

## Supplementary Material to “Biomarker Evaluation and Comparison Using the Controls as a Reference Population”

YING HUANG\*

*Fred Hutchinson Cancer Research Center Public Health Sciences,  
1100 Fairview Avenue N., M3-A410, Seattle, WA 98109, USA  
yhuang124@gmail.com*

MARGARET SULLIVAN PEPE

*Fred Hutchinson Cancer Research Center Public Health Sciences,  
1100 Fairview Avenue N., M2-B500, Seattle, WA 98109, USA*

### WEB APPENDIX A: TABLES AND FIGURES

Table 1. *P-values for comparing the percentile value distributions of benign tumor cases and ovarian cancer cases,  $n_D = 41, n_1 = 24, n_2 = 66$ . Tests comparing raw marker values between benign tumor cases and ovarian cancer cases yielded  $p < 0.0001$  for both the *t*-test and the Wilcoxon rank sum test.*

Test	Unconditional				Conditional			
	$\hat{F}$ empirical		$\hat{F}$ parametric		$\hat{F}$ empirical		$\hat{F}$ parametric	
	Asym <sup>1</sup>	Boot <sup>2</sup>	Asym	Boot	Asym	Boot	Asym	Boot
t-test	0.0009	0.0006	0.0018	0.0013	0.0005	0.0003	0.0012	0.0009
WRS <sup>3</sup>	-	< 0.0001 <sup>s</sup>	-	< 0.0001 <sup>s</sup>	< 0.0001	< 0.0001 <sup>s</sup>	< 0.0001	< 0.0001 <sup>s</sup>

<sup>1</sup>asymptotic variance

<sup>2</sup>nonparametric bootstrap variance or smoothed bootstrap variance indicated by superscript <sup>s</sup>

<sup>3</sup>Wilcoxon rank sum test

Table 2.  $P$ -values for comparing the case percentile value distributions of CA-19-9 and CA-125,  $n_{\bar{D}} = 51, n_D = 90$ .

Test Statistic	$\hat{F}$ empirical CDF		$\hat{F}$ parametric	
	Asym <sup>1</sup>	Boot <sup>2</sup>	Asym	Boot
	Marginal			
Mean Difference <sup>1</sup>	0.007	0.007	0.009	0.01
WRS <sup>3</sup>	-	< 0.0001 <sup>s</sup>	-	0.0006 <sup>s</sup>
WSR <sup>4</sup>	-	< 0.0001 <sup>s</sup>	-	0.0001 <sup>s</sup>
Sign <sup>5</sup>	-	< 0.0001 <sup>s</sup>	-	< 0.0001 <sup>s</sup>
	Covariate Adjusted			
Mean Difference	< 0.0001	< 0.0001	< 0.0001	< 0.0001
WRS	-	< 0.0001	-	< 0.0001
WSR	-	< 0.0001	-	< 0.0001
Sign	-	< 0.0001	-	< 0.0001

<sup>1</sup>asymptotic variance

<sup>2</sup>nonparametric bootstrap variance, or smoothed bootstrap variance indicated by superscript <sup>s</sup>

<sup>3</sup>Wilcoxon rank sum test statistic

<sup>4</sup>Wilcoxon signed rank test statistic

<sup>5</sup>Sign test statistic

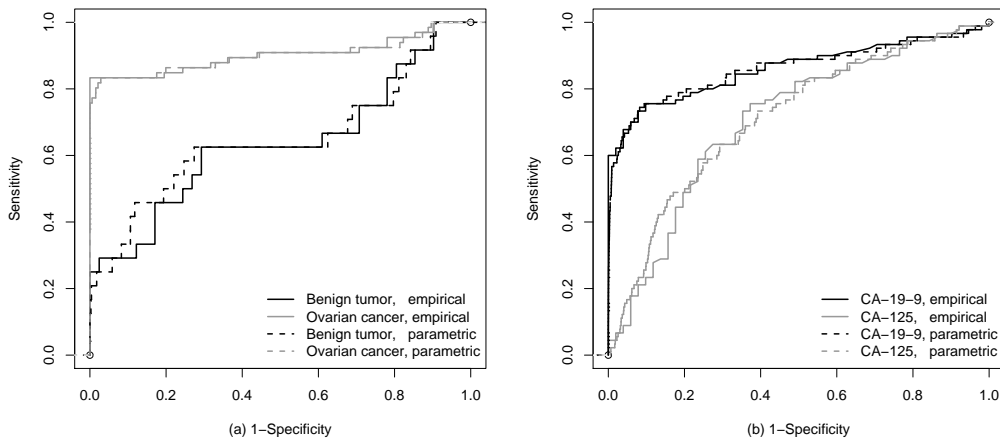


Fig. 1. A. ROC curves (control distribution is empirically or parametrically estimated) (a) for benign tumor cases and ovarian cancer cases in the ovarian cancer data, and (b) for CA-19-9 and CA-125 in the pancreatic cancer data.

## WEB APPENDIX B: PROOF OF THEOREMS

## 0.1 Proof of Theorem 1

With  $z = 1, 2$ , let  $Y_{\bar{D}i}, i = 1, \dots, n_{\bar{D}}$  and  $Y_{zj}, j = 1, \dots, n_z$  be the marker measurements in controls and the  $z^{\text{th}}$  type of cases. Let  $F$  and  $G_z$  be the corresponding marker distribution functions. When  $F$  is estimated empirically,

$$\begin{aligned}
& \sqrt{n_{\bar{D}}} (\hat{\Delta} - \Delta) \\
&= \sqrt{n_{\bar{D}}} \left\{ \frac{1}{n_1 n_{\bar{D}}} \sum_{i=1}^{n_{\bar{D}}} \sum_{j=1}^{n_1} I(Y_{\bar{D}i} \leq Y_{1j}) - \frac{1}{n_2 n_{\bar{D}}} \sum_{i=1}^{n_{\bar{D}}} \sum_{j=1}^{n_2} I(Y_{\bar{D}i} \leq Y_{2j}) - P(Y_{\bar{D}} \leq Y_1) + P(Y_{\bar{D}} \leq Y_2) \right\} \times 100 \\
&=^* \sqrt{n_{\bar{D}}} \left( \frac{1}{n_{\bar{D}}} \sum_{i=1}^{n_{\bar{D}}} [1 - G_1(Y_{\bar{D}i}) - \{1 - G_2(Y_{\bar{D}i})\}] - P(Y_{\bar{D}} \leq Y_1) + P(Y_{\bar{D}} \leq Y_2) \right) \times 100 \\
&+ \sqrt{n_{\bar{D}}} \left\{ \frac{1}{n_1} \sum_{j=1}^{n_1} F(Y_{1j}) - P(Y_{\bar{D}} \leq Y_1) \right\} \times 100 - \sqrt{n_{\bar{D}}} \left\{ \frac{1}{n_2} \sum_{j=1}^{n_2} F(Y_{2j}) - P(Y_{\bar{D}} \leq Y_2) \right\} \times 100 + o_p(1) \\
&= -\frac{1}{\sqrt{n_{\bar{D}}}} \sum_{i=1}^{n_{\bar{D}}} \{R_1(Y_{\bar{D}i}) - R_2(Y_{\bar{D}i}) - E(R_1 - R_2)\} \\
&+ \frac{1}{\sqrt{\lambda_1 n_1}} \sum_{j=1}^{n_1} \{Q_{1j} - E(Q_1)\} - \frac{1}{\sqrt{\lambda_2 n_2}} \sum_{j=1}^{n_2} \{Q_{2j} - E(Q_2)\} + o_p(1),
\end{aligned}$$

where \* can be proved with the U-statistic theory (van der Vaart, 1998).

When  $F$  is modeled parametrically, assume that  $\sqrt{n} (\hat{\theta} - \theta)$  can be represented as  $\frac{1}{\sqrt{n_{\bar{D}}}} \sum_{i=1}^{n_{\bar{D}}} \psi_i + o_p(1)$ , where  $\psi_i, i = 1, \dots, n_{\bar{D}}$  are independent identically distributed variables with  $E(\psi_i) = 0$  and  $\text{var}(\psi_i) = \Sigma(\theta)$ , we have

$$\begin{aligned}
& \sqrt{n_{\bar{D}}} (\hat{\Delta} - \Delta) \\
&= \sqrt{n_{\bar{D}}} \left\{ \frac{1}{n_1} \sum_{j=1}^{n_1} F_{\hat{\theta}}(Y_{1j}) - \frac{1}{n_2} \sum_{j=1}^{n_2} F_{\hat{\theta}}(Y_{2j}) - E_{Y_1} F_{\hat{\theta}}(Y) + E_{Y_2} F_{\hat{\theta}}(Y) \right\} \times 100 \\
&=^* \sqrt{n_{\bar{D}}} \{E_{Y_1} F_{\hat{\theta}}(Y) - E_{Y_1} F_{\theta}(Y) - E_{Y_2} F_{\hat{\theta}}(Y) + E_{Y_2} F_{\theta}(Y)\} \times 100 \\
&+ \sqrt{n_{\bar{D}}} \left\{ \frac{1}{n_1} \sum_{j=1}^{n_1} F_{\theta}(Y_{1j}) - \frac{1}{n_2} \sum_{j=1}^{n_2} F_{\theta}(Y_{2j}) - E_{Y_1} F_{\theta}(Y) + E_{Y_2} F_{\theta}(Y) \right\} \times 100 + o_p(1) \\
&= \frac{1}{n_{\bar{D}}} \sum_{i=1}^{n_{\bar{D}}} \frac{\partial \Delta_i}{\partial \theta} \psi_i + \frac{1}{\sqrt{\lambda_1 n_1}} \sum_{j=1}^{n_1} \{Q_{1j} - E(Q_1)\} - \frac{1}{\sqrt{\lambda_2 n_2}} \sum_{j=1}^{n_2} \{Q_{2j} - E(Q_2)\} + o_p(1),
\end{aligned}$$

where \* follows since  $\mathcal{F} = \{F_{\theta}(y) : \theta \in \Theta\}$  is a Donsker class (van der Vaart and Wellner, 1996, Theorem 2.7.5 on page 159).

## 0.2 Proof of Theorem 2

With  $z = 1, 2$  and  $k = 1, \dots, K$ , let  $Y_{Di}^k, i = 1, \dots, n_{\bar{D}k}$  and  $Y_{zj}^k, j = 1, \dots, n_{zk}$  be the marker measurements for controls and the  $z^{\text{th}}$  type of cases in the  $k^{\text{th}}$  covariate category. Let  $F^k$  and  $G_z^k$  be the corresponding marker distribution functions. When  $F^k$  is estimated empirically,

$$\begin{aligned}
& \sqrt{n_{\bar{D}}} (\hat{\Delta} - \Delta) \\
&= \sqrt{n_{\bar{D}}} \left\{ \frac{1}{n_1} \sum_{k=1}^K \sum_{i=1}^{n_{\bar{D}k}} \sum_{j=1}^{n_{1k}} \frac{1}{n_{\bar{D}k}} I(Y_{Di}^k \leq Y_{1j}^k) - \frac{1}{n_2} \sum_{k=1}^K \sum_{i=1}^{n_{\bar{D}k}} \sum_{j=1}^{n_{2k}} \frac{1}{n_{\bar{D}k}} I(Y_{Di}^k \leq Y_{2j}^k) \right. \\
&\quad \left. - \sum_{k=1}^K p_{1k} P(Y_D^k \leq Y_1^k) + \sum_{k=1}^K p_{2k} P(Y_D^k \leq Y_2^k) \right\} \times 100 \\
&= \sqrt{n_{\bar{D}}} \left\{ \sum_{k=1}^K \sum_{i=1}^{n_{\bar{D}k}} \frac{1}{n_{\bar{D}k}} [p_{1k} \{1 - G_1^k(Y_{Di}^k)\} - p_{2k} \{1 - G_2^k(Y_{Di}^k)\}] - \sum_{k=1}^K p_{1k} P(Y_D^k \leq Y_1^k) \right. \\
&\quad \left. + \sum_{k=1}^K p_{2k} P(Y_D^k \leq Y_2^k) \right\} \times 100 + \sqrt{n_{\bar{D}}} \left\{ \frac{1}{n_1} \sum_{k=1}^K \sum_{j=1}^{n_{1k}} F^k(Y_{1j}^k) - \sum_{k=1}^K p_{1k} P(Y_D^k \leq Y_1^k) \right\} \times 100 \\
&\quad - \sqrt{n_{\bar{D}}} \left\{ \frac{1}{n_2} \sum_{k=1}^K \sum_{j=1}^{n_{2k}} F^k(Y_{2j}^k) - \sum_{k=1}^K p_{2k} P(Y_D^k \leq Y_2^k) \right\} \times 100 + o_p(1) \\
&= \sum_{k=1}^K -\frac{1}{\sqrt{n_{\bar{D}k} p_{\bar{D}k}}} \sum_{i=1}^{n_{\bar{D}k}} \{p_{1k} R_1^k(Y_{Di}^k) - p_{2k} R_2^k(Y_{Di}^k) - p_{1k} E(R_1^k) + p_{2k} E(R_2^k)\} \\
&\quad + \frac{1}{\sqrt{\lambda_1 n_1}} \sum_{j=1}^{n_1} \{Q_{1X_j} - E(Q_{1X})\} - \frac{1}{\sqrt{\lambda_2 n_2}} \sum_{j=1}^{n_2} \{Q_{2X_j} - E(Q_{2X})\} + o_p(1).
\end{aligned}$$

This proves Theorem 2(a). Now, consider when  $F$  is estimated parametrically and suppose the parameter  $\theta$  can be represented as  $\sqrt{n_{\bar{D}}} (\hat{\theta} - \theta) = \frac{1}{\sqrt{n_{\bar{D}}}} \sum_{i=1}^{n_{\bar{D}}} \psi_i + o_p(1)$ , where  $\psi_i, i = 1, \dots, n_{\bar{D}}$  are independent identically distributed variables with  $E(\psi_i) = 0$  and  $\text{var}(\psi_i) = \Sigma(\theta)$ . Let  $X_{zj}, j = 1, \dots, n_z$  be the covariate measurements for the  $z^{\text{th}}$  type of cases, we have

$$\begin{aligned}
& \sqrt{n_{\bar{D}}} (\hat{\Delta} - \Delta) \\
&= \sqrt{n_{\bar{D}}} \left\{ \frac{1}{n_1} \sum_{j=1}^{n_1} F_{\hat{\theta}}(Y_{1j}|X_{1j}) - \frac{1}{n_2} \sum_{j=1}^{n_2} F_{\hat{\theta}}(Y_{2j}|X_{2j}) - E_{Y_1, X_1} F_{\theta}(Y|X) + E_{Y_2, X_2} F_{\theta}(Y|X) \right\} \times 100 \\
&= \sqrt{n_{\bar{D}}} \{E_{Y_1, X_1} F_{\hat{\theta}}(Y|X) - E_{Y_1, X_1} F_{\theta}(Y|X) - E_{Y_2, X_2} F_{\hat{\theta}}(Y|X) + E_{Y_2, X_2} F_{\theta}(Y|X)\} \times 100 \\
&\quad + \sqrt{n_{\bar{D}}} \left\{ \frac{1}{n_1} \sum_{j=1}^{n_1} F_{\theta}(Y_{1j}|X_{1j}) - \frac{1}{n_2} \sum_{j=1}^{n_2} F_{\theta}(Y_{2j}|X_{2j}) - E_{Y_1, X_1} F_{\theta}(Y|X) + E_{Y_2, X_2} F_{\theta}(Y|X) \right\} \times 100 + o_p(1) \\
&= \frac{1}{\sqrt{n_{\bar{D}}}} \sum_{i=1}^{n_{\bar{D}}} \frac{\partial \Delta_i}{\partial \theta} \psi_i + \frac{1}{\sqrt{\lambda_1 n_1}} \sum_{j=1}^{n_1} \{Q_{1X_j} - E(Q_{1X})\} - \frac{1}{\sqrt{\lambda_2 n_2}} \sum_{j=1}^{n_2} \{Q_{2X_j} - E(Q_{2X})\} + o_p(1),
\end{aligned}$$

where  $\{E_{Y_z, X_z}, z = 1, 2\}$  is integrated over the joint distribution of the marker  $Y$  and the covariate  $X$  in the  $z^{th}$  type of cases.

### 0.3 Proof of Theorem 3

The first result, 3(a) has been proved (DeLong *and others*, 1988). For the second result, with  $z = 1, 2$ , let  $Y_{zDj}, j = 1, \dots, n_D$  be marker  $z$  measurements in cases. Let  $F_{\theta_z}$  be the distribution function for marker  $z$  in controls and  $\theta = (\theta_1, \theta_2)$ . Suppose  $\sqrt{n_D}(\hat{\theta} - \theta)$  can be represented as  $\frac{1}{\sqrt{n_D}} \sum_{i=1}^{n_D} \psi_i + o_p(1)$ , where  $\psi_i, i = 1, \dots, n_D$  are independent identically distributed variables with  $E(\psi_i) = 0$  and  $\text{var}(\psi_i) = \Sigma(\theta)$ , we have

$$\begin{aligned}
& \sqrt{n_D}(\hat{\Delta} - \Delta) \\
&= \sqrt{n_D} \left[ \frac{1}{n_D} \sum_{j=1}^{n_D} \{F_{\hat{\theta}}(Y_{1Dj}) - F_{\hat{\theta}}(Y_{2Dj})\} - E_{Y_{1D}} F_{\theta_1}(Y) + E_{Y_{2D}} F_{\theta_2}(Y) \right] \times 100 \\
&= \sqrt{n_D} \left\{ E_{Y_{1D}} F_{\hat{\theta}_1}(Y) - E_{Y_{1D}} F_{\theta_1}(Y) - E_{Y_{2D}} F_{\hat{\theta}_2}(Y) + E_{Y_{2D}} F_{\theta_2}(Y) \right\} \times 100 \\
&+ \sqrt{n_D} \left[ \frac{1}{n_D} \sum_{j=1}^{n_D} \{F_{\theta_1}(Y_{1Dj}) - F_{\theta_2}(Y_{2Dj})\} - E_{Y_{1D}} F_{\theta_1}(Y) + E_{Y_{2D}} F_{\theta_2}(Y) \right] \times 100 + o_p(1) \\
&= \frac{1}{\sqrt{n_D}} \sum_{i=1}^{n_D} \frac{\partial \Delta_i}{\partial \theta} \psi_i + \frac{1}{\sqrt{\lambda n_D}} \sum_{j=1}^{n_D} \{Q_{1j} - Q_{2j} - E(Q_1 - Q_2)\} + o_p(1)
\end{aligned}$$

and result 3(b) follows.

### 0.4 Proof of Theorem 4

With  $z = 1, 2$ , let  $Y_{zD_i}^k, i = 1, \dots, n_{Dk}$  and  $Y_{zD_j}^k, j = 1, \dots, n_{Dk}$  be the measurements of marker  $z$  in controls and cases respectively in the  $k^{th}$  covariate category. Let  $F_z^k$  and  $G_z^k$  be the corresponding distribution functions. When  $F^k$  is estimated empirically,

$$\begin{aligned}
\sqrt{n_D}(\hat{\Delta} - \Delta) &= \sqrt{n_D} \left[ \frac{1}{n_D} \sum_{k=1}^K \sum_{i=1}^{n_{Dk}} \sum_{j=1}^{n_{Dk}} \frac{1}{n_{Dk}} \left\{ I(Y_{1D_i}^k \leq Y_{1D_j}^k) - I(Y_{2D_i}^k \leq Y_{2D_j}^k) \right\} \right. \\
&\quad \left. - \sum_{k=1}^K p_{Dk} P(Y_D^k \leq Y_{1D}^k) + \sum_{k=1}^K p_{Dk} P(Y_{2D}^k \leq Y_{2D}^k) \right] \times 100
\end{aligned}$$

$$\begin{aligned}
&= \sqrt{n_{\bar{D}}} \left( \sum_{k=1}^K \sum_{i=1}^{n_{\bar{D}k}} \frac{p_{Dk}}{n_{Dk}} [1 - G_1^k(Y_{1\bar{D}i}^k) - \{1 - G_2^k(Y_{2\bar{D}i}^k)\}] - \sum_{k=1}^K p_{Dk} P(Y_{1\bar{D}}^k \leq Y_{1D}^k) \right. \\
&\quad \left. + \sum_{k=1}^K p_{Dk} P(Y_{2\bar{D}}^k \leq Y_{2D}^k) \right) \times 100 + \sqrt{n_{\bar{D}}} \left[ \frac{1}{n_D} \sum_{k=1}^K \sum_{j=1}^{n_{Dk}} \{F_1^k(Y_{1Dj}^k) - F_2^k(Y_{2Dj}^k)\} \right. \\
&\quad \left. - \sum_{k=1}^K p_{Dk} \{P(Y_{1\bar{D}}^k \leq Y_{1D}^k) - P(Y_{2\bar{D}}^k \leq Y_{2D}^k)\} \right] \times 100 + o_p(1) \\
&= \sum_{k=1}^K -\frac{1}{\sqrt{n_{\bar{D}k}}} \sum_{i=1}^{n_{\bar{D}k}} \frac{p_{Dk}}{\sqrt{p_{\bar{D}k}}} \{R_1^k(Y_{1\bar{D}i}^k) - R_2^k(Y_{2\bar{D}i}^k) - E(R_1^k - R_2^k)\} \\
&\quad + \frac{1}{\sqrt{\lambda n_D}} \sum_{j=1}^{n_D} \{Q_{1Xj} - Q_{2Xj} - E(Q_{1X} - Q_{2X})\} + o_p(1).
\end{aligned}$$

This proves Theorem 4(a). Now consider when  $F$  is estimated parametrically. Let  $X_{\bar{D}i}, i = 1, \dots, n_{\bar{D}}$  and  $X_{Dj}, j = 1, \dots, n_D$  be the covariate measurements for controls and cases. Let  $F_{\theta_z}(\cdot|X)$  be the distribution function of marker  $z$  in controls conditional on  $X$ . Suppose  $\sqrt{n}(\hat{\theta} - \theta)$  can be represented as  $\frac{1}{\sqrt{n_{\bar{D}}}} \sum_{i=1}^{n_{\bar{D}}} \psi_i + o_p(1)$ , where  $\psi_i, i = 1, \dots, n_{\bar{D}}$  are independent identically distributed variables with  $E(\psi_i) = 0$  and  $\text{var}(\psi_i) = \Sigma(\theta)$ , we have

$$\begin{aligned}
&\sqrt{n_{\bar{D}}} (\hat{\Delta} - \Delta) \\
&= \sqrt{n_{\bar{D}}} \left[ \frac{1}{n_D} \sum_{j=1}^{n_D} \{F_{\hat{\theta}_1}(Y_{1Dj}|X_{Dj}) - F_{\hat{\theta}_2}(Y_{2Dj}|X_{Dj})\} - E_{Y_{1D}, X_D} F_{\theta_1}(Y|X) + E_{Y_{2D}, X_D} F_{\theta_2}(Y|X) \right] \times 100 \\
&= \sqrt{n_{\bar{D}}} \left\{ E_{Y_{1D}, X_D} F_{\hat{\theta}_1}(Y_{1D}|X_D) - E_{Y_{1D}, X_D} F_{1, \theta}(Y_{1D}|X_D) - E_{Y_{2D}, X_D} F_{\hat{\theta}_2}(Y_{2D}|X_D) \right. \\
&\quad \left. + E_{Y_{2D}, X_D} F_{\theta_2}(Y_{2D}|X_D) \right\} \times 100 + \sqrt{n_{\bar{D}}} \left[ \frac{1}{n_D} \sum_{j=1}^{n_D} \{F_{\theta_1}(Y_{1Dj}|X_{Dj}) - F_{\theta_2}(Y_{2Dj}|X_{Dj})\} \right. \\
&\quad \left. - E_{Y_{1D}, X_D} F_{\theta_1}(Y|X) + E_{Y_{2D}, X_D} F_{\theta_2}(Y|X) \right] \times 100 + o_p(1) \\
&= \frac{1}{\sqrt{n_{\bar{D}}}} \sum_{i=1}^{n_{\bar{D}}} \frac{\partial \Delta_i}{\partial \theta} \psi_i + \frac{1}{\sqrt{\lambda n_D}} \sum_{j=1}^{n_D} \{Q_{1Xj} - Q_{2Xj} - E(Q_{1X} - Q_{2X})\} + o_p(1).
\end{aligned}$$

#### 0.4.1 Theorem 4 (Extension Version)

Consider the setting when the two markers are adjusted for different covariates  $X_1$  and  $X_2$ . For  $z = 1, 2$ , let  $Q_{zX_z}(\hat{Q}_{zX_z})$  be the (estimated) covariate-specific percentile value for the  $z^{\text{th}}$  marker, let  $\Delta = E(Q_{1X_1}) - E(Q_{2X_2})$  and  $\hat{\Delta} = \hat{Q}_{1X_1} - \hat{Q}_{2X_2}$ . When  $X_1$  and  $X_2$  are discrete with strata  $k_1 = 1, \dots, K_1$  and  $k_2 = 1, \dots, K_2$  respectively, let  $n_{Dk_1k_2}$  and  $n_{\bar{D}k_1k_2}$  be the number of controls and cases in the intersection of the  $k_1^{\text{th}}$  stratum for covariate 1 and the  $k_2^{\text{th}}$  stratum for covariate 2. Suppose as  $n_{\bar{D}} \rightarrow \infty, n_D/n_{\bar{D}} \rightarrow \lambda \in (0, 1), n_{\bar{D}k_1k_2}/n_{\bar{D}} \rightarrow p_{\bar{D}k_1k_2} \in (0, 1), n_{Dk_1k_2}/n_D \rightarrow p_{Dk_1k_2} \in (0, 1)$ , then

$\sqrt{n_{\bar{D}}}(\hat{\Delta} - \Delta)$  converges to a mean 0 normal random variable with variance  $\sigma^2$ , where

$$(a) \quad \sigma^2 = \sum_{k_1} \sum_{k_2} \frac{\text{var} \left\{ R_1^{k_1}(Y_{1\bar{D}}^{k_1}) - R_2^{k_2}(Y_{2\bar{D}}^{k_2}) \right\}}{p_{\bar{D}k_1k_2}/p_{\bar{D}k_1k_2}^2} + \frac{\text{var}(Q_{1X_1} - Q_{2X_2})}{\lambda},$$

if for each marker, the covariate-specific reference distribution  $F(Y|X)$  is estimated empirically within each covariate category specific for the marker, where  $R_z^{k_z}(Y_{z\bar{D}}^{k_z}) = P(Y_{z\bar{D}}^{k_z} < Y_{z\bar{D}}^{k_z})$  is the percentile value for a control using his covariate-specific case distribution as the reference for the  $z^{\text{th}}$  marker in the  $k_z^{\text{th}}$  covariate category, and

$$(b) \quad \sigma^2 = \left( \frac{\partial \Delta}{\partial \theta} \right)^T \Sigma(\theta) \left( \frac{\partial \Delta}{\partial \theta} \right) + \frac{\text{var}(Q_{1X_1} - Q_{2X_2})}{\lambda}, \quad (0.1)$$

if  $F(Y|X)$  is modeled parametrically for marker  $z$  with parameter estimate  $\theta_z$ ,  $\theta = (\theta_1, \theta_2)$  and  $\Sigma(\theta)$  is the asymptotic variance of  $\sqrt{n_{\bar{D}}}(\hat{\theta} - \theta)$ . We assume that  $\Delta$  is differentiable with respect to  $\theta$  and that  $\mathcal{F} = \{F_\theta(y|x) : \theta \in \Theta\}$  is a Donsker class.

Proof follows similar steps to above.

### 0.5 Proof of Proposition 1

Let  $Y_{\bar{D}i}$ ,  $i = 1, \dots, n_{\bar{D}}$ , and  $Y_z$  be the marker measurement in controls and the  $z^{\text{th}}$  type of cases,  $z = 1, 2$ . Let  $Q_z(\hat{Q}_z)$  be the corresponding (estimated) case percentile value. We have

$$\begin{aligned} Q_1 \stackrel{d}{=} Q_2 &\Rightarrow P\{F(Y_1) \leq t\} = P\{F(Y_2) \leq t\} \quad \forall t \in (0, 1) \\ &\Rightarrow P\{Y_1 \leq F^{-1}(t)\} = P\{Y_2 \leq F^{-1}(t)\} \\ &\Rightarrow Y_1 \stackrel{d}{=} Y_2 \\ &\Rightarrow \hat{Q}_1|Y_{\bar{D}i} \stackrel{d}{=} \hat{Q}_2|Y_{\bar{D}i}, i = 1, \dots, n_{\bar{D}} \quad \text{and} \quad \hat{Q}_1 \stackrel{d}{=} \hat{Q}_2. \end{aligned}$$

### 0.6 Proof of Proposition 2

Let  $Y_{\bar{D}}$  and  $Y_D$  be the marker measurements for controls and cases respectively,  $Y_{\bar{D}} \sim F$ ,  $Y_D \sim G$ . Let  $\hat{Q}$  be the estimated case percentile value. When  $\hat{F}$  is the empirical CDF, for  $m \in (0, 1, \dots, n_{\bar{D}})$ ,

$$\begin{aligned} P(n_{\bar{D}}\hat{Q}/100 = m) &= P\left\{n_{\bar{D}}\hat{F}(Y_D) = m\right\} = P\left\{\sum_{i=1}^{n_{\bar{D}}} I(Y_{\bar{D}} \leq Y_D) = m\right\} \\ &= \frac{n_{\bar{D}}!}{(n_{\bar{D}} - m)!m!} \int_{-\infty}^{\infty} F^m(y) \{1 - F(y)\}^{n_{\bar{D}} - m} dG(y) \\ &= \frac{n_{\bar{D}}!}{(n_{\bar{D}} - m)!m!} \int_0^1 w^m (1 - w)^{n_{\bar{D}} - m} dG\{F^{-1}(w)\}, \end{aligned}$$

which is a function of ROC. Therefore, when  $\text{ROC}_1(t) = \text{ROC}_2(t) \forall t$ , i.e. when  $Q_1 \stackrel{d}{=} Q_2$ ,  $\hat{Q}_1$  and  $\hat{Q}_2$  have the same distribution.

## 0.7 Proof of Proposition 3

With  $z = 1, 2$ , let  $\hat{Q}_{zi}, i = 1, \dots, n_D$  be the estimated case percentile values for the  $z^{\text{th}}$  marker. Because of the exchangeability among cases, the correlation between  $(\hat{Q}_{1i}, \hat{Q}_{2i})$  and  $(\hat{Q}_{1j}, \hat{Q}_{2j})$  are the same for each pair of  $i, j$ . In another words,  $(\hat{Q}_{11}, \hat{Q}_{21}), \dots, (\hat{Q}_{1n_D}, \hat{Q}_{2n_D})$  are correlated samples from a bivariate distribution with exchangeable correlation. Under  $H_0 : Q_1 \stackrel{d}{=} Q_2$ ,  $\hat{Q}_1$  and  $\hat{Q}_2$  have the same distribution when  $F$  is estimated empirically (Proposition 2). Consequently  $U_j = \hat{Q}_{1j} - \hat{Q}_{2j}, j = 1, \dots, n_D$  are correlated (exchangeably) samples from a symmetric univariate distribution  $H$ . Due to the complete exchangeability among those  $U_j$ 's, if  $H$  is continuous,  $Z_j = P(U_j > 0)$  is Bernoulli random variable with  $P(Z_j = 1) = 0.5$ , thus  $E(S) = \sum_{j=1}^{n_D} Z_j$  is equal to  $1/2$ . In addition,  $Z_j$  is independent of  $r(|U_j|)$ , where  $r(|U_j|)$  is the rank of  $|U_j|$  in the sample, thus,  $E(T) = \sum_{j=1}^{n_D} E\{Z_j r(|U_j|)\} = \sum_{j=1}^{n_D} \frac{1}{2} E\{r(|U_j|)\} = (n_D + 1)/4$ .

## 0.8 Proof of Proposition 4

It is equivalent to prove that the Mann-Whitney test statistic (MW) of the two groups of estimated percentile values has mean 0.5 when  $F$  is estimated empirically. We have

$$\begin{aligned} MW &= \frac{1}{n_D} \sum_{i=1}^{n_D} \sum_{j=1}^{n_D} \left\{ I(\hat{Q}_{1i} < \hat{Q}_{2j}) + \frac{I(\hat{Q}_{1i} = \hat{Q}_{2j})}{2} \right\} \\ E(MW) &= \frac{1}{n_D} \sum_{i=1}^{n_D} \sum_{j=1}^{n_D} \left\{ P(\hat{Q}_{1i} < \hat{Q}_{2j}) + \frac{1}{2} P(\hat{Q}_{1i} = \hat{Q}_{2j}) \right\} \\ &= P(\hat{Q}_{1i} < \hat{Q}_{2j}) + \frac{P(\hat{Q}_{1i} = \hat{Q}_{2j})}{2}. \end{aligned}$$

As proved in Proposition 2, when  $\text{ROC}_1(t) = \text{ROC}_2(t) \forall t$ , i.e. when  $Q_1 \stackrel{d}{=} Q_2$ ,  $\hat{Q}_1$  and  $\hat{Q}_2$  have the same distribution when  $\hat{F}$  is estimated empirically. Therefore,  $E(MW) = 0.5$ .

## REFERENCES

- DELONG, E. R., DELONG, D. M. AND CLARKE-PEARSON, D. L. (1988). Comparing the areas under two or more correlated receiver operating characteristic curves: a nonparametric approach. *Biometrics*, **44**(3), 837-845.
- VAN DER VAART, A. W. (1998) *Asymptotic Statistics*. Cambridge University Press.
- VAN DER VAART, A. W. AND WELLNER, J. A. (1996) *Weak Convergence and Empirical processes*. Springer-Verlag, New York.