which converges in probability to $\mathbb{E}[RZe^{-Z^T\gamma_0}] + O_p(1/\sqrt{N})$. We can then rewrite this expectation as $\mathbb{E}[Ze^{-Z^T\gamma_0}\pi(Z;\gamma_0)^{-1}|R=1]\pi^*$ thanks to Lemma 1. Then $\sqrt{n}(\hat{\gamma} - \gamma_0)$ is asymptotically normal by Proposition 1 and consequently the second term in (10) is asymptotically normal by Slutsky. Note that the first two terms in (10) are asymptotically independent because they arise from different samples, and thus we have proved (8). The variance in (9) follows by an application of Lemma 1. $\qquad\square$

PROPOSITION 3. Consider the conditions of Proposition 1 to be satisfied, and Assumptions A.1–A.4 to hold. Then

$$V_{\pi^*}^{-1/2}\sqrt{n}(\hat{\pi}^* - \pi^*) \xrightarrow{d} \mathcal{N}(0,1)$$

as $N \to \infty$ where

$$V_{\pi^*} := (\pi^*)^2 \left[ \mathbb{E}\left[ \frac{\pi^* - \pi(Z;\gamma_0)}{\pi(Z;\gamma_0)^2} \middle| R=1 \right] + \kappa(\pi^*)^2 W^T \Sigma^v W \right].$$

with $W := \mathbb{E}[e^{-Z^T\gamma_0}\pi(Z;\gamma_0)^{-1}Z|R=1]$

PROOF. In order to show asymptotic normality, we can first rewrite $\hat{\pi}^* - \pi^*$ as

$$(11) \quad \frac{\sum_{i=1}^{N} R_i}{N} - \pi^* + \pi^*\left(1 - \frac{\hat{N}}{N}\right) + \left(\frac{N}{\hat{N}} - 1\right)\left(\frac{\sum_{i=1}^{N} R_i}{N} - \pi^*\right)$$

$$+ \pi^*\left(\frac{N}{\hat{N}} - 1\right)\left(1 - \frac{\hat{N}}{N}\right)$$

The third term in (11) is $O_p\left(\frac{1}{\sqrt{n}}\max\left\{\frac{1}{\sqrt{n}}, \frac{1}{\sqrt{n^v}}\right\}\right)$ thanks to Proposition 2, the weak law of law of large numbers, and Slutsky. The last term is $O_p\left(\max\left\{\frac{1}{n}, \frac{1}{n^v}\right\}\right)$ by Proposition 2. Thus, $\sqrt{n}(\hat{\pi}^* - \pi^*)$ is equal to

$$(12) \quad \sqrt{n}\left(\frac{\sum_{i=1}^{N} R_i}{N} - \pi^*\right) + \sqrt{n}\pi^*\left(1 - \frac{\hat{N}}{N}\right) + O_p\left(\max\left\{\frac{1}{\sqrt{n}}, \frac{\sqrt{n}}{n^v}\right\}\right).$$

We can plug in the expansion of $\sqrt{n}(1 - \hat{N}/N)$ in (10) to rewrite (12) as

$$(13) \quad \frac{\sqrt{n}}{N}\sum_{i=1}^{N} R_i\left(1 - \frac{\pi^*}{\pi(Z_i;\gamma_0)}\right) + \pi^*\sqrt{n^v}(\hat{\gamma} - \gamma_0)\sqrt{\kappa}\frac{1}{N}\sum_{i=1}^{N} R_i e^{-Z_i^T\gamma_0}Z_i$$

$$+ O_p\left(\max\left\{\frac{1}{\sqrt{n}}, \frac{\sqrt{n}}{n^v}\right\}\right) + \sqrt{\kappa}O_p(\sqrt{n^v}\|\hat{\gamma} - \gamma_0\|^2)$$

Note that the first term is a sum of i.i.d. random variables bounded by Assumption A.4 and thus converges in distribution to $\mathcal{N}(0, (\pi^*)^2(\mathbb{E}[\pi(Z;\gamma_0)^{-1}]\pi^* - 1))$. Then asymptotic normality of $\sqrt{n}(1 - \hat{N}/N)$ follows from analogous arguments as those in the proof of Proposition 2. $\qquad\square$

PROPOSITION 4. Consider the conditions of Proposition 1 to be satisfied, and Assumptions A.1–A.4 to hold. Then

$$V_{q^*}^{-1/2}\sqrt{n}(\hat{q}^* - q^*) \xrightarrow{d} \mathcal{N}(0,1)$$

as $N \to \infty$ where

$$V_{q^*} = \pi^* q^* \left[ \pi^* \mathbb{E} \left[ \frac{q^* - \alpha^* \pi(Z; \gamma_0)}{\pi(Z; \gamma_0)^2} \middle| R = 1 \right] + 1 - \alpha^* + \pi^* \kappa W^T \Sigma^v W \right]$$

with $W := \mathbb{E}[e^{-Z^T \gamma_0} \pi(Z; \gamma_0)^{-1} Z | R = 1]$, and $\alpha^* := \mathbb{E}[A | R = 1]$.

PROOF. This proof follows from the same set of arguments as the proof of Propositions 2 and 3, hence it is omitted. □

PROPOSITION 5. Consider the conditions of Proposition 1 and Assumptions A.1–A.4 to hold. Let $\theta_0 \in \text{Int}(\Theta)$ be defined by the moment condition (1) and $\hat{\theta} \in \Theta$ be the estimator that solves the estimating equation (2). Then

$$\Sigma^{-1/2} \sqrt{n} (\hat{\theta} - \theta_0) \xrightarrow{d} \mathcal{N}(0, I_d)$$

as $n \to \infty$ with $\Sigma := J_\theta^{-1} \Xi J_\theta^{-1}$ and

$$\Xi := \mathbb{E} \left[ h(A, Z, X; \theta_0, \gamma_0) h(A, Z, X; \theta_0, \gamma_0)^T | R = 1 \right] + \kappa J_\gamma \Sigma^v J_\gamma^T,$$

$$J_\theta := \nabla_\theta G(\theta_0, \gamma_0) = \mathbb{E}[q(X; \theta_0)(1 - q(X; \theta_0)) \pi(Z; \gamma_0)^{-1} X X^T | R = 1],$$

$$J_\gamma := \nabla_\gamma G(\theta_0, \gamma_0) = \mathbb{E}[q(X; \theta_0) e^{-\gamma_0^T Z} X Z^T | R = 1].$$

PROOF. To show that $\hat{\theta}$ is consistent for $\theta_0$ and asymptotically normal, we can verify the following three Assumptions from Yuan and Jennrich (1998):

**C.1** $\hat{G}_N(\theta_0, \hat{\gamma}) \xrightarrow{p} 0$ as $N \to \infty$.
**C.2** There exists a neighborhood $\Theta$ of $\theta_0$ on which with probability one all $\hat{G}_N(\theta, \hat{\gamma})$ are continuously differentiable and $\nabla_\theta \hat{G}_N(\theta, \hat{\gamma})$ converge uniformly to a nonstochastic limit that is nonsigular at $\theta_0$.
**C.3** $\Xi^{-1/2} \sqrt{n} \hat{G}_N(\theta_0, \hat{\gamma}) \xrightarrow{d} \mathcal{N}(0, I_{d_x})$ as $N \to \infty$ for some matrix $\Xi$.

Since C.3 implies C.1, we only need to prove C.2 and C.3.

We first show that C.3 holds. Consider the following Taylor expansion of $\sqrt{n} \hat{G}_N(\theta_0, \hat{\gamma})$

$$\frac{\sqrt{n}}{N} \sum_{i=1}^N R_i h(A_i, Z_i, X_i; \theta_0, \gamma_0) - \frac{1}{N} \sum_{i=1}^N R_i q(X_i; \theta_0) e^{-\gamma_0^T Z_i} X_i Z_i^T \frac{\sqrt{n}}{\sqrt{n^v}} \sqrt{n^v} (\hat{\gamma} - \gamma_0)$$

$$+ \frac{1}{N} \sum_{i=1}^N R_i q(X_i; \theta_0) e^{-\tilde{\gamma}^T Z_i} X_i \frac{\sqrt{n}}{n^v} \left( \sqrt{n^v} (\hat{\gamma} - \gamma_0)^T Z_i \right)^2.$$

where $\tilde{\gamma}$ is a convex combination of $\hat{\gamma}$ and $\gamma_0$. The first term is a i.i.d. sum whose terms have bounded moments by A.4. Thus, it is asymptotically normal by the central limit theorem. The term $N^{-1} \sum_{i=1}^N R_i q(X_i; \theta_0) e^{-\gamma_0^T Z_i} X_i Z_i^T$ is an i.i.d. average formed by terms that have finite moments by A.4, so it converges in probability to $J_\gamma := \mathbb{E}[R q(X; \theta_0) e^{-\gamma_0^T Z} X Z^T]$ by the weak law of large numbers. Thus, the second term is asymptotically normal by Proposition 1 and Slutsky. Using similar arguments as in the proof of Proposition 2, we can show that the third term is $o_p(1)$. The first two terms are asymptotically independent because they are arise from separate samples, hence it follows that C.3 is verified.

To show that C.2 holds, it suffices to show that for any random vector $\theta_N \in \Theta$ converging in probability to $\theta_0$, $\hat{G}_N(\theta_N, \hat{\gamma}) := \nabla_\theta \hat{G}_N(\theta, \hat{\gamma}) \xrightarrow{p} J_\theta$ for some nonstochastic function $J_\theta :=$

$J(\theta_0)$ (Theorem 1 in Iséki (1957)). Consider the following Taylor expansion of $\hat{\dot{G}}_N(\theta_N, \hat{\gamma})$

$$(14) \quad \frac{1}{N}\sum_{i=1}^{N} R_i q(X_i; \theta_N)(q(X_i; \theta_N) - 1)(1 + e^{-Z_i^T \gamma_0})X_i X_i^T + (\hat{\gamma} - \gamma_0)^T \nabla_\gamma \hat{\dot{G}}_N(\theta_N, \tilde{\gamma})$$

where $\tilde{\gamma}$ is a convex combination of $\hat{\gamma}$ and $\gamma_0$. The first term in (14) can be rewritten as

$$\frac{1}{N}\sum_{i=1}^{N} R_i q(X_i; \theta_0)(q(X_i; \theta_0) - 1)(1 + e^{-Z_i^T \gamma_0})X_i X_i^T$$

$$+ \frac{1}{N}\sum_{i=1}^{N} R_i \left[ q(X_i; \theta_N)(q(X_i; \theta_N) - 1) - q(X_i; \theta_0)(q(X_i; \theta_0) - 1) \right] (1 + e^{-Z_i^T \gamma_0})X_i X_i^T$$

where the first term converges to $J_\theta := \mathbb{E}[RXX^T e^{-X^T\theta_0}/\pi(Z; \gamma)]$ by the weak law of large numbers. The second term is $o_p(1)$ by Cauchy-Schwarz and A.4; the upper bound can be derived using a similar strategy as in expression (7) of the proof of Proposition 1. For the second term in (14), we have that

$$(15) \quad (\hat{\gamma} - \gamma_0)^T \nabla_\gamma \hat{\dot{G}}_N(\theta_N, \tilde{\gamma})$$

$$= (\hat{\gamma} - \gamma_0)^T \frac{1}{N}\sum_{i=1}^{N} R_i Z_i q(X_i; \theta_N)(1 - q(X_i; \theta_N))e^{-Z_i^T \tilde{\gamma}}X_i X_i^T$$

where each element of the $d_x \times d_x$ matrix can be upper bounded by

$$\|\hat{\gamma} - \gamma_0\| M^2 e^{\sqrt{d}M \sup_{\gamma \in \Gamma}\|\gamma\|}\sqrt{d_z}.$$

Together with Proposition 1, this implies each of the elements in (15) is $o_p(1)$. It follows that C.2 holds true.

Under Assumptions C.1, C.2, and C.3, the result follows by an application of Theorem 4 in Yuan and Jennrich (1998) and our Lemma 1. $\square$

Finally, we turn to the result on generalized estimation equations (GEEs) given by Proposition 6. In the proof, we will use $\alpha_0$, which is such that for each $1 \leq i \leq N$, $K \geq 2$, and $1 \leq k < j \leq K_i$,

$$\alpha_0 := \frac{\mathbb{E}[\mathbf{A}_{ik}\mathbf{A}_{ij}|\mathbf{X}_i] - \mathbb{E}[\mathbf{A}_{ij}|\mathbf{X}_i]\mathbb{E}[\mathbf{A}_{ik}|\mathbf{X}_i]}{\sqrt{\text{Var}(\mathbf{A}_{ik}|\mathbf{X}_i)}\sqrt{\text{Var}(\mathbf{A}_{ij}|\mathbf{X}_i)}}$$

where $\mathbb{E}[\mathbf{A}_{ij}|\mathbf{X}_i] = q(\mathbf{X}_i; \theta_0)$ and $\text{Var}(\mathbf{A}_{ik}|\mathbf{X}_i) = q(\mathbf{X}_i; \theta_0)(1 - q(\mathbf{X}_i; \theta_0))$ by A.5.

PROPOSITION 6. Assume that the conditions of Proposition 1 and A.1–A.5 hold. Assume that the entries of $W(\mathbf{X}, \theta, \alpha)^{-1}$ and their derivatives are continuous. Let $\hat{\theta}$ be the estimate of $\theta$ obtained by solving the estimating equation

$$\hat{G}_N(\theta, \hat{\alpha}, \hat{\gamma}) := \frac{1}{N}\sum_{i=1}^{N} R_i \mathbf{X}_i D_i(\theta)W_i(\theta, \hat{\alpha})^{-1}\left(\mathbf{A}_i - \frac{\mathbf{q}_i(\theta)}{\pi(Z_i; \hat{\gamma})}\right) = 0.$$

Let $\hat{\alpha}$ be an estimator of $\alpha_0$ such that $\hat{\alpha} - \alpha_0 = O_p(1/\sqrt{N})$. Then

$$\Sigma^{-1/2}\sqrt{n}(\hat{\theta} - \theta_0) \overset{d}{\to} \mathcal{N}(0, I_{d_x})$$

as $N \to \infty$. $\Sigma := J_\theta^{-1}(\Xi)J_\theta^{-1}$ with

$$\Xi := \mathbb{E}[h(\mathbf{X}, Z, \mathbf{A}; \theta_0, \alpha_0)h(\mathbf{X}, Z, \mathbf{A}; \theta_0, \alpha_0)^T | R = 1] + \kappa J_\gamma \Sigma^v J_\gamma^T$$

$$J_\theta := \mathbb{E}\left[\mathbf{X}D(\mathbf{X}; \theta_0)W(\mathbf{X}; \theta_0, \alpha_0)^{-1}\nabla_\theta \mathbf{q}(\theta_0)\pi(Z; \gamma_0)^{-1} | R = 1\right]$$

$$J_\gamma := \mathbb{E}\left[\mathbf{X}D(\mathbf{X}; \theta_0)W(\mathbf{X}; \theta_0, \alpha_0)^{-1}\mathbf{q}(\theta_0)Ze^{-Z^T\gamma_0} | R = 1\right]$$

where $h(\mathbf{X}, Z, \mathbf{A}; \theta_0, \alpha_0) := \mathbf{X}D(\mathbf{X}; \theta_0)W(\mathbf{X}; \theta_0, \alpha_0)^{-1}(\mathbf{A} - \mathbf{q}(\theta_0)\pi(Z; \gamma_0)^{-1})$, $\mathbf{q}(\theta_0) := (q(X_1; \theta_0), \ldots, q(X_K; \theta_0))^T$ with $X_k$ for $1 \leq k \leq K$ being the $k^{th}$ column of $\mathbf{X}$.

PROOF. To prove consistency of $\hat{\theta}$ for $\theta_0$ and its asymptotic normality, we will use the results of Yuan and Jennrich (1998) which rely on the following three conditions.

**C.1** $\hat{G}_N(\theta_0, \hat{\alpha}, \hat{\gamma}) \xrightarrow{p} 0$ as $N \to \infty$.

**C.2** There exists a neighborhood $\Theta$ of $\theta_0$ on which with probability one $\nabla_\theta \hat{G}_N(\theta, \hat{\alpha}, \hat{\gamma})$ is continuously differentiable and its derivatives converge uniformly to a nonstochastic limit that is nonsigular at $\theta_0$.

**C.3** $\Xi^{-1/2}\sqrt{n}\hat{G}_N(\theta_0, \hat{\alpha}, \hat{\gamma}) \xrightarrow{d} \mathcal{N}(0, I_{d_x})$ for some positive definite matrix $\Xi$.

To show that C.2 holds, it suffices to show that for any random vector $\theta_N \in \Theta$ converging in probability to $\theta_0$, $\hat{\tilde{G}}_N(\theta, \hat{\alpha}, \hat{\gamma}) := \hat{\tilde{G}}_N(\theta_N, \hat{\alpha}, \hat{\gamma}) \xrightarrow{p} J(\theta_0)$ as $N \to \infty$ for some nonstochastic limit $J_\theta := J(\theta_0)$. First note that $\nabla_\theta \hat{G}_N(\theta, \alpha, \gamma)$ represents the mean of $N$ i.i.d. observations. For $i = 1, \ldots, N$, the $i^{th}$ observation is finite by A.4 and the fact that $\max_{i,j=1}^{N_i} |W_i(\theta, \alpha)^{-1}|$ is bounded. To simplify the presentation, let us rewrite $\hat{G}_N(\theta, \alpha, \gamma) = \hat{G}_N^D(\theta, \alpha, \gamma) + \hat{G}_N^W(\theta, \alpha, \gamma) + \hat{G}_N^q(\theta, \alpha, \gamma)$ where

$$\hat{G}_N^D(\theta, \alpha, \gamma) := \frac{1}{N}\sum_{i=1}^N \mathbf{X}_i[\nabla_\theta D_i(\theta)]W_i(\theta, \alpha)^{-1}\left(\mathbf{A}_i - \mathbf{q}_i(\theta)\frac{1}{\pi(Z_i; \gamma)}\right)$$

$$\hat{G}_N^W(\theta, \alpha, \gamma) := \frac{1}{N}\sum_{i=1}^N \mathbf{X}_i D_i(\theta)[\nabla_\theta W_i(\theta, \alpha)^{-1}]\left(\mathbf{A}_i - \mathbf{q}_i(\theta)\frac{1}{\pi(Z_i; \gamma)}\right)$$

$$\hat{G}_N^q(\theta, \alpha, \gamma) := -\frac{1}{N}\sum_{i=1}^N R_i \mathbf{X}_i D_i(\theta)W_i(\theta, \alpha)^{-1}\frac{1}{\pi(Z_i; \gamma)}[\nabla_\theta \mathbf{q}_i(\theta)]$$

Consider the following Taylor expansion of $\hat{\tilde{G}}_N(\theta_N, \hat{\alpha}, \hat{\gamma})$:

$$(16) \quad \hat{G}_N(\theta_N, \alpha_0, \gamma_0) + (\hat{\alpha} - \alpha_0)\nabla_\alpha \hat{\tilde{G}}_N(\theta_N, \tilde{\alpha}, \gamma_0)$$

$$+ (\hat{\gamma} - \gamma_0)^T \nabla_\gamma \hat{\tilde{G}}_N(\theta_N, \alpha_0, \tilde{\gamma}) + (\hat{\gamma} - \gamma_0)^T \nabla_\alpha \nabla_\gamma \hat{\tilde{G}}_N(\theta_N, \tilde{\alpha}, \tilde{\gamma})(\hat{\alpha} - \alpha_0)$$

where $\tilde{\gamma}$ is a convex combination of $\hat{\gamma}$ and $\gamma_0$, while $\tilde{\alpha}$ is a convex combination of $\hat{\alpha}$ and $\alpha_0$.

The first term in (16) can be rewritten as

$$(17) \quad \hat{\tilde{G}}_N(\theta_0, \alpha_0, \gamma_0) + \hat{\tilde{G}}_N(\theta_N, \alpha_0, \gamma_0) - \hat{\tilde{G}}_N(\theta_0, \alpha_0, \gamma_0).$$

The first term in (17) is an average of $N$ terms that are i.i.d. and finite by A.1, A.4, and the boundedness of $W_i(\theta, \alpha)^{-1}$. Since the model for the mean is corectly specified by A.5, the first term converges in probability $-\mathbb{E}[RD(\mathbf{X}; \theta_0)W(\mathbf{X}; \theta_0, \alpha_0)^{-1}\nabla_\theta \mathbf{q}(\theta_0)\pi(Z; \gamma_0)^{-1}]$ by the weak law of large numbers and iterated expectations.

We can then show that the difference between the remaining two terms in (17) converges to 0 in probability. Let $\bar{w} := \max_{\theta \in \theta, \alpha \in [-1,1]} \max_{i,j,k} |(W_i(\theta, \alpha)^{-1})_{jk}|$. We have that $\hat{G}_N^D(\theta_N, \alpha_0, \gamma_0)$ is equal to

$$\frac{1}{N} \sum_{i=1}^N \mathbf{X}_i \nabla_\theta D_i(\theta_N) W_i(\theta_N, \alpha_0)^{-1} \left( A_i - \frac{q_i(\theta_N)}{\pi(Z_i; \gamma_0)} \right)$$

$$< \frac{1}{N} \sum_{i=1}^N M^2 \bar{w} u_{d_x} u_{K_i}^T I_{K_i} u_{K_i} u_{K_i}^T \left[ \left( A_i - \frac{q_i(\theta_0)}{\pi(Z_i; \gamma_0)} \right) + \frac{1}{\pi(Z_i; \gamma_0)} (q_i(\theta_0) - q_i(\theta_N)) \right] u_{d_x}^T$$

where $u_n := (1, \ldots, 1)^T$ has length $n \in \mathbb{Z}^+$. The inequality follows from A.4 and the boundedness of $W(\theta_N, \alpha)^{-1}$. Then $N^{-1} \sum_{i=1}^N (A_i - q_i(\theta_0)/\pi(Z_i; \theta_0))$ is an average of i.i.d. terms that are finite and thus converges in proability to 0 by the weak law of large numbers and A.5. By A.1, A.4, and Cauchy-Schwarz, we have the following upper bound

$$\frac{1}{N} \sum_{i=1}^N \frac{1}{\pi(Z_i; \gamma_0)} (q_i(\theta_0 - q_i(\theta_N))) < \frac{1}{\epsilon} e^{\|\theta_0\| M \sqrt{d}} \left( e^{\|\theta_N - \theta_0\| M \sqrt{d}} - 1 \right)$$

where the RHS converges to 0 in probability by the weak law of large numbers and the continuous mapping theorem. It follows that $\hat{G}_N^D(\theta_N, \alpha_0, \gamma_0)$ converges to 0 in probability. Using similar arguments, we can show that $\hat{G}_N^W(\theta_N, \alpha_0, \gamma_0)$ converges in probability to 0. Clearly, $\hat{G}_N^D(\theta_0, \alpha_0, \gamma_0)$ and $\hat{G}_N^W(\theta_0, \alpha_0, \gamma_0)$ converge to 0 in probability by the weak law of large numbers. Next, we need to show that $\hat{G}_N^q(\theta_N, \alpha, \gamma) - \hat{G}_N^q(\theta_0, \alpha, \gamma)$ converges to 0 in probability. This difference can be rewritten as

$$\frac{1}{N} \sum_{i=1}^N \frac{R_i}{\pi(Z_i; \gamma_0)} \mathbf{X}_i \left[ D_i(\theta_0) W_i(\theta_0, \alpha_0)^{-1} \nabla \mathbf{q}(\theta_0) - D_i(\theta_N) W(\theta_N, \alpha_0)^{-1} \nabla \mathbf{q}(\theta_N) \right].$$

By A.4, the convergence boils down to showing that for all $i = 1, \ldots, N$ and $1 \leq k, j \leq K_i$,

$$(W_i(\theta_0, \alpha_0)^{-1})_{jk} q(X_{ik}; \theta_0)(1 - q(X_{ik}; \theta_0)) q(X_{ij}; \theta_0)(1 - q(X_{ij}; \theta_0))$$
$$- (W_i(\theta_N, \alpha_0)^{-1})_{jk} q(X_{ik}; \theta_N)(1 - q(X_{ik}; \theta_N)) q(X_{ij}; \theta_N)(1 - q(X_{ij}; \theta_N))$$

converges in probability to 0. This can be upper bounded by

$$(18) \quad e^{\theta_0^T(X_{ik} + X_{ij})}(1 + e^{\theta_0^T X_{ik}})(1 + e^{\theta_0^T X_{ij}})(W_i(\theta_0, \alpha_0)^{-1})_{kj}$$

$$\left| \frac{e^{(\theta_N^T(X_{ik} + X_{ij})}}{e^{(\theta_0^T(X_{ik} + X_{ij})}} \frac{1 + e^{\theta_N^T X_{ik}}}{1 + e^{\theta_0^T X_{ik}}} \frac{1 + e^{\theta_N^T X_{ij}}}{1 + e^{\theta_0^T X_{ij}}} \frac{(W_i(\theta_N, \alpha_0)^{-1})_{kj}}{(W_i(\theta_0, \alpha_0)^{-1})_{kj}} - 1 \right|$$

$$< e^{2\|\theta_0\| \sqrt{d} M}(1 + e^{\|\theta_0\| \sqrt{d} M})^2 \left| e^{4\|\theta_N - \theta_0\| \sqrt{d} M}(1 + o_p(1)) - 1 \right|$$

where the inequality follows from Cauchy-Schwarz, the fact that $\|X\| \leq \sqrt{d} \|X\|_\infty < M$, and the continuous mapping theorem thanks to the fact that $W_i$ has continuous derivatives and $\theta_N \to \theta_0$ as $N \to \infty$. Since $\theta_N \xrightarrow{p} \theta_0$ as $N \to \infty$, the RHS in (18) is $o_p(1)$. Thus, we can conclude that the difference between the second and third terms in (17) converges in probability to 0.

Next, we turn again to the Taylor expansion of $\hat{\dot{G}}_N(\theta_N, \hat{\alpha}, \hat{\gamma})$ in (16). Note that $\nabla_\alpha \hat{\dot{G}}_N(\theta_N, \tilde{\alpha}, \gamma_0)$ is bounded and $\hat{\alpha} - \alpha_0 = o_p(1)$, so their product converges to 0 in probability by Slutsky. The third term $\nabla_\gamma \hat{\dot{G}}_N(\theta_N, \alpha_0, \tilde{\gamma})$ is also bounded and, by Proposition 1, $\hat{\gamma} - \gamma_0 \xrightarrow{p} 0$ as $N \to \infty$, thus we can appply Cauchy-Schwarz to show that the product converges to 0 in probability. The fourth term in the RHS of (16) can be shown to converge in probability to 0 using analogous arguments. It follows that C.2 holds.

Since C.3 implies C.1, we only need to show that C.3 holds. To show that C.3 holds, consider the following Taylor expansion of $\sqrt{n}\hat{G}(\theta_0, \hat{\alpha}, \hat{\gamma})$:

$$
(19) \quad \sqrt{n}\Bigg[\hat{G}_N(\theta_0, \alpha_0, \gamma_0) + (\hat{\alpha} - \alpha_0)\nabla_\alpha \hat{G}_N(\theta_0, \alpha_0, \gamma_0) + (\hat{\gamma} - \gamma_0)^T \nabla_\gamma \hat{G}_N(\theta_0, \alpha_0, \gamma_0)
$$
$$
+ (\hat{\alpha} - \alpha_0)^2 \nabla_\alpha^2 \hat{G}_N(\theta_0, \tilde{\alpha}, \gamma_0) + (\hat{\alpha} - \alpha_0)(\hat{\gamma} - \gamma_0)^T
$$
$$
\nabla_\alpha \nabla_\gamma \hat{G}_N(\theta_0, \alpha_0, \gamma_0) + (\hat{\gamma} - \gamma_0)^T \nabla_\gamma^2 \hat{G}_N(\theta_0, \alpha, \tilde{\gamma})(\hat{\gamma} - \gamma_0)
$$
$$
+ (\hat{\alpha} - \alpha_0)^2 (\hat{\gamma} - \gamma_0)^T \nabla_\alpha^2 \nabla_\gamma \hat{G}_N(\theta_0, \tilde{\alpha}, \gamma_0) + (\hat{\alpha} - \alpha_0)(\hat{\gamma} - \gamma_0)\nabla_\alpha \nabla_\gamma^2 \hat{G}_N(\theta_0, \alpha_0, \tilde{\gamma})(\hat{\gamma} - \gamma_0)
$$
$$
+ (\hat{\alpha} - \alpha_0)^2 (\hat{\gamma} - \gamma_0)^T \nabla_\alpha^2 \nabla_\gamma^2 \hat{G}_N(\theta_0, \alpha, \tilde{\gamma})(\hat{\gamma} - \gamma_0)\Bigg].
$$

The first term in (19) is composed by $N$ i.i.d. bounded random variables and thus, by the central limit theorem,

$$
\sqrt{n}\hat{G}_N(\theta_0, \alpha_0, \gamma_0) \xrightarrow{d} \mathcal{N}(0, \Xi)
$$

as $N \to \infty$ where $\Xi := \text{Var}(G(\theta_0, \alpha_0, \gamma_0)$.

For the second term in (19), $\sqrt{n}(\hat{\alpha} - \alpha_0) = O_p(1)$ by Assumption, while $\nabla_\alpha \hat{G}_N(\theta_0, \alpha_0, \gamma_0)$ converges in probability to 0 by the weak law of large numbers and A.5. By the continuous mapping theorem, their product then converges in probability to 0.

For the third term,

$$
\sqrt{n}(\hat{\gamma} - \gamma_0)^T \nabla_\gamma \hat{G}_N(\theta_0, \alpha_0, \gamma_0) = \sqrt{n^v}(\hat{\gamma} - \gamma_0)^T \sqrt{\kappa}\nabla_\gamma \hat{G}_N(\theta_0, \alpha_0, \gamma_0)
$$

where $\nabla_\gamma \hat{G}_N(\theta_0, \alpha_0, \gamma_0) \xrightarrow{p} J_\gamma := \mathbb{E}[\nabla_\gamma G(\theta_0, \alpha_0, \gamma_0)]$ as $n^v \to \infty$ by the weak law of large numbers. In addition, $\sqrt{n^v}(\hat{\gamma} - \gamma_0) \xrightarrow{d} \mathcal{N}(0, \Sigma^v)$ by propositon 1. Thus, the term converges in distribution by Slutsky to $\mathcal{N}(0, J_\gamma \Sigma^v J_\gamma)$.

For the fourth term, $\nabla_\alpha^2 \hat{G}_N(\theta_0, \tilde{\alpha}, \gamma_0)$ is bounded and $\sqrt{n}(\hat{\alpha} - \alpha_0) = O_p(1)$, so their product converges in probability to 0 as $N \to \infty$. Using similar arguments, it is easy to see that all remaining terms converge to 0 in probability as well. Thus C.3 is satisied.

The result of the Proposition follows from Theorem 4 of Yuan and Jennrich (1998).

$\square$

## REFERENCES

ANDREWS, D. W. and MONAHAN, J. C. (1992). An improved heteroskedasticity and autocorrelation consistent covariance matrix estimator. *Econometrica: Journal of the Econometric Society* 953–966.

AVAKAME, E. F., FYFE, J. J. and MCCOY, C. (1999). "Did you call the police? What did they do?" An empirical assessment of Black's theory of mobilization of law. *Justice Quarterly* 16 765–792.

AZUR, M. J., STUART, E. A., FRANGAKIS, C. and LEAF, P. J. (2011). Multiple imputation by chained equations: what is it and how does it work? *International journal of methods in psychiatric research* 20 40–49.

BACHMAN, R. (1998). The factors related to rape reporting behavior and arrest: New evidence from the National Crime Victimization Survey. *Criminal justice and behavior* 25 8–29.

BARNETT-RYAN, C., LANGTON, L. and PLANTY, M. (2014). The nation's two crime measures, 2014. *US Department of Justice, Washington, DC*.

BASU, D. (2011). An essay on the logical foundations of survey sampling, part one. In *Selected Works of Debabrata Basu* 167–206. Springer.

BAUMER, E. P. (2002). Neighborhood disadvantage and police notification by victims of violence. *Criminology* **40** 579–616.

BAUMER, E. P. and LAURITSEN, J. L. (2010). Reporting crime to the police, 1973–2005: A multivariate analysis of long-term trends in the National Crime Survey (NCS) and National Crime Victimization Survey (NCVS). *Criminology* **48** 131–185.

BECK, A. J. and BLUMSTEIN, A. (2018). Racial disproportionality in US state prisons: Accounting for the effects of racial and ethnic differences in criminal involvement, arrests, sentencing, and time served. *Journal of Quantitative Criminology* **34** 853–883.

BERK, R., BUJA, A., BROWN, L., GEORGE, E., KUCHIBHOTLA, A. K., SU, W. and ZHAO, L. (2019). Assumption lean regression. *The American Statistician*.

BINDER, D. A. (1983). On the variances of asymptotically normal estimators from complex surveys. *International Statistical Review/Revue Internationale de Statistique* 279–292.

BLUMSTEIN, A. and COHEN, J. (1979). Estimation of individual crime rates from arrest records. *J. Crim. L. & Criminology* **70** 561.

BLUMSTEIN, A. and COHEN, J. (1987). Characterizing criminal careers. *Science* **237** 985–991.

BLUMSTEIN, A. et al. (1986). *Criminal Careers and" Career Criminals,"* **2**. National Academies.

BLUMSTEIN, A., COHEN, J., PIQUERO, A. R. and VISHER, C. A. (2010). Linking the crime and arrest processes to measure variations in individual arrest risk per crime (Q). *Journal of Quantitative Criminology* **26** 533–548.

BÖHNING, D. and VAN DER HEIJDEN, P. G. (2009). A covariate adjustment for zero-truncated approaches to estimating the size of hidden and elusive populations. *The Annals of Applied Statistics* **3** 595–610.

BRAME, R., FAGAN, J., PIQUERO, A. R., SCHUBERT, C. A. and STEINBERG, L. (2004). Criminal careers of serious delinquents in two cities. *Youth Violence and Juvenile Justice* **2** 256–272.

BREIMAN, L. (2001). Random forests. *Machine learning* **45** 5–32.

BUIL-GIL, D., MEDINA, J. and SHLOMO, N. (2021). Measuring the dark figure of crime in geographic areas: Small area estimation from the crime survey for England and Wales. *The British journal of criminology* **61** 364–388.

BUIL-GIL, D., MORETTI, A. and LANGTON, S. H. (2021). The accuracy of crime statistics: Assessing the impact of police data bias on geographic crime analysis. *Journal of Experimental Criminology* 1–27.

BUJA, A., BROWN, L., BERK, R., GEORGE, E., PITKIN, E., TRASKIN, M., ZHANG, K. and ZHAO, L. (2019a). Models as approximations I: Consequences illustrated with linear regression. *Statistical Science* **34** 523–544.

BUJA, A., BROWN, L., KUCHIBHOTLA, A. K., BERK, R., GEORGE, E., ZHAO, L. et al. (2019b). Models as Approximations II: A Model-Free Theory of Parametric Regression. *Statistical Science* **34** 545–565.

BYRD, J. and LIPTON, Z. (2019). What is the effect of importance weighting in deep learning? In *International Conference on Machine Learning* 872–881. PMLR.

CERNAT, A., BUIL-GIL, D., PINA-SÁNCHEZ, J., MURRIÀ-SANGENÍS, M. et al. (2021). Estimating crime in place: Moving beyond residence location.

D'ALESSIO, S. J. and STOLZENBERG, L. (2003). Race and the probability of arrest. *Social forces* **81** 1381–1397.

DUGAN, L. (2003). Domestic violence legislation: Exploring its impact on the likelihood of domestic violence, police involvement, and arrest. *Criminology & Public Policy* **2** 283–312.

FISHER, B. S., DAIGLE, L. E., CULLEN, F. T. and TURNER, M. G. (2003). Reporting sexual victimization to the police and others: Results from a national-level study of college women. *Criminal justice and behavior* **30** 6–38.

FITZMAURICE, G. M., LAIRD, N. M. and ROTNITZKY, A. G. (1993). Regression models for discrete longitudinal responses. *Statistical Science* 284–299.

FITZMAURICE, G., DAVIDIAN, M., VERBEKE, G. and MOLENBERGHS, G. (2008). *Longitudinal data analysis*. CRC press.

FOGLIATO, R., XIANG, A., LIPTON, Z., NAGIN, D. and CHOULDECHOVA, A. (2021). On the Validity of Arrest as a Proxy for Offense: Race and the Likelihood of Arrest for Violent Crimes. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*. *AIES '21* 100–111. Association for Computing Machinery, New York, NY, USA. https://doi.org/10.1145/3461702.3462538

FULLER, W. A. (2011). *Sampling statistics* **560**. John Wiley & Sons.

GRAHAM, J. W., OLCHOWSKI, A. E. and GILREATH, T. D. (2007). How many imputations are really needed? Some practical clarifications of multiple imputation theory. *Prevention science* **8** 206–213.

HECKMAN, J. J. (1979). Sample selection bias as a specification error. *Econometrica: Journal of the econometric society* 153–161.

HORVITZ, D. G. and THOMPSON, D. J. (1952). A generalization of sampling without replacement from a finite universe. *Journal of the American statistical Association* **47** 663–685.

HUBBARD, A. E., AHERN, J., FLEISCHER, N. L., VAN DER LAAN, M., SATARIANO, S. A., JEWELL, N., BRUCKNER, T. and SATARIANO, W. A. (2010). To GEE or not to GEE: comparing population average and mixed models for estimating the associations between neighborhood risk factors and health. *Epidemiology* 467–474.

HUGGINS, R. (1989). On the statistical analysis of capture experiments. *Biometrika* **76** 133–140.

ISÉKI, K. (1957). A theorem on continuous convergence. *Proceedings of the Japan Academy* **33** 355–356.

KANG, J. D. and SCHAFER, J. L. (2007). Demystifying double robustness: A comparison of alternative strategies for estimating a population mean from incomplete data. *Statistical science* **22** 523–539.

KOCHEL, T. R., WILSON, D. B. and MASTROFSKI, S. D. (2011). EFFECT OF SUSPECT RACE ON OFFICERS' ARREST DECISIONS. *Criminology* **49** 473–512.

LANTZ, B. and WENGER, M. R. (2019). The co-offender as counterfactual: A quasi-experimental within-partnership approach to the examination of the relationship between race and arrest. *Journal of experimental criminology* 1–24.

LEE, S.-M. and CHAO, A. (1994). Estimating population size via sample coverage for closed capture-recapture models. *Biometrics* 88–97.

LIANG, K.-Y. and ZEGER, S. L. (1986). Longitudinal data analysis using generalized linear models. *Biometrika* **73** 13–22.

LITTLE, R. J. and RUBIN, D. B. (2019). *Statistical analysis with missing data* **793**. John Wiley & Sons.

LOEFFLER, C. E., HYATT, J. and RIDGEWAY, G. (2019). Measuring self-reported wrongful convictions among prisoners. *Journal of Quantitative Criminology* **35** 259–286.

LOHR, S. L. (2007). Comment: Struggles with Survey Weighting and Regression Modeling. *Statistical Science* **22** 175 – 178. https://doi.org/10.1214/088342307000000159

LUMLEY, T. and SCOTT, A. (2017). Fitting regression models to survey data. *Statistical Science* 265–278.

LYTLE, D. J. (2014). The effects of suspect characteristics on arrest: A meta-analysis. *Journal of Criminal Justice* **42** 589–597.

MORGAN, R. E. and TRUMAN, J. (2021). Criminal victimization, 2020. *Washington, DC: National Crime Victimization Survey, Bureau of Justice Statistics. Retrieved Jan* **4** 2022.

MORGAN, R. E., OF JUSTICE STATISTICS (BJS), B., OF JUSTICE, U. D., OF JUSTICE PROGRAMS, O. and OF AMERICA, U. S. (2017). Race and hispanic origin of victims and offenders, 2012-15. *Victims and Offenders* **2012** 15.

NAGIN, D. S. (2013). Deterrence in the twenty-first century. *Crime and justice* **42** 199–263.

NEWEY, W. K. and MCFADDEN, D. (1994). Large sample estimation and hypothesis testing. *Handbook of econometrics* **4** 2111–2245.

PETERSEN, C. G. J. (1896). The yearly immigration of young plaice in the Limfjord from the German sea. *Rept. Danish Biol. Sta.* **6** 1–48.

PIQUERO, A. R. and BRAME, R. W. (2008). Assessing the race–crime and ethnicity–crime relationship in a sample of serious adolescent delinquents. *Crime & Delinquency* **54** 390–422.

POLLEY, E. C. and VAN DER LAAN, M. J. (2010). Super learner in prediction.

POPE, C. E. and SNYDER, H. N. (2003). *Race as a factor in juvenile arrests*. Citeseer.

RACINE, J. and LI, Q. (2004). Nonparametric estimation of regression functions with both categorical and continuous data. *Journal of Econometrics* **119** 99–130.

RENNISON, C. M. (2010). An investigation of reporting violence to the police: A focus on Hispanic victims. *Journal of Criminal Justice* **38** 390–399.

RICHARDSON, R., SCHULTZ, J. and CRAWFORD, K. (2019). Dirty data, bad predictions: How civil rights violations impact police data, predictive policing systems, and justice. *New York University Law Review Online, Forthcoming*.

ROBERTS, A. and LYONS, C. J. (2009). Victim-offender racial dyads and clearance of lethal and nonlethal assault. *Journal of research in crime and delinquency* **46** 301–326.

ROBERTS, A. and LYONS, C. J. (2011). Hispanic victims and homicide clearance by arrest. *Homicide Studies* **15** 48–73.

RUBIN-BLEUER, S. and KRATINA, I. S. (2005). On the two-phase framework for joint model and design-based inference. *The Annals of Statistics* **33** 2789–2810.

SÄRNDAL, C.-E., SWENSSON, B. and WRETMAN, J. (2003). *Model assisted survey sampling*. Springer Science & Business Media.

SKOGAN, W. G. (1974). The validity of official crime statistics: An empirical investigation. *Social Science Quarterly* 25–38.

SKOGAN, W. G. (1977). Dimensions of the dark figure of unreported crime. *Crime & Delinquency* **23** 41–50.

SOLON, G., HAIDER, S. J. and WOOLDRIDGE, J. M. (2015). What are we weighting for? *Journal of Human resources* **50** 301–316.

STEFFENSMEIER, D., FELDMEYER, B., HARRIS, C. T. and ULMER, J. T. (2011). Reassessing trends in black violent crime, 1980–2008: Sorting out the "Hispanic effect" in Uniform Crime Reports arrests, National Crime Victimization Survey offender estimates, and US prisoner counts. *Criminology* **49** 197–251.

SUGIYAMA, M., KRAULEDAT, M. and MÜLLER, K.-R. (2007). Covariate shift adaptation by importance weighted cross validation. *Journal of Machine Learning Research* **8**.

THOMPSON, M. (1997). *Theory of sample surveys* **74**. CRC Press.

TIBSHIRANI, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)* **58** 267–288.

UNITED STATES DEPARTMENT OF JUSTICE, F. B. O. I. (2008a). National Incident-Based Reporting System, 2006. https://doi.org/10.3886/ICPSR22407.v1

UNITED STATES DEPARTMENT OF JUSTICE, F. B. O. I. (2008b). Uniform Crime Reporting Program Data [United States]: Police Employee (LEOKA) Data, 2006. https://doi.org/10.3886/ICPSR22402.v1

UNITED STATES DEPARTMENT OF JUSTICE, F. B. O. I. (2009a). National Incident-Based Reporting System, 2007. https://doi.org/10.3886/ICPSR25113.v1

UNITED STATES DEPARTMENT OF JUSTICE, F. B. O. I. (2009b). Uniform Crime Reporting Program Data [United States]: Police Employee (LEOKA) Data, 2007. https://doi.org/10.3886/ICPSR25104.v1

UNITED STATES DEPARTMENT OF JUSTICE, F. B. O. I. (2010a). National Incident-Based Reporting System, 2008. https://doi.org/10.3886/ICPSR27647.v1

UNITED STATES DEPARTMENT OF JUSTICE, F. B. O. I. (2010b). Uniform Crime Reporting Program Data [United States]: Police Employee (LEOKA) Data, 2008. https://doi.org/10.3886/ICPSR27646.v1

UNITED STATES DEPARTMENT OF JUSTICE, F. B. O. I. (2011a). Uniform Crime Reporting: National Incident-Based Reporting System, 2009. https://doi.org/10.3886/ICPSR30770.v1

UNITED STATES DEPARTMENT OF JUSTICE, F. B. O. I. (2011b). Uniform Crime Reporting Program Data [United States]: Police Employee (LEOKA) Data, 2009. https://doi.org/10.3886/ICPSR30765.v1

UNITED STATES DEPARTMENT OF JUSTICE, F. B. O. I. (2012a). Uniform Crime Reporting: National Incident-Based Reporting System, 2010. https://doi.org/10.3886/ICPSR33530.v1

UNITED STATES DEPARTMENT OF JUSTICE, F. B. O. I. (2012b). Uniform Crime Reporting Program Data: Police Employee (LEOKA) Data, 2010. https://doi.org/10.3886/ICPSR33525.v1

UNITED STATES DEPARTMENT OF JUSTICE, F. B. O. I. (2013a). Uniform Crime Reporting Program Data: National Incident-Based Reporting System, 2011. https://doi.org/10.3886/ICPSR34585.v1

UNITED STATES DEPARTMENT OF JUSTICE, F. B. O. I. (2013b). Uniform Crime Reporting Program Data: Police Employee (LEOKA) Data, 2011. https://doi.org/10.3886/ICPSR34584.v1

UNITED STATES DEPARTMENT OF JUSTICE, F. B. O. I. (2014a). Uniform Crime Reporting Program Data: National Incident-Based Reporting System, 2012. https://doi.org/10.3886/ICPSR35035.v1

UNITED STATES DEPARTMENT OF JUSTICE, F. B. O. I. (2014b). Uniform Crime Reporting Program Data: Police Employee (LEOKA) Data, 2012. https://doi.org/10.3886/ICPSR35020.v1

UNITED STATES DEPARTMENT OF JUSTICE, F. B. O. I. (2015a). Uniform Crime Reporting Program Data: National Incident-Based Reporting System, 2013. https://doi.org/10.3886/ICPSR36120.v2

UNITED STATES DEPARTMENT OF JUSTICE, F. B. O. I. (2015b). Uniform Crime Reporting Program Data: Police Employee (LEOKA) Data, 2013. https://doi.org/10.3886/ICPSR36119.v1

UNITED STATES DEPARTMENT OF JUSTICE, F. B. O. I. (2016a). Uniform Crime Reporting Program Data: National Incident-Based Reporting System, 2014. https://doi.org/10.3886/ICPSR36398.v1

UNITED STATES DEPARTMENT OF JUSTICE, F. B. O. I. (2016b). Uniform Crime Reporting Program Data: Police Employee (LEOKA) Data, 2014. https://doi.org/10.3886/ICPSR36395.v1

UNITED STATES DEPARTMENT OF JUSTICE, B. O. J. S. (2017a). National Crime Victimization Survey, 2016. Technical Documentation.

UNITED STATES DEPARTMENT OF JUSTICE, F. B. O. I. (2017b). Uniform Crime Reporting Program Data: National Incident-Based Reporting System, 2015. https://doi.org/10.3886/ICPSR36795.v1

UNITED STATES DEPARTMENT OF JUSTICE, F. B. O. I. (2017c). Uniform Crime Reporting Program Data: Police Employee (LEOKA) Data, 2015. https://doi.org/10.3886/ICPSR36791.v1

UNITED STATES DEPARTMENT OF JUSTICE, F. B. O. I. F. (2019). 2019 National Incident-Based Reporting System User Manual.

UNITED STATES DEPARTMENT OF JUSTICE, B. O. J. S. (2021). National Crime Victimization Survey, Concatenated File, [United States], 1992-2020. https://doi.org/10.3886/ICPSR38136.v1

VAN DER HEIJDEN, P. G., BUSTAMI, R., CRUYFF, M. J., ENGBERSEN, G. and VAN HOUWELINGEN, H. C. (2003). Point and interval estimation of the population size using the truncated Poisson regression model. *Statistical Modelling* **3** 305–322.

VAN DER LAAN, M. J., POLLEY, E. C. and HUBBARD, A. E. (2007). Super learner. *Statistical applications in genetics and molecular biology* **6**.

VAN DER VAART, A. W. (2000). *Asymptotic statistics* **3**. Cambridge university press.

WHITE, H. (2014). *Asymptotic theory for econometricians*. Academic press.

XIE, M. and BAUMER, E. P. (2019a). Neighborhood immigrant concentration and violent crime reporting to the police: A multilevel analysis of data from the National Crime Victimization Survey. *Criminology* **57** 237–267.

XIE, M. and BAUMER, E. P. (2019b). Crime victims' decisions to call the police: Past research and new directions. *Annual Review of Criminology*.

XIE, M. and LAURITSEN, J. L. (2012). Racial context and crime reporting: A test of Black's stratification hypothesis. *Journal of quantitative criminology* **28** 265–293.

XIE, M. and LYNCH, J. P. (2017). The effects of arrest, reporting to the police, and victim services on intimate partner violence. *Journal of research in crime and delinquency* **54** 338–378.

YUAN, K.-H. and JENNRICH, R. I. (1998). Asymptotics of estimating equations under natural conditions. *Journal of Multivariate Analysis* **65** 245–260.