

Harvard University
Harvard University Biostatistics Working Paper Series

Year 2006

Paper 36

Regression Analysis for the Partial Area Under
the ROC Curve

Tianxi Cai*

Lori E. Dodd†

*Harvard University, tcai@hsph.harvard.edu

†National Institutes of Health, National Cancer Institute, ld135k@nih.gov

This working paper is hosted by The Berkeley Electronic Press (bepress) and may not be commercially reproduced without the permission of the copyright holder.

<http://biostats.bepress.com/harvardbiostat/paper36>

Copyright ©2006 by the authors.

Regression Analysis for the Partial Area Under the ROC Curve

TIANXI CAI

Department of Biostatistics, Harvard University, Boston, MA 02115

LORI E. DODD

Biometric Research Branch, National Cancer Institute, Rockville, MD 20892

SUMMARY

Performance evaluation of any classification method is fundamental to its acceptance in practice. Evaluation should consider the dependence of a classifier's accuracy on relevant covariates in addition to its overall accuracy. When developing a classifier with a continuous output that allocates units into one of two groups, receiver operating characteristic (ROC) curve analysis is appropriate. The partial area under the ROC curve (pAUC) is a summary measure of the ROC curve used to make statistical inference when only a region of the ROC space is of interest. We propose a new pAUC regression method to evaluate covariate effects on the diagnostic accuracy. We provide asymptotic distribution theory and procedures for making statistical inference that allows for correlated observations. Graphical methods and goodness-of-fit statistics for model checking are also developed. Simulation studies demonstrate that the large-sample theory provides reasonable inference in small samples and the new estimator is considerably more efficient than the estimator proposed by Dodd and Pepe (2003a). Application to an analysis of prostate-specific antigen (PSA), a biomarker for early detection of prostate cancer, demonstrates the utility of the method in practice.

Key words : Diagnostic Accuracy, Generalized Linear Model, Model Checking.

1. Introduction

Binary classification is a relevant undertaking in a wide variety of statistical fields. Algorithms such as support vector machines and neural networks have been applied, for example, to detect automobile insurance claim fraud (Viaene et al., 2002) or to predict peptide binding (Brusic et al., 1998). In the medical field, a multitude of medical tests, such as biomarkers and imaging modalities have been developed to screen and diagnose disease, as well to predict outcome and monitor response to therapy. Rigorous evaluation of any classification method is a prerequisite to its wide-spread use. A method must be shown to be accurate and factors influencing the accuracy of a method must be adequately understood.

Accuracy may be summarized by the percent of correct classifications. However, a more refined analysis of accuracy considers the false positive error and the false negative error separately, as each has a unique associated cost. For a continuous outcome variable, Y , let $Y \geq c$ denote a positive classification. Throughout, the two states are referred to as “diseased” and “disease-free”, however more general terminology could be used. Additionally, the term “test” refers generally to the continuous output of a classifier, such as a biomarker or a neural network result. The true positive rate (TPR), is defined as $S_D(c) \equiv P(Y \geq c \mid \text{diseased})$, while the false positive rate (FPR), is defined as $S_{\bar{D}}(c) \equiv P(Y \geq c \mid \text{disease-free})$. The receiver operating characteristic (ROC) curve plots $\{(S_{\bar{D}}(c), S_D(c)), c \in (-\infty, \infty)\}$, or, equivalently, $\{(u, \text{ROC}(u)), u \in (0, 1)\}$. The curve describes the inherent capacity of the test in discriminating the two states, without linking the test to any specific positivity criterion.

A single summary index is useful as a descriptive of overall test performance and for hypothesis testing. The most common summary index of the ROC curve is the area under the

curve (AUC) (Bamber, 1975; DeLong, DeLong and Clarke-Pearson, 1988). The AUC can be interpreted as the probability that a randomly selected case with disease will be regarded with greater suspicion than a randomly selected disease-free case. Often, interest does not lie in the entire range of FPRs, and consequently, only part of the area under the ROC curve is relevant. For example, very low false positive rates such as $FPR \leq 0.05$ have been advocated in settings such as cancer screening (Baker and Pinsky, 2001) and hence analysis should be restricted to the portion of the curve corresponding to $FPR \leq 0.05$. Alternatively, a restricted region of TPRs may be of interest (Jiang, et al., 1996). Noting that a definition with respect to TPRs is straightforward, we consider the partial AUC (pAUC) for a range of FPRs, without loss of generality, say $FPR \in (0, u]$, for some $u \leq 1$. The pAUC is given as $pAUC(u) = \int_0^u ROC(u)du$ (McClish, 1989; Thompson and Zucchini, 1989), which has a value of u when a test is perfect and of $u^2/2$ when a test is uninformative. Another reason to analyze the pAUC rather than the entire AUC is that a summary of the entire ROC curve fails to consider the plot as a composite of different segments with different diagnostic implications (Dwyer, 1996). This is particularly important if prominent differences between ROC plots in specific regions are muted or reversed when the total area is considered.

Methods for estimating and comparing pAUCs are available (McClish, 1989; Wieand et al., 1989; Zhang et al, 2002; Pepe, 2003; Dodd and Pepe, 2003a). Generalizations of these methods to regression modeling assists with further characterization of a classifier. As an example, consider PSA, a biomarker for prostate cancer. Since a biomarker that detects cancer prior to the onset of clinical symptoms is of clinical interest, a model of PSA accuracy with a covariate representing the time prior to clinical diagnosis is of interest. This will provide

information about by how much PSA advances diagnosis. In addition, if there is a relationship between PSA accuracy and age, a model that includes age as a covariate might identify ages for targeting PSA screening programs.

Two approaches to the pAUC regression analysis have been proposed (Thompson and Zucchini, 1989; Dodd and Pepe, 2003a). The method proposed by Thompson and Zucchini (1989) does not accommodate continuous covariates and is not applicable to many types of data. Dodd and Pepe (2003a) present a more flexible pAUC regression method, however, they do not provide theoretical justification for their estimator and rely on bootstrap to estimate the variance. Furthermore, their models require making unnecessary assumptions. Specifically, they model the pAUC comparing test results of diseased subjects, Y_D , with continuous covariate \mathbf{Z}_1 to test results of disease-free, $Y_{\bar{D}}$, with continuous covariate \mathbf{Z}_0 as:

$$\text{pAUC}_{\mathbf{Z}_1, \mathbf{Z}_0}(u) = \eta \{ \beta_0 + \beta_1^T \mathbf{Z}_1 + \beta_2^T (\mathbf{Z}_1 - \mathbf{Z}_0) \} \quad (1)$$

for a given link function $\eta : (-\infty, \infty) \rightarrow [0, u]$. This formulation requires modeling the effect of $\mathbf{Z}_1 - \mathbf{Z}_0$, the difference between the covariate levels in the two populations in addition to the quantity of interest $\beta_1^T \mathbf{Z}_1$. However, when assessing the test accuracy adjusting for covariates, the interest only lies in comparing the distribution of Y_D and $Y_{\bar{D}}$ among subjects with matched common covariates. Thus β_2 is not of scientific interest and (1) imposes unnecessary modelling.

In this article, we propose to model the covariate specific pAUC assuming (1) only when $\mathbf{Z}_0 = \mathbf{Z}_1$. Our estimation approach is based on the concept of *placement values* (Hanley and Hajian-Tilaki, 1997; Pepe and Cai, 2004), defined as particular standardizations of the raw measurements relative to the reference populations. In section 2, we introduce placement

values and illustrate how they can be used to estimate pAUC when there is no covariate. In section 3, we propose a marginal regression model for the pAUC and derive inference procedures for the regression parameters allowing for clustered data. Simulation studies in section 4 suggest that the new approach performs well. Furthermore, the new estimator, while always being more robust, is considerably more efficient than the Dodd and Pepe (2003a) estimator. To examine whether the specified regression model is appropriate for the data, in section 5, we present both graphical procedures as well as goodness of fit testing statistics for model checking. Section 6 gives results from the application of the proposed method to a PSA dataset. Some discussion is provided in section 7.

2. Placement Values and pAUC Estimation

As in Pepe and Cai (2004), we choose the disease-free population as the reference population and define the placement value for Y_D as $U_D \equiv S_{\bar{D}}(Y_D)$. Then U_D quantifies the degree of separation between the two populations. Moreover,

$$P(U_D \leq u) = P\{S_{\bar{D}}(Y_D) \leq u\} = P\{Y_D \geq S_{\bar{D}}^{-1}(u)\} = \text{ROC}(u), \quad \text{and}$$

$$E(U_D) = \int_0^1 u \, d\text{ROC}(u) = 1 - \int_0^1 \text{ROC}(u) \, du = 1 - \text{AUC}.$$

DeLong et al. (1988) and Hanley and Hajian-Tilaki (1997) interpreted the nonparametric estimate of the AUC as one minus the sample mean of the empirically estimated placement values. Placement values have been used recently to make inference about ROC regression models (Pepe and Cai, 2004; Cai, 2004). Here, we propose to make inference about the pAUC based on *truncated* placement values.

We first illustrate our proposal by constructing a non-parametric estimator for the pAUC

in the absence of covariates. Suppose we have N_D data records for n_D diseased subjects $\{Y_{Dik}, k = 1, \dots, K_i, i = 1, \dots, n_D\}$ and $N_{\bar{D}}$ data records for $n_{\bar{D}}$ disease-free subjects $\{Y_{\bar{D}jl}, l = 1, \dots, K_j, j = n_D + 1, \dots, n_D + n_{\bar{D}}\}$, where $N_D = \sum_{i=1}^{n_D} K_i$ and $N_{\bar{D}} = \sum_{j=n_D+1}^{n_D+n_{\bar{D}}} K_j$. Each subject may have more than one data record in the analysis and these records could be correlated, but we assume that K_i and K_j are relatively small with respect to n_D and $n_{\bar{D}}$. We also assume that $S_{\bar{D}}(y) \equiv P(Y_{\bar{D}jl} \geq y) = P(Y_{\bar{D}j'l'} \geq y)$ and $S_D(y) \equiv P(Y_{Dik} \geq y) = P(Y_{Di'k'} \geq y)$. Then, the placement value for Y_{Dik} is $U_{Dik} \equiv S_{\bar{D}}(Y_{Dik})$. Without loss of generality let $u_a = 0$. Further, let $U_{Dik}^{(u)} \equiv \min(U_{Dik}, u)$ denote the truncated placement value and $\hat{U}_{Dik}^{(u)} \equiv \min(\hat{U}_{Dik}, u)$ be the empirical estimator of $U_{Dik}^{(u)}$, where $\hat{U}_{Dik} = \hat{S}_{\bar{D}}(Y_{Dik})$ and $\hat{S}_{\bar{D}}(y) = N_{\bar{D}}^{-1} \sum_{j=n_D+1}^{n_D+n_{\bar{D}}} \sum_{l=1}^{K_j} I(Y_{\bar{D}jl} \geq y)$.

Using integration by parts, we find that the marginal mean of the truncated placement values relates to the pAUC through

$$E(U_{Dik}^{(u)}) = \int_0^u \{1 - \text{ROC}(v)\} dv = u - \text{pAUC}(u).$$

This motivates us to estimate the pAUC(u) with

$$\widehat{\text{pAUC}}(u) = u - \frac{1}{N_D} \sum_{i=1}^{n_D} \sum_{k=1}^{K_i} \hat{U}_{Dik}^{(u)} = \frac{1}{N_D} \sum_{i=1}^{n_D} \sum_{k=1}^{K_i} \hat{V}_{Dik}^{(u)},$$

where $\hat{V}_{Dik}^{(u)} = u - \hat{U}_{Dik}^{(u)}$. When $K_i = K_j = 1$, this estimator is equivalent to the non-parametric estimate proposed by Dodd and Pepe (2003a). Since Dodd and Pepe (2003a) did not provide large sample theory for $\widehat{\text{pAUC}}(u)$, we show in appendix A the consistency of $\widehat{\text{pAUC}}(u)$ and that the distribution of $n_D^{\frac{1}{2}} \{\widehat{\text{pAUC}}(u) - \text{pAUC}(u)\}$ is approximately $N(0, \hat{\sigma}^2)$ accounting for within cluster correlation, where $\hat{\sigma}^2 = n_D^{-1} \sum_{i=1}^{n_D} \hat{P}_{Di}^2 + n_{\bar{D}}^{-1} \sum_{j=n_D+1}^{n_D+n_{\bar{D}}} \hat{P}_{\bar{D}j}^2$, $\hat{P}_{Di} = \frac{n_D}{N_D} \sum_k \hat{V}_{Dik}^{(u)} - \widehat{\text{pAUC}}(u)$, and $\hat{P}_{\bar{D}j} = \frac{(n_D/n_{\bar{D}})^{\frac{1}{2}}}{N_D} \sum_l \sum_{i,k} I(\hat{U}_{Dik} \leq u) \{\hat{U}_{Dik} - I(Y_{\bar{D}jl} \geq Y_{Dik})\}$.

3. Partial AUC Regression

Next, we use truncated placement values to develop an estimating equation for pAUC regression models. Let $\mathbf{X}_{ik} = (\mathbf{Z}_{Dik}, \mathbf{Z}_{ik})$ denote the covariates associated with Y_{Dik} and \mathbf{Z}_{jl} be the covariates associated with Y_{Djl} . Covariates denoted by \mathbf{Z} are relevant to both diseased and disease-free subjects. Examples might include the subject's age or the type of biomarker represented by Y (see Pepe 2003, chapter 6). Covariates denoted by \mathbf{Z}_D are specific to subjects with disease, but not applicable to disease-free subjects. Examples include severity of disease and timing of biomarker measurement prior to onset of clinical symptoms. In the presence of disease subject specific covariates, \mathbf{Z}_D , one would be interested in comparing the distribution of Y among those diseased subjects with covariates $\mathbf{X} = (\mathbf{Z}, \mathbf{Z}_D)$, to the distribution of Y among those disease-free subjects with covariates \mathbf{Z} .

We assume a marginal model for the covariate specific pAUC:

$$\int_0^u \text{ROC}_{\mathbf{X}_{ik}}(v) dv \equiv \text{pAUC}_{\mathbf{X}_{ik}}(u) = \eta(\boldsymbol{\beta}_0^\top \vec{\mathbf{X}}_{ik}), \quad (2)$$

where $\text{ROC}_{\mathbf{X}}(v) = P\{Y_D \geq S_{\bar{D}, \mathbf{Z}}^{-1}(v) \mid \mathbf{X} = (\mathbf{Z}_D, \mathbf{Z})\}$, $S_{\bar{D}, \mathbf{Z}}(y) = P(Y_{\bar{D}jl} \geq y \mid \mathbf{Z}_{jl} = \mathbf{Z})$, and $\vec{\mathbf{X}}_{ik} = (1, \mathbf{X}_{ik})$. To estimate $\boldsymbol{\beta}_0$, we define the placement value for the test result Y_{Dik} with covariate \mathbf{X}_{ik} as $U_{Dik} \equiv S_{\bar{D}, \mathbf{Z}_{ik}}(Y_{Dik})$. It is straightforward to show that $E\{\min(U_{Dik}, u) \mid \mathbf{X}_{ik}\} = \int_0^u \{1 - \text{ROC}_{\mathbf{X}_{ik}}(v)\} dv = u - \text{pAUC}_{\mathbf{X}_{ik}}(u)$. This motivates us to estimate $\boldsymbol{\beta}_0$ by solving

$$\frac{1}{N_D} \sum_{i=1}^{n_D} \sum_{k=1}^{K_i} w(\vec{\mathbf{X}}_{ik}) \vec{\mathbf{X}}_{ik} \left\{ \widehat{V}_{Dik}^{(u)} - \eta(\boldsymbol{\beta}^\top \vec{\mathbf{X}}_{ik}) \right\} = 0, \quad (3)$$

where $w(\cdot)$ is a given positive weight function, $\widehat{V}_{Dik}^{(u)} = u - \min(\widehat{U}_{Dik}, u)$, $\widehat{U}_{Dik} = \widehat{S}_{\bar{D}, \mathbf{Z}_{ik}}(Y_{Dik})$ and $\widehat{S}_{\bar{D}, \mathbf{Z}}(y)$ is a consistent estimate of $S_{\bar{D}, \mathbf{Z}}(y)$. If the covariates \mathbf{Z} are discrete, $S_{\bar{D}, \mathbf{Z}}(y)$ can be estimated non-parametrically within covariate specific subsets. When continuous covariates

are included, we recommend semi-parametric regression models for $S_{\bar{D}, \mathbf{Z}}(y)$. For example, one could assume a flexible semi-parametric location-scale model (Pepe, 1998; Heagerty and Pepe, 1999). Other types of semi-parametric models such as linear transformation models (Han, 1997; Cai, Wei and Wilcox, 2000) could also be considered. We do not assume any specific model for $S_{\bar{D}, \mathbf{Z}}(y)$, but require that the resulting estimator of $S_{\bar{D}, \mathbf{Z}}(y)$ is $n_{\bar{D}}^{\frac{1}{2}}$ -consistent. We note that the Dodd and Pepe estimator also requires semi-parametric regression models for the conditional quantile of $Y_{\bar{D}}$ (Dodd and Pepe, 2002, page 620).

Let $\hat{\beta}$ denote the solution to (3). We show in appendix B that $\hat{\beta}$ is unique and consistent. To obtain interval estimates of specific components of β_0 , we also show in appendix B that accounting for the correlation within each subject, $n_{\bar{D}}^{\frac{1}{2}}(\hat{\beta} - \beta_0)$ is asymptotically equivalent to a sum of independent terms indexed by subjects:

$$n_{\bar{D}}^{\frac{1}{2}}(\hat{\beta} - \beta_0) \approx \mathbb{A}^{-1} \left\{ n_{\bar{D}}^{-\frac{1}{2}} \sum_{i=1}^{n_{\bar{D}}} \mathfrak{B}_{\bar{D}i} + n_{\bar{D}}^{-\frac{1}{2}} \sum_{j=n_{\bar{D}}+1}^{n_{\bar{D}}+n_{\bar{D}}} \mathfrak{B}_{\bar{D}j} \right\},$$

where $\mathbb{A} = E\{\dot{\eta}(\beta_0^T \vec{\mathbf{X}}_{ik}) \vec{\mathbf{X}}_{ik}^{\otimes 2}\}$, $\dot{\eta}(x) = d\eta(x)/dx$, $\mathfrak{B}_{\bar{D}i} = K_{\bar{D}}^{-1} \sum_{k=1}^{K_i} w(\mathbf{X}_{ik}) \vec{\mathbf{X}}_{ik} \{V_{\bar{D}ik}^{(u)} - \eta(\beta_0^T \vec{\mathbf{X}}_{ik})\}$, $\mathfrak{B}_{\bar{D}j}$ is the limit of $\frac{p_{\bar{D}j}^{\frac{1}{2}}}{N_{\bar{D}}} \sum_{i=1}^{n_{\bar{D}}} \sum_{k=1}^{K_i} w(\mathbf{X}_{ik}) \vec{\mathbf{X}}_{ik} \int_0^u I_{\bar{D}j}(v; \mathbf{Z}_{ik}) d\text{ROC}_{\mathbf{X}_{ik}}(v)$ and $I_{\bar{D}j}(v, \mathbf{Z})$ is defined in appendix B. It follows from the multivariate central limit theorem that the distribution of $n_{\bar{D}}^{\frac{1}{2}}(\hat{\beta} - \beta_0)$ can be approximated by $N(0, \Sigma)$. Σ can be consistently estimated by

$$\hat{\mathbb{A}}^{-1} \left\{ n_{\bar{D}}^{-1} \sum_{i=1}^{n_{\bar{D}}} \hat{\mathfrak{B}}_{\bar{D}i} \hat{\mathfrak{B}}_{\bar{D}i}^T + n_{\bar{D}}^{-1} \sum_{j=n_{\bar{D}}+1}^{n_{\bar{D}}+n_{\bar{D}}} \hat{\mathfrak{B}}_{\bar{D}j} \hat{\mathfrak{B}}_{\bar{D}j}^T \right\} \hat{\mathbb{A}}^{-1},$$

where $\hat{\mathbb{A}} = \frac{1}{N_{\bar{D}}} \sum_{i=1}^{n_{\bar{D}}} \sum_{k=1}^{K_i} \dot{\eta}(\hat{\beta}^T \vec{\mathbf{X}}_{ik}) \vec{\mathbf{X}}_{ik} \vec{\mathbf{X}}_{ik}^T$, $\hat{\mathfrak{B}}_{\bar{D}i} = K_{\bar{D}}^{-1} \sum_{k=1}^{K_i} w(\mathbf{X}_{ik}) \vec{\mathbf{X}}_{ik} \{V_{\bar{D}ik}^{(u)} - \eta(\hat{\beta}^T \vec{\mathbf{X}}_{ik})\}$, $\hat{\mathfrak{B}}_{\bar{D}j} = \frac{p_{\bar{D}j}^{\frac{1}{2}}}{N_{\bar{D}}} \sum_{i=1}^{n_{\bar{D}}} \sum_{k=1}^{K_i} w(\mathbf{X}_{ik}) \vec{\mathbf{X}}_{ik} I(\hat{U}_{\bar{D}ik} \leq u) \hat{I}_{\bar{D}j}(\hat{U}_{\bar{D}ik}; \mathbf{Z}_{ik})$, and $\hat{I}_{\bar{D}j}(v, \mathbf{Z})$ is obtained by replacing all the theoretical quantities in $I_{\bar{D}j}(v, \mathbf{Z})$ by their empirical counterparts.

4. Model Checking Procedures

The proposed inference procedures require the specification of the link function $\eta(\cdot)$. Here, we present a graphical method as well as statistical tests to assess whether model (2), with a given link function $\eta(\cdot)$, is appropriate for the data. Noting that $\text{pAUC}_{\mathbf{X}_{ik}}(u)$ is the conditional mean of $V_{\text{Dik}}^{(u)}$, we define the residuals for fitting model (2) as $\hat{e}_{ik} = \hat{V}_{\text{Dik}}^{(u)} - \eta(\hat{\boldsymbol{\beta}}^\top \vec{\mathbf{X}}_{ik})$. To examine the appropriateness of model (2), we first check the functional form for each component of the covariate \mathbf{X} . For $q = 1, \dots, p$, we consider the following moving sum of the \hat{e}_{ik} 's over the $X_{ik}^{(q)}$:

$$\bar{W}_q(x; b) = \frac{1}{N_D} \sum_{i=1}^{n_D} \sum_{k=1}^{K_i} I(x - b < X_{ik}^{(q)} \leq x) \hat{e}_{ik} \quad (4)$$

for a pre-specified positive block size b , where $X_{ik}^{(q)}$ is the q th element of \mathbf{X}_{ik} . Moving sum of residuals was proposed by Lin, Wei & Ying (2002) to test the goodness of fit for generalized linear models. When $b = \infty$, (4) corresponds to the partial residual process considered by Su and Wei (1991).

Under H_0 that model (2) holds, $\bar{W}_q(x; b)$ is expected to fluctuate around 0. To obtain the large sample distribution of $\bar{W}_q(x; b)$, let $\mathcal{L} = (\mathcal{L}_1, \dots, \mathcal{L}_{n_D+n_D})$ be a random sample from the standard normal distribution which is independent of the data. Define

$$\begin{aligned} n_D^{\frac{1}{2}} \widehat{W}_q(x; b) &= n_D^{\frac{1}{2}} \sum_{i=1}^{n_D} \widehat{W}_{\text{Dqi}}(x) \mathcal{L}_i + n_D^{\frac{1}{2}} \sum_{j=n_D+1}^{n_D+n_D} \widehat{W}_{\text{Dqj}}(x) \mathcal{L}_j, \\ \widehat{W}_{\text{Dqi}}(x; b) &= K_D^{-1} \sum_{k=1}^{K_i} I(x - b < X_{ik}^{(q)} \leq x) \left\{ \hat{V}_{\text{Dik}}^{(u)} - \eta(\hat{\boldsymbol{\beta}}^\top \vec{\mathbf{X}}_{ik}) \right\} + \hat{\mathbf{R}}_q(x; b)^\top \hat{\mathbf{A}}(\hat{\boldsymbol{\beta}})^{-1} \hat{\boldsymbol{\beta}}_{\text{Di}}, \\ \widehat{W}_{\text{Dqj}}(x; b) &= \frac{p_{10}}{N_D} \sum_{i=1}^{n_D} \sum_{k=1}^{K_i} I(x - b < X_{ik}^{(q)} \leq x) I(\hat{U}_{\text{Dik}} \leq u) \hat{I}_{\text{Dj}}(\hat{U}_{\text{Dik}}, \mathbf{Z}_{ik}) + \hat{\mathbf{R}}_q(x; b)^\top \hat{\mathbf{A}}^{-1} \hat{\boldsymbol{\beta}}_{\text{Dj}}, \end{aligned}$$

and $\hat{\mathbf{R}}_q(x; b) = \frac{1}{N_D} \sum_{i,k} I(x - b < X_{ik}^{(q)} \leq x) \eta(\hat{\boldsymbol{\beta}}_0^\top \vec{\mathbf{X}}_{ik}) \vec{\mathbf{X}}_{ik}$. In appendix C, we show that under H_0 , the conditional distribution of $n_D^{\frac{1}{2}} \widehat{W}_q(x; b)$ given the data is the same in the limit as the

unconditional distribution of $n_D^{\frac{1}{2}} \bar{W}_q(x; b)$. To approximate the null distribution of $W_q(x; b)$, we simulate a number of realizations from $n_D^{\frac{1}{2}} \widehat{W}_q(x; b)$ by repeatedly generating the normal samples of \mathcal{L} while fixing the data at their observed values. To assess how unusual the observed process $\bar{W}_q(x; b)$ is under H_0 , one may plot $\bar{W}_q(x; b)$ along with a few realizations from $\widehat{W}_q(x; b)$ and supplement the graphical display with an estimated p-value from a supremum-type test statistic $S_q = \sup_x |\bar{W}_q(x; b)|$. An unusually large observed value s_q would suggest improper specification of the functional form of X_q . In practice, the p-value, $P(S_q \geq s_q)$, can be approximated by $P(\widehat{S}_q \geq s_q)$, where $\widehat{S}_q = \sup_x |\widehat{W}_q(x; b)|$. We estimate $P(\widehat{S}_q \geq s_q)$ by generating a large number \mathcal{J} , say $\mathcal{J} = 5000$, of realizations from $\widehat{W}_q(\cdot; b)$.

To assess the linearity of the model given in (2) and more generally the link function $\eta(\cdot)$, we consider the moving sum of residuals over the fitted values:

$$\bar{W}_\eta(x; b) = n_D^{\frac{1}{2}} \sum_{i=1}^{n_D} \sum_{k=1}^{K_i} I(x - b < \widehat{\beta}^\top \vec{\mathbf{X}}_{ik} < x) \widehat{e}_{ik}.$$

The null distribution of $\bar{W}_\eta(x; b)$ can be approximated by the conditional distribution of $\widehat{W}_\eta(x; b)$, which is obtained from $\widehat{W}_q(x; b)$ by replacing $I(x - b < X_{ik}^{(q)} \leq x)$ with $I(x - b < \widehat{\beta}^\top \vec{\mathbf{X}}_{ik} \leq x)$. As noted in Lin et al. (2002), although S_η is referred to as the link function test, anomalies in \bar{W}_η may reflect mis-specification of the link function, of the functional form of the response variable or of the linear predictor.

5. Simulation Studies

5.1 Asymptotic Inference in Finite Samples

To evaluate the finite sample performance of the method, first we examine the variance estimator for $\widehat{\text{pAUC}}(u)$ when there is no covariate. We simulate Y_D from $N(10, 1.5^2)$ and $Y_{\bar{D}}$

from $N(9, 1)$. The induced ROC curve has a partial area of 0.0726 for $\text{FPR} \leq 0.2$. The results, summarized in Table 1, show that the standard error estimates based on large sample approximation are close to the true sampling standard errors. In addition, for confidence intervals, the empirical coverage probabilities are close to their nominal counterparts.

Next, we examine the validity of the large sample approximations in the regression setting for making inference in finite sample sizes. We simulate data from the following models:

$$Y_{D_i} = 10 + 1.3Z_i - \epsilon_{D_i}, \quad \text{for } i = 1, \dots, n_D, \quad (5)$$

$$Y_{\bar{D}_j} = 9 + 0.5Z_j - \epsilon_{\bar{D}_j}, \quad \text{for } j = 1, \dots, n_{\bar{D}}, \quad (6)$$

where Z is generated from $\text{Uniform}(0, C)$. We first set $C = 1$ and generate $\epsilon_{D_i} \sim N(0, 1.5^2)$ and $\epsilon_{\bar{D}_j} \sim N(0, 1)$. The induced pAUC model is:

$$\text{pAUC}_z(u) = \eta_u(1 + 0.8z), \quad \text{where } \eta_u = \int_0^u \Phi \left\{ \frac{x + \Phi^{-1}(v)}{1.5} \right\} dv.$$

We refer to this as the normal-normal model. We choose $u = 0.2$ and fit the data with $\text{pAUC}_z(u) = \eta_u(\beta_0 + \beta_1 z)$. To estimate the FPR conditional on covariates, we use a semi-parametric location model (Heagerty and Pepe, 1999): $S_{D,Z}(y) = S_0(y - \gamma Z)$, where γ and S_0 are unspecified. In Table 2(a), we present the bias, the sampling standard error, average of the standard error estimates and the coverage probability of the 95% confidence intervals for β_0 and β_1 . The standard error estimates are close to the true sampling standard errors. In addition, the empirical coverage probabilities are close to their nominal counterparts.

In another study, we also use models (5) and (6), but simulate ϵ_{D_i} and $\epsilon_{\bar{D}_j}$ from extreme value distributions and Z from $\text{Uniform}(0, 2)$. The corresponding link function η_u is then

$$\eta_u(x) = u - \frac{1 - \exp \{ -(1 - u)^{1 + \exp(x)} \}}{1 + \exp(x)}.$$

The results for $u = 0.2$, summarized in Table 2(b), also show that the asymptotic approximations behave reasonably in finite samples.

5.2 Comparison with Existing Method

To compare the proposed method to the Dodd and Pepe (2003a) approach, we simulate data from models (5) and (6) with ϵ_{Di} and $\epsilon_{\bar{D}j}$ generated from zero-mean normal distributions and extreme value distributions. For each simulated data, we obtain point estimates of β_0 and β_1 with the proposed approach by fitting $\text{pAUC}_z(u) = \eta_u(\beta_0 + \beta_1 z)$, and with Dodd and Pepe (2003a) by fitting $\text{pAUC}_{z_D, z_{\bar{D}}}(u) = \eta_u\{\beta_0 + \beta_1 z_D + \beta_2(z_D - z_{\bar{D}})\}$. The results in Table 3 show that even though the new approach uses a more robust model, the new estimator is more efficient than the Dodd and Pepe (2003a) estimator. At sample sizes of $n_{\bar{D}} = 400$ and $n_D = 100$, the empirical efficiency of the Dodd and Pepe (2003a) method relative to the new method is 57% for β_0 and 51% for β_1 when $\epsilon_{\bar{D}} \sim N(0, 1)$ and $\epsilon_D \sim N(0, 1.5^2)$. When ϵ_D and $\epsilon_{\bar{D}}$ are generated from the extreme value distribution, the relative efficiency is 51% for β_0 and 43% for β_1 .

5.3 Mis-specified Link Function

To examine the properties of the estimator under a mis-specified link function, we simulate data from models (5) and (6) with $Z \sim \text{Uniform}(0, 12)$, and fit the data to the model:

$$\text{pAUC}_z(u) = u\Phi(\beta_0 + \beta_1 z). \quad (7)$$

We generate $\epsilon_{\bar{D}}$ from a standard normal. For ϵ_D , we consider two scenarios, 1) $N(0, 1.5^2)$, and 2) a mixture of $N(2, 3^2)$ with probability 0.3 and $N(7, 1)$ with probability 0.7. To explore how far away from (7) the true underlying link functions are, we examine the linearity of $\Phi^{-1}\{\eta_u(x)/u\}$ in x , where η_u is the true link function. In Figure 1, we can see that (7) is a fair

approximation for the first setting, especially for $x \leq 8$, but not so for the second setting.

As shown in Table 4, the predicted pAUC based on the linear model in (7) has little bias in the first setting, but the bias is substantial in the second setting. To improve the approximation, we instead fit a quadratic spline model for the covariate effect:

$$\text{pAUC}_z(u) = u\Phi \left\{ \beta_0 + \beta_1 z + \beta_2 z^2 + \sum_{k=1}^K b_k (z - \kappa_k)_+^2 \right\} \quad (8)$$

where $x_+ \equiv \max(0, x)$ and $\kappa_1, \dots, \kappa_K$ are the pre-specified knots. In this study, we use 3 knots at 3, 6, and 9. The results, also presented in in Table 4, suggest that the spline model (8) is rather robust with respect to the mis-specification of the link function.

6. Example : Early detection of prostate cancer with PSA

PSA levels in serum are used to screen men for prostate cancer. However considerable controversy exists as to its value. A longitudinal case-control study of PSA as a screening marker for prostate cancer was nested within the Beta-Carotene and Retinol Efficacy Trial, in an effort to evaluate the accuracy of PSA, prior to onset of clinical symptoms, in diagnosing prostate cancer (Thornquist et al, 1993; Etzioni et al, 1999). As part of the protocol, serum was drawn and stored periodically from study participants. 88 subjects developed prostate cancer during the study and their serum samples were analyzed for PSA levels. An age-matched set of 88 control subjects also had their stored serum samples analyzed for PSA levels. The median number of PSA measurements per subject is 4 and the median time interval between two consecutive measurements is 1 year.

Among subjects that develop cancer it is likely that PSA measured closer to the time of onset of clinical symptoms is more predictive of disease than measures taken earlier in time.

Additionally, increasing age is associated with increasing serum PSA level and could affect the discriminatory capacity of PSA. To understand the time and age effect on PSA accuracy, we consider a pAUC model with a covariate T , defined as the time (in years) between the onset of symptoms and the time at which the serum sample was drawn, and an additional covariate $z = \text{age at measurement (in years)}$. We choose the upper bound of FPR as 0.02 which was considered in Baker (2000) for PSA screening and fit the following model

$$\text{pAUC}_{z,T}(0.02) = 0.02\Phi(\beta_0 + \beta_z z + \beta_t T), \quad (9)$$

to the data. Using our approach, the estimate of β_t is -0.091 per year with standard error 0.053 and the coefficient for age, β_z , is estimated as 0.0005 per year of age with standard error 0.020. The negative coefficient for T implies that discrimination improves as T decreases, i.e., when PSA is measured closer to diagnosis. The coefficient for age is almost 0 (p-value = 0.79) suggesting that discrimination is about the same in younger men as in older men.

To examine whether model (9) is appropriate for the data, we consider \bar{W}_z, \bar{W}_T for checking the linearity in specific covariate effects and \bar{W}_η for checking the link function. Figures 2(a)-(c) display the observed processes $\bar{W}_z, \bar{W}_T, \bar{W}_\eta$ along with realizations of $\widehat{W}_z, \widehat{W}_T$ and \widehat{W}_η . The p-values based on the sup-statistics with $\mathcal{J} = 5000$ are 0.38 for the linearity in age, 0.0085 for the linearity in T and is 0.026 for the link function. Thus, the linearity assumption in the time effect is problematic. This motivates us to consider the following model:

$$\text{pAUC}_{z,T}(0.02) = 0.02\Phi(\beta_0 + \beta_z z + \beta_T T + \beta_{T^2} T^2 + \beta_{T^3} T^3) \quad (10)$$

to allow for a non-linear time effect. The resulting estimate of the age effect is $\widehat{\beta}_z = 0.005$ (s.e. = 0.021). The estimated time effects in model (10) are $\widehat{\beta}_T = -0.59$ (s.e. = 0.14), $\widehat{\beta}_{T^2} = -0.10$ (s.e.

$= 0.032$) and $\hat{\beta}_{T,3} = -0.0053$ (s.e. = 0.0022). We apply the model checking procedure again for model (10). The residual plots, shown in Figure 2 (d) – (f), along with the p-values (0.42 for S_z , 0.52 for S_T and 0.50 for S_η) indicate that the revised model is reasonable.

Figure 3 displays the estimated pAUCs and their 95% confidence bands for patients who are 60 years old at different times before clinical diagnosis. For example, when $T = 2$ years, the estimated pAUC(0.02) is 0.0071 (s.e. = 0.0014). Therefore, if we define the restricted reference population as all 60-year old disease-free men with PSA value exceeding its corresponding 98th percentile, there is a $0.0071/0.02 = 36\%$ chance that a randomly selected 60-year old man with cancer whose PSA is measured at 2 years prior to diagnosis is higher than that of a man randomly selected from the restricted reference population. This probability can also be viewed as the average TPR over the range of $FPR \leq 0.02$. Thus the average TPR fluctuates around 30% when $T \geq 3$ years and then improves more quickly to 57% when T decreases to 6 months. This indicates that PSA may not be accurate for detecting prostate cancer early. To fully understand the predictive accuracy of PSA, one needs to further evaluate the positive and negative predictive values of PSA which may be assessed through prospective studies.

7. Remarks

This paper provides an alternative pAUC regression method to Dodd and Pepe (2003a). Advantages of the proposed method include large-sample theory, improved efficiency and model checking procedures. When $u = 1$, the proposed estimator also provides an alternative to the AUC regression approach developed by Dodd and Pepe (2003b). The proposed inference procedure also accounts for possible within-cluster correlation. The model checking procedures are based on a simulation technique that has a minimal computational burden

relative to other re-sampling methods such as the bootstrap. This offers a formal goodness of fit method that is not available with existing AUC and pAUC regression methods. Additional simulation studies indicated that the proposed tests have proper sizes at least when $\min(n_{\bar{D}}, n_D) \geq 200$. The power of the tests would depend on the degree of the model misspecification and remains to be investigated. When applied to the PSA example, the procedure indicated an important non-linearity, which resulted in a revised and better-fitting model. The validity of the proposed inference procedure also requires the correct specification of the FPR model. Goodness of fit for typical FPR models such as the semi-parametric location scale model may be examined based on existing methods such as the procedures proposed in Lin, Wei and Ying (2002) and Cai and Zheng (2005).

Although the focus here was on a general regression model, the method is easily adapted to compare accuracies of two tests, as considered by Wieand et al. (1989). It is straightforward to extend our procedures to make inference about the difference of two pAUCs for both paired and unpaired data. With a single covariate indicating test type, one can create a model based on (2) to examine the difference in the accuracy of two tests. The resulting estimator is equivalent to the estimator proposed by Wieand et al. (1989) when $K_{\bar{D}j} = K_{Di} = 1$.

Appendix

A. Large Sample Properties of $\widehat{\text{pAUC}}(u)$

For technical reasons, we assume that potentially every diseased subject has $\mathcal{K} = \max(K_1, \dots, K_{n_D})$ records and the n_D sets of random vectors $\{\bar{\mathbf{Y}}_{Di}\}$ or $\{(\mathbf{Y}_{Di}, \bar{\mathbf{X}}_i)\}$ with covariates, are independent and identically distributed, where $\bar{\mathbf{Y}}_{Di} = (Y_{Di1}, \dots, Y_{DiK_D})$ and $\bar{\mathbf{X}}_i = (\mathbf{X}_{i1}, \dots, \mathbf{X}_{iK_D})$. Al-

though not every subject with disease has \mathcal{K} records, the presence or absence of individual records in a cluster does not depend on the observations. Corresponding assumptions are made for observations from disease-free subjects.

We assume that $\text{ROC}_{\mathbf{X}}(\cdot)$ is continuously differentiable. The uniform consistency of $\widehat{S}_{\bar{D}}(\cdot)$ and the uniform law of large numbers (Pollard, 1990) ensure the consistency of $\widehat{\text{pAUC}}(u)$. It remains to show the large sample distribution of $\widehat{\text{pAUC}}(u)$. To this end, let $\widehat{\mathcal{I}}_{\bar{D}}(u) = S_{\bar{D}}(\widehat{S}_{\bar{D}}^{-1}(u))$ and $\widetilde{\text{pAUC}}(u) = \frac{1}{N_{\bar{D}}} \sum_{i,k} V_{\text{Dik}}^{(u)}$. We note that

$$\begin{aligned} n_{\bar{D}}^{\frac{1}{2}} \left\{ \widehat{\text{pAUC}}(u) - \text{pAUC}(u) \right\} &= n_{\bar{D}}^{\frac{1}{2}} \left\{ \widetilde{\text{pAUC}}(\widehat{\mathcal{I}}_{\bar{D}}(u)) - \widetilde{\text{pAUC}}(u) \right\} + n_{\bar{D}}^{\frac{1}{2}} \left\{ \widetilde{\text{pAUC}}(u) - \text{pAUC}(u) \right\} \\ &\quad + \frac{n_{\bar{D}}^{\frac{1}{2}}}{N_{\bar{D}}} \sum_{i,k} I(U_{\text{Dik}} \leq \widehat{\mathcal{I}}_{\bar{D}}(u)) \left\{ u - \widehat{\mathcal{I}}_{\bar{D}}(u) - \widehat{U}_{\text{Dik}} + U_{\text{Dik}} \right\}. \end{aligned}$$

It has been shown that $\sup_u |\widehat{\mathcal{I}}_{\bar{D}}(u) - u| \rightarrow 0$ and $n_{\bar{D}}^{\frac{1}{2}} \{\widehat{\mathcal{I}}_{\bar{D}}(u) - u\}$ is asymptotically equivalent to $n_{\bar{D}}^{-\frac{1}{2}} \sum_{j=n_{\bar{D}}+1}^{n_{\bar{D}}+n_{\bar{D}}} I_{\bar{D}j}(u)$, where $I_{\bar{D}j}(u) = \sum_{l=1}^{K_j} \{u - I(Y_{\bar{D}jl} \geq S_{\bar{D}}^{-1}(u))\}$ (Cai and Pepe, 2002). This, coupled with the equicontinuity of the process $n_{\bar{D}}^{\frac{1}{2}} \{\widetilde{\text{pAUC}}(u) - \text{pAUC}(u)\}$, ensures that

$$\begin{aligned} n_{\bar{D}}^{\frac{1}{2}} \left\{ \widehat{\text{pAUC}}(u) - \text{pAUC}(u) \right\} &\approx n_{\bar{D}}^{\frac{1}{2}} \left\{ \text{pAUC}(\widehat{\mathcal{I}}_{\bar{D}}(u)) - \text{pAUC}(u) \right\} + n_{\bar{D}}^{\frac{1}{2}} \left\{ \widetilde{\text{pAUC}}(u) - \text{pAUC}(u) \right\} \\ &\quad + \frac{n_{\bar{D}}^{\frac{1}{2}}}{N_{\bar{D}}} \sum_{i,k} I(U_{\text{Dik}} \leq \widehat{\mathcal{I}}_{\bar{D}}(u)) \left\{ u - \widehat{\mathcal{I}}_{\bar{D}}(u) - \widehat{U}_{\text{Dik}} + U_{\text{Dik}} \right\}. \end{aligned}$$

It follows from a Taylor series expansion, the convergence of $\frac{1}{N_{\bar{D}}} \sum_{i,k} I(U_{\text{Dik}} \leq u) \rightarrow \text{ROC}(u)$ and the equicontinuity of the process $n_{\bar{D}}^{\frac{1}{2}} \{\widehat{\mathcal{I}}_{\bar{D}}(u) - u\}$ that

$$\begin{aligned} n_{\bar{D}}^{\frac{1}{2}} \left\{ \widehat{\text{pAUC}}(u) - \text{pAUC}(u) \right\} &\approx n_{\bar{D}}^{\frac{1}{2}} \text{ROC}(u) \left\{ \widehat{\mathcal{I}}_{\bar{D}}(u) - u \right\} + n_{\bar{D}}^{\frac{1}{2}} \left\{ \widetilde{\text{pAUC}}(u) - \text{pAUC}(u) \right\} \\ &\quad - p_{10}^{\frac{1}{2}} \int_0^u n_{\bar{D}}^{\frac{1}{2}} \left\{ \widehat{\mathcal{I}}_{\bar{D}}^{-1}(v) - v \right\} d\text{ROC}(v) - n_{\bar{D}}^{\frac{1}{2}} \text{ROC}(u) \left\{ \widehat{\mathcal{I}}_{\bar{D}}(u) - u \right\} \\ &\approx n_{\bar{D}}^{-\frac{1}{2}} \sum_{i=1}^{n_{\bar{D}}} \mathcal{P}_{\text{Di}} + n_{\bar{D}}^{-\frac{1}{2}} \sum_{j=n_{\bar{D}}+1}^{n_{\bar{D}}+n_{\bar{D}}} \mathcal{P}_{\bar{D}j}, \end{aligned}$$

where $\mathcal{P}_{D_i} = K_D^{-1} \sum_{k=1}^{K_i} V_{D_{ik}}^{(u)} - \text{pAUC}(u)$, $\mathcal{P}_{\bar{D}_j} = p_{10}^{\frac{1}{2}} \int_0^u I_{\bar{D}_j}(v) d\text{ROC}(v)$ and K_D is the limit of N_D/n_D . It follows from the central limit theorem that $n_D^{\frac{1}{2}} \{\widehat{\text{pAUC}}(u) - \text{pAUC}(u)\}$ converges in distribution to a zero-mean normal with variance σ^2 , where σ^2 is the limit of $n_D^{-1} \sum_{i=1}^{n_D} \mathcal{P}_{D_i}^2 + n_{\bar{D}}^{-1} \sum_{j=n_D+1}^{n_D+n_{\bar{D}}} \mathcal{P}_{\bar{D}_j}^2$. A consistent estimate of σ^2 is $\hat{\sigma}^2$ which is obtained by replacing all the theoretical quantities in $n_D^{-1} \sum_{i=1}^{n_D} \mathcal{P}_{D_i}^2 + n_{\bar{D}}^{-1} \sum_{j=n_D+1}^{n_D+n_{\bar{D}}} \mathcal{P}_{\bar{D}_j}^2$ by their empirical counterparts.

B. Large Sample Properties of $\hat{\beta}$

To show the existence and uniqueness of $\hat{\beta}$, we assume that the covariates $\mathbf{X} = (\mathbf{Z}, \mathbf{Z}_D)$ are bounded, the estimators of $S_{\bar{D}, \mathbf{Z}}(y)$ are uniformly consistent and $n_{\bar{D}}^{\frac{1}{2}} \{\hat{S}_{\bar{D}, \mathbf{Z}}(y) - S_{\bar{D}, \mathbf{Z}}(y)\}$ converges weakly to a Gaussian process uniformly in y and \mathbf{Z} . Without loss of generality, we also assume that $n_{\bar{D}}^{\frac{1}{2}} \{\hat{\mathcal{I}}_{\bar{D}, \mathbf{Z}}(u) - u\}$ can be approximated by a sum of independent terms:

$$\sup_{u, \mathbf{Z}} \left| n_{\bar{D}}^{\frac{1}{2}} \left\{ \hat{\mathcal{I}}_{\bar{D}, \mathbf{Z}}(u) - u \right\} - n_{\bar{D}}^{-\frac{1}{2}} \sum_{j=n_D+1}^{n_D+n_{\bar{D}}} I_{\bar{D}_j}(u, \mathbf{Z}) \right| \rightarrow 0 \quad (11)$$

in probability, where $\hat{\mathcal{I}}_{\bar{D}, \mathbf{Z}}(u) = S_{\bar{D}}(\hat{S}_{\bar{D}, \mathbf{Z}}^{-1}(u))$. Let $\bar{\mathbf{V}}(\beta)$ denote the left hand side of (3). It is easy to see that $\frac{\partial \bar{\mathbf{V}}(\beta)}{\partial \beta} = \hat{\mathbb{A}}(\beta)$, where $\hat{\mathbb{A}}(\beta) = \frac{1}{N_D} \sum_{i,k} \dot{\eta}(\beta^\top \bar{\mathbf{X}}_{ik}) \bar{\mathbf{X}}_{ik}^{\otimes 2}$, which is nonnegative definite. Furthermore, $\hat{\mathbb{A}}(\beta_0) \rightarrow \mathbb{A}$. When $\bar{\mathbf{X}}_{ik}$ is non-degenerate, \mathbb{A} is positive definite. Now, since $\bar{\mathbf{V}}(\beta_0) \rightarrow 0$, by the standard inverse function theorem, there exists a unique solution $\hat{\beta}$ to the equation $\bar{\mathbf{V}}(\beta)$ in a neighborhood of β_0 . This, coupled with the nonnegativity of $\hat{\mathbb{A}}(\beta)$, ensures the uniqueness of the root of $\bar{\mathbf{V}}(\beta) = 0$ in the entire domain of β asymptotically. The above proof also implies that $\hat{\beta}$ is strongly consistent.

By the consistency of $\hat{\beta}$ and a Taylor series expansion of $\bar{\mathbf{V}}(\hat{\beta})$ around β_0 , we obtain

$$n_D^{\frac{1}{2}} (\hat{\beta} - \beta) \approx \mathbb{A}^{-1} n_D^{\frac{1}{2}} \bar{\mathbf{V}}(\beta_0). \quad (12)$$

Define $V_{Dik}^{(u)} = u - \min(u, U_{Dik})$, $e_{ik} = V_{Dik}^{(u)} - \text{pAUC}_{\mathbf{X}_{ik}}(u)$, $\tilde{\mathbf{V}}(u) = \frac{1}{N_D} \sum_{i,k} w(\mathbf{X}_{ik}) \vec{\mathbf{X}}_{ik} e_{ik}$ and $\bar{\mathbf{V}}_1 = \frac{1}{N_D} \sum_{i,k} w(\mathbf{X}_{ik}) \vec{\mathbf{X}}_{ik} (\hat{V}_{Dik}^{(u)} - V_{Dik}^{(u)})$. Then $\bar{\mathbf{V}}(\beta_0) = \tilde{\mathbf{V}}(u) + \bar{\mathbf{V}}_1$. We first show the large sample approximation for $n_D^{\frac{1}{2}} \bar{\mathbf{V}}_1$. We note that

$$\begin{aligned} n_D^{\frac{1}{2}} \bar{\mathbf{V}}_1 &= \frac{n_D^{\frac{1}{2}}}{N_D} \sum_{i,k} w(\mathbf{X}_{ik}) \vec{\mathbf{X}}_{ik} \left[I(U_{Dik} \leq \hat{\mathcal{I}}_{\bar{\mathbf{D}}, \mathbf{Z}_{ik}}(u)) \left\{ \hat{\mathcal{I}}_{\bar{\mathbf{D}}, \mathbf{Z}_{ik}}(u) - U_{Dik} \right\} - V_{Dik}^{(u)} \right] \\ &\quad + \frac{n_D^{\frac{1}{2}}}{N_D} \sum_{i,k} w(\mathbf{X}_{ik}) \vec{\mathbf{X}}_{ik} I(U_{Dik} \leq \hat{\mathcal{I}}_{\bar{\mathbf{D}}, \mathbf{Z}_{ik}}(u)) \left\{ u - \hat{\mathcal{I}}_{\bar{\mathbf{D}}, \mathbf{Z}_{ik}}(u) - \hat{U}_{Dik} + U_{Dik} \right\}. \end{aligned}$$

It follows from the equicontinuity of $n_D^{\frac{1}{2}} \tilde{\mathbf{V}}(\cdot)$ and the uniform consistency of $\hat{S}_{\bar{\mathbf{D}}, \mathbf{Z}}(\cdot)$ that

$$\begin{aligned} n_D^{\frac{1}{2}} \bar{\mathbf{V}}_1 &\approx \frac{n_D^{\frac{1}{2}}}{N_D} \sum_{i,k} w(\mathbf{X}_{ik}) \vec{\mathbf{X}}_{ik} \left\{ \text{ROC}_{\mathbf{X}_{ik}}(u) - I(U_{Dik} \leq u) \right\} \left\{ \hat{\mathcal{I}}_{\bar{\mathbf{D}}, \mathbf{Z}_{ik}}(u) - u \right\} \\ &\quad - \frac{n_D^{\frac{1}{2}}}{N_D} \sum_{i,k} w(\mathbf{X}_{ik}) \vec{\mathbf{X}}_{ik} I(U_{Dik} \leq u) (\hat{U}_{Dik} - U_{Dik}). \end{aligned} \quad (13)$$

Since $n_D^{\frac{1}{2}} \{ \hat{\mathcal{I}}_{\bar{\mathbf{D}}, \mathbf{Z}}(u) - u \}$ converges weakly to a Gaussian process, using the strong law of large numbers and the strong representation theorem (Pollard, 1990), one can show that (13) $\rightarrow 0$ in probability. Therefore,

$$n_D^{\frac{1}{2}} \bar{\mathbf{V}}_1 \approx \frac{n_D^{\frac{1}{2}}}{N_D} \sum_{i,k} w(\mathbf{X}_{ik}) \vec{\mathbf{X}}_{ik} \int_0^u \left\{ \hat{\mathcal{I}}_{\bar{\mathbf{D}}, \mathbf{Z}_{ik}}(u) - u \right\} d\text{ROC}_{\mathbf{X}_{ik}}(u).$$

This, coupled with (11) and (12), implies that $n_D^{\frac{1}{2}} (\hat{\beta} - \beta_0) \approx \mathbb{A}^{-1} \{ n_D^{-\frac{1}{2}} \sum_{i=1}^{n_D} \mathfrak{B}_{Di} + n_D^{-\frac{1}{2}} \sum_{j=n_D+1}^{n_D+n_D} \mathfrak{B}_{Dj} \}$.

C. Large Sample Distribution of $\bar{W}(\mathbf{x})$ Under Model (2)

Let $I_{ik}^{(q)}(x; b)$ denote $I(x - b < X_{ik}^{(q)} \leq x)$. By the consistency of $\hat{\beta}$ and the Taylor series expansion, uniformly in x , we have

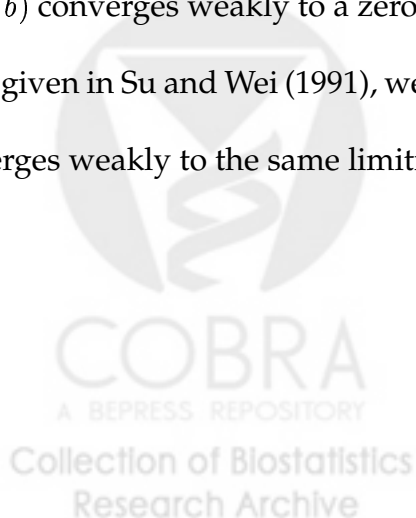
$$n_D^{\frac{1}{2}} \bar{W}_q(x; b) \approx \frac{n_D^{\frac{1}{2}}}{N_D} \sum_{i,k} I_{ik}^{(q)}(x; b) (\hat{V}_{Dik}^{(u)} - V_{Dik}^{(u)}) + \frac{n_D^{\frac{1}{2}}}{N_D} \sum_{i,k} I_{ik}^{(q)}(x; b) e_{ik} - \hat{\mathbf{R}}_q(x; b)^\top (\hat{\beta} - \beta_0).$$

Furthermore, the uniform law of large numbers (Pollard, 1990) implies that $\widehat{\mathbf{R}}_q(x)$ converges, uniformly in x , to a non-random function, $\mathbf{R}_q(x)$. Using arguments similar to those given in Appendix B and the large sample properties of $n_D^{\frac{1}{2}}(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)$, one can show that

$$n_D^{\frac{1}{2}}\bar{W}_q(x; b) \approx n_D^{\frac{1}{2}} \sum_{i=1}^{n_D} \left\{ W_{q_{D_1i}}(x; b) + W_{q_{D_2i}}(x; b) \right\} + n_D^{-\frac{1}{2}} \sum_{j=n_D+1}^{n_D+n_{\bar{D}}} \left\{ W_{q_{\bar{D}_1j}}(x; b) + W_{q_{\bar{D}_2j}}(x; b) \right\},$$

where $W_{q_{D_1i}}(x; b) = K_D^{-1} \sum_{k=1}^{K_i} I_{ik}^{(q)}(x; b) e_{ik}$, $W_{q_{D_2i}}(x; b) = R_q(x; b)^\top \mathbb{A}^{-1} \mathfrak{B}_{D_i}$, $W_{q_{\bar{D}_1j}}(x; b)$ is the limit of $\frac{p_{j0}}{N_D} \sum_{i,k} I_{ik}^{(q)}(x; b) \int_0^u I_{D_j}(v, \mathbf{Z}_{ik}) d\text{ROC}_{\mathbf{x}_{ik}}(v)$, and $W_{q_{\bar{D}_2j}}(x; b) = R_q(x; b)^\top \mathbb{A}^{-1} \mathfrak{B}_{\bar{D}_j}$.

For fixed x , $n_D^{-\frac{1}{2}} \sum_{i=1}^{n_D} \{W_{q_{D_1i}}(x; b) + W_{q_{D_2i}}(x; b)\}$ and $n_D^{-\frac{1}{2}} \sum_{j=n_D+1}^{n_D+n_{\bar{D}}} \{W_{q_{\bar{D}_1j}}(x; b) + W_{q_{\bar{D}_2j}}(x; b)\}$ are essentially sums of independent and identically distributed zero-mean random variables. It follows from the multivariate central limit theorem that $W_q(x; b)$ converges in finite dimensional distributions to a zero-mean Gaussian process. Since $R_q(x; b)^\top \mathbb{A}^{-1}$ is non-random and $n_D^{-\frac{1}{2}} \sum_{i=1}^{n_D} \mathfrak{B}_{D_i} + n_D^{-\frac{1}{2}} \sum_{j=n_D+1}^{n_D+n_{\bar{D}}} \mathfrak{B}_{\bar{D}_j}$ does not involve x , $n_D^{-\frac{1}{2}} \sum_{i=1}^{n_D} W_{q_{D_2i}}(x; b) + n_D^{-\frac{1}{2}} \sum_{j=n_D+1}^{n_D+n_{\bar{D}}} W_{q_{\bar{D}_2j}}(x; b)$ is tight. Now, both $W_{q_{D_1i}}(\mathbf{x})$ and $W_{q_{\bar{D}_1j}}(x; b)$ are uniformly bounded monotone functions, which are clearly manageable (Pollard, 1990, p38). It follows from the functional central limit theorem (Pollard, 1990, p53) that $n_D^{-\frac{1}{2}} \sum_{i=1}^{n_D} W_{q_{D_1i}}(x; b) + n_D^{-\frac{1}{2}} \sum_{j=n_D+1}^{n_D+n_{\bar{D}}} W_{q_{\bar{D}_1j}}(x; b)$ is tight. Hence, $W_q(x; b)$ converges weakly to a zero-mean Gaussian process. Appealing arguments similar to those given in Su and Wei (1991), we have that, conditional on the data, the process $n_D^{\frac{1}{2}}\widehat{W}_q(x; b)$ converges weakly to the same limiting Gaussian process as that of $n_D^{\frac{1}{2}}\bar{W}_q(x; b)$.



REFERENCES

- Baker, S. G. (2000). Identifying combinations of cancer markers for further study as triggers of early intervention. *Biometrics* **56**, 1082–7.
- Baker, S. G. and Pinsky, P. F. (2001). A proposed design and analysis for comparing digital and analog mammography: Special receiver operating characteristic methods for cancer screening. *J. Amer. Statist. Assoc.* **96**, 421–28.
- Bamber, D. (1975). The area above the ordinal dominance graph and the area below the receiver operating characteristic graph. *J. Mathematical Psychology* **12**, 387–415.
- Brusic, V., Rudy, G., Honeyman, M., Hammer, J. and Harrison, L. (1998). Prediction of MHC class II-binding peptides using an evolutionary algorithm and artificial neural network. *Bioinformatics* **14**, 121–30.
- Cai, T. (2004). Semiparametric ROC regression analysis with placement values. *Biostatistics* **5**, 45–60.
- Cai, T. and Pepe, M. S. (2002). Semiparametric receiver operating characteristic analysis to evaluate biomarkers for disease. *J. Amer. Statist. Assoc.* **97**, 1099–107.
- Cai, T., Wei, L. J. and Wilcox, M. (2000). Semiparametric regression analysis for clustered failure time data. *Biometrika* **87**, 867–78.
- Cai, T. and Zheng, Y. (2005). Model checking for roc regression analysis. Technical report, Harvard University.
- 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 ap-

- proach. *Biometrics* **44**, 837–45.
- Dodd, L. and Pepe, M. S. (2003a). Partial AUC estimation and regression. *Biometrics* **59**, 614–23.
- Dodd, L. and Pepe, M. S. (2003b). Semi-parametric regression for the area under the receiver operating characteristic curve. *J. Amer. Statist. Assoc.* **98**, 409–17.
- Dwyer (1996). In pursuit of a piece of the ROC. *Radiology* **201**, 621–5.
- Etzioni, R., Pepe, M., Longton, G., Hu, C. and Goodman, G. (1999). Incorporating the time dimension in receiver operating characteristic curves: A case study of prostate cancer. *Medical Decision Making* **19**, 242–51.
- Han, A. K. (1987). A non-parametric analysis of transformations. *J. Econometrics* **35**, 191–209.
- Hanely, J. A. and Hajian-Tilaki, K. O. (1997). Sampling variability of nonparametric estimate of the areas under receiver operating characteristic curves: An update. *Academic Radiology* **4**, 49–58.
- Heagerty, P. J. and Pepe, M. S. (1999). Semiparametric estimation of regression quantiles with application to standardizing weight for height and age in US children. *Appl. Statist.* **48**, 533–51.
- Jiang, Y. L., Metz, C. E. and Nishikawa, R. M. (1996). A receiver operating characteristic partial area index for highly sensitive diagnostic tests. *Radiology* **201**, 745–50.
- Lin, D. Y., Wei, L. J. and Ying, Z. (2002). Model-checking techniques based on cumulative residuals. *Biometrics* **58**, 1–12.
- McClish, R. J. (1989). Analyzing a portion of the ROC curve. *Medical Decision Making* **9**, 190–5.
- Pepe, M. S. (1998). Three approaches to regression analysis of receiver operating characteristic

- curves for continuous test results. *Biometrics* **54**, 124–35.
- Pepe, M. S. (2003). *The Statistical Evaluation of Medical Tests for Classification and Prediction*. Oxford University Press, United Kingdom.
- Pepe, M. S. and Cai, T. (2004). The analysis of placement values for evaluating discriminatory measures. *Biometrics* **60**, 528–35.
- Pollard, D. (1990). *Empirical Processes: Theory and Applications*. Hayward, CA: Institute of Mathematical Statistics.
- Su, J. Q. and Wei, L. J. (1991). A lack-of-fit test for the mean function in a generalized linear model. *J. Amer. Statist. Assoc.* **86**, 420–6.
- Thompson, M. L. and Zucchini, W. (1989). On the statistical analysis of ROC curves. *Statist. in Medicine* **8**, 1277–90.
- Thornquist, M. D., Omenn, G. S., Goodman, G. E. and et al (1993). Statistical design and monitoring of the carotene and retinol efficacy trial. *Controlled Clinical Trials* **14**, 308–24.
- Viaene, S., Derrig, R. A., Baesens, B. and Dedene, G. (2002). A comparison of state-of-the-art classification techniques for expert automobile insurance claim fraud detection. *J. Risk and Insurance* **69**, 372–421.
- Wieand, S., Gail, M. H., James, B. R. and James, K. L. (1989). A family of nonparametric statistics for comparing diagnostic markers with paired or unpaired data. *Biometrika* **76**, 585–92.
- Zhang, D. D., Zhou, X.-H., Freeman, Daniel H., J. and Freeman, J. L. (2002). A non-parametric method for the comparison of partial areas under ROC curves and its application to large health care data sets. *Statist. in Medicine* **21**, 701–15.

Table 1. *The Bias, Sampling Standard Error (SSE), sample average of the Estimated Standard Errors (ESE) and empirical Coverage Probability (CovP) of the 95% confidence interval for \widehat{pAUC} . Results are based on 1000 simulated datasets.*

| | $n_D = 100$ | | | | $n_D = 200$ | | | |
|-------------|-------------|------|------|------|-------------|------|------|------|
| | Bias | SSE | ESE | CovP | Bias | SSE | ESE | CovP |
| $n_D = 100$ | .0005 | .011 | .011 | .946 | .0005 | .010 | .010 | .945 |
| $n_D = 200$ | .0007 | .010 | .010 | .940 | .0002 | .008 | .008 | .944 |



Table 2. Bias, sampling standard error (SSE), average of the estimated standard error estimator (ESE), and the coverage probability (CovP) of the 95% confidence interval. Each entry is based on 1000 simulation samples.

(a) $N(0, 1)$ versus $N(0, 1.5^2)$

| $(n_{\bar{D}}, n_D)$ | β_0 | | | | β_1 | | | |
|----------------------|-----------|------|------|------|-----------|------|------|------|
| | Bias | SSE | ESE | CovP | Bias | SSE | ESE | CovP |
| (100, 100) | .008 | .430 | .431 | .950 | .044 | .714 | .724 | .960 |
| (100, 200) | .023 | .332 | .347 | .960 | .000 | .551 | .576 | .957 |
| (200, 100) | .021 | .376 | .389 | .963 | -.009 | .643 | .663 | .959 |
| (400, 100) | -.002 | .369 | .367 | .949 | .020 | .635 | .632 | .955 |

(b) Extreme Value versus Extreme Value

| $(n_{\bar{D}}, n_D)$ | β_0 | | | | β_1 | | | |
|----------------------|-----------|------|------|------|-----------|------|------|------|
| | Bias | SSE | ESE | CovP | Bias | SSE | ESE | CovP |
| (100, 100) | .031 | .453 | .472 | .950 | .007 | .323 | .349 | .967 |
| (100, 200) | .020 | .408 | .419 | .947 | .011 | .290 | .304 | .963 |
| (200, 100) | .002 | .385 | .391 | .955 | .008 | .277 | .294 | .966 |
| (400, 100) | -.010 | .334 | .343 | .951 | .017 | .257 | .263 | .954 |

Table 3. Estimates of β_0 and β_1 compared with their respective actual values $\beta_0 = 1$ and $\beta_1 = 0.8$, based on the Dodd & Pepe approach (D&P) and on the New approach (New). Results are based on 1000 simulated datasets.

(a) $N(0, 1)$ versus $N(0, 1.5^2)$

| $(n_{\bar{D}}, n_D)$ | Bias | | | | Mean Squared Error | | | |
|----------------------|------------------------|-------------------------|------------------------|-------------------------|------------------------|-------------------------|------------------------|-------------------------|
| | β_0^{New} | $\beta_0^{\text{D\&P}}$ | β_1^{New} | $\beta_1^{\text{D\&P}}$ | β_0^{New} | $\beta_0^{\text{D\&P}}$ | β_1^{New} | $\beta_1^{\text{D\&P}}$ |
| (100, 100) | .008 | -.004 | .044 | .028 | .185 | .623 | .512 | 2.182 |
| (100, 200) | .023 | .015 | .000 | -.017 | .111 | .559 | .304 | 1.994 |
| (200, 100) | .021 | .004 | -.009 | -.010 | .148 | .327 | .414 | 1.103 |
| (400, 100) | -.002 | -.002 | .020 | .013 | .136 | .237 | .404 | .787 |

(b) Extreme Value versus Extreme Value

| $(n_{\bar{D}}, n_D)$ | Bias | | | | Mean Squared Error | | | |
|----------------------|------------------------|-------------------------|------------------------|-------------------------|------------------------|-------------------------|------------------------|-------------------------|
| | β_0^{New} | $\beta_0^{\text{D\&P}}$ | β_1^{New} | $\beta_1^{\text{D\&P}}$ | β_0^{New} | $\beta_0^{\text{D\&P}}$ | β_1^{New} | $\beta_1^{\text{D\&P}}$ |
| (100, 100) | .031 | -.045 | .007 | .064 | .206 | .678 | .104 | .542 |
| (100, 200) | .020 | -.010 | .011 | .023 | .167 | .540 | .084 | .435 |
| (200, 100) | .002 | -.002 | .008 | .010 | .148 | .342 | .077 | .251 |
| (400, 100) | -.010 | -.032 | .017 | .029 | .111 | .217 | .066 | .153 |

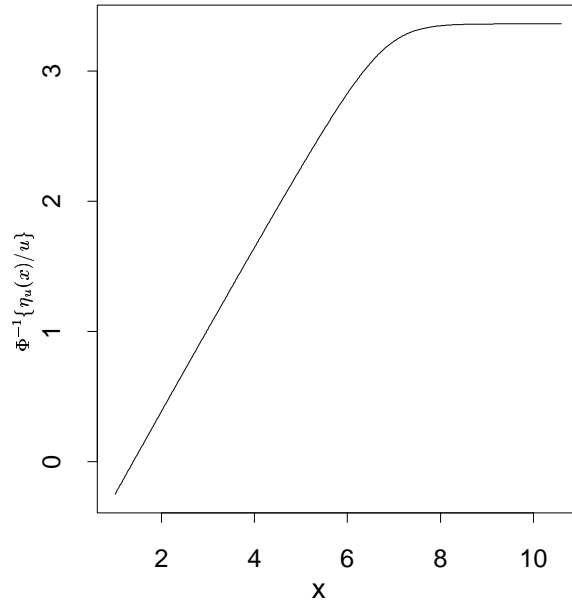
Table 4. Bias and mean squared error (MSE) of the predicted pAUC. For each dataset, we fit with two models: $pAUC_z(0.2) = 0.2\Phi(\beta_0 + \beta_1 z)$ (Linear) and $pAUC_z(0.2) = 0.2\Phi\{\beta_0 + \beta_1^T \mathcal{R}(z)\}$ (Spline).

The results are based on 1000 simulated datasets with sample size $n_D = n_{\bar{D}} = 200$.

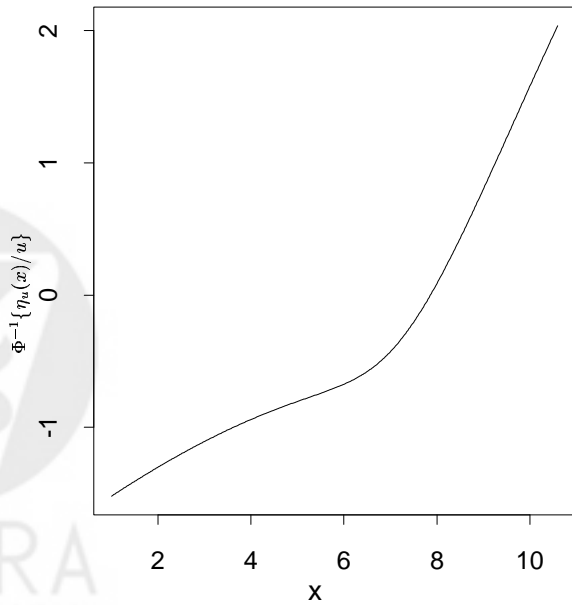
| z | True $pAUC_z(0.2)$ | $\epsilon_D \sim N(0, 1.5^2)$ | | | | $\epsilon_D \sim \text{Normal Mixture}$ | | | |
|----|-----------------------|-------------------------------|---------|--------|--------|---|---------|--------|--------|
| | | Bias | | MSE | | Bias | | MSE | |
| | | Linear | Spline | Linear | Spline | Linear | Spline | Linear | Spline |
| 2 | 2.4E-2 | 5.1E-4 | -8.8E-4 | 1.1E-4 | 2.4E-4 | -9.8E-3 | 1.5E-4 | 1.3E-4 | 1.3E-4 |
| 4 | 3.6E-2 | 9.1E-5 | 3.5E-4 | 1.6E-5 | 3.5E-5 | -3.0E-3 | -3.1E-4 | 7.0E-5 | 1.3E-4 |
| 6 | 4.8E-2 | -1.1E-4 | -8.2E-4 | 6.8E-7 | 1.3E-5 | 1.7E-2 | -6.0E-4 | 3.6E-4 | 1.3E-4 |
| 8 | 8.0E-2 | 4.2E-5 | 8.1E-5 | 1.3E-8 | 1.0E-8 | 2.6E-2 | 1.9E-3 | 7.5E-4 | 2.3E-4 |
| 10 | 1.6E-1 | 7.3E-5 | 8.0E-5 | 5.7E-7 | 6.1E-5 | -1.3E-2 | 1.6E-4 | 2.4E-4 | 2.0E-4 |



Figure 1. Plot of function $\Phi^{-1}\{\eta_u(x)/u\}$ for $u = 0.2$.



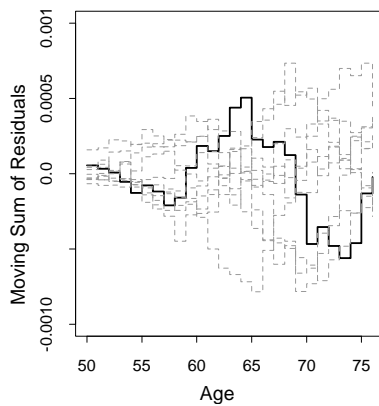
(a) Normal



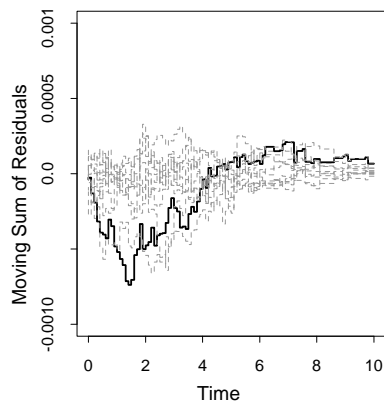
(b) Mixture

Figure 2. Plot of moving sums of residuals: (a) & (d) for testing linear age effect with $b = 10$ (interquartile range of age); (b) & (e) for testing linear time effect with $b = 3$; (c) & (f) for testing the linearity of the model with $b = 1$; The observed pattern is shown by the thick solid curve, and 10 simulated realizations under the null are shown by the dotted curve.

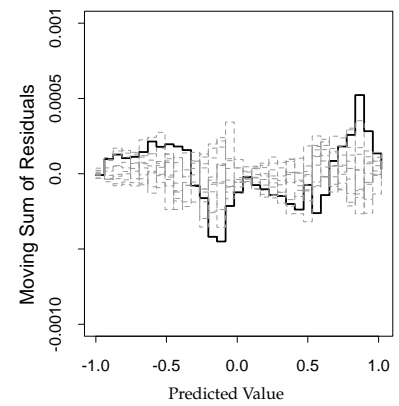
(a) – (c) : Linear covariate effect model



(a) p-value = 0.38

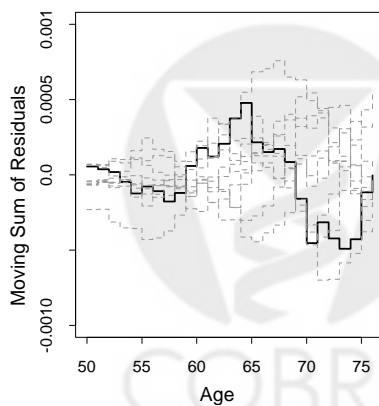


(b) p-value = 0.0085

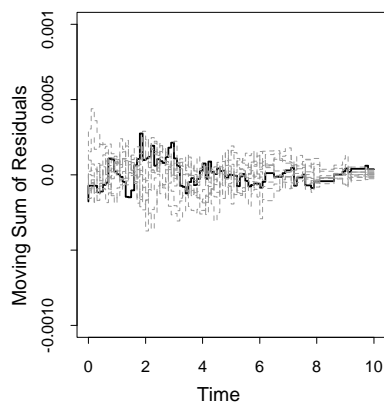


(c) p-value = 0.026

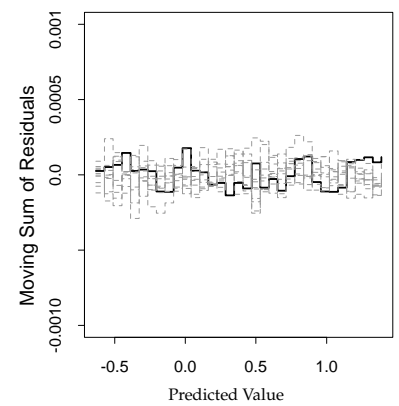
(d) – (f) : Cubic time effect and linear age effect model



(d) p-value = 0.42



(e) p-value = 0.52



(f) p-value = 0.50

Figure 3. Predicted pAUC for PSA as a biomarker of prostate cancer in 60 year old men. Shown also are their 95% confidence intervals.

