

Partial AUC Estimation and Regression

Lori E. Dodd^{1,*} and Margaret S. Pepe^{2,**}

¹Biometric Research Branch, National Cancer Institute, 6130 Executive Blvd, MSC 7434,
Rockville, Maryland 20892, U.S.A.

²Department of Biostatistics, University of Washington, Box 357232, Seattle, Washington 98195, U.S.A.

**email:* doddl@mail.nih.gov

***email:* mspepe@u.washington.edu

SUMMARY. Accurate diagnosis of disease is a critical part of health care. New diagnostic and screening tests must be evaluated based on their abilities to discriminate diseased from nondiseased states. The partial area under the receiver operating characteristic (ROC) curve is a measure of diagnostic test accuracy. We present an interpretation of the partial area under the curve (AUC), which gives rise to a nonparametric estimator. This estimator is more robust than existing estimators, which make parametric assumptions. We show that the robustness is gained with only a moderate loss in efficiency. We describe a regression modeling framework for making inference about covariate effects on the partial AUC. Such models can refine knowledge about test accuracy. Model parameters can be estimated using binary regression methods. We use the regression framework to compare two prostate-specific antigen biomarkers and to evaluate the dependence of biomarker accuracy on the time prior to clinical diagnosis of prostate cancer.

KEY WORDS: Diagnostic testing; Mann-Whitney U-statistic; Regression; Receiver operating characteristic curve.

1. Introduction

Screening and diagnostic tests are familiar and ever-evolving tools of modern medicine. Populations of healthy individuals characterized as at-risk are commonly screened for diseases such as cancer and heart disease. Early detection by screening is considered essential to alleviate disease burden, and considerable resources have been devoted to developing new screening tests. New diagnostic tests that are less invasive, less expensive, and more accurate than existing procedures are sought for diagnosis of many conditions. Technologies that measure gene and protein expression, as well as new imaging procedures, all hold promise in this regard. Prior to widespread application, however, rigorous evaluation of test accuracy and of factors that effect it is compulsory.

Inherent in the analysis of screening and diagnostic tests are costs and benefits, both monetary and nonmonetary, associated with true-positive and false-positive diagnoses. Consider a continuous test result Y for which $Y > c$ indicates a positive test result, and let D and \bar{D} denote diseased and nondiseased states, respectively. The true-positive rate at a threshold c , $\text{TPR}(c)$, is defined as $P(Y > c | D) \equiv S_D(c)$. The corresponding false-positive rate, $\text{FPR}(c)$, is $P(Y > c | \bar{D}) \equiv S_{\bar{D}}(c)$. Costs and benefits are associated with any given $\{\text{FPR}(c), \text{TPR}(c)\}$ pair. The receiver operating characteristic (ROC) curve plots $\{\text{FPR}(c), \text{TPR}(c)\}$ for all possible thresholds c , and provides a visual description of the trade-offs between TPRs and FPRs as one changes the threshold stringency (Figure 1). We can write the ROC curve as a function of $t = S_{\bar{D}}(c)$ as follows: $\text{ROC}(t) = S_D\{S_{\bar{D}}^{-1}(t)\}$. An uninformative test is represented

by a straight line from the $(0, 0)$ vertex to $(1, 1)$, while a curve pulled closer towards the $(0, 1)$ vertex indicates a better-performing test.

Frequently, the best threshold is not known when a test is under evaluation, and it may vary depending on the setting in which the test is implemented. A summary measure that aggregates performance information across a range of possible thresholds is desirable. The area under the ROC curve (AUC), defined as $\int_0^1 \text{ROC}(t) dt$, summarizes across all thresholds, and is the most commonly used measure of diagnostic accuracy for quantitative tests. However, the AUC summarizes test performance over regions of the ROC space in which one would never operate, i.e., for $\{\text{FPR}(c), \text{TPR}(c)\}$ values of no practical interest. In population screening, large monetary costs result from high false-positive rates; hence the region of the curve corresponding to low false-positive rates is of primary interest. In diagnostic testing, it is critical to maintain a high TPR in order not to miss detecting subjects with disease. In this case, interest is in the region of the ROC curve corresponding only to acceptably high TPRs. In this article, we consider a summary index for the ROC curve restricted to a clinically relevant range of false-positive rates. The partial AUC is

$$\text{AUC}(t_0, t_1) = \int_{t_0}^{t_1} \text{ROC}(t) dt, \quad (1)$$

where the interval (t_0, t_1) denotes the false-positive rates of interest. The analogue that restricts to a range of true-positive rates will also be discussed. Selecting the interval (t_0, t_1) is an

important practical issue. The choice depends on the particular setting and should depend on the costs of a false-positive diagnosis, as well as the benefits of a true positive. Baker (2000) develops a “utility” function to specify a target partial ROC region. Obuchowski (1997) uses decision analysis to associate patient outcome with desirable diagnostic accuracy values. Such methods could be adapted, for example, to determine a maximum allowable false-positive rate or the lowest desirable true-positive rate. This is a complex area requiring input from health services and economic specialists.

Although the partial AUC has been proposed before (McClish, 1989; Thompson and Zucchini, 1989; Jiang, Metz, and Nishikawa, 1996) and has gained popularity, particularly in screening research (Baker and Pinsky, 2001), little attention has been devoted to statistical inference about it. We provide a probabilistic interpretation for the partial AUC that gives rise to a novel nonparametric estimator. Through simulation studies, the estimator is compared with the existing estimators that are all based on parametric assumptions. We show that the increased robustness of the nonparametric estimator is gained at the expense of a moderate loss in efficiency. Then we present a regression framework for evaluating covariate effects on the partial AUC. The approach extends a method developed recently for regression analysis of the full AUC (Dodd and Pepe, 2003). Since the partial AUC has more appeal in many practical settings, this represents an important generalization.

The work is motivated by a study of prostate-specific antigen (PSA), a serum biomarker. PSA is a screening tool for prostate cancer that has been the focus of much research. Important questions to consider when evaluating this biomarker are, which form of PSA is best (total or ratio of free to total PSA), and by how long does this test advance the lead time, or time prior to clinical diagnosis of disease? In this article, we show that a parsimonious regression model of the partial AUC with PSA type, and lead time as covariates, assists with this type of evaluation.

2. Partial AUC

The partial AUC is simply the area under the ROC curve between t_0 and t_1 (Figure 1). With an uninformative test, $\text{TPR}(c) = \text{FPR}(c)$ for all c , and the partial AUC is the area of a trapezoid equal to $1/2(t_1 + t_0)(t_1 - t_0)$. For a perfect test, $\text{ROC}(t) = 1$ for all $t \in (0, 1)$, and the partial AUC is the area of the rectangle with height 1 and base $t_1 - t_0$.

2.1 Interpretations of the Partial AUC

Assume Y^D and $Y^{\bar{D}}$ are continuous random variables with survivor functions S_D and $S_{\bar{D}}$, respectively. Let $F = 1 - S$. The partial AUC is the joint probability that $Y^D > Y^{\bar{D}}$ and that $Y^{\bar{D}}$ falls within the range of clinically relevant quantiles. To see this, observe that:

$$\begin{aligned} \text{AUC}(t_0, t_1) &= \int_{t_0}^{t_1} \text{ROC}(t) dt = \int_{t_0}^{t_1} S_D\{S_{\bar{D}}^{-1}(t)\} dt \\ &= \int_{S_{\bar{D}}^{-1}(t_1)}^{S_{\bar{D}}^{-1}(t_0)} S_D(y) dF_{\bar{D}}(y) \end{aligned}$$

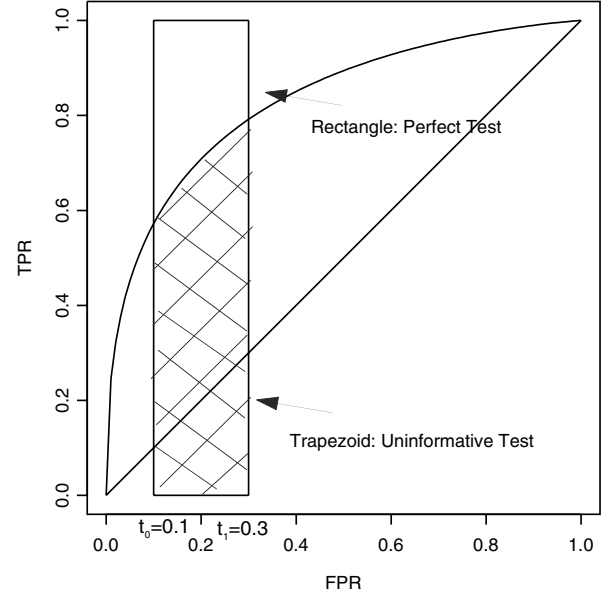


Figure 1. Illustration of an ROC curve and its partial AUC ($t_0 = 0.1$, $t_1 = 0.3$).

$$\begin{aligned} &= \int_{S_{\bar{D}}^{-1}(t_1)}^{S_{\bar{D}}^{-1}(t_0)} S_D(y) f_{\bar{D}}(y) dy \\ &= P\{Y^D > Y^{\bar{D}}, Y^{\bar{D}} \in (S_{\bar{D}}^{-1}(t_1), S_{\bar{D}}^{-1}(t_0))\} \quad (2) \end{aligned}$$

To simplify notation let $q_0 = S_{\bar{D}}^{-1}(t_0)$ and $q_1 = S_{\bar{D}}^{-1}(t_1)$. The partial AUC can also be written as a weighted conditional expectation $\text{AUC}(t_0, t_1) = (t_1 - t_0)P\{Y^D > Y^{\bar{D}} | Y^{\bar{D}} \in (q_1, q_0)\}$.

A second interpretation arises from the concept of placement values (Pepe, 2003). The term $S_D(y^{\bar{D}})$ is the *placement* of a given nondiseased test result, $Y^{\bar{D}} = y^{\bar{D}}$, in the survivor function of results of diseased. For a good test, the nondiseased observations fall in the lower tail of the distribution of diseased. Hence, larger placement values indicate a better test. Hanley and Hajian-Tilaki (1997) showed a connection between average placement values and the full AUC. We extend their result to the partial AUC here, as it provides an interesting interpretation of the partial AUC, as well as a simpler estimator:

$$\begin{aligned} &P\{Y^D > Y^{\bar{D}}, Y^{\bar{D}} \in (q_1, q_0) | Y^{\bar{D}} = y^{\bar{D}}\} \\ &= P\{Y^D > y^{\bar{D}}, y^{\bar{D}} \in (q_1, q_0)\} \\ &= I\{y^{\bar{D}} \in (q_1, q_0)\}P(Y^D > y^{\bar{D}}) \\ &= I\{y^{\bar{D}} \in (q_1, q_0)\}S_D(y^{\bar{D}}) \equiv \text{Pl}^D(y^{\bar{D}}) \quad (3) \end{aligned}$$

We refer to $\text{Pl}^D(y^{\bar{D}})$ as a “restricted placement value” (with respect to the distribution of results from the diseased population). Note that a full placement value is defined as in (3), except there is no restriction to $Y^{\bar{D}} \in (q_0, q_1)$. The restricted placement is weighted so that only those values that fall within the relevant quantiles are considered. Note that $E\{\text{Pl}^D(Y^{\bar{D}})\} = \text{AUC}(t_0, t_1)$. Thus, the partial AUC can be thought of as an average of restricted placement values. It

is straightforward to show that the partial AUC is also the expected restricted placement value from conditioning on observations from the diseased population.

One may wish to rescale the partial AUC, especially in regression analysis (see Section 5). We define the partial AUC odds as

$$\begin{aligned} & \text{AUC}(t_0, t_1) / \{(t_1 - t_0) - \text{AUC}(t_0, t_1)\} \\ &= P\{Y^D > Y^{\bar{D}} \mid Y^{\bar{D}} \in (q_1, q_0)\} / \\ & P\{Y^D < Y^{\bar{D}} \mid Y^{\bar{D}} \in (q_1, q_0)\}. \end{aligned} \quad (4)$$

This is the ratio of the probability of a correct ordering of a randomly selected diseased and nondiseased test result to the probability of an incorrect ordering, with both probabilities *conditional* on the test result of nondiseased arising from the region of interest. These odds have value of $(t_1 + t_0) / \{2 - (t_1 + t_0)\}$ when a test is uninformative, and of infinity for a perfect test.

2.2 Restricting the True-Positive Range

The partial AUC just described restricts the ROC region of interest to false-positive rates that take values in (t_0, t_1) . In some settings, one may wish to restrict to a range of true-positive rates (Jiang et al., 1996). By transforming the ROC curve to a plot of $\{S_D(c), 1 - S_{\bar{D}}(c)\}$, interpretations of a partial AUC corresponding to true-positive rates in an interval are easily obtained. The curve is no longer the classic ROC curve. It is simply a 270° rotation of Figure 1. We refer to this curve as the specificity-ROC curve (ROC_{spe}), since specificity is plotted on the y-axis. For $u = S_D(c)$, this curve is described by $\text{ROC}_{\text{spe}}(u) = 1 - S_{\bar{D}}\{S_D^{-1}(u)\} = F_{\bar{D}}\{S_D^{-1}(u)\}$. The partial AUC for a range of true positives, denoted $\text{AUC}^{\text{TP}}(u_0, u_1)$, is defined as

$$\begin{aligned} \int_{u_0}^{u_1} \text{ROC}_{\text{spe}}(u) du &= \int_{S_D^{-1}(u_0)}^{S_D^{-1}(u_1)} \{1 - S_{\bar{D}}(y)\} dS_D(y) \\ &= P[Y^D > Y^{\bar{D}}, Y^D \in \{S_D^{-1}(u_1), S_D^{-1}(u_0)\}]. \end{aligned} \quad (5)$$

Note that $\int_0^1 S_D(y) dS_D(y) = \int_0^1 \{1 - S_{\bar{D}}(y)\} dS_D(y)$. Therefore, the two representations give the same full AUC. However, the quantiles that define the partial AUC are $\{S_D^{-1}(t_1), S_D^{-1}(t_0)\}$, and are derived from the nondiseased distribution for the classic ROC curve. For the case just presented, the quantiles are $\{S_D^{-1}(u_1), S_D^{-1}(u_0)\}$, and arise from the distribution of tests of diseased subjects. With this exception, the interpretations are the same. Thus, we focus on the classic ROC curve, recognizing that methods developed here easily extend to the case in which restriction of true-positive rates is required.

3. Estimators

3.1 Parametric Estimators

Before proposing our nonparametric estimator, we briefly describe existing estimators. McClish (1989) describes what we refer to as the normal-distributions partial AUC estimator (NDE), which assumes that the test results for diseased and nondiseased populations follow normal distributions with different means and variances. Maximum likelihood

estimates of these parameters provide the ROC curve estimate, and numerical integration of it gives the partial AUC estimator: $\widehat{\text{AUC}}(t_0, t_1) = \int_{t_0}^{t_1} \Phi\{\hat{a} + \hat{b}\Phi^{-1}(t)\} dt$. Here Φ denotes the standard normal cumulative distribution function (CDF), $\hat{a} = (\hat{\mu}_D - \hat{\mu}_{\bar{D}}) / \hat{\sigma}_D$ and $\hat{b} = \hat{\sigma}_{\bar{D}} / \hat{\sigma}_D$, where μ and σ^2 denote mean and variance, respectively.

Another approach is to parameterize the ROC curve directly. The most common form is the binormal ROC curve, $\text{ROC}(t) = \Phi\{a + b\Phi^{-1}(t)\}$, for which the partial AUC is given by the same formula above. However, because this approach does not parameterize the distributions of test results, but only the ROC curve that describes the *relationship* between their distributions, it stipulates far weaker assumptions. Two methods have been proposed for estimating parameters a and b of the binormal ROC curve (Metz, Herman, and Shen, 1998; Pepe, 2000). Both methods are distribution-free in that they are functions only of the ranks of the data.

3.2 Proposed Nonparametric Estimator

Denote $V_{ij}^{q_0, q_1} = I\{Y_i^D > Y_j^{\bar{D}}, Y_j^{\bar{D}} \in (q_1, q_0)\}$ and observe that $E(V_{ij}^{q_0, q_1}) = P\{Y^D > Y^{\bar{D}}, Y^{\bar{D}} \in (q_1, q_0)\} = \text{AUC}(t_0, t_1)$, according to the interpretation of equation (1). This suggests the following method-of-moments estimator:

$$\widehat{\text{AUC}}(t_0, t_1) = (n_D n_{\bar{D}})^{-1} \sum_{ij} V_{ij}^{q_0, q_1}, \quad (6)$$

where n_D and $n_{\bar{D}}$ denote the number of observations from the population of diseased and nondiseased, respectively.

In some circumstances, the quantiles (q_0, q_1) will be known. In others, they will not, in which case empirical quantile estimates are substituted. If the empirical quantile value does not coincide precisely with the desired value, as may happen with small sample sizes, values are linearly interpolated. Observe that, when $t_0 = 0$ and $t_1 = 1$, $\widehat{\text{AUC}}(t_0, t_1) = (n_D n_{\bar{D}})^{-1} \sum_{ij} I(Y_i^D > Y_j^{\bar{D}})$. Hence, for the full AUC the estimator reduces to the Mann-Whitney U-statistic, the classic nonparametric AUC estimator. Further, note that this results in the same value as the area calculated from the empirical ROC curve, using the trapezoidal rule when there are no ties in the data.

The same estimator is derived by consideration of empirical restricted placement values. Let \hat{S} denote an empirical survivor function and write the empirical placement value corresponding to an observation from a disease-free subject, $Y^{\bar{D}}$, as:

$$\begin{aligned} \widehat{\text{Pl}}^D(Y_j^{\bar{D}}) &= I\{Y_j^{\bar{D}} \in (q_1, q_0)\} \hat{S}_D(Y_j^{\bar{D}}) \\ &= I\{Y_j^{\bar{D}} \in (q_1, q_0)\} (1/n_D) \sum_i I(Y_i^D > Y_j^{\bar{D}}). \end{aligned} \quad (7)$$

The sample average is

$$\begin{aligned} (1/n_{\bar{D}}) \sum_j \widehat{\text{Pl}}^D(Y_j^{\bar{D}}) \\ = (n_D n_{\bar{D}})^{-1} \sum_j \sum_i I\{Y_i^D > Y_j^{\bar{D}}, Y_j^{\bar{D}} \in (q_1, q_0)\}. \end{aligned}$$

Thus, the nonparametric estimator is an average of restricted placement values within the disease reference distribution. Likewise, the estimator can be written as the average of the placement values within the nondiseased reference distribution. This generalizes the corresponding results for the full AUC given in DeLong, DeLong, and Clarke-Pearson (1988) and Hanley and Hajian-Tilaki (1997). Calculations using placement values are considerably faster computationally, and are used in the simulations described next.

Asymptotic distribution theory is nonstandard because the binary indicators, $V_{ij}^{q_0, q_1}$, are cross-correlated. It can be shown that the projection $(n_D n_{\bar{D}})^{-1} \sum_{ij} \text{PI}^{\bar{D}}(Y_i^D) + \text{PI}^D(Y_j^{\bar{D}})$ is asymptotically equivalent to (6). Since this projection is a sum of independent terms, standard theory provides consistency and asymptotic normality results. For complete details, refer to Dodd (2001). We recommend using the bootstrap to obtain variance estimates.

4. Simulations

4.1 Small Sample Performance of the Estimator

Three different ROC models were simulated to evaluate performance of the proposed method of inference. The models simulated are a normal-distributions model, a proportional-hazards ROC model, and an extreme-value ROC model. The proportional-hazards ROC model assumes that the hazard function for the disease test result distribution is proportional to that of the nondisease distribution. If the ratio of hazards is r , then $\text{ROC}(t) = t^r$. The extreme-value ROC model involves two parameters (e, f), and is of the form $\text{ROC}(t) = \exp[-(1/e)\exp\{-f\Phi^{-1}(t)\}]$ (Cai and Pepe, 2003). Extensive simulations of these models described in Dodd (2001) showed that the nonparametric method produced estimates with little bias and that confidence intervals using the bootstrap standard error estimator and normal quantiles provided good coverage probability.

We present results for a normal-distributions model that assumes $Y^D \sim N(1.5, 1.44)$ and $Y^{\bar{D}} \sim N(0, 1)$ (see Figure 1). Partial AUCs are considered for $(t_0, t_1) = \{(0, 0.1), (0, 0.2), (0.1, 0.2), (0.1, 0.3)\}$ when quantiles are both known and estimated empirically. Sample sizes were generated with both equal and unequal numbers in each group, such that $(n_D, n_{\bar{D}}) = \{(10, 50), (10, 100), (50, 50), (50, 100), (100, 100), (100, 200), (100, 300)\}$. All resulted in estimates with little bias (Table 1). The largest amount of bias (5.4%) occurs for $n_D = 10$ and $n_{\bar{D}} = 50$, for $(t_0, t_1) = (0.1, 0.2)$, when known quantiles are used. Bias becomes negligible as sample size increases and as more of the curve is integrated.

In most cases, bootstrapped standard errors, computed with 200 bootstrap samples, slightly overestimate the truth when the quantiles are estimated. In the case when quantiles are known, the bootstrapped standard error estimator is reasonably unbiased. Consequently, confidence interval coverage tends to be above the nominal level with estimated quantile values, but close to the nominal level with the known threshold. Surprisingly, when empirical quantiles are substituted, there is an increase in efficiency compared to the case when the true quantile is used. The true standard errors associated with $\text{AUC}\{t_0(\hat{q}), t_1(\hat{q})\}$ are consistently smaller than those for $\text{AUC}\{t_0(q), t_1(q)\}$. The reason for this is unclear to us, but it is

a real phenomenon observed throughout our simulation studies. Simulations suggest that the empirical quantile estimator is negatively correlated with the partial AUC estimator, hence reducing the variance. However, explicit characterization of this correlation is difficult to obtain.

4.2 Robustness

To investigate if the nonparametric estimator gains robustness over other estimators, we considered settings where the ROC curve deviated from the classic binormal form. We generate an ROC model in which test results arise from mixtures of distributions. Assume that $Y^D \sim N(0, 1)$, but that $Y^{\bar{D}}$ is a mixture of two normal distributions, f_{Z_1} and f_{Z_2} , with $Z_1 \sim N(\mu_{D,1}, \sigma_{D,1}^2)$ and $Z_2 \sim N(\mu_{D,2}, \sigma_{D,2}^2)$, such that the probability density of $Y^{\bar{D}}$ is $p f_{Z_1}(y) + (1-p) f_{Z_2}(y)$, for a mixing proportion $p \in (0, 1)$. The resulting ROC curve is simply a mixture of ROC curves weighted by the appropriate mixing probabilities and can be expressed as $\text{ROC}(t) = p\Phi\{a_1 + b_1\Phi^{-1}(t)\} + (1-p)\Phi\{a_2 + b_2\Phi^{-1}(t)\}$, where $a_i = (\mu_{D,i} - \mu_{\bar{D},i})/(\sigma_{D,i})$ and $b_i = \sigma_{\bar{D},i}/\sigma_{D,i}$ for $i = 1, 2$. We set $\mu_{D,1} = 5, \sigma_{D,1}^2 = 1.2, \mu_{D,2} = 0, \sigma_{D,2}^2 = 1$, and $p = 0.3$. The true ROC curve is the solid line shown in Figure 2. We use $t_0 = 0$ and three separate maximum false-positive rates, $t_1 = 0.05, 0.1, 1.0$. Very low false-positive rates, such as $t_1 = 0.05$, have been advocated in settings such as cancer screening (Baker and Pinsky, 2001).

Estimates of ROC curves from the NDE, the Pepe estimator (PPE) and the Metz estimator (MZE) are also given in Figure 2. Recall that the normal-distributions estimator assumes that test results are normally distributed in diseased and nondiseased populations. Using sample means and variances, the ROC curve is calculated as $\widehat{\text{ROC}}(t) = \Phi\{\hat{a} + \hat{b}\Phi^{-1}(t)\}$ with $\hat{a} = (\hat{\mu}_D - \hat{\mu}_{\bar{D}})/\hat{\sigma}_D$ and $\hat{b} = \hat{\sigma}_{\bar{D}}/\hat{\sigma}_D$. Clearly, all of these curves fail to accurately represent the true

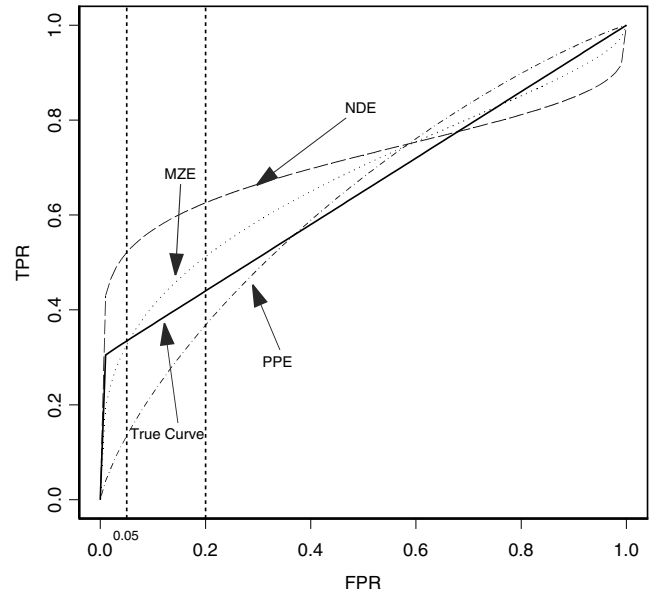


Figure 2. Estimates of the ROC based on data from the mixture model. Shown are curves that use the average parameter estimates.

Table 1
Bias and coverage probability of the nonparametric method with known and estimated quantiles

n_D	$n_{\bar{D}}$	$\widehat{\text{AUC}}\{t_0(\widehat{q}), t_1(\widehat{q})\}$			$\widehat{\text{AUC}}\{t_0(q), t_1(q)\}$		
		Bias	$\widehat{\text{SE}}(\text{SE}) \times 100$	CP	Bias	$\widehat{\text{SE}}(\text{SE}) \times 100$	CP
AUC(0,0.1) = 0.042							
10	50	3.28%	1.578(1.600)	0.95	0.71%	2.315(2.344)	0.88
10	100	2.11%	1.443(1.492)	0.95	0.11%	1.971(1.905)	0.93
50	50	4.30%	0.109(0.101)	0.96	1.90%	0.200(0.192)	0.90
50	100	2.60%	0.085(0.087)	0.94	<0.01%	0.152(0.146)	0.93
100	100	1.50%	0.075(0.072)	0.95	-0.20%	0.138(0.139)	0.94
100	200	0.20%	0.06(0.061)	0.94	0.60%	0.102(0.103)	0.95
AUC(0,0.2) = 0.107							
10	50	1.20%	2.957(2.978)	0.96	2.98%	4.198(4.388)	0.92
10	100	0.58%	2.761(2.848)	0.94	0.14%	3.557(3.347)	0.93
50	50	1.50%	0.185(0.172)	0.97	1.60%	0.347(0.338)	0.92
50	100	1.20%	0.152(0.155)	0.94	1.80%	0.269(0.252)	0.92
100	100	0.80%	0.128(0.127)	0.94	-1.40%	0.224(0.237)	0.96
100	200	0.00%	0.107(0.109)	0.94	-0.20%	0.178(0.178)	0.95
AUC(0.1,0.2) = 0.065							
10	50	-0.17%	1.929(1.593)	0.99	5.40%	3.119(3.403)	0.88
10	100	-0.43%	1.616(1.530)	0.95	-0.23%	2.479(2.569)	0.92
50	50	-0.30%	0.136(0.086)	1.00	1.40%	0.284(0.280)	0.92
50	100	0.40%	0.094(0.080)	0.98	3.00%	0.217(0.203)	0.92
100	100	0.30%	0.082(0.065)	0.99	-2.10%	0.191(0.197)	0.95
100	200	-0.20%	0.061(0.055)	0.96	-0.80%	0.147(0.144)	0.94
AUC(0.1,0.3) = 0.140							
10	50	-0.27%	3.113(2.877)	0.96	2.01%	4.871(4.915)	0.92
10	100	-0.30%	2.820(2.774)	0.94	-0.11%	3.939(4.037)	0.94
50	50	-0.30%	0.194(0.151)	0.99	1.40%	0.416(0.414)	0.95
50	100	0.10%	0.151(0.142)	0.96	1.90%	0.301(0.305)	0.94
100	100	0.30%	0.126(0.114)	0.97	-2.10%	0.281(0.293)	0.95
100	200	-0.10%	0.103(0.099)	0.95	0.40%	0.218(0.216)	0.95

Note: $\widehat{\text{SE}}$ is bootstrapped standard error ($\times 100$) from 200 bootstrap samples, SE is true standard error ($\times 100$), and CP is coverage probability for 95% confidence intervals using $\widehat{\text{SE}}$ with a normal quantile. Results represent 1000 realizations of normal-distributions model.

curve, which is not of the binormal form. They parameterize the ROC curve incorrectly, and hence produce biased partial AUC estimates (Table 2). The largest bias was observed with the normal-distributions estimator. In contrast to this estimator and those of Pepe and Metz, the nonparametric estimators produce estimates with very small bias. Observe that the amount of bias decreases for both the Metz and the Pepe estimators as more of the curve is integrated. This is consistent with results showing that the binormal model estimators produce AUC estimates that are robust to departures from this model (Hanley, 1996).

Additionally, note that the more parametric methods model and estimate the ROC over the entire (0, 1) range and then integrate the relevant portion to determine the partial AUC. In contrast, the proposed method *directly* estimates the partial AUC over the false-positive rate range of interest. A more robust approach might be to model over (t_0, t_1) , which can be accommodated by the Pepe approach. When the estimates are computed while restricting to (0, 0.05) and (0, 0.2), however, there is still bias with the Pepe method. The mean estimates restricting to the corresponding intervals are $\widehat{\text{AUC}}(0, 0.05) = 0.013$ (-18% bias) and $\widehat{\text{AUC}}(0, 0.2) = 0.044$ (-40% bias).

4.3 Efficiency

We compare the efficiency of the methods for estimating the partial AUC under the normal-distributions model. In this setting, the normal-distributions estimator is more efficient than the others, as it is the maximum likelihood estimator and makes use of the information that the data are normally distributed. The results presented in Figure 3 indicate that the nonparametric estimator with known quantile, $\text{NPE}(q)$, is the least efficient of the methods. In agreement with Table 1, we observe that the nonparametric estimator is more efficient when the quantile is estimated than when the known value is used. The distribution-free binormal estimators are more efficient than the nonparametric estimator, but have reduced efficiency relative to the approach that models the test results. Although very similar, the Pepe estimator appears to be slightly more efficient than the Metz estimator.

The efficiency of all methods approaches that of the normal-distributions estimator as $t \rightarrow 1$. Recall that when $t = 1$, the nonparametric estimator is the Mann-Whitney U-statistic. The normal-distributions estimator is simply a transformation of the standardized differences in means, i.e., of the Z-statistic. Hence, the comparison of the variance

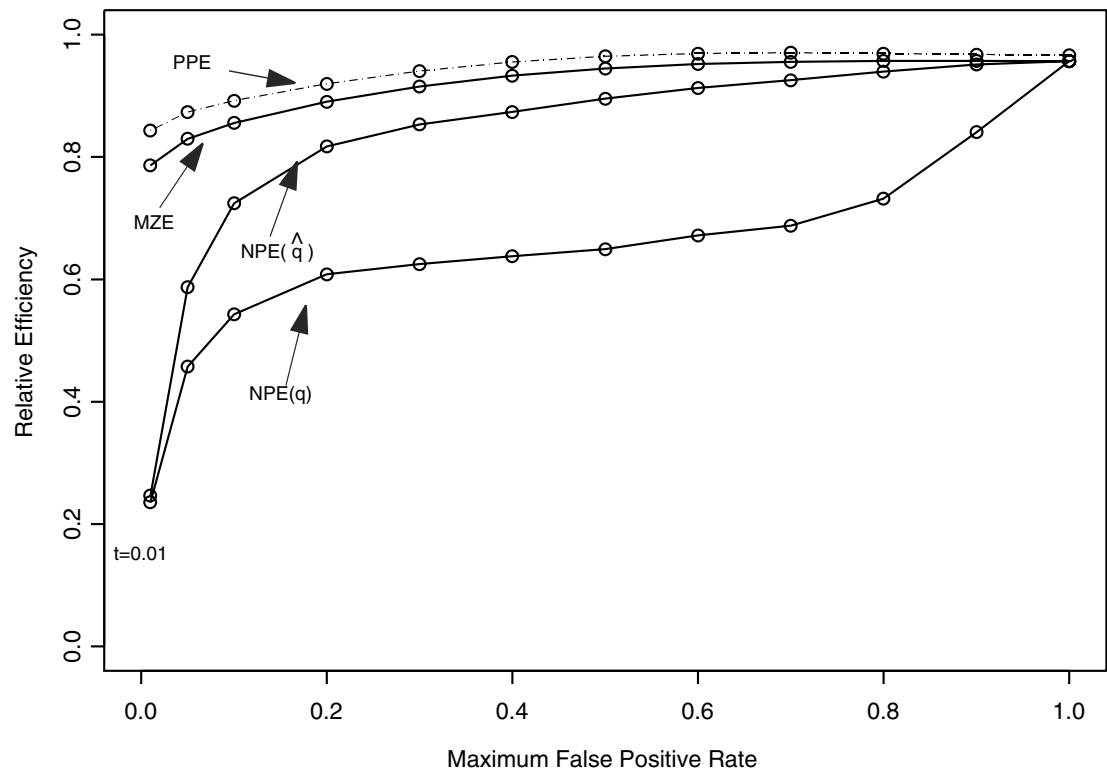


Figure 3. Efficiency of estimators relative to the maximum likelihood normal-distributions estimator. Shown are results with sample sizes of 200 per group and 2000 replicates per scenario.

of the nonparametric estimator with the variance of the normal-distributions estimator is akin to the comparison of the efficiency of the Mann-Whitney U-statistic to that of the Z-statistic. It is well known that the asymptotic relative ef-

iciency of the Mann-Whitney U-statistic to the t-statistic is 0.95 (Lehmann, 1997). Thus it is no surprise that, in this study, the relative efficiency is 0.94 at $t = 1$.

4.4 Recommendations

The normal-distributions estimator, while most efficient, produces unacceptably biased estimates. Estimators that parameterize the ROC curve are similarly not robust. The nonparametric estimator provides substantially more robust estimation. The added robustness is at the expense of moderate losses in efficiency. Hence, we recommend use of the nonparametric estimator. Indeed, for partial areas with $t > 0.1$, it has similar efficiency to the Pepe and Metz methods.

5. Regression Analysis

Accuracy may depend on information other than the test result itself. For example, accuracy may depend on how close in time a subject is to clinical diagnosis. Patient characteristics, such as age or family history of disease, may also be relevant determinants of accuracy. This information may guide decisions about what populations are most likely to benefit from testing or for what populations test innovations are needed. We extend regression methodology we previously developed for the AUC to a regression methodology for the partial AUC summary index of test accuracy. For a more extensive discussion of modeling approaches for the AUC and details about fitting, refer to Dodd and Pepe (2003). Here, we elaborate only on aspects unique to the partial AUC.

Table 2		
Bias in partial AUC estimates from ROC mixture model.		
Shown are results with 200 samples per group and 1000 simulations of the model.		
Method	Estimate	Percent bias
AUC(0, 0.05) = 0.0158		
Nonparametric estimator (\hat{q})	0.016	1
Nonparametric estimator (\hat{q})	0.016	<1
Pepe estimator	0.003	-81
Metz estimator	0.012	22
Normal-distributions estimator	0.023	46
AUC(0, 0.2) = 0.074		
Nonparametric estimator (\hat{q})	0.075	1
Nonparametric estimator (\hat{q})	0.074	<1
Pepe estimator	0.043	-41
Metz estimator	0.078	5
Normal-distributions estimator	0.110	49
AUC(0, 1) = 0.649		
Nonparametric estimator	0.651	<1
Pepe estimator	0.634	-2
Metz estimator	0.675	4
Normal-distributions estimator	0.709	9

5.1 Models

To assess the effect of covariates on test accuracy, we propose the following partial AUC regression models. Consider a vector of covariates X . Define the covariate-specific partial AUC as $\text{AUC}_X(t_0, t_1) = P\{Y^D > Y^{\bar{D}}, Y^D \in (q_0, q_1) \mid X\}$. For a specified link function g , the general model is given by

$$\text{AUC}_X(t_0, t_1) = g(X^T \beta) \quad (8)$$

Possible link functions include the logit or probit forms. However, since $\text{AUC}(t_0, t_1)$ has an upper bound of $(t_1 - t_0)$, a generalization of the logit that incorporates this constraint is appropriate: $g^{-1}(u, t_0, t_1) = \log\{u/(t_1 - t_0 - u)\}$. When this link is used, an interpretation that corresponds to the partial AUC odds defined earlier follows. Consider the model

$$\log[\text{AUC}_X(t_0, t_1)/\{(t_1 - t_0) - \text{AUC}_X(t_0, t_1)\}] = \beta_0 + \beta_1 X. \quad (9)$$

It follows that e^{β_1} is a ratio of partial AUC odds for $X + 1$ to X . When $e^{\beta_1} > 1$, the partial AUC odds are an increasing function of X , and accuracy increases with X . If $e^{\beta_1} < 1$, the partial AUC odds are a decreasing function of X . Refer to Dodd and Pepe (2003) for details about model specification.

5.2 Estimating Function

Consider the indicators $V_{ij}^{(q_0, q_1)}$ as defined in (6), but now we condition on the covariates X . They have mean $E\{V_{ij}^{(q_0, q_1)} \mid X\} = \text{AUC}_X(t_0, t_1)$, which suggests that binary regression methods can be used to estimate partial AUC model parameters with the following estimating equation:

$$V_{n_D, n_{\bar{D}}}(\beta) = \sum_i^{n_D} \sum_j^{n_{\bar{D}}} (\partial \theta_X / \partial \beta) \nu^{-1}(\theta_X) \left(V_{ij}^{(q_0, q_1)} - \theta_X \right) = 0. \quad (10)$$

Here, $\theta_X = \text{AUC}_X(t_0, t_1)$, $\nu^{-1}(\theta_X)$ is the variance of $V^{(q_0, q_1)}$, and $(\partial \theta_X / \partial \beta)$ is a $(p \times 1)$ vector of the partial derivatives of θ_X with respect to the model parameters β . This resembles the classic estimating equation for binary regression. However, the binary indicators are cross-correlated in the sense that, for a given i , the set of binary variables $\{V_{ij}, j = 1, \dots, n_{\bar{D}}\}$ are correlated because they are a function of Y_i^D . Consistency and asymptotic normality for parameter estimates from a similar estimating equation is developed in Dodd and Pepe (2003). The same theory applies here, at least when the quantiles (q_0, q_1) are known. See Dodd (2001) for a full exposition.

5.3 Comparing Two Tests

The proposed estimating equation gives rise to a standard approach when a model to compare two tests is of interest. Consider the following model: $\log[\text{AUC}_X(t_0, t_1)/\{t_1 - t_0 - \text{AUC}_X(t_0, t_1)\}] = \beta_0 + \beta_1 X$, where X is an indicator of test type. The scorelike test of $\beta_1 = 0$, based on the estimating equation in (10), reduces to $\widehat{\text{AUC}}_{X=1}(t_0, t_1) - \widehat{\text{AUC}}_{X=0}(t_0, t_1)$, where $\widehat{\text{AUC}}(t_0, t_1)$ is defined as in (6). This statistic is a member of the family of statistics proposed by Wieand et al. (1989) that takes integrated differences in weighted ROC curves to compare two tests. The Wieand et al. approach does not require bootstrapping for inference, as does ours when the quantiles, q_0 and q_1 , are unknown. The current approach, on the other hand, provides a more general regression framework.

5.4 Implementation

To obtain parameter estimates, algorithms developed for binary regression are used. First, however, the quantiles (q_0, q_1) must be specified. If known, they are simply substituted into (10). If quantiles are unknown and do not depend on covariates, they can be estimated empirically from the nondisease data at hand. If the unknown quantiles depend on covariates, a quantile regression model may be specified (see Heagerty and Pepe, 1999). Then, the binary indicators based on pairs of diseased and nondiseased observations are computed. When covariates are categorical and there are sufficient observations of diseased and nondiseased within each category, all possible pairs of diseased and nondiseased *within a given category* are created. That is, the n_D^k and $n_{\bar{D}}^k$ observations corresponding to the category or level of covariate k are selected and we calculate $V_{ij}^{(q_0, q_1)}$ at each X for $k = 1, \dots, K$. Note that the estimating function is modified to indicate the sum over k . When covariates are continuous, pairing of diseased and nondiseased observations $(Y_i^D, Y_j^{\bar{D}})$ may be undertaken for all possible pairs, although we prefer to pair those observations only with covariate values near one another. Finally, if covariates are unique to the diseased group, as with the "time prior to diagnosis" covariate in the PSA example that follows, the covariate does not restrict the pairing. Pairing is based only on covariates common to D and \bar{D} . We refer to the article on AUC regression by Dodd and Pepe (2003) for a detailed discussion of pairing in the presence of covariates. The same considerations are relevant here.

Logistic or probit regression estimation routines in standard statistical packages can be used to solve the estimating equations. Note that standard errors reported from these packages will not be valid, even with a robust sandwich variance estimator, because of the crossed-correlation structure. We use the bootstrap to calculate the standard errors for parameter estimates. Bootstrap samples are taken by sampling subjects as the primary unit.

6. Prostate-Specific Antigen Analysis

The data analyzed here are taken from a retrospective sampling of stored serum from the α -Tocopherol and β -Carotene Study (ATBC), described by Heinonen et al. (1998). Although the primary goal of this study was to evaluate the effect of dietary supplementation of α -tocopherol and β -carotene on lung cancer risk, the development of prostate cancer was also recorded. Additionally, serum samples were collected and stored both at baseline and three years later. For those 240 subjects who were diagnosed with prostate cancer during the eight-year follow-up period, their serum samples were retrospectively evaluated for prediagnostic levels of prostate-specific antigen. Age-matched serum samples for 237 non-prostate-diagnosed subjects were selected for comparative purposes. Two ways of quantifying PSA were considered, the total PSA in serum (denoted by "total") and the ratio of free to total PSA in serum (denoted by "ratio"). Etzioni et al. (1999) previously compared these two measures using a different dataset. In addition to comparing total to ratio PSA, we examine the effect on accuracy of time from serum sampling to clinical diagnosis. One would expect that PSA levels taken close to the time of clinical time of diagnosis would be more predictive. Let "test" denote an indicator that takes a value of one for total PSA and zero for PSA ratio.

Table 3
PSA partial AUC regression model parameter estimates

Coefficient	Estimate	Std. err.	95% CI	AUC(0,0.4) Odds ratio
Intercept	0.99	0.29	(0.43, 1.55)	2.70
PSA measure (1 for total, 0 for ratio)	0.98	0.54	(−0.09, 2.04)	2.67
Time prior to diagnosis (per year) test/time	0.05	0.05	(−0.04, 0.14)	1.05
Interaction	0.16	0.08	(0.01, 0.31)	1.17

The term “time” denotes years prior to clinical diagnosis at which PSA was measured, so that -7 indicates 7 years before diagnosis and 0 indicates sampling concurrent with diagnosis.

Studies of the false-positive rates of PSA at the standard threshold of 4.0 ng/mL vary widely, with values ranging from 0.10 to 0.70 reported in the literature (Barry, 2001; Tanguay et al., 2002). We consider the following partial AUC regression model from $t_0 = 0$ to the midpoint of positive rates, i.e., $t_1 = 0.4$:

$$\begin{aligned} \log[\text{AUC}(0, 0.4) / \{0.4 - \text{AUC}(0, 0.4)\}] \\ = \beta_0 + \beta_1 \text{test} + \beta_2 \text{time} + \beta_3 \text{test} * \text{time}. \end{aligned} \quad (11)$$

The “time” covariate is irrelevant in the disease-free group, hence the $(1 - t_1)$ th empirical quantile within each test type was substituted for q_1 . Refer to Table 3 for parameter estimates. Bootstrap sampling was conducted with 200 replicates using case-control sampling with subjects as the primary unit.

In accordance with Etzioni’s results, total PSA appears to be the better marker for prostate cancer. At time = 0, the partial AUC odds ratio is 2.67, which, at 5% ($p = 0.07$), is almost statistically significant. As expected, PSA accuracy improves when subjects are measured at times closer to clinical diagnosis of prostate cancer. Although this is true for both measures, the interaction term indicates that the time effect is different for the two measures. There is a 17% greater increase in partial AUC odds for each year for total PSA relative to ratio PSA. For total PSA, the partial AUC odds increase by about 23% for each year closer to diagnosis, while for ratio PSA the partial AUC, odds only increase by 5% for each year. Stated another way, the relative performance of the measures changes with time. The relative odds, estimated as 2.67 at diagnosis, is 1.4 at four years prior to diagnosis (Figure 4). Figure 4 also plots the fitted model. Empirical partial AUCs are averaged within a time window, where intervals are selected so that there are sufficient observations per window,

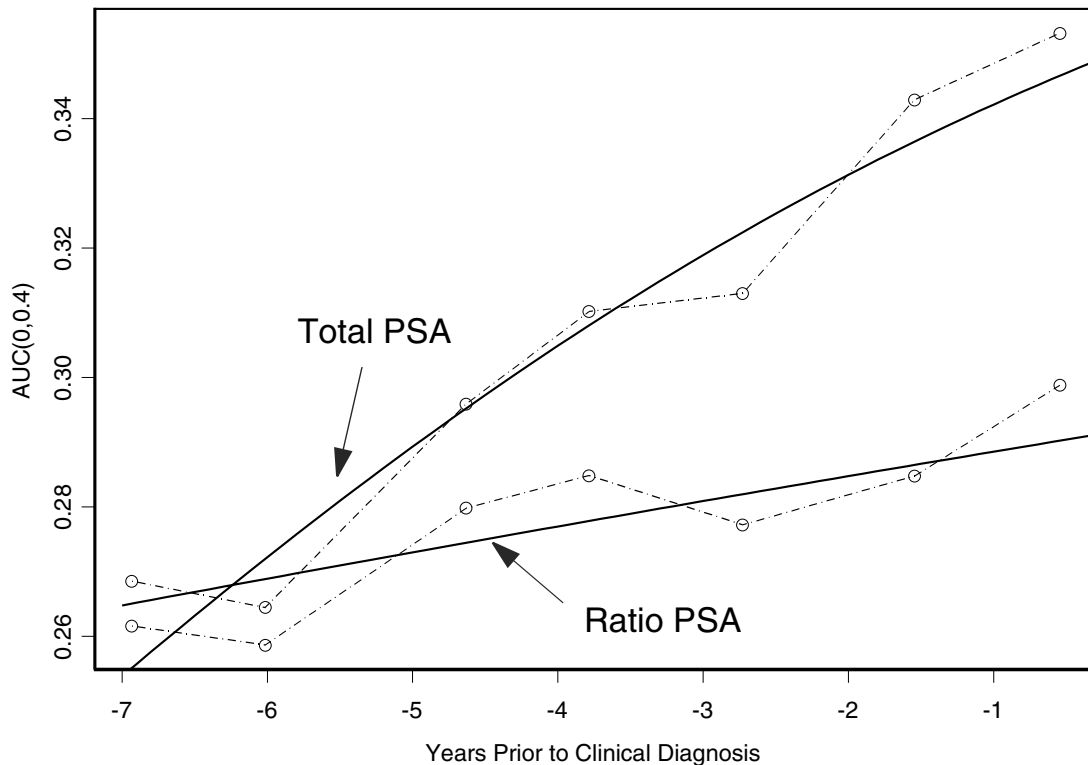


Figure 4. AUC(0, 0.4) model fit. Solid lines represent fitted model. Dotted lines represent empirical values averaged over a time window.

while still avoiding wide time intervals. The model appears to give a reasonable fit, as the fitted and empirical lines are fairly close.

7. Discussion

The partial AUC summarizes test accuracy over a clinically relevant region of the ROC curve. This summary measure can be interpreted as a joint probability that $Y^D > Y^{\bar{D}}$ and that Y^D lies within the quantiles corresponding to the relevant false-positive region. Similarly, the measure can be thought of as the expected restricted placement value. These related interpretations give rise to a nonparametric partial AUC estimator that is more robust than existing estimators, and loses only moderate efficiency relative to them. Simulation studies demonstrate reasonable performance of the approach under a range of models. Interestingly, the estimator is more efficient when estimated quantiles are substituted, as opposed to when their true values are used. When true quantiles are known, one may prefer to estimate them, for a gain in efficiency.

We extend the approach to make inference about partial AUC regression models. One could also use the derived variables approach, as proposed by Thompson and Zucchini (1989) or modify the jackknifed-AUC approach of Dorfman, Berbaum, and Metz (1992) for regression modeling. These methods have been shown to be less efficient in previous studies (Dodd and Pepe, 2003). Further, they are not sufficiently flexible for the range of models of scientific relevance. For example, neither of these methods could be applied to the PSA analysis presented, because the time prior to diagnosis covariate is continuous, and these other methods are restricted to discrete covariate types.

Lastly, we emphasize that, although the partial AUC estimator is a more clinically relevant summary measure of accuracy, the choice of the appropriate restricted region may be arguable. More research is necessary to provide guidance for determining the relevant region. Such a method would inevitably depend on information about the costs and benefits associated with true- and false-positive diagnoses.

ACKNOWLEDGEMENTS

Research supported by grants T32 HL07183 and R01 GM54438 from the National Institutes of Health. The authors acknowledge Phil Taylor for providing the PSA dataset.

RÉSUMÉ

Un diagnostic précis d'une maladie est un élément essentiel des soins. De nouveaux diagnostics et des tests de dépistage doivent être évalués en se basant sur leurs capacités à discriminer les états 'malade' et 'non-malade.' L'aire partielle sous la courbe ROC (APSC) est une mesure de la précision d'un test diagnostique. Nous présentons une interprétation de l'APSC, qui aboutit à un estimateur non-paramétrique. Cet estimateur est plus robuste que certains estimateurs existant et qui font des hypothèses paramétriques. Nous montrons que la robustesse est améliorée avec une perte modérée dans l'efficacité. Nous décrivons une modélisation par régression pour faire des inférences de l'effets des co-variables sur l'APSC. De tels modèles peuvent affiner la connaissance sur la précision des tests. Les paramètres du modèle peuvent être estimés par des méthodes de régression

binaires. Nous utilisons la régression pour comparer deux marqueurs spécifiques des antigènes prostatiques et pour évaluer la précision des marqueurs au cours du temps, antérieurement au diagnostic clinique du cancer de la prostate.

REFERENCES

- Baker, S. (2000). Identifying combinations of cancer markers for further study as triggers of early intervention. *Biometrics* **56**, 1082–1087.
- Baker, S. and Pinsky, P. (2001). A proposed design and analysis for comparing digital and analog mammography: Special receiver operating characteristic methods for cancer screening. *Journal of the American Statistical Association* **96**, 421–428.
- Barry, M. (2001). Prostate-specific-antigen testing for early diagnosis of prostate cancer. *New England Journal of Medicine* **344**, 1373–1377.
- Cai, T. and Pepe, M. (2003). Semi-parametric ROC analysis to evaluate biomarkers for disease. *Journal of the American Statistical Association* **97**, 1099–1107.
- DeLong, E. R., DeLong, D., and Clarke-Pearson, D. (1988). Comparing the areas under two or more correlated receiver operating characteristic curves: A nonparametric approach. *Biometrics* **44**, 837–845.
- Dodd, L. (2001). Regression methods for areas and partial areas under the receiver operating characteristic curve. Ph.D. thesis, University of Washington, Seattle.
- Dodd, L. and Pepe, M. (2003). Semi-parametric regression for the area under the receiver operating characteristic curve. *Journal of the American Statistical Association* **98**, 409–417.
- Dorfman, D., Berbaum, K., and Metz, C. (1992). Receiver operating characteristic analysis: Generalization to the population of readers and patients with the jackknife method. *Investigative Radiology* **27**, 723–731.
- Etzioni, R., Pepe, M., Longton, G., Hu, C., and Goodman, G. (1999). Incorporating the time dimension in receiver operating characteristic curves. *Medical Decision Making* **19**, 242–251.
- Hanley, J. (1996). The use of the “binormal” model for parametric ROC analysis of quantitative diagnostic tests. *Statistics of Medicine* **15**, 1575–1585.
- Hanley, J. and Hajian-Tilaki, K. (1997). Sampling variability of nonparametric estimates of the areas under receiver operating characteristic curves: An update. *Academic Radiology* **17**, 49–58.
- Heagerty, P. and Pepe, M. (1999). Semiparametric estimation of regression quantiles with application to standardizing weight for height and age in US children. *Applied Statistics* **48**, 533–551.
- Heinonen, O., Albanes, D., Vitarmo, J., and Taylor, P. (1998). Prostate cancer and supplementation with α -tocopherol and β -carotene: Incidence and mortality in a controlled trial. *Journal of the National Cancer Institute* **90**, 440–447.
- Jiang, Y., Metz, C., and Nishikawa, R. (1996). A receiver operating characteristic partial area index for highly sensitive diagnostic tests. *Radiology* **201**, 745–750.
- Lehmann, E. (1997). *Testing Statistical Hypotheses*, Chapter 6, 314–321. New York: Springer.

- McClish, D. (1989). Analyzing a portion of the ROC curve. *Medical Decision Making* **9**, 190–195.
- Metz, C., Herman, B., and Shen, J.-H. (1998). Maximum likelihood estimation of receiver operating characteristic (ROC) curves from continuously-distributed data. *Statistics in Medicine* **17**, 1033–1053.
- Obuchowski, N. (1997). Testing for equivalence of diagnostic tests. *American Journal of Roentgenology* **168**, 13–17.
- Pepe, M. (2000). An interpretation for the ROC curve and inference using GLM procedures. *Biometrics* **56**, 352–359.
- Pepe, M. (2003). *Statistical Evaluation of Medical Tests for Classification and Prediction*. Oxford: Oxford University Press.
- Tanguay, S., Begin, L., Elhilali, H., Karakiewicz, P., and Aprikian, A. (2002). Comparative evaluation of total PSA, free/total PSA, and complexed PSA in prostate cancer detection. *Adult Urology* **59**, 261–265.
- Thompson, M. L. and Zucchini, W. (1989). On the statistical analysis of ROC curves. *Statistics in Medicine* **8**, 1277–1290.
- Wieand, S., Gail, M., James, B., and James, K. (1989). A family of nonparametric statistics for comparing diagnostic markers with paired or unpaired data. *Biometrika* **76**, 585–592.

Received July 2002. Revised March 2003.

Accepted March 2003.