# Supporting information for 'Covariate adjustment in continuous biomarker assessment' by

Ziyi Li, Yijian Huang, Dattatraya Patil, and Matin G. Sanda

## S1. Simulation Study with Discrete Covariates

We use the same setting as the simulation study in Janes and Pepe (2009), where covariate $Z$ is a discrete scalar. As discussed in the introduction section of main manuscript, our point estimator reduces to the nonparametric estimator considered by Janes and Pepe (2009) when covariate takes finite values. However, our inference procedure, especially the sample-based variance estimator, is different from the proposal in Janes and Pepe (2009) and is the focus of the evaluation here.

The biomarker from cases $M_1$ and controls $M_0$ follows normal distribution conditional on $Z$: $M_0 \sim N(0,1)$ and $M_1 \sim N(0.9,1)$ if $Z = 0$, $M_0 \sim N(0.2,1)$ and $M_1 \sim N(0.9,1)$ if $Z = 1$. We consider both specificity at controlled sensitivity level and sensitivity at controlled specificity level in this study, since specificity at controlled senstivity level is of interest for this paper whereas the reverse way allows us to directly compare with the performance in Janes and Pepe (2009). The true covariate-adjusted specificity is 0.20, 0.31, 0.48, and 0.60 under controlled sensitivity level 0.95, 0.90, 0.85 and 0.80. The true covariate-adjusted sensitivity is 0.21, 0.33, 0.50, 0.80 under controlled specificity level 0.95, 0.90, 0.85 and 0.80. We implement both sample-based standard error and bootstrap-derived standard error. With mean and SE of estimators, we construct Wald-type confidence intervals as well as logit-transformed confidence interval (Pepe et al. 2003, page 102). Janes and Pepe (2009) reported that logit-transformed confidence interval can improve coverage when controlled specificity is close to 0 or 1.

Table S1 and S2 present the simulation results from this setting. Our standard error estimators are close to the standard deviation obtained in both tables, indicating that our proposed inference procedures are effective. Comparing Table S1 and the results in Janes and Pepe (2009), our sample-based inference procedure achieves better covarage rate, especially when controlled specificity is 0.95. Although Table S2 focus on estimating specificity under controlled sensitivity and Janes and Pepe (2009) is the reverse way, the evaluation metrics such as percentage bias and confidence interval coverage rates are comparable between the two studies. We again find that our proposed method has similar or even higher coverage rate comparing with Janes and Pepe (2009).

## S2. Proof of Theorems

*Proof of Theorem 1.* We first establish the consistency and asymptotic normality of $\widehat{\boldsymbol{\beta}}$. These results for quantile regression have been established by, for example, Koenker (2005, section 4.1.1 and theorem 4.1) under fixed design. Although we consider random design, among other assumptions, similar arguments follow through to give the consistency of $\widehat{\boldsymbol{\beta}}$ and

$$n_1^{1/2}(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0) \quad \overset{d}{\to} \quad N\left(0, \ \rho_0(1-\rho_0)\boldsymbol{D}_1^{-1}\boldsymbol{D}_0\boldsymbol{D}_1^{-1}\right), \tag{1}$$

Table S1: Results under the simulation setting with discrete covariate values. Sensitivity $\rho_0$ under controlled specificity $\phi_0$ is estimated and presented.

| $n_1 = n_0$ | Bias | SD | Sample-based | | | Boostrap-based | | |
|---|---|---|---|---|---|---|---|---|
| | | | SE | Cov | LCov | SE | Cov | LCov |
| $\phi_0 = 0.95,\ \rho_0 = 0.21$ | | | | | | | | |
| 100 | 310 | 885 | 1190 | 93.00 | 92.66 | 905 | 92.02 | 90.48 |
| 200 | 166 | 642 | 795 | 91.52 | 91.74 | 671 | 93.78 | 92.94 |
| 500 | 93 | 416 | 467 | 91.72 | 91.52 | 429 | 93.56 | 93.10 |
| 1000 | 44 | 293 | 324 | 93.16 | 93.02 | 304 | 93.88 | 93.74 |
| $\phi_0 = 0.90,\ \rho_0 = 0.33$ | | | | | | | | |
| 100 | 200 | 920 | 1126 | 92.50 | 93.30 | 962 | 93.78 | 94.20 |
| 200 | 95 | 670 | 770 | 92.20 | 92.98 | 691 | 94.04 | 94.38 |
| 500 | 49 | 424 | 466 | 92.86 | 92.76 | 438 | 93.88 | 93.96 |
| 1000 | 19 | 300 | 323 | 93.68 | 93.48 | 309 | 94.12 | 94.60 |
| $\phi_0 = 0.85,\ \rho_0 = 0.50$ | | | | | | | | |
| 100 | 108 | 872 | 1017 | 92.42 | 93.84 | 919 | 93.98 | 95.80 |
| 200 | 39 | 640 | 706 | 92.58 | 93.46 | 657 | 93.96 | 94.88 |
| 500 | 20 | 400 | 434 | 93.46 | 93.62 | 416 | 93.74 | 94.54 |
| 1000 | 2 | 284 | 300 | 94.00 | 94.20 | 291 | 94.34 | 94.56 |
| $\phi_0 = 0.50,\ \rho_0 = 0.80$ | | | | | | | | |
| 100 | -15 | 607 | 683 | 93.70 | 95.22 | 640 | 95.10 | 96.12 |
| 200 | -13 | 421 | 469 | 94.32 | 95.34 | 444 | 94.70 | 95.80 |
| 500 | 0 | 270 | 286 | 94.24 | 94.94 | 276 | 94.82 | 95.04 |
| 1000 | -16 | 192 | 200 | 95.04 | 95.04 | 195 | 95.10 | 94.96 |

Bias, $(\widehat{\rho}-\rho)\times 10^4$; SD, standard deviation $(\times 10^4)$; SE, mean standard error $(\times 10^4)$; Cov (%) and LCov (%), coverage rates of 95% confidence interval and logit-transformed confidence interval.

Table S2: Results under the simulation setting with discrete covariate values. Specificity $\rho_0$ under controlled sensitivity $\phi_0$ is estimated and presented.

| $n_1 = n_0$ | Bias | SD | Sample-based | | | Boostrap-based | | |
|---|---|---|---|---|---|---|---|---|
| | | | SE | Cov | LCov | SE | Cov | LCov |
| $\rho_0 = 0.95,\ \phi_0 = 0.19$ | | | | | | | | |
| 100 | 202 | 819 | 1071 | 90.76 | 93.08 | 839 | 93.12 | 91.46 |
| 200 | 94 | 570 | 724 | 91.04 | 92.30 | 608 | 94.94 | 94.58 |
| 500 | 32 | 372 | 417 | 90.28 | 91.70 | 382 | 94.20 | 94.16 |
| 1000 | 22 | 261 | 287 | 91.36 | 91.74 | 270 | 94.58 | 94.66 |
| $\rho_0 = 0.90,\ \phi_0 = 0.30$ | | | | | | | | |
| 100 | 144 | 871 | 1117 | 90.86 | 92.34 | 909 | 94.26 | 95.28 |
| 200 | 42 | 612 | 730 | 90.94 | 92.02 | 651 | 94.68 | 95.26 |
| 500 | 27 | 389 | 439 | 91.76 | 92.28 | 410 | 94.80 | 94.98 |
| 1000 | 6 | 282 | 303 | 92.08 | 92.36 | 290 | 93.98 | 94.30 |
| $\rho_0 = 0.85,\ \phi_0 = 0.47$ | | | | | | | | |
| 100 | 42 | 886 | 1061 | 91.10 | 92.36 | 912 | 94.18 | 95.66 |
| 200 | 16 | 632 | 709 | 90.88 | 91.74 | 647 | 93.84 | 94.68 |
| 500 | 10 | 398 | 436 | 92.48 | 92.68 | 411 | 94.32 | 94.66 |
| 1000 | -6 | 286 | 301 | 92.58 | 92.78 | 288 | 94.22 | 94.32 |
| $\rho_0 = 0.80,\ \phi_0 = 0.59$ | | | | | | | | |
| 100 | 1 | 832 | 971 | 91.22 | 92.28 | 863 | 94.48 | 95.84 |
| 200 | -8 | 588 | 660 | 91.40 | 92.04 | 612 | 94.52 | 95.50 |
| 500 | 0 | 372 | 402 | 92.62 | 92.90 | 384 | 94.72 | 95.12 |
| 1000 | -3 | 265 | 280 | 93.46 | 93.42 | 270 | 94.74 | 94.76 |

Bias, $(\widehat{\phi} - \phi) \times 10^4$; SD, standard deviation $(\times 10^4)$; SE, mean standard error $(\times 10^4)$; Cov (%) and LCov (%), coverage rates of 95% confidence interval and logit-transformed confidence interval.

where $\boldsymbol{D}_0 = E\widetilde{\boldsymbol{Z}}_1^{\otimes 2}$ and $\boldsymbol{D}_1 = E\{F_1^{'}(\widetilde{\boldsymbol{Z}}_1^T\boldsymbol{\beta}_0)\widetilde{\boldsymbol{Z}}_1^{\otimes 2}\}$.

Next we turn to $\widehat{\phi}$. By Condition 4a and the consistency of $\widehat{\boldsymbol{\beta}}$, the consistency of $\widehat{\phi}$ can be easily established. For asymptotic normality, we have

$$n_0^{1/2}(\widehat{\phi} - \phi_0) = n_0^{-1/2}\sum_{i=1}^{n_0}\{I(M_{0i} \leq \widetilde{\boldsymbol{Z}}_{0i}^T\widehat{\boldsymbol{\beta}}) - \Pr(M_0 \leq \widetilde{\boldsymbol{Z}}_0^T\boldsymbol{\beta}_0)\}$$

$$= n_0^{-1/2}\sum_{i=1}^{n_0}\{I(M_{0i} \leq \widetilde{\boldsymbol{Z}}_{0i}^T\widehat{\boldsymbol{\beta}}) - \Pr(M_0 \leq \widetilde{\boldsymbol{Z}}_0^T\boldsymbol{\beta}|\boldsymbol{\beta} = \widehat{\boldsymbol{\beta}})\}$$

$$+ n_0^{1/2}\{\Pr(M_0 \leq \widetilde{\boldsymbol{Z}}_0^T\boldsymbol{\beta}|\boldsymbol{\beta} = \widehat{\boldsymbol{\beta}}) - \Pr(M_0 \leq \widetilde{\boldsymbol{Z}}_0^T\boldsymbol{\beta}_0)\}$$

$$\equiv A_n(\widehat{\boldsymbol{\beta}}) + B_n.$$

Since $F_0(t; \boldsymbol{z})$ is differentiable at $\widetilde{\boldsymbol{Z}}_0^T\boldsymbol{\beta}_0$, in light of (1), Delta method leads to

$$B_n \overset{d}{\to} N\left(0,\ c\rho_0(1 - \rho_0)\boldsymbol{D}_2^T\boldsymbol{D}_1^{-1}\boldsymbol{D}_0\boldsymbol{D}_1^{-1}\boldsymbol{D}_2\right), \tag{2}$$

where $\boldsymbol{D}_2 = E\{F_0^{'}(\widetilde{\boldsymbol{Z}}_0^T\boldsymbol{\beta}_0)\widetilde{\boldsymbol{Z}}_0\}$. Meanwhile, $A_n(\widehat{\boldsymbol{\beta}})$ can be written as

$$A_n(\boldsymbol{\beta}_0) + \{A_n(\widehat{\boldsymbol{\beta}}) - A_n(\boldsymbol{\beta}_0)\}, \tag{3}$$

where $A_n(\boldsymbol{\beta}_0) = n_0^{-1/2}\sum_{i=1}^{n_0}\{I(M_{0i} \leq \widetilde{\boldsymbol{Z}}_{0i}^T\boldsymbol{\beta}_0) - \phi_0\}$. By central limit theorem,

$$A_n(\boldsymbol{\beta}_0) \overset{d}{\to} N\left(0, \phi_0(1 - \phi_0)\right). \tag{4}$$

On the other hand,

$$E[\{A_n(\widehat{\boldsymbol{\beta}}) - A_n(\boldsymbol{\beta}_0)\}^2] = n_0^{-1}\sum_{i=1}^{n_0}E\{I(M_{0i} \leq \widetilde{\boldsymbol{Z}}_{0i}^T\widehat{\boldsymbol{\beta}}) - I(M_{0i} \leq \widetilde{\boldsymbol{Z}}_{0i}^T\boldsymbol{\beta}_0)$$

$$- \Pr(M_0 \leq \widetilde{\boldsymbol{Z}}_0^T\widehat{\boldsymbol{\beta}}|\widehat{\boldsymbol{\beta}}) + \Pr(M_0 \leq \widetilde{\boldsymbol{Z}}_0^T\boldsymbol{\beta}_0)\}^2$$

$$= E\left(E\left[\{I(M_0 \leq \widetilde{\boldsymbol{Z}}_0^T\widehat{\boldsymbol{\beta}}) - I(M_0 \leq \widetilde{\boldsymbol{Z}}_0^T\boldsymbol{\beta}_0)\right.\right.$$

$$\left.\left. - F_0(\widetilde{\boldsymbol{Z}}_0^T\widehat{\boldsymbol{\beta}}) - F_0(\widetilde{\boldsymbol{Z}}_0^T\boldsymbol{\beta}_0)\}^2|\widetilde{\boldsymbol{Z}}_0^T\widehat{\boldsymbol{\beta}}\right]\right)$$

$$\leq E|I(M_0 \leq \widetilde{\boldsymbol{Z}}_0^T\widehat{\boldsymbol{\beta}}) - I(M_0 \leq \widetilde{\boldsymbol{Z}}_0^T\boldsymbol{\beta}_0)|$$

$$\leq E|F_{0|\boldsymbol{z}}(\widetilde{\boldsymbol{Z}}_0^T\widehat{\boldsymbol{\beta}}) - F_{0|\boldsymbol{z}}(\widetilde{\boldsymbol{Z}}_0^T\boldsymbol{\beta}_0)|.$$

By Markov's inequality, $A_n(\widehat{\boldsymbol{\beta}}) - A_n(\boldsymbol{\beta}_0) \overset{d}{\to} 0$.

Together with (2) and (4), Slutsky's theorem yields the result. $\square$

*Proof of Theorem 2.* Start with the cases, and write

$$\Psi_n(\boldsymbol{\beta}, \rho) = n_1^{-1}\sum_{i=1}^{n}\widetilde{\boldsymbol{Z}}_{1i}\{I(M_{1i} > \widetilde{\boldsymbol{Z}}_{1i}^T\boldsymbol{\beta}) - \rho\},$$

$$\Psi(\boldsymbol{\beta}, \rho) = E\left[\widetilde{\boldsymbol{Z}}_1\{I(M_1 > \widetilde{\boldsymbol{Z}}_1^T\boldsymbol{\beta}) - \rho\}\right].$$

It is known that $\{I(M_1 > \widetilde{\boldsymbol{Z}}_1^T\boldsymbol{\beta}) : \boldsymbol{\beta} \in \mathbb{R}^p\}$ is Donsker (e.g. Kosorok 2007, lemma 9.12). Furthermore, $\widetilde{\boldsymbol{Z}}_1$ is bounded by Condition 2. By permanence property of the Donsker class, $\{\widetilde{\boldsymbol{Z}}_1 I(M_1 > \widetilde{\boldsymbol{Z}}_1^T\boldsymbol{\beta}) : \boldsymbol{\beta} \in \mathbb{R}^p\}$ is Donsker. Since Donsker implies Glivenko-Cantelli, it follows that, almost surely

$$\sup_{\boldsymbol{\beta}, \rho \in [\rho_1, \rho_2]} ||\Psi_n(\boldsymbol{\beta}, \rho) - \Psi(\boldsymbol{\beta}, \rho)|| = o(1).$$

4

Thus, $||\Psi\{\widehat{\boldsymbol{\beta}}(\rho),\rho\}|| \le ||\Psi_n\{\widehat{\boldsymbol{\beta}}(\rho),\rho\}|| + ||\Psi_n\{\widehat{\boldsymbol{\beta}}(\rho),\rho\} - \Psi\{\widehat{\boldsymbol{\beta}}(\rho),\rho\}||$ leads to, almost surely,

$$\sup_{\rho\in[\rho_1,\rho_2]} ||\Psi\{\widehat{\boldsymbol{\beta}}(\rho),\rho\}|| = o(1).$$

It remains to be shown that, for any $\epsilon > 0$, there exists $\delta > 0$ such that if $\sup_{\rho\in[\rho_1,\rho_2]} ||\Psi\{\boldsymbol{\beta}(\rho),\rho\}|| <$ $\delta$, then $\sup_{\rho\in[\rho_1,\rho_2]} ||\boldsymbol{\beta}(\rho) - \boldsymbol{\beta}_0(\rho)|| < \epsilon$. Suppose that this is not true. Thus, for each $\delta > 0$, there exists $(\zeta,\nu)$ such that $||\Psi(\zeta,\nu) - \Psi\{\boldsymbol{\beta}_0(\nu),\nu\}|| < \delta$ and $||\zeta - \boldsymbol{\beta}_0(\nu)|| > c$ for some constant $c > 0$. Then, there exists a subsequence of $(\zeta,\nu)$ that converges to, say, $(\zeta_0,\nu_0)$, which implies that $\zeta_0 \ne \boldsymbol{\beta}_0(\nu_0)$ also solves $\Psi(\boldsymbol{\beta},\nu_0)$. This contradicts the fact that $\boldsymbol{\beta}_0(\rho)$ is the unique solution of $\Psi(\beta,\rho)$ for all $\rho \in [\rho_1,\rho_2]$, as guaranteed by Condition 3 and 4a. Therefore,

$$\sup_{\rho\in[\rho_1,\rho_2]} ||\widehat{\boldsymbol{\beta}}(\rho) - \boldsymbol{\beta}_0(\rho)|| = o(1)$$

almost surely.

In light of the above Donsker result, for given $\rho$, $n_1^{1/2}\{\Psi_n(\boldsymbol{\beta},\rho) - \Psi(\boldsymbol{\beta},\rho)\}$ converges weakly to a Gaussian process. Under Conditions 2 and 4b, $n_1^{1/2}\{\Psi_n(\boldsymbol{\beta},\rho) - \Psi(\boldsymbol{\beta},\rho)\}$ is asymptotically uniformly equicontinuous in probability using arguments similar to Huang (2017), appendix. Thus, for any positive sequence $d_n = o(1)$,

$$\sup_{||\boldsymbol{\beta}-\boldsymbol{\beta}'||<d_n,\ \rho\in[\rho_1,\rho_2]} n_1^{1/2}||\Psi_n(\boldsymbol{\beta},\rho) - \Psi_n(\boldsymbol{\beta}',\rho) - \Psi(\boldsymbol{\beta},\rho) + \Psi(\boldsymbol{\beta}',\rho)|| = o_p(1);$$

note that the above expression does not actually involve $\rho$. Therefore,

$$\sup_{\rho\in[\rho_1,\rho_2]} ||\Psi_n\{\boldsymbol{\beta}_0(\rho),\rho\} + \Psi\{\widehat{\boldsymbol{\beta}}(\rho),\rho\}|| = o_p(n_1^{-1/2}).$$

Under Condition 4b, by component-wise Taylor expansion, one can show that, almost surely,

$$\sup_{\rho\in[\rho_1,\rho_2]} \frac{||\Psi\{\widehat{\boldsymbol{\beta}}(\rho),\rho\} + E[\widetilde{\boldsymbol{Z}}_1^{\otimes 2} f_1\{\widetilde{\boldsymbol{Z}}_1^T\boldsymbol{\beta}_0(\rho)\widetilde{\boldsymbol{Z}}_1\}]\{\widehat{\boldsymbol{\beta}}(\rho) - \boldsymbol{\beta}_0(\rho)\}||}{||\widehat{\boldsymbol{\beta}}(\rho) - \boldsymbol{\beta}_0(\rho)||} = o(1).$$

Thus,

$$n_1^{1/2}\{\widehat{\boldsymbol{\beta}}(\rho) - \boldsymbol{\beta}_0(\rho)\} = n_1^{1/2}\big(E[\widetilde{\boldsymbol{Z}}_1^{\otimes 2} f_1\{\widetilde{\boldsymbol{Z}}_1^T\boldsymbol{\beta}_0(\rho)\widetilde{\boldsymbol{Z}}_1\}]\big)^{-1}\Psi_n\{\boldsymbol{\beta}_0(\rho),\rho\} + o_p(1),$$

uniformly in $\rho \in [\rho_1,\rho_2]$. Therefore, $n_1^{1/2}\{\widehat{\boldsymbol{\beta}}(\cdot) - \boldsymbol{\beta}(\cdot)\}$ over $[\rho_1,\rho_2]$ converges weakly to a Gaussian process.

Now, we turn to the controls. Write $\Gamma_n(\boldsymbol{\beta}) = n_0^{-1}\sum_{j=1}^{n_0} I(M_{0j} \le \widetilde{\boldsymbol{Z}}_{0j}^T\boldsymbol{\beta})$ and $\Gamma(\boldsymbol{\beta}) = \Pr(M_0 \le \widetilde{\boldsymbol{Z}}_0^T\boldsymbol{\beta})$. Similar arguments as above give

$$\sup_{\boldsymbol{\beta}} |\Gamma_n(\boldsymbol{\beta}) - \Gamma(\boldsymbol{\beta})| = o(1).$$

Thus,

$$\sup_{\rho\in[\rho_1,\rho_2]} |\widehat{\phi}(\rho) - \phi_0(\rho)| \le \sup_{\rho\in[\rho_1,\rho_2]} |\Gamma\{\widehat{\boldsymbol{\beta}}(\rho)\} - \Gamma\{\boldsymbol{\beta}_0(\rho)\}| + o(1) = o(1)$$

almost surely, given the strong consistency of $\widehat{\boldsymbol{\beta}}(\cdot)$ and the continuity of $\Gamma(\boldsymbol{\beta})$. To establish the weak convergence of $\widehat{\phi}(\rho)$, one can show that, for any positive sequence $d_n = o(1)$,

$$\sup_{||\boldsymbol{\beta}-\boldsymbol{\beta}'||<d_n} n_0^{1/2}|\Gamma_n(\boldsymbol{\beta}) - \Gamma_n(\boldsymbol{\beta}') - \Gamma(\boldsymbol{\beta}) + \Gamma(\boldsymbol{\beta}')| = o_p(1),$$

using similar arguments as for the cases. Therefore

$$\begin{aligned}
n_0^{1/2}\{\widehat{\phi}(\rho) - \phi_0(\rho)\} &= n_0^{1/2}\big[\Gamma_n\{\widehat{\boldsymbol{\beta}}(\rho)\} - \Gamma\{\boldsymbol{\beta}_0(\rho)\}\big] \\
&= n_0^{1/2}\big[\Gamma_n\{\boldsymbol{\beta}_0(\rho)\} - \Gamma\{\boldsymbol{\beta}_0(\rho)\}\big] + n_0^{1/2}\big[\Gamma\{\widehat{\boldsymbol{\beta}}(\rho)\} - \Gamma\{\boldsymbol{\beta}_0(\rho)\}\big] + o_p(1) \\
&= n_0^{1/2}\big[\Gamma_n\{\boldsymbol{\beta}_0(\rho)\} - \Gamma\{\boldsymbol{\beta}_0(\rho)\}\big] + n_0^{1/2}\Gamma'\{\boldsymbol{\beta}_0(\rho)\}\{\widehat{\boldsymbol{\beta}}(\rho) - \boldsymbol{\beta}_0(\rho)\} + o_p(1)
\end{aligned}$$

uniformly in $\rho \in [\rho_1,\rho_2]$. Then, the weak convergence of $\widehat{\phi}(\rho)$ follows. $\square$

## S3. Details about two monotonization methods

To recover the monotonicity of the constructed ROC curves, we applied two monotonization methods based on Huang (2017). Using the notations in the main manuscript Section 2, for linear dynamic regression model with covariate $\boldsymbol{z}$ and coefficient $\boldsymbol{\beta}$, the quantile regression model is

$$F_1^{-1}(t; \boldsymbol{z}) = \widetilde{\boldsymbol{z}}^T \boldsymbol{\beta}(t),$$

where $\widetilde{\boldsymbol{z}} = (1, \boldsymbol{z})$. Denote the estimator of $\boldsymbol{\beta}(t)$ by $\widehat{\boldsymbol{\beta}}(t)$ and the estimator for specificity by $\widehat{\phi}(\cdot)$. Note that both $\widehat{\boldsymbol{\beta}}(t)$ and $\widehat{\phi}(\cdot)$ are piece-wise constant and thus we can identify a countable set of breakpoints. Given a starting quantile point $\tau$, let

$$\max \left[ t : t < \tau, \sup_{\boldsymbol{z} \in \mathcal{Z}_s} \{ \widetilde{\boldsymbol{z}}^T \widehat{\boldsymbol{\beta}}(t) - \widetilde{\boldsymbol{z}}^T \widehat{\boldsymbol{\beta}}(\tau) \} \leq 0 \right]$$

be the left nearest monotonicity-respecting neighbor and

$$\min \left[ t : t > \tau, \inf_{\boldsymbol{z} \in \mathcal{Z}_s} \{ \widetilde{\boldsymbol{z}}^T \widehat{\boldsymbol{\beta}}(t) - \widetilde{\boldsymbol{z}}^T \widehat{\boldsymbol{\beta}}(\tau) \} \geq 0 \right]$$

be the right nearest monotonicity respecting neighbor. We denote the collection of all these points, including the starting point $\tau$, by $\mathcal{M}$. Huang (2017) proposed to adopt an adaptive interpolation method to connect the original estimator $\widehat{\boldsymbol{\beta}}(\cdot)$ linearly between adjacent points in the break point set $\mathcal{M}$. Denote the monotonicity-respecting estimator by $\widetilde{\boldsymbol{\beta}}(\cdot)$. For any $t$ between two adjacent points in $\mathcal{M}$, say $\tau_1 < t < \tau_2$, $\widetilde{\boldsymbol{\beta}}(t)$ is constructed by

$$\widetilde{\boldsymbol{\beta}}(t) = \frac{\tau_2 - t}{\tau_2 - \tau_1} \widehat{\boldsymbol{\beta}}(\tau_1) + \frac{t - \tau_1}{\tau_2 - \tau_1} \widehat{\boldsymbol{\beta}}(\tau_2).$$

For $t < min(\mathcal{M})$ we set $\widetilde{\boldsymbol{\beta}}(t) = \widehat{\boldsymbol{\beta}}(min\mathcal{M})$, and for $t > max(\mathcal{M})$ we set $\widetilde{\boldsymbol{\beta}}(t) = \widehat{\boldsymbol{\beta}}(max\mathcal{M})$. The regression-based monotonization method directly applies the above approach on the coefficient estimator $\widehat{\boldsymbol{\beta}}$ of our quantile regression model. The ROC-based monotonization method uses regular quantile regression estimator $\widehat{\boldsymbol{\beta}}(\cdot)$ to obtain the estimated specificities $\widehat{\phi}(\cdot)$, and then applies the adaptive interpolation approach on $\widehat{\phi}(\cdot)$ to obtain $\widetilde{\phi}(\cdot)$.
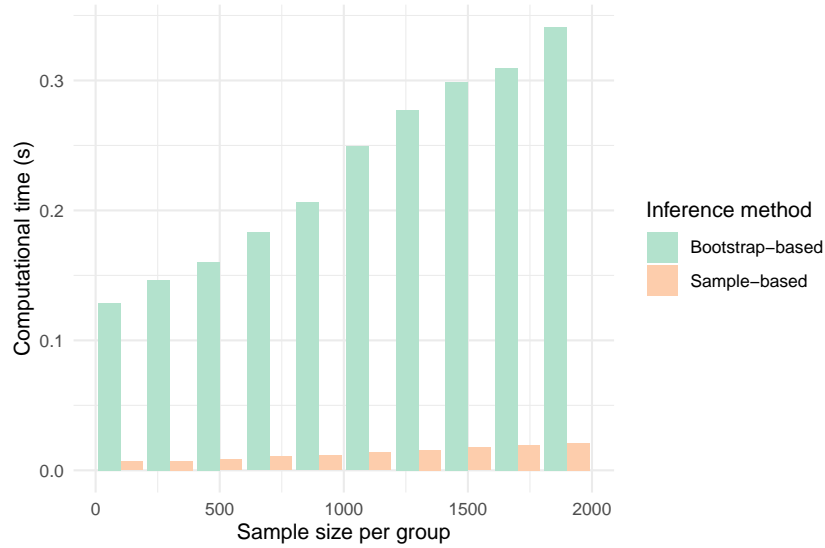
Figure S1: Summarization of the computation time using two inference methods with different sample sizes.

# References

Huang, Y. (2017). Restoration of monotonicity respecting in dynamic regression. Journal of the American Statistical Association 112(518), 613–622.

Janes, H. and M. S. Pepe (2009). Adjusting for covariate effects on classification accuracy using the covariate-adjusted receiver operating characteristic curve. Biometrika 96(2), 371–382.

Kosorok, M. R. (2007). Introduction to empirical processes and semiparametric inference. Springer Science & Business Media.

Pepe, M. S. et al. (2003). The statistical evaluation of medical tests for classification and prediction. Medicine.