

The Population and Personalized Areas Under the Receiving Operating Characteristic Curve

Haben Michael¹ and Lu Tian²

¹Department of Mathematics and Statistics, University of Massachusetts

²Department of Biomedical Data Science, Stanford University

ABSTRACT: We consider two generalizations of the area under the receiving operating characteristic curve (“AUC”), a popular measure of discrimination, to accommodate clustered data. We describe situations in which the two cluster AUCs diverge and other situations in which they coincide. Differences are described using concrete models and visualizations, while quantitative results are used to relate the two generalizations. Procedures for joint estimation and inference are also presented, along with a simulation study. We apply the results to data collected on urban policing behavior.

KEYWORDS: AUC, Confounding, Clustered data, Simpson’s paradox

1 Introduction

The AUC is a widely used measure of how well a scalar predictor discriminates between two outcomes. As a population parameter, the AUC is the probability that the value of a randomly sampled predictor from one of the outcome classes is less than an independently sampled predictor from the other outcome class. There are several ways to generalize the AUC to accommodate clustered data. What we refer to as the “population AUC” appears to be the most commonly studied. The population AUC evaluates the predictor’s typical effect on an entire population, as further discussed below.

While the population AUC is an important part of understanding the usefulness of a predictor, the medical field has lately focused on personalizing treatment. For example, in 2018 the National Academy of Medicine concluded: “The individuality of the patient should be at the core of every treatment decision. One-size-fits-all approaches to treating medical

conditions are inadequate; instead, treatments should be tailored to individuals based on heterogeneity of clinical characteristics and their personal preferences.”

We examine a “personalized AUC” in conjunction with the population AUC. These two evaluations may give different accounts of the usefulness of a marker. In the extreme case, the phenomenon known as Simpson’s paradox may occur: The personalized AUC may be nearly uninformative while the population AUC is nearly perfectly predictive, or vice versa. Modern accounts of Simpson’s paradox, working in the framework of causal inference, delineate situations in which the personalized AUC is appropriate, and other situations in which the population AUC is appropriate.

Previous Literature. Obuchowski (1997) proposes a nonparametric, asymptotic estimator for the variance of an estimator for the population AUC. We give an alternate derivation here. We also clarify the statistical model and target of inference. Rosner and Grove (1999) give a formula for the finite-sample variance under certain distributional assumptions. Lee and Dehling (2005) discuss the asymptotic behavior of generalized U-statistics with clustered data, a class which includes the estimator for the population AUC discussed in Obuchowski (1997) and below. Liu et al. (2005) suggest a bootstrap approach, provided the data follow a generalized linear model. Toledano (2003) summarizes techniques in use for analyzing clustered ROC curves. Pearl (2014) gives an overview of Simpson’s paradox, an effect illustrated by the examples in Section 3 below, from the standpoint of causal inference. Michael et al. (2019) analyzes population and personalized versions of the ROC curve which may, in principle, be used to obtain estimates of the population and personalized AUCs discussed here under certain distributional assumptions. The analysis here is nonparametric and further avoids the inefficiency introduced by first estimating the entire ROC curve, which may not be feasible for many smaller data sets.

Organization of the article. In Section 2 we introduce the AUC and the two generalizations to clustered data considered here, the population and personalized AUCs. In Section 3 we discuss several examples of data-generating processes that highlight the differences be-

tween the two AUCs. In Section 4 we describe two results available under the assumption of independence between the cluster sizes and the values of the predictor of interest. In Section 5 we give the asymptotic joint distribution of estimators for the two AUCs. The performance of these estimators are then analyzed using synthetic data generated under two models in Section 6. Section 7 contains an application of the methods to data on urban policing patterns. In Section 8 we conclude and give directions for future work. Software implementing the methods used in the paper and code for replicating the tables and figures are publicly available at the first author’s website.

2 Setting and Notation

Let $\psi : \mathbb{R}^2 \rightarrow \mathbb{R}$ denote the function $(x, y) \mapsto \{x < y\} + \frac{1}{2}\{x = y\}$, using $\{\cdot\}$ to denote the indicator function. Given independent draws X and Y from two distributions F_X and F_Y , the AUC is defined as

$$\theta = E(\psi(X, Y)) = P(X < Y) + \frac{1}{2}P(X = Y).$$

Given samples X_1, X_2, \dots, X_M , IID as F_X and Y_1, Y_2, \dots, Y_N , IID as F_Y , an unbiased estimator of the AUC of F_X and F_Y is

$$\frac{1}{MN} \sum_{i=1}^M \sum_{j=1}^N \psi(X_i, Y_j).$$

The function ψ is referred to as the “kernel.” The AUC is often used to evaluate how effectively the data distinguish the two distributions. The AUC is close to 1/2 when the distinction is poor, and equals 1/2 in the extreme case that $F_X = F_Y$. The AUC is close to 1 when the distinction is better. In this extreme, there is a number $c \in \mathbb{R}$ such that always $X < c$ and $Y > c$, and then $\theta = 1$. We informally refer to the two classes given by the two distributions as “control” and “case,” and the scalar predictor as “marker.” Switching the

observations designated “control” and “case” reflects the AUC across $1/2$, $AUC \mapsto 1 - AUC$, so $|AUC - \frac{1}{2}|$ is often of greater interest than the AUC itself.

We extend the AUC to accommodate 1) vectors of case and control observations and 2) dependence between case and control observations. Examples of data of this type are:

1. The predictors are longitudinal measurements of tumor antigens (CEA, CA15-3, TPS), and the outcomes are progression or non-progression of breast cancer (Emir et al., 2000).
2. The predictors are longitudinal measurements of levels of vascular endothelial growth factor and a soluble fragment of Cytokeratin 19, and the outcomes are progression or non-progression of non-small cell lung cancer (Wu and Wang, 2011).
3. The predictors are longitudinal measurements of an HIV positive patient’s CD4 counts, and the outcome is “blip” status, a binary measurement representing a transient spike in viral load (Michael et al., 2019).

Let (X, Y, M, N) be a random vector with joint distribution P such that X and Y are sequences and M and N are counting numbers.

$$\begin{aligned} (X, Y, M, N) &\sim P \\ X = (X_1, X_2, \dots) &\in \mathbb{R}^N, Y = (Y_1, Y_2, \dots) \in \mathbb{R}^N \\ M, N &\in 1, 2, 3, \dots, E(M) < \infty, E(N) < \infty. \end{aligned} \tag{1}$$

Informally, we regard X and Y as vectors of length M and N , ignoring the rest of the sequences. The formulation (1) lets us avoid working with vectors of variable length.

Extend the AUC kernel $\psi(\cdot, \cdot)$ to vector arguments as

$$\psi(x, y) = \psi((x_1, \dots, x_m), (y_1, \dots, y_n)) = \sum_{i=1}^m \sum_{j=1}^n \left(\{x_i < y_j\} + \frac{1}{2}\{x_i = y_j\} \right). \tag{2}$$

We define the personalized AUC as

$$\theta_{11}(P) = E\left(\frac{\psi(X, Y)}{MN}\right). \quad (3)$$

With (X_1, Y_1, M_1, N_1) and (X_2, Y_2, M_2, N_2) , being two independent draws from P , we define the population AUC as

$$\theta_{12}(P) = \frac{E\psi(X_1, Y_2)}{E(M_1)E(N_2)} \quad (4)$$

$$(X_1, Y_1, M_1, N_1), (X_2, Y_2, M_2, N_2) \stackrel{\text{iid}}{\sim} P.$$

The personalized AUC may be undefined if M or N can take the value 0 with positive probability, which is the reason for restricting them to counting numbers. The population AUC may still be well-defined and some analyses do allow $M = 0$ or $N = 0$ (Obuchowski, 1997). In applications where $M = 0$ or $N = 0$ is possible, our analysis is therefore conditional on $M > 0, N > 0$, a sub-population in which all clusters have at least 1 case and 1 control observation.

For estimation, suppose a sample is given,

$$(X_1, Y_1, M_1, N_1), \dots, (X_I, Y_I, M_I, N_I) \stackrel{\text{iid}}{\sim} P.$$

An unbiased estimator of θ_{11} is

$$\hat{\theta}_{11} = \frac{1}{I} \sum_{i=1}^I \frac{\psi(X_i, Y_i)}{M_i N_i}.$$

A consistent estimator of θ_{12} is

$$\hat{\theta}_{12} = \frac{\sum_i \sum_{i \neq j} \psi(X_i, Y_j)}{\sum_i M_i \sum_i N_i}. \quad (5)$$

Letting P_I denote the empirical distribution of the sample, $\hat{\theta}_{11} = \theta_{11}(P_I)$, while $\hat{\theta}_{12} =$

$\theta_{12}(P_I) + O(I)$ (the estimator $\theta_{12}(P_I)$ is discussed at (7) below).

Both the population and personalized AUC, like the usual AUC, are bounded between 0 and 1, $\frac{1}{2}$ represents poor discrimination, and distance from $\frac{1}{2}$ represents increasing discrimination. However, they describe distinct measures of discrimination. It is possible for one to be informative and therefore far from $1/2$, while the other is non-informative, or close to $1/2$. Whereas the personalized AUC is the average AUC of a typical cluster, the population AUC is, setting aside ties in the data, the probability that a typical control observation in the population is less than a typical case observation. The following proposition makes this description precise. The consistency of $\hat{\theta}_{12}$ follows from Corollary 8.

Proposition 1. *1. Let $(X_1, Y_1, M_1, N_1), \dots, (X_I, Y_I, M_I, N_I)$, be a random sample of size I IID according to P . Let P_I be the joint distribution of independent random selections from among the elements of X_1, \dots, X_I , and Y_1, \dots, Y_I , and let $(\xi_I, \eta_I) \sim P_I$. Then $\theta(P_I) = Pr(\xi_I < \eta_I) + \frac{1}{2}Pr(\xi_I = \eta_I) \rightarrow \theta_{12}(P)$ as $I \rightarrow \infty$.*

2. Let ξ follow the distribution of $X_{1i} \mid M = m$ with probability $P(M = m)/E(M)$, $i = 1, \dots, m, m = 1, 2, \dots$, and, independently, let η follow the distribution of $Y_{2j} \mid N = n$ with probability $P(N = n)/E(N)$, $j = 1, \dots, n, n = 1, 2, \dots$. Then $\theta_{12}(P) = Pr(\xi < \eta) + \frac{1}{2}Pr(\xi = \eta)$.

The definition of the population AUC (4) allows for dependence between (M, N) and (X, Y) in capturing a population-level AUC in the sense of Proposition 1. Practical reasons to avoid assuming $(X, Y) \perp\!\!\!\perp (M, N)$ include informative censoring, imbalanced designs, and confounding by indication; Further examples are given in Benhin et al. (2005) and Bugni et al. (2022). As an alternative definition of the population AUC, consider

$$\theta'_{12} = E \left(\frac{\psi(X_1, Y_2)}{M_1 N_2} \right). \quad (6)$$

This parameter is formally a closer counterpart to the personalized AUC (3), but does not take into account different cluster sizes, with a small cluster contributing as much as a large

cluster. This estimator would not therefore represent discrimination at the population level, except in case $(X, Y) \perp\!\!\!\perp (M, N)$.

Similar to the population AUC estimator (5), Obuchowski (1997) presents the estimator

$$\frac{\sum_i \sum_{j \neq i} \psi(X_i, Y_j)}{\sum_i M_i \sum_i N_i} = \hat{\theta}_{12} + \frac{\sum_i \psi(X_i, Y_i)}{\sum_i M_i \sum_i N_i}. \quad (7)$$

This estimator may be obtained as $\theta_{12}(P_I)$, where P_I is the empirical distribution given a sample of size I . It differs from $\hat{\theta}_{12}$ only in including the diagonal terms, an asymptotically negligible $O(1/I)$ bias. The definition (4) was chosen in part as the probability limit of (7). Though Obuchowski (1997) does not enunciate a clear statistical model, the analysis of (7) rather than the simpler (6) perhaps suggests that Obuchowski (1997) too contemplates $(X, Y) \not\perp\!\!\!\perp (M, N)$.

The population AUC, which appears more prominently in past research, may lay a claim to being the more natural generalization of the usual AUC since it equals the usual AUC when $M = N = 1$. Below we argue that in general the population and personalized AUCs are both important, complementary tools in evaluating an estimator. In the other direction, we give inequalities that may be used in some situations to relate the two cluster AUCs.

3 Examples

3.1 Random effects model

We illustrate the population and personalized AUCs and their differences using a generic random effects model with a location shift parameter. We show that the location shift can be used to control the within-cluster informativity of the observations, thereby controlling the personalized AUC, while separately the random effect can be used to control informativity across clusters, controlling the population AUC. Real data illustrating these contrasts, in less dramatic fashion than the artificial examples constructed here, are presented in Section

7.

Let the distribution of (X, Y, M, N) given M, N be

$$\begin{aligned} X | M, N &\sim Z(M, N) + \xi_i^x, i = 1, \dots, M \\ Y | M, N &\sim Z(M, N) + \xi_j^y + \Delta, j = 1, \dots, N \end{aligned} \quad (8)$$

Here, $\Delta > 0$ is a non-random location shift between the control and case values, Z is a random, cluster-level effect, and $\xi_i^x, \xi_j^y, i = 1, \dots, M, j = 1, \dots, N$, are IID individual effects. The within-cluster dependence is induced by Z . The individual effects ξ_i^x, ξ_j^y are assumed to be independent of (M, N) , but Z is not assumed to be so. To keep things simple, we assume continuous densities are available, and so $\psi(x, y) = \{x < y\}$.

The personalized AUC is

$$\begin{aligned} \theta_{11} &= E \left(\frac{\psi_{11}}{M_1 N_1} \right) = E \left(\frac{1}{M_1 N_1} \sum_{i=1}^{M_1} \sum_{j=1}^{N_1} \{X_{1i} < Y_{1j}\} \right) \\ &= E \left(\frac{1}{M_1 N_1} \sum_{i=1}^{M_1} \sum_{j=1}^{N_1} \{Z_1 + \xi_i^x < Z_1 + \xi_j^y + \Delta\} \right) \\ &= E \left(\frac{1}{M_1 N_1} \sum_{i=1}^{M_1} \sum_{j=1}^{N_1} P(\xi_i^x - \xi_j^y < \Delta | M_1, N_1) \right) \\ &= P(\xi_1 - \xi_2 < \Delta). \end{aligned} \quad (9)$$

Lemma 9 was used to pull the conditional expectation inside the double sum.

The population AUC is

$$\begin{aligned}
\theta_{12} &= \frac{1}{E(M)E(N)} E \left(\sum_{i=1}^{M_1} \sum_{j=1}^{N_2} \{X_{1i} < Y_{2j}\} \right) \\
&= \frac{1}{E(M)E(N)} E \left(\sum_{i=1}^{M_1} \sum_{j=1}^{N_2} P(Z_1 + \xi^x < Z_2 + \xi^y + \Delta \mid M_1, N_1, M_2, N_2) \right) \\
&= \frac{1}{E(M)E(N)} E (M_1 N_2 P(Z_1 + \xi^x < Z_2 + \xi^y + \Delta \mid M_1, N_1, M_2, N_2)) \\
&= E \left(\frac{M_1 N_2}{E(M)E(N)} \{Z_1 - Z_2 + (\xi^x - \xi^y) < \Delta\} \right) \tag{10}
\end{aligned}$$

The last expression is a covariance-like term lying between 0 and 1.

Informative personalized AUC, uninformative population AUC

From (9), $\theta_{11} \rightarrow 1$ as $\Delta \rightarrow \infty$.

Let Z be independent of $(M, N, \vec{\xi}^x, \vec{\xi}^y)$. Then $\theta_{12} = P(Z_1 + \xi^x - (Z_2 + \xi^y) < \Delta)$. As a difference of two IID random variables, $Z_1 + \xi^x - (Z_2 + \xi^y)$ is symmetric about 0, and $\theta_{12} = P(Z_1 + \xi^x - (Z_2 + \xi^y) < \Delta) = 1 - P(Z_1 + \xi^x - (Z_2 + \xi^y) \geq \Delta) = 1 - 1/2P(|Z_1 + \xi^x - (Z_2 + \xi^y)| \geq \Delta)$. As $P(|Z_1 + \xi^x - (Z_2 + \xi^y)| \geq \Delta) \geq 1 - 2\Delta|f_{Z+\xi}|_\infty \geq 1 - 2\Delta|f_Z|_\infty$, a sufficient condition for $\theta_{12} \rightarrow 1/2$ is $|f_Z|_\infty \rightarrow 0$. For example, suppose Z belongs to a scale family, $f_Z = f_{Z_0}(Z/\sqrt{\text{Var}(Z)})/\sqrt{\text{Var}(Z)}$ for a fixed density f_{Z_0} , $|f_{Z_0}|_\infty < \infty$, and $\text{Var}(Z) \rightarrow \infty$.

Therefore, for $\Delta = E(Y_{11}) - E(X_{11})$ large enough, θ_{11} is arbitrarily close to 1, while for any fixed Δ , for $\text{Var}(Z)$ large enough, θ_{12} may approach 1/2.

Informative population AUC, uninformative personalized AUC

From (9), $\theta_{11} \rightarrow 1/2$ as $\Delta \rightarrow 0$, $\xi^x - \xi^y$ being symmetric about 0.

The covariance-like term (10) may approach 1 when there is a large negative covariance between M, N , and $Z_1 - Z_2$, i.e., a large negative covariance between M and Z or large positive covariance between N and Z , or both. Suppose (1) $\text{Corr}(M_1 N_2, Z_1 - Z_2)$ is close to -1 , (2) $\text{Var}(M_1 N_2)$ is large, (3) $\Delta \downarrow 0$, (4) the counts M, N are bounded, and (5) $\text{Var}(\xi)$ is small. For Δ sufficiently small, $P(Z_1 - Z_2 + (\xi^x - \xi^y) < \Delta) > 1/2 - \epsilon$, since $Z_1 - Z_2 + (\xi^x - \xi^y)$ is a symmetric RV. As $\text{Var}(M_1 N_2)$ maxes out, by (4) $M_1 N_2$ approximates a balanced two-

point distribution, so $P(M_1N_2 > 2E(M_1N_2) - \epsilon) > 1/2 - \epsilon$ and $P(M_1N_2 < \epsilon) > 1/2 - \epsilon$. As $\text{Corr}(M_1N_2, Z_1 - Z_2) \rightarrow -1$ and $\text{Var}(\xi) \rightarrow 0$, $\text{Corr}(M_1N_2, \{Z_1 - Z_2 + (\xi^x - \xi^y)\}) \rightarrow -1$. As $\text{Corr}(M_1N_2, Z_1 - Z_2 + (\xi^x - \xi^y))$ approaches perfect negative linearity, $P(\{Z_1 - Z_2 + (\xi^x - \xi^y) < \Delta\} \cap \{M_1N_2 > 2E(M_1N_2) - \epsilon\}) > P(M_1N_2 > 2E(M_1N_2) - \epsilon) - \epsilon > 1/2 - 2\epsilon$. Therefore,

$$\begin{aligned}
\theta_{12} &= E \left(\frac{M_1N_2}{E(M)E(N)} \{Z_1 - Z_2 + (\xi^x - \xi^y) < \Delta\} \right) \\
&> \frac{1}{E(M)E(N)} P(\{Z_1 - Z_2 + (\xi^x - \xi^y) < \Delta\} \cap \{M_1N_2 > 2E(M_1N_2) - \epsilon\}) \\
&\quad \times (2E(M_1N_2) - \epsilon) \\
&> \frac{1}{E(M)E(N)} (1/2 - 2\epsilon)(2E(M_1N_2) - \epsilon) \\
&= 1 - o(1).
\end{aligned}$$

Figure 1 presents a simulation using Gaussian data to demonstrate the discussed differences between the population and personalized AUCs. The model is an example of the random effects model (8) and is discussed further in Section 6. Though a large location shift can push the personalized AUC close to 1, large inter-cluster variance relative to intra-cluster variance keeps the population AUC uninformative. Similarly, if the number of case observations relative to control is positively associated with the observation values, the population AUC may approach 1 irrespective of the personalized AUC.

3.2 Binary response model

Models for case and control data are often given by specifying the status conditional on the marker, rather than vice versa as in Example 3.1. Let σ denote a monotone link such as the probit or logistic function. Fixing the cluster size $M + N = k$, let continuous cluster effects Z , continuous within-cluster effects ξ , and within-cluster status indicators D specify

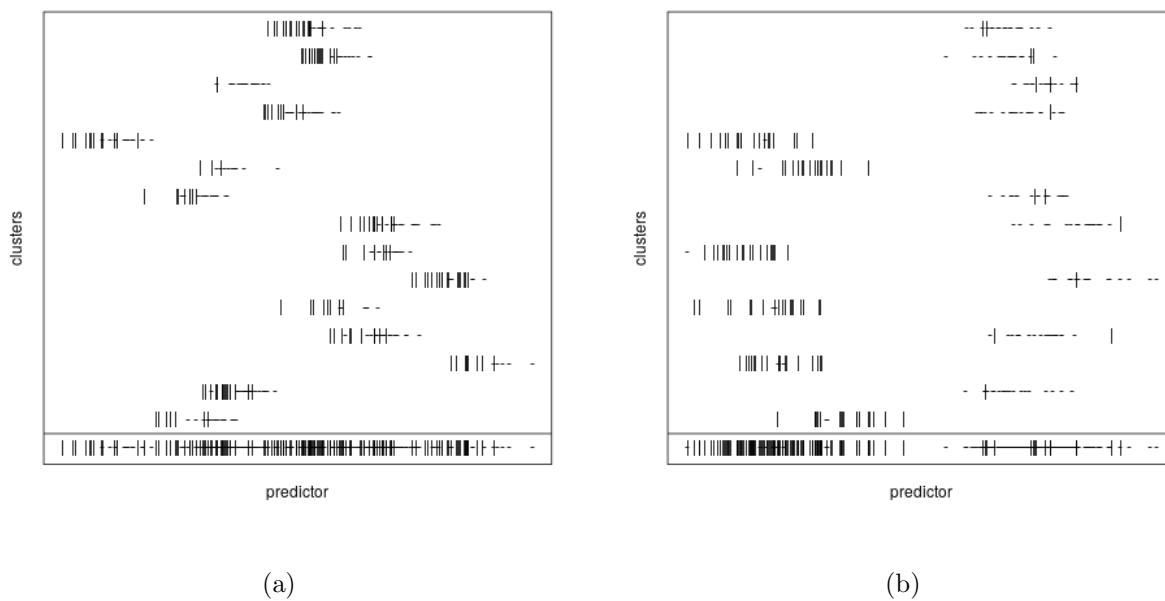


Figure 1: Two visualizations contrasting the personalized and population AUCs. Each gives rug plots of fifteen clusters of data, each cluster sampled IID according to a bivariate normal model, with the unclustered data combined at the bottom. Case observations are represented with “-” and control observations with “|”. On the left, the personalized AUC is informative and the population AUC uninformative. The reverse situation is presented on the right.

the distribution of a cluster as follows:

$$\begin{aligned}\vec{\xi} &= (\xi_1, \dots, \xi_k) \text{ IID} \\ Z &\perp\!\!\!\perp (\xi_1, \dots, \xi_k) \\ B_i &= Z + \xi_i, i = 1, \dots, k \\ D_i | \vec{Z}, \vec{\xi} &\sim \text{bernoulli with parameter } \sigma(\beta_0 Z + \beta_1 \xi_i), i = 1, \dots, k \\ M &= \sum_{i=1}^k (1 - D_i), \quad N = \sum_{i=1}^k D_i.\end{aligned}$$

The control and case observations in a cluster, X_i and Y_i , are then those B_i such that $D_i = 0$ and $D_i = 1$, respectively. Here the cluster allocations M and N and the markers \vec{B} can be dependent, both being functions of Z and $\vec{\xi}$, though they are conditionally independent given the statuses \vec{D} .

Suppose first that $\beta_0 > 0$ and $\beta_1 = 0$, so $P(D_i = 1 | Z, \vec{\xi}) = \sigma(\beta_0 Z)$. The population AUC is

$$\begin{aligned}\theta_{12} &= \frac{1}{E(M)E(N)} E \left(\sum_{i=1}^k \sum_{j=1}^k \{B_{1i} < B_{2j}\} \{D_{1i} = 0\} \{D_{2j} = 1\} \right) \\ &= \frac{1}{E(M)E(N)} E \left(\sum_{i=1}^k \sum_{j=1}^k P(B_{1i} < B_{2j} | D_{1i} = 0, D_{2j} = 1) \{D_{1i} = 0\} \{D_{2j} = 1\} \right) \\ &= P(B_{11} < B_{21} | D_{11} = 0, D_{21} = 1) \\ &= P(Z_{11} - Z_{21} < \xi_{21} - \xi_{11} | D_{11} = 0, D_{21} = 1).\end{aligned}$$

Since $Z_{11} | D_{11} = 0$ is stochastically less than $Z_{21} | D_{21} = 1$, with the difference increasing in β_0 , and since the ξ s are independent of the Z s and D s, the last line is $> \frac{1}{2}$, with the difference increasing in β_0 .

On the other hand, since $\beta_1 = 0$ implies $\vec{\xi} \perp\!\!\!\perp (\vec{D}, M, N)$, the personalized AUC is

$$\begin{aligned} \theta_{11} &= E \left(\frac{1}{MN} \sum_{i=1}^k \sum_{j=1}^k \{B_{1i} < B_{1j}\} \{D_{i1} = 0 \text{ and } D_{ij} = 1\} \mid M > 0, N > 0 \right) \\ &= E \left(\frac{1}{MN} \sum_{i=1}^k \sum_{j=1}^k \{\xi_{1i} < \xi_{1j}\} \{D_{i1} = 0 \text{ and } D_{ij} = 1\} \mid M > 0, N > 0 \right) \\ &= P(\xi_{11} < \xi_{12}) = 1/2. \end{aligned}$$

Two possible instances of the model:

- (a) The cluster effect Z represents a genuine signal of disease status D , such as viral load does for HIV status, and ξ represents non-systematic measurement error on instruments measuring Z . In this case, the population AUC better matches expectations of an AUC measurement than the personalized AUC. The biomarker B isn't completely uninformative, as θ_{11} suggests.
- (b) The cluster effect Z is a subject's dose of a possibly ineffective drug, and larger doses are administered to sicker patients. The subject-specific measurements ξ represent non-systematic measurement error again. Here the association between the marker and disease status implied by the population AUC is spurious, and may or may not be of value to the analyst. It is possible that the personalized AUC, which does not convey any association, is preferable.

Reversing the roles of the cluster-level effect Z and within-cluster effects ξ , suppose $\beta_0 = 0$ and $\beta_1 > 0$, so that $\theta_{12} \approx 1/2$ and $\theta_{11} > 1/2$. Two instances of this second model are:

- (c) The markers B are measurements on a patient, and D denotes the presence of a disease that depends little or not at all on a baseline measure Z but is indicated by the deviations ξ from the baseline. As a second example, the markers B are post-test measurements on a population that has been stratified by pre-test measurement Z . The subject effects $\xi_i = B_i - Z$ represent the difference between post-test and pre-

test measurements, and the status indicators D represent an effective or ineffective intervention. Here the personalized AUC probably carries the correct interpretation.

- (d) A population clustered along any given dimension Z , and, analogous to (b), uptake of a possibly ineffective drug is confounded by indication. That is, sicker individuals, those for which D_i is more likely to be 1, take higher doses ξ_i of the drug. Here again a causal analysis would suggest the population AUC as less misleading than the personalized AUC, though a non-causal analysis, e.g., an intention-to-treat analysis, may point to the personalized AUC.

Simpson’s paradox, understood broadly, refers to situations where data is clustered and exhibits a consistent trend at each cluster, but exhibits a contrary trend when the unclustered data is analyzed. The examples in Section 3.1 are instances of this phenomenon. The individual and population AUCs are clustered and unclustered analyses that can yield opposite conclusions about the quality of the predictor. Contemporary analyses of Simpson’s paradox show the importance of considering both the individual and population AUCs.

Working in the framework of causal inference, Pearl (2014) argues that the paradox arises from the subtle relationship between causal intervention and statistical conditioning. Human judgments, which align more closely with causal relations, may be contradicted by one of the analyses when it represents a non-causal association. Resolution of the paradox therefore amounts to formally identifying which of the two analyses represents causal relationships, if either. The correct analysis in any given situation, whether the clustered or unclustered analysis, requires information about the underlying causal relationships between the treatment, outcome, and clustering variable.

4 Simplifications when $(X, Y) \perp\!\!\!\perp (M, N)$

Under some conditions, the cluster AUC parameters θ_{12} and θ_{11} may simplify to the $M = N = 1$ case. The examples given in Section 3 are of this sort. The exchangeable cluster

structure enables the simplification.

Proposition 2. *Given $(X, Y, M, N) \sim P$, suppose that $E(\psi(X_{1k}, Y_{1l}) \mid M, N)$ and $E(\psi(X_{1k}, Y_{2l}) \mid M, N)$ do not depend on k, l . Then $\theta_{11}(P) = E\psi(X_{11}, Y_{11})$ and $\theta_{12}(P) = E\psi(X_{11}, Y_{21})$.*

In order for $\hat{\theta}_{12} \rightarrow 1$ while $\hat{\theta}_{11} \not\rightarrow 1$ in the random effects model discussed in Section 3, it was necessary that $(X, Y) \not\perp\!\!\!\perp (M, N)$. Theorem 3 below bounds θ_{12} by θ_{11} in one situation where $(X, Y) \perp\!\!\!\perp (M, N)$, namely, when M and N are each constant.

Theorem 3. *Let $(X, Y, M, N) \sim P$ be given as in (1). Assume further that $M = m$ and $N = n$ are constant. Then*

$$\frac{1}{2} \left(\theta_{11} + \frac{\sum_{k,l} P(X_{1k} = Y_{1l})}{2mn} \right)^2 \leq \theta_{12} \leq 1 - \frac{1}{2} \left(1 - \theta_{11} + \frac{\sum_{k,l} P(X_{1k} = Y_{1l})}{2mn} \right)^2.$$

The theorem follows from the lemma,

Lemma 4. *Given a pair of scalar random variables (X, Y) with joint distribution P , let $P_{\perp\!\!\!\perp}$ be the product measure of the marginals, i.e., for all real a, b ,*

$$P_{\perp\!\!\!\perp}(\{x < a\} \cap \{y < b\}) = P(\{x < a\})P(\{y < b\}).$$

Then

$$\frac{1}{2}(P(X < Y) + P(X = Y))^2 \leq P_{\perp\!\!\!\perp}(X < Y) + \frac{1}{2}P_{\perp\!\!\!\perp}(X = Y) \leq 1 - \frac{1}{2}(1 - P(X < Y))^2.$$

Let the random vector (X, Y, M, N) follow P with constant $M = N = 1$, so that P may be regarded as the joint distribution of (X, Y) , assumed continuous. Then the conclusion of the Lemma is

$$\frac{1}{2}(\theta_{11}(P))^2 \leq \theta_{12}(P) \leq 1 - \frac{1}{2}(1 - \theta_{11}(P))^2, \quad (11)$$

equivalently,

$$1 - \sqrt{2(1 - \theta_{12})} \leq \theta_{11} \leq \sqrt{2\theta_{12}}.$$

When the personalized AUC is completely uninformative, $\theta_{11} = 1/2$, the informativity of the population AUC is limited, $1/8 \leq \theta_{12} \leq 7/8$. However, when the population AUC is completely uninformative, $\theta_{12} = 1/2$, the above bounds on the personalized AUC, which are tight, are vacuous, $0 \leq \theta_{11} \leq 1$. Situations as described in Section 3, where the population AUC $\rightarrow 1$ while the personalized AUC $\rightarrow 1/2$, appear to require some dependence between M, N and X, Y .

5 Asymptotic Distribution of $(\theta_{12}, \theta_{11})$

Theorem 5 gives the asymptotic joint distribution of the individual and population AUCs. It is stated in somewhat greater generality for any square-integrable kernel, not just the AUC kernel (2). The proof is the same for any random variables M, N , such that $EM \neq 0, EN \neq 0, EM^{-2} < \infty, EN^{-2} < \infty$, i.e., M and N need not be the lengths of X and Y . Let V denote the space of finite sequences.

Theorem 5. *Let $\psi : V \times V \rightarrow \mathbb{R}$, $(X, Y, M, N) \sim P$ with $(X, Y) \in V \times V$, $\psi \in L^2(P)$, and let M and N be counting numbers > 0 with finite means. Given a sample $W_i = (X_i, Y_i, M_i, N_i), i = 1, \dots, I$, from P ,*

$$\sqrt{I}(\hat{\theta}_{12} - \theta_{12}, \hat{\theta}_{11} - \theta_{11}) \rightsquigarrow \mathcal{N}(0, \Sigma)$$

with

$$\begin{aligned}\Sigma_{11} &= \lim_{I \rightarrow \infty} I \text{Var}(\hat{\theta}_{12}) = E \left(\frac{E(\psi_{12} | W_1) + E(\psi_{21} | W_1)}{E(M)E(N)} - \theta_{12} \left(\frac{M_1}{E(M)} + \frac{N_1}{E(N)} \right) \right)^2 \\ \Sigma_{22} &= \lim_{I \rightarrow \infty} I \text{Var}(\hat{\theta}_{11}) = \text{Var}(\psi_{11}/(M_1 N_1)) \\ \Sigma_{12} &= \lim_{I \rightarrow \infty} I \text{Cov}(\hat{\theta}_{12}, \hat{\theta}_{11}) = \theta_{12} E \left(\frac{\psi_{11}}{M_1 N_1} \left(\frac{\psi_{12} + \psi_{21}}{E\psi_{12}} - \frac{M_1}{E(M)} - \frac{N_1}{E(N)} \right) \right)\end{aligned}$$

Corollary 6. *Under the assumptions of Theorem 5, let $(X_1, Y_1, M_1, N_1), \dots, (X_I, Y_I, M_I, N_I)$, be IID according to P . For $1 \leq i \leq I$ define*

$$\begin{aligned}\psi_{i.} &= I^{-1} \sum_{j=1}^I \psi(X_i, Y_j), \\ \psi_{.i} &= I^{-1} \sum_{j=1}^I \psi(X_j, Y_i), \\ \phi_i &= \frac{\psi(X_i, Y_i)}{M_i N_i},\end{aligned}$$

and analogously for $M.$, $N.$, and $\psi..$. The asymptotic covariance matrix Σ of $(\hat{\theta}_{12}, \hat{\theta}_{11})$ may be consistently estimated by $\hat{\Sigma}$ given by:

$$\begin{aligned}\hat{\Sigma}_{11} &= \frac{1}{I-1} \sum_{i=1}^I \left(\frac{\psi_{i.} + \psi_{.i}}{M_i N_i} - \hat{\theta}_{12} \left(\frac{M_i}{M.} + \frac{N_i}{N.} \right) \right)^2 \\ \hat{\Sigma}_{22} &= \frac{1}{I-1} \sum_{i=1}^I (\phi_i - \phi.)^2 \\ \hat{\Sigma}_{12} &= \frac{1}{I} \sum_{i=1}^I \left(\frac{\phi_i}{\phi.} \left(\frac{\psi_{i.} + \psi_{.i}}{\psi..} - \frac{M_i}{M.} - \frac{N_i}{N.} \right) \right)\end{aligned}$$

Proof. See Sen (1960) for convergence results for random variables like $\psi_{i.}$ and $\psi_{.i}$. □

The estimator $\hat{\Sigma}_{11}$ of the asymptotic variance of $\hat{\theta}_{12}$ is the same as given by Obuchowski (1997), derived by a different method. The finite-sample performance of this estimator is examined in Section 6.

6 Simulation

We examine estimation and inference on the population and personalized AUCs jointly. Many of the choices and parameters follow the simulation in Obuchowski (1997) examining what is here referred to as the population AUC. Key differences include: 1) In our model $M > 0, N > 0$, to ensure that the personalized AUC is well-defined; 2) Whereas Obuchowski (1997) take $I = 100$, we take the number of clusters to be $I = 60$ in the coverage simulation, Section 6.2, and $I = 10$ in the power simulation, Section 6.3.

6.1 Data models

To generate (M, N) , first a preliminary number $\bar{M} + \bar{N}$ of combined case and control observations belonging in a sample is randomly selected from among $k \in \{2, 3, 4, 5\}$. Next, to obtain the allocation to case and control observations, $\bar{M} + \bar{N}$ normal variables are sampled with unit variance and common pairwise correlation $\rho_{MN} \in \{0, 0.1, 0.4, 0.8\}$. A preliminary number \bar{M} of control observations is taken to be those greater than 0, and the remainder the preliminary number \bar{N} of case observations. Finally, 1 is added to each to obtain the final number of control and case observations, $M = \bar{M} + 1, N = \bar{N} + 1$. The greater the correlation ρ_{MN} , the greater the imbalance between case and control observations within the clusters.

Two related models were considered for $(X, Y) \mid (M, N)$.

Bivariate normal model A popular parametric model for the AUC is the “binormal” model, where the case and control observations are taken to be jointly Gaussian (Hanley, 1988). Following Obuchowski (1997) we extend this model to accommodate clustered data by modeling the observations as multivariate normal vectors with an exchangeable correlation

structure.

$$(X, Y) \mid (M, N) \sim \mathcal{N}_{M+N} \left(\begin{pmatrix} 0 \cdot \mathbb{1}_M \\ \Delta \cdot \mathbb{1}_N \end{pmatrix}, \rho \mathbb{1}_{M+N} \mathbb{1}_{M+N}^T + (1 - \rho) Id_{M+N} \right) \quad (12)$$

That is, the case and control observations of a given cluster all have unit variance and share the same pairwise correlation ρ , all the case observations have mean $\Delta > 0$, and all the control observations mean 0. The bivariate normal model is in fact an example of the random effect model described in Section 3, though the random effect is not given explicitly in (12). As the impact of the random effect discussed there is only to change the intra-cluster correlation or mean in (12), it is redundant to the usual multivariate normal parameters. Moreover, further parameters such as for a non-zero control mean $E(X_{11})$ or non-unit variances $\text{Var}(X_{11})$ and $\text{Var}(Y_{11})$ are redundant for our purpose of modeling AUCs.

Using Proposition 2,

$$\begin{aligned} \theta_{12}(P) &= \Phi \left(\frac{\Delta}{\sqrt{2}} \right) \\ \theta_{11}(P) &= \Phi \left(\frac{\Delta}{\sqrt{2(1-\rho)}} \right) \end{aligned} \quad (13)$$

The formulas (13) show that $\theta_{11} > \theta_{12}$ and further that θ_{12} and θ_{11} are simultaneously $> 1/2$, $= 1/2$, or $< 1/2$. We give two benefits. The first is that $(\theta_{12}, \theta_{11})$ can be restricted without loss of generality to $[1/2, 1] \times [1/2, 1]$, switching control and case labels if necessary. The pair $(\theta_{12}, \theta_{11})$ may then serve as a parameterization of the bivariate normal model (12), solving for Δ and ρ in (13). The second involves testing. Though AUCs are often compared by magnitude, e.g., $H_0 : AUC_1 - AUC_2 > 0$, one is usually interested in the discrimination, i.e., $|AUC_1 - 1/2|$ versus $|AUC_2 - 1/2|$. With the latter view, the hypothesis $H_0 : AUC_1 - AUC_2 > 0$ is ambiguous, indicating that AUC_1 is more discriminating than AUC_2 when both are greater than $1/2$, but less discriminating if both are less than $1/2$. A further complication when comparing AUCs in general, which will not be solved by switching

the class designations, is that one AUC may be greater than $1/2$ and the other less. These complications are avoided in the bivariate normal model for the personalized and population AUCs. A test of $\theta_{12} = \theta_{11}$ versus $\theta_{12} < \theta_{11}$ is also a test of discrimination, $|\theta_{12} - 1/2| = |\theta_{11} - 1/2|$ versus $|\theta_{12} - 1/2| < |\theta_{11} - 1/2|$.

Censored bivariate normal model

We also examine the bivariate normal model under censoring, a mixed discrete-continuous distribution. Let $a > 0$, let $(\bar{X}, \bar{Y}) | (M, N)$ be sampled as in (12), and let

$$\begin{aligned} (X, Y) | (M, N) = & (-a\{\bar{X} \leq -a\} + \bar{X}\{-a < \bar{X} < a\} + a\{\bar{X} \geq a\}, \\ & -a\{\bar{Y} \leq -a\} + \bar{Y}\{-a < \bar{Y} < a\} + a\{\bar{Y} \geq a\}). \end{aligned} \quad (14)$$

That is, observations (\bar{X}, \bar{Y}) are generated as in the bivariate normal model (12), and the values are then clipped to $\pm a$. This type of data-generating process is used by Obuchowski (1997) to model radiologists' scores, which lie on a 0–100% scale and often accumulate at 0% and 100%.

Let

$$(X_{11}, Y_{11}) \sim \mathcal{N}_2 \left(\begin{pmatrix} 0 \\ \Delta \end{pmatrix}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix} \right).$$

Again using Proposition 2 to reduce to the $M = N = 1$ case,

$$\begin{aligned} \theta_{12}(P) = & - \int_{-a}^a \Phi(x - \Delta) \phi(x) dx + \frac{1}{2} (\Phi(a) - \Phi(a - \Delta) - \Phi(-a - \Delta) \\ & + \Phi(a) (\Phi(-a - \Delta) + \Phi(a - \Delta)) + 1) \\ \theta_{11}(P) = & \int_{-a}^a \int_x^a f_{X_{11}, Y_{11}}(x, y) dy dx + \Phi(-a) + 1 - \Phi(a - \Delta) - \frac{1}{2} P(X_{11} < -a, Y_{11} < -a) \\ & - \frac{1}{2} P(X_{11} > a, Y_{11} > a) - P(X_{11} < -a, Y_{11} > a). \end{aligned}$$

Due to the censoring, the AUCs may be bounded below 1 in this model, regardless of the magnitude of the location shift between the underlying control and case observations. As

$\Delta \rightarrow \infty$, θ_{12} and θ_{11} both tend to $\frac{1}{2}(1 + \Phi(a))$.

6.2 Coverage

The parameters Δ and ρ were set to correspond to a population AUC of $\theta_{12} \in \{0.7, 0.8\}$ and personalized AUCs of $\theta_{11} \in \{0.7, 0.8, 0.9\}$ with $\theta_{11} \geq \theta_{12}$. For each setting of $\rho_{MN}, \theta_{12}, \theta_{11}$, 1,000 replicates of size $I = 60$ were sampled and used to form a confidence ellipse for $(\theta_{12}, \theta_{11})$. Specifically, with $\hat{\theta}_{12}, \hat{\theta}_{11}$ computed as in Section 2 and Σ as in Theorem 5, under P ,

$$\left| \Sigma^{-1/2} \begin{pmatrix} \theta_{12} \\ \theta_{11} \end{pmatrix} - \begin{pmatrix} \hat{\theta}_{12} \\ \hat{\theta}_{11} \end{pmatrix} \right|^2 \quad (15)$$

has a chi-squared distribution with 2 degrees of freedom. If q is an upper α quantile of this distribution, then

$$\left\{ \begin{pmatrix} x \\ y \end{pmatrix} : \left| \Sigma^{-1/2} \begin{pmatrix} x \\ y \end{pmatrix} - \begin{pmatrix} \hat{\theta}_{12} \\ \hat{\theta}_{11} \end{pmatrix} \right|^2 < q \right\}$$

is a level $1 - \alpha$ confidence region for $(\theta_{12}, \theta_{11})$, which then covers $(\theta_{12}, \theta_{11})$ when (15) is $< q$. In the simulation, we substitute for Σ the asymptotic approximation $\hat{\Sigma}$ given in Corollary 6. Results are presented in Table 1. The bias is on the order of a hundredth at this sample size, and the coverage is generally close to .95. There is some degradation in the coverage as $(\theta_{12}, \theta_{11})$ approach $(1, 1)$.

6.3 Power

We examine the power of testing the null hypothesis $H_0 : \theta_{12} = \theta_{11}$ using the proposed variance estimators under the bivariate normal model (12). Restricting to $\rho > 0$ in (12), the set of alternatives to $H_0 : 1/2 < \theta_{12} = \theta_{11}$ is $H_A : 1/2 < \theta_{12} < \theta_{11}$, i.e., where the personalized AUC is more discriminating than the population AUC.

| parameters | | | coverage | bias x 1000 | | | | |
|---------------|---------------|-------------|----------|---------------|---------------|---------------|---------------|---------------|
| θ_{12} | θ_{11} | ρ_{MN} | | θ_{12} | θ_{11} | Σ_{11} | Σ_{12} | Σ_{22} |
| 0.70 | 0.70 | 0.00 | 0.93 | 2.09 | 0.94 | -5.66 | -6.38 | -4.61 |
| 0.70 | 0.70 | 0.10 | 0.94 | 2.18 | 1.22 | 1.94 | 1.56 | 2.46 |
| 0.70 | 0.70 | 0.40 | 0.94 | 1.14 | 1.07 | 6.41 | 5.58 | 4.11 |
| 0.70 | 0.70 | 0.80 | 0.93 | 2.53 | -1.24 | 2.95 | 0.59 | 0.50 |
| 0.70 | 0.80 | 0.00 | 0.93 | 1.16 | -0.76 | 3.59 | 5.23 | 2.04 |
| 0.70 | 0.80 | 0.10 | 0.93 | 1.93 | 0.41 | 5.65 | 5.82 | -1.17 |
| 0.70 | 0.80 | 0.40 | 0.93 | -0.21 | 0.75 | -3.59 | -1.25 | -0.97 |
| 0.70 | 0.80 | 0.80 | 0.93 | -0.95 | -2.43 | -9.59 | -2.95 | -0.72 |
| 0.80 | 0.80 | 0.00 | 0.93 | 0.69 | -0.33 | -1.49 | -1.51 | -0.79 |
| 0.80 | 0.80 | 0.10 | 0.93 | 2.21 | 0.36 | 4.94 | 3.34 | 2.84 |
| 0.80 | 0.80 | 0.40 | 0.92 | 0.59 | -0.42 | 0.90 | 2.44 | 3.47 |
| 0.80 | 0.80 | 0.80 | 0.93 | 1.02 | -0.22 | 4.23 | 4.06 | 5.33 |
| 0.80 | 0.90 | 0.00 | 0.90 | 2.39 | 1.87 | -1.05 | -0.45 | -2.89 |
| 0.80 | 0.90 | 0.10 | 0.91 | 0.46 | -0.28 | 0.83 | 1.15 | 0.33 |
| 0.80 | 0.90 | 0.40 | 0.92 | -0.05 | -0.95 | 8.86 | 3.31 | 1.96 |
| 0.80 | 0.90 | 0.80 | 0.91 | 3.20 | -1.50 | -0.94 | -0.08 | -0.12 |

(a) Binormal model (12)

| parameters | | | coverage | bias x 1000 | | | | |
|---------------|---------------|-------------|----------|---------------|---------------|---------------|---------------|---------------|
| θ_{12} | θ_{11} | ρ_{MN} | | θ_{12} | θ_{11} | Σ_{11} | Σ_{12} | Σ_{22} |
| 0.70 | 0.70 | 0.00 | 0.94 | 0.51 | -0.08 | 4.92 | 2.03 | 2.29 |
| 0.70 | 0.70 | 0.10 | 0.94 | -3.41 | -3.31 | 9.87 | 7.38 | 8.83 |
| 0.70 | 0.70 | 0.40 | 0.95 | 1.18 | 1.87 | 1.48 | -5.06 | -2.10 |
| 0.70 | 0.70 | 0.80 | 0.93 | 4.67 | 1.30 | -2.15 | -1.17 | -1.80 |
| 0.70 | 0.80 | 0.00 | 0.92 | 1.68 | 0.84 | -1.95 | -2.88 | -5.02 |
| 0.70 | 0.80 | 0.10 | 0.93 | 2.60 | 0.53 | 2.16 | 2.91 | 3.38 |
| 0.70 | 0.80 | 0.40 | 0.93 | 1.16 | -0.54 | -2.68 | -0.33 | 0.03 |
| 0.70 | 0.80 | 0.80 | 0.93 | 5.34 | 1.74 | 1.25 | -3.48 | -8.08 |
| 0.80 | 0.80 | 0.00 | 0.94 | 1.59 | 0.00 | 0.44 | -0.96 | 0.82 |
| 0.80 | 0.80 | 0.10 | 0.94 | 0.64 | -1.45 | 1.13 | -0.25 | 0.33 |
| 0.80 | 0.80 | 0.40 | 0.93 | 2.03 | -1.19 | 5.15 | 3.03 | 2.46 |
| 0.80 | 0.80 | 0.80 | 0.93 | 1.26 | 0.04 | -2.46 | -1.73 | -0.16 |
| 0.80 | 0.90 | 0.00 | 0.92 | 1.25 | -8.08 | 2.21 | -1.19 | -4.51 |
| 0.80 | 0.90 | 0.10 | 0.92 | 1.81 | -7.24 | -0.98 | 0.64 | 4.38 |
| 0.80 | 0.90 | 0.40 | 0.92 | 1.76 | -6.56 | -1.60 | -3.79 | -3.01 |
| 0.80 | 0.90 | 0.80 | 0.91 | 3.27 | -7.55 | 1.44 | 1.89 | 3.81 |

(b) Binormal model with censoring (14)

Table 1: The results of a simulation examining the coverage of a nominal 95% confidence ellipse obtained using the asymptotic estimator given in Section 5. For θ_{11} and θ_{12} , the bias is computed as the mean difference between the estimates and the known true values. For the elements of the covariance matrix Σ_{ij} , the bias is the mean difference between the estimates given by Theorem 5 and the empirical covariance.

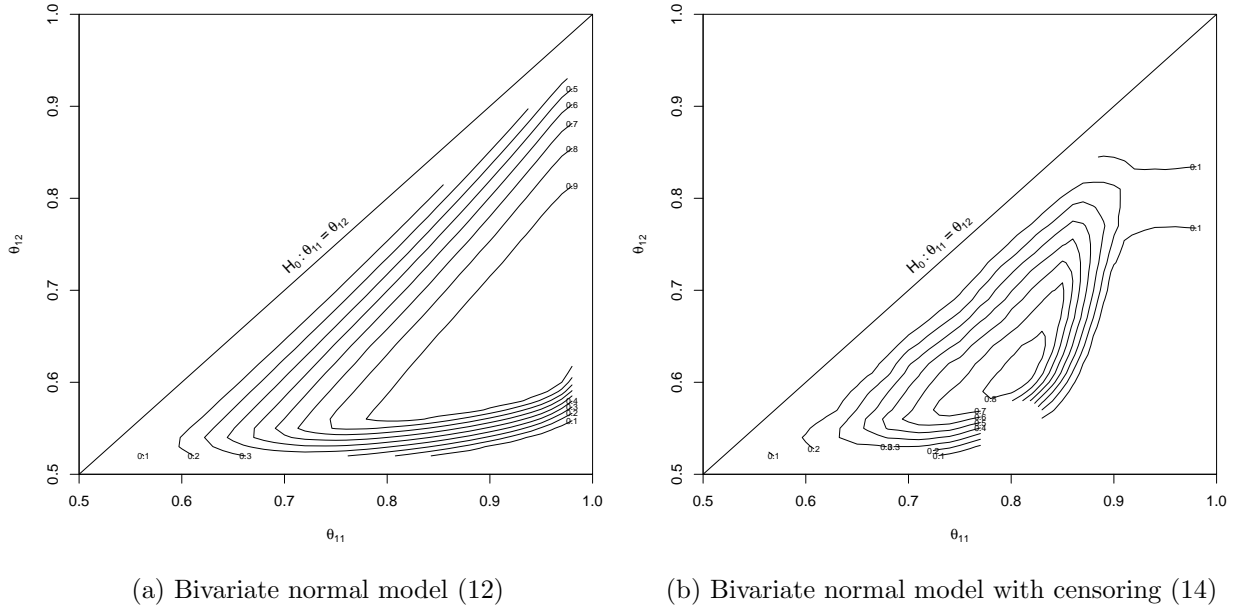


Figure 2: Empirical power function of the test of $H_0 : \theta_{12} = \theta_{11}$ versus $\theta_{12} < \theta_{11}$ using the asymptotic estimator given in Section 5. In the bivariate normal model with or without censoring, the null is equivalent to $H_0 : |\theta_{12} - 1/2| = |\theta_{11} - 1/2|$, equal informativity.

The data is generated under (12) using $(\theta_{12}, \theta_{11})$ selected from points randomly and uniformly selected in $[\frac{1}{2}, 1] \times [\frac{1}{2}, 1], \{\theta_{11} \geq \theta_{12}\}$. Estimates $\hat{\theta}_{12}, \hat{\theta}_{11}$, and $\hat{\Sigma}$ were then obtained as described above. The test is carried out by testing the significance of the z-statistic

$$(\hat{\theta}_{12} - \hat{\theta}_{11}) / \sqrt{c^t \hat{\Sigma} c}$$

where the contrast vector c is $(1, -1)^t$.

The observed power functions are plotted in Fig. 2. The number of clusters was chosen to be $I = 10$, few relative to the setting in Obuchowski (1997) , since the qualitative behavior of the power surface appears clearer with fewer clusters.

7 Data analysis

We examine data on police behavior and give 3 analyses leading to 3 different relationships between the population and personalized AUCs: the population AUC 1) significantly more than, 2) significantly less than, and 3) not significantly different from the personalized AUC.

The data consists of Terry stops in New York City and Boston. As a legal concept, a Terry stop is a policing procedure whereby an officer briefly detains an individual based on a reasonable suspicion that a crime has been committed, which is a lower evidentiary bar than required to arrest the individual. Terry stops are colloquially referred to as “stop and frisks” though the suspect need not be frisked or searched. The analysis here focuses on the relationship between the duration of the stop and race of the suspect. We cluster the stops according to precinct, in the case of NYC, and according to the officer conducting the stop, in the case of Boston. There is an extensive literature examining the relationship between race and Terry stops. Duration of the stop in particular is examined in, e.g., Ridgeway (2006), clustering at the precinct level in, e.g., Goel et al. (2016), and clustering at the officer level in, e.g., Ridgeway and MacDonald (2009).

The NYC data consists of measurements on 54,587 stops carried out between 2017 and 2021. The Boston data consists of 6,591 stops carried out between 2019 and 2021. The stop durations range between 0 minutes and 1–2 hours, with modes at multiples of 5 minutes, and 15 minutes being the most commonly recorded duration. Measurement error is inevitable; see Section 3.2, (a)—(d), for two contrasting illustrations of how it may affect the interpretation of the results here. While data is available for years prior to the cutoffs used here, key covariates used in the analysis were either missing or coded differently in the earlier data. So that the personalized AUC could be estimated, the data was further restricted to those clusters with at least 1 control and 1 case observation, where the interpretation of “control” or “case” depends on the racial classification under analysis below. The final number of clusters and cluster sizes are given in Table 3.

The racial classifications we consider are Black, White, and Hispanic; see Table 3 for

breakdowns. In the first two analyses below, race is considered apart from Hispanic ethnicity, i.e., Black and White is taken to include Black Hispanic and White Hispanic, somewhat uncommon in analyses of policing. In the third the more common treatment of Black and White as exclusive of Hispanics is considered.

1. $\theta_{12} < \theta_{11}$. With Black race as the binary classification, the AUC analysis looks for a difference in location between the distribution of stop durations of non-Black (“control”) and Black (“case”) suspects. For the NYC data, the population AUC estimate is $\hat{\theta}_{12} = 0.46$ with 95% CI 0.45—0.47, significantly different from the null value of $1/2$. The personalized AUC estimate is $\hat{\theta}_{11} = 0.50$ with a 95% CI 0.47—0.53. A test of equality $H_0 : \theta_{12} = \theta_{11}$ against $\theta_{12} < \theta_{11}$ returns a p-value of .05%. The Boston data is similar. The population AUC estimate is 0.46 [0.42, 0.50] and the personalized AUC estimate is 0.52 [0.46, 0.58]. A test of equality $H_0 : \theta_{12} = \theta_{11}$ against $\theta_{12} < \theta_{11}$ returns the p-value .91%. Confidence ellipses are plotted in Figure 3. The data recalls the situation depicted in Fig. 1b, though of course the difference between the two AUCs is less dramatic here than in the artificial example constructed there.
2. $\theta_{11} < \theta_{12}$. We next consider differences in duration of stop between non-White (“control”) or White (“case”) suspect status. As Table 2 indicates, the vast majority of suspects are either Black or White, when those categories are taken inclusive of Hispanics, so one might expect that the analysis for non-White/White status to be nearly the same as the analysis for Black/non-Black status, therefore simply reversing the direction of the results just given, i.e., reflecting the AUCs across $1/2$. That expectation largely holds for the NYC data, where the population and personalized AUCs are 0.53 [0.52, 0.54] and 0.50 [0.48, 0.53], and the population AUC remains the only one of the two significantly different from the null value $1/2$. For the Boston estimates, however, the personalized AUC, 0.46 [0.40, 0.53], is more informative than the population AUC, 0.52 [0.48, 0.55], with the test of equality versus $\theta_{11} < \theta_{12}$ returning a p-value of 2.5%. This analysis therefore corresponds to the situation in Fig. 1a.

| | NYC | | | Boston | | |
|--------------------|-----------------------|-------|-------|-----------------------|-------|-------|
| group | mean duration (SD) | count | freq. | mean duration (SD) | count | freq. |
| Asian | 14.24 (21.16) | 1139 | 0.02 | 25.00 (24.22) | 53 | 0.01 |
| Black Hispanic | 11.01 (17.12) | 4675 | 0.09 | 15.28 (18.73) | 391 | 0.06 |
| Black non-Hispanic | 10.99 (16.78) | 31588 | 0.58 | 19.06 (28.93) | 3448 | 0.55 |
| White Hispanic | 11.21 (15.15) | 11486 | 0.21 | 15.63 (15.96) | 578 | 0.09 |
| White non-Hispanic | 12.85 (16.18) | 4854 | 0.09 | 21.74 (33.01) | 1760 | 0.28 |
| other | 11.84 (17.70) | 261 | 0.00 | 20.89 (23.90) | 93 | 0.01 |

Table 2: Summary estimates on the duration of Terry stops by racial group.

3. No significant difference between θ_{12} and θ_{11} . Finally, we consider duration of the stop between non-Hispanic (“control”) and Hispanic (“case”) suspects. For both the NYC and Boston data, neither the population AUC nor personalized AUC is significantly different from the null value $1/2$, and the test of equality of the two AUCs fails to reject. As a second example, in Boston, whether one takes the case status to be non-Hispanic Black or non-Hispanic White, the two AUCs are statistically indistinguishable from each other and each is indistinguishable from the null value $1/2$.

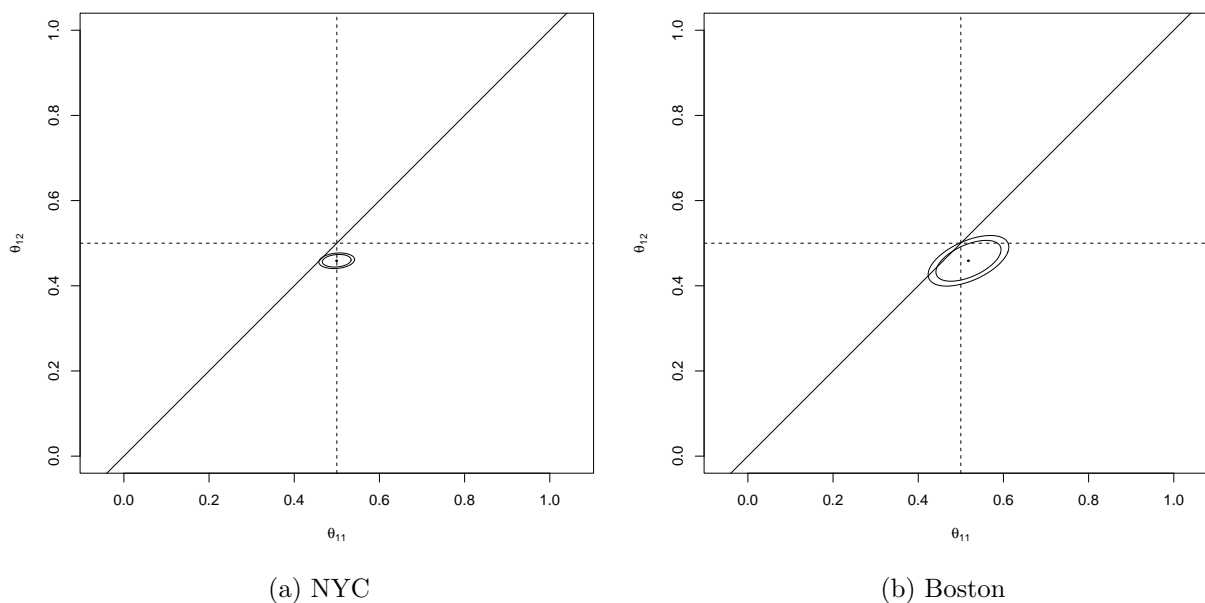
The decision to cluster at the officer or precinct level, as opposed to, say, the time of day of the stop, age of the suspect, or other partition of the data, is in part arbitrary. For the application of the definitions and results given in the previous sections, the decision amounts to the idealization that the officers’ or precincts’ data are drawn independently from a universe of officers or precinct Terry stop data. At the same time, many current analyses, such as cited above, besides this IID assumption further impose modeling assumptions such as linear random effects or logistic links. The approach here has the advantage of being otherwise nonparametric.

8 Discussion

We have compared and contrasted two generalizations of the AUC to accommodate clustered, paired data. Straightforward extensions include allowing for multiple dependent AUCs, clus-

| case group | data set | I | ΣM_i | ΣN_i | θ_{12} | θ_{11} | $H_0 : \theta_{12} = \theta_{11}$ |
|--------------------|----------|-----|--------------|--------------|-------------------|-------------------|-----------------------------------|
| Black | NYC | 187 | 17698 | 36152 | 0.46 [0.45, 0.47] | 0.50 [0.47, 0.53] | 0.00 |
| | Boston | 112 | 418 | 585 | 0.46 [0.42, 0.50] | 0.52 [0.46, 0.58] | 0.02 |
| Black non-Hispanic | NYC | 185 | 22348 | 31490 | 0.47 [0.46, 0.48] | 0.51 [0.48, 0.53] | 0.01 |
| | Boston | 117 | 464 | 569 | 0.48 [0.44, 0.51] | 0.50 [0.44, 0.56] | 0.30 |
| Black Hispanic | NYC | 154 | 48847 | 4672 | 0.48 [0.47, 0.49] | 0.49 [0.47, 0.52] | 0.42 |
| | Boston | 41 | 494 | 62 | 0.44 [0.37, 0.51] | 0.49 [0.40, 0.59] | 0.09 |
| White | NYC | 185 | 37547 | 16298 | 0.53 [0.52, 0.54] | 0.50 [0.48, 0.53] | 0.04 |
| | Boston | 109 | 614 | 385 | 0.52 [0.48, 0.55] | 0.46 [0.40, 0.53] | 0.05 |
| White non-Hispanic | NYC | 148 | 48327 | 4838 | 0.56 [0.55, 0.58] | 0.52 [0.49, 0.55] | 0.00 |
| | Boston | 106 | 631 | 324 | 0.52 [0.47, 0.56] | 0.49 [0.43, 0.56] | 0.39 |
| White Hispanic | NYC | 176 | 42333 | 11463 | 0.51 [0.50, 0.52] | 0.49 [0.47, 0.52] | 0.30 |
| | Boston | 62 | 631 | 89 | 0.48 [0.41, 0.55] | 0.47 [0.39, 0.56] | 0.81 |
| Hispanic | NYC | 180 | 37693 | 16125 | 0.50 [0.49, 0.51] | 0.49 [0.46, 0.52] | 0.41 |
| | Boston | 85 | 706 | 151 | 0.46 [0.41, 0.50] | 0.48 [0.41, 0.55] | 0.51 |

Table 3: Estimates of the population and personalized AUCs of the duration of Terry stops by racial group.

Figure 3: Level 95% and 99% Confidence ellipses for the estimates of $(\theta_{11}, \theta_{12})$ for duration of Terry stop by non-Black/Black status.

ters that are only exchangeable or otherwise fall short of being IID, and covariate-adjusted AUCs. A more delicate extension would allow for estimation of the personalized AUC when some clusters have no control or no case observations. As the personalized AUC is not currently defined for such clusters either the definition would need to be re-worked or a model would need to be introduced for the missing values corresponding to those clusters. No major changes would be required of the analysis under a strong enough assumption such as ignorability, i.e., the assumption that the behavior of the personalized AUC (or the pair) is the same on $M > 1, N > 1$ as on the entire population.

The authors wish to thank Prof. Maria Cuellar for helpful consultation regarding the data analysis, and an anonymous reviewer for contributing the substance of Prop. 1(2).

Haben Michael
University of Massachusetts
hmichael@math.umass.edu

Lu Tian
Stanford University
lutian@stanford.edu

References

- Bamber, D. (1975). The area above the ordinal dominance graph and the area below the receiver operating characteristic graph. *Journal of mathematical psychology*, 12(4):387–415.
- Benhin, E., Rao, J., and Scott, A. (2005). Mean estimating equation approach to analysing cluster-correlated data with nonignorable cluster sizes. *Biometrika*, 92(2):435–450.
- Bugni, F., Canay, I., Shaikh, A., and Tabord-Meehan, M. (2022). Inference for cluster randomized experiments with non-ignorable cluster sizes. *arXiv preprint ArXiv:2204.08356*.
- Emir, B., Wieand, S., Jung, S.-H., and Ying, Z. (2000). Comparison of diagnostic markers with repeated measurements: a non-parametric ROC curve approach. *Statistics in Medicine*, 19(4):511–523.
- Goel, S., Rao, J. M., and Shroff, R. (2016). Precinct or prejudice? Understanding racial disparities in new york city’s stop-and-frisk policy. *The Annals of Applied Statistics*, 10(1):365–394.

- Hanley, J. A. (1988). The robustness of the “binormal” assumptions used in fitting ROC curves. *Medical decision making*, 8(3):197–203.
- Lee, A. J. (2019). *U-statistics: Theory and Practice*. Routledge.
- Lee, M.-L. T. and Dehling, H. G. (2005). Generalized two-sample u-statistics for clustered data. *Statistica Neerlandica*, 59(3):313–323.
- Liu, H., Li, G., Cumberland, W. G., Wu, T., et al. (2005). Testing statistical significance of the area under a receiving operating characteristics curve for repeated measures design with bootstrapping. *Journal of Data Science*, 3(3):257–278.
- Michael, H., Tian, L., and Ghebremichael, M. (2019). The ROC curve for regularly measured longitudinal biomarkers. *Biostatistics*, 20(3):433–451.
- Obuchowski, N. A. (1997). Nonparametric analysis of clustered ROC curve data. *Biometrics*, pages 567–578.
- Pearl, J. (2014). Comment: Understanding Simpson’s paradox. *The American Statistician*, 68(1):8–13.
- Ridgeway, G. (2006). Assessing the effect of race bias in post-traffic stop outcomes using propensity scores. *Journal of Quantitative Criminology*, 22(1):1–29.
- Ridgeway, G. and MacDonald, J. M. (2009). Doubly robust internal benchmarking and false discovery rates for detecting racial bias in police stops. *Journal of the American Statistical Association*, 104(486):661–668.
- Rosner, B. and Grove, D. (1999). Use of the Mann–Whitney U-test for clustered data. *Statistics in medicine*, 18(11):1387–1400.
- Sen, P. K. (1960). On some convergence properties of U-tatistics. *Calcutta Statistical Association Bulletin*, 10(1-2):1–18.
- Toledano, A. Y. (2003). Three methods for analysing correlated roc curves: a comparison in real data sets from multi-reader, multi-case studies with a factorial design. *Statistics in medicine*, 22(18):2919–2933.
- Wu, Y. and Wang, X. (2011). Optimal weight in estimating and comparing areas under the receiver operating characteristic curve using longitudinal data. *Biometrical journal*, 53(5):764–778.

Deferred proofs are given below, beginning with a few short technical results relied on by the main results. To save space, the value of a bivariate kernel $\psi(X_i, Y_j)$ is abbreviated as ψ_{ij} in some places.

The following lemma gives a convergence result for a two-sample U -statistic with kernel of degree $(1, 1)$ where the data is paired. The corresponding definitions and result for independent samples are given in, e.g., Lee (2019). Let V denote the space of finite sequences.

Lemma 7. *Given a sample $(X_0, Y_0), (X_1, Y_1), \dots, (X_I, Y_I)$ on $V \times V$ IID according to P and a function $\psi : V \times V \rightarrow \mathbb{R}$ in $L^2(P)$, define*

$$U_I = (I)_2^{-1} \sum_{\substack{1 \leq i, j \leq I \\ i \neq j}} \psi(X_i, Y_j), \quad V_I = I^{-2} \sum_{1 \leq i, j \leq I} \psi(X_i, Y_j),$$

and

$$\hat{U}_I = I^{-1} \sum_{i=1}^I (E(\psi(X_i, Y_0) \mid X_i, Y_i) + E(\psi(X_0, Y_i) \mid X_i, Y_i)) - 2E\psi(X_1, Y_2).$$

Then

$$E(U_I - EU_I - \hat{U}_I)^2 = O(I^{-2}) \text{ and } E(V_I - EV_I - \hat{U}_I)^2 = O(I^{-2}).$$

Proof of Lemma 7. Define

$$\bar{\psi}_{ij} = \psi(X_i, Y_j) - E(\psi(X_i, Y_0) \mid X_i, Y_i) - E(\psi(X_0, Y_j) \mid X_j, Y_j) + E\psi(X_1, Y_2).$$

Then, for $i \neq j$, $E(\bar{\psi}_{ij} \mid (X_i, Y_i)) = E(\bar{\psi}_{ij} \mid (X_j, Y_j)) = 0$, implying

$$\begin{aligned} E(U_I - EU_I - \hat{U}_I)^2 &= E \left((I)_2^{-1} \sum_{i \neq j} \bar{\psi}_{ij} \right)^2 \\ &= (I)_2^{-2} \sum_{i \neq j} E \bar{\psi}_{ij}^2 + O(I^{-2}) \\ &= O(I^{-2}). \end{aligned}$$

For the second equation,

$$\begin{aligned} E(U_I - EU_I - V_I + EV_I)^2 &= I^{-2} E \left((I)_2^{-1} \sum_{i \neq j} \left(\psi_{ij} + \psi_{ii} - \frac{I}{I-1} E\psi_{11} + E\psi_{12} \right) \right)^2 \\ &\leq I^{-2} \left((I)_2^{-1} \sum_{i \neq j} E \left(\psi_{ij} + \psi_{ii} - \frac{I}{I-1} E\psi_{11} + E\psi_{12} \right)^2 \right) \\ &= O(I^{-2}). \end{aligned}$$

□

Corollary 8. *With the same setup as Lemma 7, $U_I - EU_I \rightarrow 0$ a.s. and $\sqrt{I}(U_I - EU_I)/\sqrt{\text{Var}(U_I)} \rightarrow \mathcal{N}(0, 1)$ in distribution.*

Proof of Corollary 8. By Lemma 7, $U_I - EU_I \rightarrow \hat{U}_I$ a.s. and $\sqrt{I}(U_I - EU_I - \hat{U}_I) \rightarrow 0$ in quadratic mean, and \hat{U}_I is an IID sum subject to the usual LLN and CLT. \square

Proof of Proposition 1. 1. By the LLN $I^2/(\sum_i M_i \sum_i N_i) \rightarrow 1/(E(M)E(N))$ almost surely and by Corollary 8 $\sum_{i,j} \psi_{ij}/I^2 \rightarrow E\psi_{12}$ almost surely. Conditioning on the sample,

$$\begin{aligned} E\psi(\xi_I, \eta_I) &= E(E(\psi(\xi_I, \eta_I) \mid (X_1, Y_1, M_1, N_1), \dots, (X_I, Y_I, M_I, N_I))) \\ &= E\left(\frac{\sum_{1 \leq i, j \leq I} \sum_{1 \leq k \leq M_i, 1 \leq l \leq N_j} \psi(X_{ik}, Y_{jl})}{\sum_{i=1}^I M_i \sum_{i=1}^I N_i}\right) \\ &= E\left(\frac{\sum_{1 \leq i, j \leq I} \psi_{ij}}{\sum_{i=1}^I M_i \sum_{i=1}^I N_i}\right) \rightarrow \frac{E\psi_{12}}{E(M)E(N)} = \theta_{12}. \end{aligned}$$

The limit is justified since $\sum_{i,j} \psi_{i,j}/(\sum_i M_i \sum_i N_i) \leq 1$.

2.

$$\begin{aligned} P(\xi < \eta) + \frac{1}{2}P(\xi = \eta) &= E\psi(\xi, \eta) \\ &= \sum_{m=1}^{\infty} \sum_{i=1}^m \sum_{n=1}^{\infty} \sum_{j=1}^n E(\psi(X_{1i}, Y_{2j}) \mid M = m, N = n) \frac{P(M = m)P(N = n)}{(EM)(EN)} \\ &= \sum_{m=1}^{\infty} \sum_{n=1}^{\infty} E\left(\sum_{i=1}^m \sum_{j=1}^n (\psi(X_{1i}, Y_{2j}) \mid M = m, N = n)\right) \frac{P(M = m)P(N = n)}{(EM)(EN)} \\ &= \frac{1}{(EM)(EN)} E\left(\sum_{i=1}^M \sum_{j=1}^N \psi(X_{1i}, Y_{2j})\right). \end{aligned}$$

\square

Proof of Proposition 2.

$$\begin{aligned} \theta_{11}(P) &= E\left(\frac{\sum_{k=1}^M \sum_{l=1}^N \psi(X_{1k}, Y_{1l})}{MN}\right) \\ &= E\left(\frac{1}{MN} E\left(\sum_{k=1}^M \sum_{l=1}^N \psi(X_{1k}, Y_{1l}) \mid M, N\right)\right) \\ &= E\left(\frac{1}{MN} MNE(\psi(X_{11}, Y_{11}) \mid M, N)\right) = E\psi(X_{11}, Y_{11}). \end{aligned}$$

Similar to the above,

$$\begin{aligned}\theta_{12}(P) &= \frac{E\left(\sum_{k=1}^{M_1} \sum_{l=1}^{N_2} \psi(X_{1k}, Y_{2l})\right)}{E(M)E(N)} \\ &= \frac{E(M)E(N)E\psi(X_{11}, Y_{21})}{E(M)E(N)} = E\psi(X_{11}, Y_{21}).\end{aligned}$$

□

Lemma 9. *Given integrable random variables M, V, X_1, X_2, \dots , such that $M \in \{1, 2, \dots\}$ and $\sum_{i=1}^{\infty} E(|X_i|; M \geq i) < \infty$,*

$$E\left(\sum_{i=1}^M X_i \middle| M, V\right) = \sum_{i=1}^M E(X_i \mid M, V)$$

Proof of Lemma 9.

$$\begin{aligned}E\left(\sum_{i=1}^M X_i \middle| M, V\right) &= E\left(\sum_{m=1}^{\infty} \{M = m\} \sum_{i=1}^m X_i \middle| M, V\right) \\ &= \sum_{m=1}^{\infty} E\left(\{M = m\} \sum_{i=1}^m X_i \middle| M, V\right) \\ &= \sum_{m=1}^{\infty} \sum_{i=1}^m \{M = m\} E(X_i \mid M, V) \\ &= \sum_{i=1}^M E(X_i \mid M, V),\end{aligned}$$

the interchange in the second equality allowed since $E\left|\sum_{i=1}^M X_i\right| \leq \sum_{i=1}^{\infty} E(|X_i|; M \geq i) < \infty$. □

Proof of Lemma 4. We introduce the bound in a simple case. Each cluster contributes just one control and one case observation each, and their joint distribution P is supported on finitely many points in the plane:

$$\begin{aligned}P &= \sum_{i=1}^B p_i \delta_{(x_i, y_i)} \\ (x_i, y_i) &\in \mathbb{R}^2 \text{ and } 0 \leq p_i \leq 1, i = 1, \dots, B \\ p_1 + \dots + p_B &= 1.\end{aligned}$$

For this simple example, assume further that all the x_i and y_i are distinct, so $\psi(x, y) = \{x < y\}$.

The personalized AUC is

$$\theta_{11}(P) = P(X < Y) = \sum_{i: x_i < y_i} p_i.$$

The population AUC depends on the product of the marginals of X and Y , say, P_{\perp} ,

$$\theta_{12}(P) = P_{\perp}(X < Y).$$

Since all the x -coordinates of the support points are distinct, the marginal distribution of X is simply $P_{\perp}(X = x) = \sum_i p_i \delta_{x_i}(x)$. Similarly, $P_{\perp}(Y = y) = \sum_i p_i \delta_{y_i}(y)$. The product measure is therefore a weighted sum of B^2 atoms, $P_{\perp}(X = x, Y = y) = \sum_{i,j} p_i p_j \delta_{(x_i, y_j)}(x, y)$. We give a lower bound for the population AUC $P_{\perp}(X < Y)$. An atom of P lying in $\{x < y\}$ of mass p contributes p^2 to the mass given by $P_{\perp}(X < Y)$. Each pair of atoms of P lying in $\{x < y\}$ of mass p and q contributes, beyond p^2 and q^2 , at least pq and possibly $2pq$ to the mass given by $P_{\perp}(X < Y)$. See Figure 4. Therefore

$$\begin{aligned} \theta_{12}(P) = P_{\perp}(X < Y) &\geq \sum_{i: x_i < y_i} p_i^2 + \sum_{i: x_i < y_i} \sum_{\substack{j: x_j < y_j \\ i < j}} p_i p_j \\ &= \frac{1}{2} \left(\sum_{i: x_i < y_i} p_i \right)^2 + \frac{1}{2} \sum_{i: x_i < y_i} p_i^2 \\ &\geq \frac{1}{2} \left(\sum_{i: x_i < y_i} p_i \right)^2 + \frac{1}{2|\{i : x_i < y_i\}|} \left(\sum_{i: x_i < y_i} p_i \right)^2 \\ &= \frac{1}{2} (1 + |\{i : x_i < y_i\}|^{-1}) \theta_{11}(P)^2. \end{aligned}$$

The first inequality is tight when each pair i, j such that $x_i < y_i$ and $x_j < y_j$ contributes exactly $p_i p_j$, i.e., when the square given by x_i, x_j and y_i, y_j has exactly one corner in $\{x < y\}$, so that $y_i - x_i < x_j - x_i$ whenever $x_i < x_j$. The second inequality is Cauchy-Schwarz, and is tight when all the atoms in $\{x < y\}$ have the same mass.

By symmetry,

$$P_{\perp}(X > Y) \geq \frac{1}{2} (1 + |\{i : x_i > y_i\}|^{-1}) P(X > Y)^2,$$

leading to an upper bound

$$\theta_{12} \leq 1 - \frac{1}{2} (1 + |\{i : x_i > y_i\}|^{-1}) (1 - \theta_{11})^2.$$

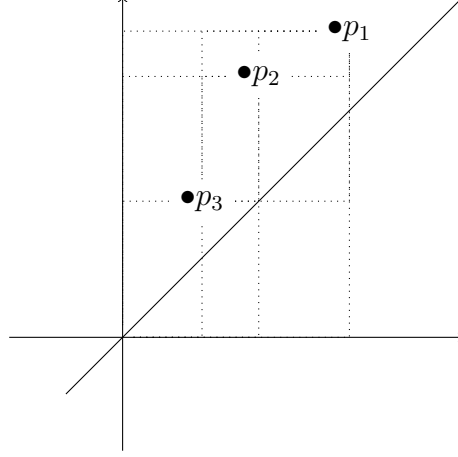


Figure 4: The case $M = N = 1$ and finitely supported (X, Y) . The distance between the atoms p_1 and p_2 is small relative to their distances to the line $x = y$, so they contribute $(p_1 + p_2)^2$ to the mass of $\{x < y\}$ under product of the marginals. The distance between p_1 and p_3 is relatively large, so they contribute only $(p_1 + p_3)^2 - p_1 p_3$.

Combining these bounds,

$$\frac{1}{2}\theta_{11}^2 \leq \theta_{12} \leq 1 - \frac{1}{2}(1 - \theta_{11})^2,$$

or equivalently,

$$1 - \sqrt{2(1 - \theta_{12})} \leq \theta_{11} \leq \sqrt{2\theta_{12}}.$$

Define for $n \in \mathbb{N}$ approximations to θ_{11} and θ_{12} by

$$\begin{aligned} A_{ij}^{(n)} &= \left\{ (x, y) : \frac{i}{2^n} \leq x < \frac{i+1}{2^n}, \frac{j}{2^n} \leq y < \frac{j+1}{2^n} \right\}, \quad -2^{2n} \leq i, j < 2^{2n} - 1 \\ \theta_{11}^{(n)} &= \sum_{i=-2^{2n}}^{2^{2n}-1} \sum_{j=i+1}^{2^{2n}-1} P(A_{ij}^{(n)}) + \frac{1}{2} \sum_{i=-2^{2n}}^{2^{2n}-1} P(A_{ii}^{(n)}) \\ \theta_{12}^{(n)} &= \sum_{i=-2^{2n}}^{2^{2n}-1} \sum_{j=i+1}^{2^{2n}-1} P_{\perp}(A_{ij}^{(n)}) + \frac{1}{2} \sum_{i=-2^{2n}}^{2^{2n}-1} P_{\perp}(A_{ii}^{(n)}). \end{aligned}$$

Since $\bigcup_n \bigcup_i \bigcup_{j>i} A_{ij}^{(n)} = \{x < y\}$ and $\bigcap_n \bigcup_i A_{ii}^{(n)} = \{x = y\}$, by continuity of measure $\theta_{11}^{(n)} \rightarrow \theta_{11}$ and $\theta_{12}^{(n)} \rightarrow \theta_{12}$. Therefore, it is enough to establish the inequality (11) for $\theta_{11}^{(n)}$ and $\theta_{12}^{(n)}$.

Fixing n ,

$$\begin{aligned}
& \sum_{i=-2^{2n}}^{2^{2n-2}} \sum_{j=i+1}^{2^{2n-1}} P_{\perp}(A_{ij}^{(n)}) = \sum_{i=-2^{2n}}^{2^{2n-2}} \sum_{j=i+1}^{2^{2n-1}} P_{\perp}(A_{ij}^{(n)}) \\
& = \sum_{i=-2^{2n}}^{2^{2n-2}} \sum_{j=i+1}^{2^{2n-1}} P_{\perp}\left(\frac{i}{2^n} \leq x < \frac{i+1}{2^n}\right) P_{\perp}\left(\frac{j}{2^n} \leq y < \frac{j+1}{2^n}\right) \\
& \geq \sum_{i=-2^{2n}}^{2^{2n-2}} \sum_{j=i+1}^{2^{2n-1}} (P(A_{ii}^{(n)}) + \sum_{k=i+1}^{2^{2n-1}} P(A_{ik}^{(n)}))(P(A_{jj}^{(n)}) + \sum_{l=-2^{2n}}^{j-1} P(A_{lj}^{(n)})) \\
& = \sum_{i=-2^{2n}}^{2^{2n-2}} \sum_{j=i+1}^{2^{2n-1}} \left(\sum_{k=i+1}^{2^{2n-1}} P(A_{ik}^{(n)}) \sum_{l=-2^{2n}}^{j-1} P(A_{lj}^{(n)}) + P(A_{ii}^{(n)}) \sum_{l=-2^{2n}}^{j-1} P(A_{lj}^{(n)}) \right. \\
& \quad \left. + P(A_{jj}^{(n)}) \sum_{k=i+1}^{2^{2n-1}} P(A_{ik}^{(n)}) + P(A_{ii}^{(n)})P(A_{jj}^{(n)}) \right).
\end{aligned}$$

We lower bound the first three terms in parentheses.

First term:

$$\begin{aligned}
& \sum_{i=-2^{2n-2}}^{2^{2n-2}-2} \sum_{j=i+1}^{2^{2n-1}-1} \sum_{k=i+1}^{2^{2n-1}-1} P(A_{ik}^{(n)}) \sum_{l=-2^{2n}}^{j-1} P(A_{lj}^{(n)}) \\
&= \sum_{i=-2^{2n}}^{2^{2n-2}-2} \sum_{k=i+1}^{2^{2n-1}-1} P(A_{ik}^{(n)}) \sum_{j=i+1}^{2^{2n-1}-1} \sum_{l=-2^{2n}}^{j-1} P(A_{lj}^{(n)}) \\
&\geq \sum_{i=-2^{2n}}^{2^{2n-2}-2} \sum_{k=i+1}^{2^{2n-1}-1} P(A_{ik}^{(n)}) \sum_{j=i+1}^{2^{2n-1}-1} \sum_{l=i}^{j-1} P(A_{lj}^{(n)}) \\
&= \sum_{i=-2^{2n}}^{2^{2n-2}-2} \sum_{k=i+1}^{2^{2n-1}-1} P(A_{ik}^{(n)}) \sum_{l=i}^{2^{2n-2}-2} \sum_{j=l+1}^{2^{2n-1}-1} P(A_{lj}^{(n)}) \\
&= \sum_{i=-2^{2n}}^{2^{2n-2}-2} \sum_{k=i+1}^{2^{2n-1}-1} P(A_{ik}^{(n)}) \sum_{j=i+1}^{2^{2n-1}-1} P(A_{ij}^{(n)}) + \sum_{i=-2^{2n}}^{2^{2n-2}-2} \sum_{k=i+1}^{2^{2n-1}-1} P(A_{ik}^{(n)}) \sum_{l=i+1}^{2^{2n-2}-2} \sum_{j=l+1}^{2^{2n-1}-1} P(A_{lj}^{(n)}) \\
&\geq \sum_{i=-2^{2n}}^{2^{2n-2}-2} \sum_{j=i+1}^{2^{2n-1}-1} P(A_{ij}^{(n)})^2 + \sum_{i=-2^{2n}}^{2^{2n-2}-2} \sum_{k=i+1}^{2^{2n-2}-2} \sum_{j=k+1}^{2^{2n-1}-1} P(A_{ij}^{(n)})P(A_{ik}^{(n)}) + \sum_{i=-2^{2n}}^{2^{2n-2}-2} \sum_{k=i+1}^{2^{2n-1}-1} P(A_{ik}^{(n)}) \sum_{l=i+1}^{2^{2n-2}-2} \sum_{j=l+1}^{2^{2n-1}-1} P(A_{lj}^{(n)}) \\
&= \sum_{\substack{i \neq k \text{ or } j \neq l \\ j > i \text{ and } l > k}} P(A_{ij}^{(n)})P(A_{kl}^{(n)}) + \sum_{i=-2^{2n}}^{2^{2n-2}-2} \sum_{j=i+1}^{2^{2n-1}-1} P(A_{ij}^{(n)})^2 \\
&= \frac{1}{2} \left(\sum_{i=-2^{2n}}^{2^{2n-2}-2} \sum_{j=i+1}^{2^{2n-1}-1} P(A_{ij}^{(n)}) \right)^2 + \frac{1}{2} \sum_{i=-2^{2n}}^{2^{2n-2}-2} \sum_{j=i+1}^{2^{2n-1}-1} P(A_{ij}^{(n)})^2.
\end{aligned}$$

Middle two terms:

$$\begin{aligned}
& \sum_{i=-2^{2n}}^{2^{2n-2}-2} \sum_{j=i+1}^{2^{2n-1}-1} \left(P(A_{ii}^{(n)}) \sum_{l=-2^{2n}}^{j-1} P(A_{lj}^{(n)}) + P(A_{jj}^{(n)}) \sum_{k=i+1}^{2^{2n-1}-1} P(A_{ik}^{(n)}) \right) \\
&= \sum_{i=-2^{2n}}^{2^{2n-2}-2} P(A_{ii}^{(n)}) \sum_{l=i}^{2^{2n-2}-2} \sum_{j=l+1}^{2^{2n-1}-1} P(A_{lj}^{(n)}) + \sum_{j=-2^{2n}+1}^{2^{2n-1}-1} P(A_{jj}^{(n)}) \sum_{i=-2^{2n}}^{j-1} \sum_{k=i+1}^{2^{2n-1}-1} P(A_{ik}^{(n)}) \\
&= \sum_{i=-2^{2n}}^{2^{2n-2}-2} P(A_{ii}^{(n)}) \sum_{l=i}^{2^{2n-2}-2} \sum_{j=l+1}^{2^{2n-1}-1} P(A_{lj}^{(n)}) + \sum_{i=-2^{2n}+1}^{2^{2n-1}-1} P(A_{ii}^{(n)}) \sum_{l=-2^{2n}}^{i-1} \sum_{j=l+1}^{2^{2n-1}-1} P(A_{lj}^{(n)}) \\
&= \left(\sum_{i=-2^{2n}}^{2^{2n-1}-1} P(A_{ii}^{(n)}) \right) \left(\sum_{l=-2^{2n}}^{2^{2n-2}-2} \sum_{j=l+1}^{2^{2n-1}-1} P(A_{lj}^{(n)}) \right).
\end{aligned}$$

The second-to-last equality is just renaming indices.

With these lower bounds,

$$\begin{aligned}
\theta_{12}^{(n)} &= \sum_{i=-2^{2n}}^{2^{2n}-1} \sum_{j=i+1}^{2^{2n}-1} P_{\perp}(A_{ij}^{(n)}) + \frac{1}{2} \sum_{i=-2^{2n}}^{2^{2n}-1} P_{\perp}(A_{ii}^{(n)}) \\
&\geq \frac{1}{2} \left(\sum_{i=-2^{2n}}^{2^{2n}-2} \sum_{j=i+1}^{2^{2n}-1} P(A_{ij}^{(n)}) \right)^2 + \left(\sum_{i=-2^{2n}}^{2^{2n}-1} P(A_{ii}^{(n)}) \right) \left(\sum_{l=-2^{2n}}^{2^{2n}-2} \sum_{j=l+1}^{2^{2n}-1} P(A_{lj}^{(n)}) \right) + \\
&\quad \sum_{i=-2^{2n}}^{2^{2n}-2} \sum_{j=i+1}^{2^{2n}-1} P(A_{ii}^{(n)})P(A_{jj}^{(n)}) + \frac{1}{2} \sum_{i=-2^{2n}}^{2^{2n}-1} P(A_{ii}^{(n)})^2 \\
&= \frac{1}{2} \left(\sum_{i=-2^{2n}}^{2^{2n}-2} \sum_{j=i+1}^{2^{2n}-1} P(A_{ij}^{(n)}) + \sum_{i=-2^{2n}}^{2^{2n}-1} P(A_{ii}^{(n)}) \right)^2 \\
&= \frac{1}{2} \left(\theta_{11}^{(n)} + \frac{1}{2} \sum_{i=-2^{2n}}^{2^{2n}-1} P(A_{ii}^{(n)}) \right)^2 \\
&= \frac{1}{2} \left(\theta_{11}^{(n)} + \frac{1}{2} P(X = Y) \right)^2 + o(1).
\end{aligned}$$

The upper bound then follows by the same symmetry argument as given in Section 4. \square

Proof of Theorem 3. With

$$\theta_{11} = \frac{1}{mn} E(\psi_{11}) = \frac{1}{mn} \sum_{i,j} (P(X_{1i} < Y_{1j}) + \frac{1}{2} P(X_{1i} = Y_{1j}))$$

Lemma 4 gives

$$\begin{aligned}
\theta_{12} &= \frac{1}{mn} E(\psi_{12}) = \frac{1}{mn} \sum_{i,j} (P(X_{1i} < Y_{2j}) + \frac{1}{2} P(X_{1i} = Y_{2j})) \\
&\geq \frac{1}{mn} \sum_{i,j} \frac{1}{2} (P(X_{1i} < Y_{1j}) + P(X_{1i} = Y_{1j}))^2 \\
&\geq \frac{1}{2} \left(\frac{1}{mn} \sum_{i,j} (P(X_{1i} < Y_{1j}) + P(X_{1i} = Y_{1j})) \right)^2 \\
&= \frac{1}{2} \left(\theta_{11} + \frac{1}{2mn} \sum_{i,j} P(X_{1i} = Y_{1j}) \right)^2.
\end{aligned}$$

The second inequality is Jensen's inequality, which is tight when the pairwise AUCs are all equal. The other bound follows similarly. \square

Proof of Theorem 5. By Lemma 7,

$$\sqrt{I} \left(\frac{(I)_2^{-1} \sum_{i \neq j} \psi_{ij} - E\psi_{12}}{\text{sd}(\sqrt{I}(I)_2^{-1} \sum_{i \neq j} \psi_{ij})}, \frac{I^{-2} \sum_{i,j} M_i N_j - E(M)E(N)}{\text{sd}(I^{-3/2} \sum_{i,j} M_i N_j)}, \frac{I^{-1} \sum_i \psi_{ii}/(M_i N_i) - E(\psi_{11}/M_1 N_1)}{\text{sd}(\psi_{11}/M_1 N_1)} \right)$$

converges to

$$I^{-1/2} \sum_{i=1}^I \left(\frac{E(\psi_{i0} | W_i) + E(\psi_{0i} | W_i) - 2E\psi_{12}}{\text{sd}(E(\psi_{10} | W_1) + E(\psi_{01} | W_1))}, \frac{M_i E(N) + N_i E(M) - 2E(M)E(N)}{\text{sd}(M_1 E(N) + N_1 E(M))}, \frac{\psi_{ii}/(M_i N_i) - E(\psi_{11}/M_1 N_1)}{\text{sd}(\psi_{11}/M_1 N_1)} \right)$$

in mean-square. The latter is an IID sum with finite covariance matrix and is asymptotically normal by the usual CLT. Applying the delta method with the function $(x, y, z) \mapsto (x/y, z)$, with derivative

$$\begin{pmatrix} 1/y & -x/y^2 & 0 \\ 0 & 0 & 1 \end{pmatrix} \Big|_{(x,y)=(\theta_{12}, E(M)E(N))}$$

for $y \neq 0$, i.e., $E(M) \neq 0, E(N) \neq 0$, gives the asymptotic normality and stated asymptotic covariance matrix of $(\theta_{11}, \theta_{12})$. □